

Variational Information Maximization in Stochastic Environments

Felix Agakov



Doctor of Philosophy
Institute for Adaptive and Neural Computation
School of Informatics
University of Edinburgh
2005

Abstract

Information maximization is a common framework of unsupervised learning, which may be used for extracting informative representations \mathbf{y} of the observed patterns \mathbf{x} . The key idea there is to maximize mutual information (MI), which is a formal measure of coding efficiency. Unfortunately, exact maximization of MI is computationally tractable only in a few special cases; more generally, approximations need to be considered. Here we describe a family of variational lower bounds on mutual information which gives rise to a formal and theoretically rigorous approach to information maximization in large-scale stochastic channels. We hope that the results presented in this work are potentially interesting for maximizing mutual information from several perspectives. First of all, our method optimizes a proper lower bound, rather than a surrogate objective criterion or an approximation of MI (which may only be accurate under specific asymptotic assumptions, and weak or even undefined when the assumptions are violated). Secondly, the flexibility of the choice of the variational distribution makes it possible to generalize and improve simple bounds on MI. For example, we may introduce tractable *auxiliary variational* bounds on MI, which may be used to improve on any simple generic approach without altering properties of the original channel. Thirdly, the suggested variational framework is typically simpler than standard variational approaches to maximizing the conditional likelihood in stochastic autoencoder models, while it leads to the same fixed points in its simplest formulation; this gives rise to more efficient optimization procedures. Finally, in some cases the variational framework results in optimization procedures which only require local computations, which may be particularly attractive from the neuro-biological perspective. Possibly the most important contribution of this work is a rigorous and general framework for maximizing the mutual information in intrinsically intractable channels. We show that it gives rise to simple, stable, and easily generalizable optimization procedures, which outperform and supersede many of the common approximate information-maximizing techniques. We demonstrate our results by considering clustering, dimensionality reduction, and binary stochastic coding problems, and discuss a link to approximate statistical inference.

Acknowledgements

I thank my teachers for their support and patience. My family for their support and love. And Mikhail Khodorkovsky for his honor and courage.

Declaration

I declare that this thesis was composed by myself, that the work contained herein is my own except where explicitly stated otherwise in the text, and that this work has not been submitted for any other degree or professional qualification except as specified.

(Felix Agakov)

Table of Contents

1	Introduction	1
1.1	Background	1
1.2	Maximum Likelihood and Maximum Mutual Information	5
1.3	Information-Maximization in Noisy Channels	7
1.3.1	Information Content, Entropy, and Mutual Information	8
1.4	Computational Problems of the Exact Formulation	13
1.4.1	Tractable Special Cases	13
1.5	Common Approximations of Mutual Information	14
1.5.1	<i>As-if</i> Gaussian Approximation of $I(\mathbf{x}, \mathbf{y})$	14
1.5.2	Other Approximations of $I(\mathbf{x}, \mathbf{y})$	15
1.6	Thesis Overview	17
2	Variational Information Maximization	20
2.1	Generic Lower Bound on Mutual Information	21
2.1.1	Definition of a Simple Lower Bound on Mutual Information	21
2.1.2	The Variational IM Algorithm	22
2.1.3	Reconstruction of the Source Patterns	25
2.1.4	Posterior Approximations	26
2.2	Tractable Choices of Variational Decoders	27
2.2.1	Structured Decoders	28
2.2.2	Gaussian Decoders and the Link to Linsker’s Criterion	32
2.3	An Auxiliary Variational Lower Bound on Mutual Information	34
2.3.1	Representations	34
2.3.2	An Auxiliary Variational Lower Bound on $I(\mathbf{x}, \mathbf{y})$	36
2.4	Summary	38
3	Likelihood, Conditional Likelihood, and Information Maximization	40
3.1	Introduction	40
3.2	Information Maximization and Maximum Likelihood	42
3.2.1	Variational Information Maximization and Exact Likelihood Training	43
3.2.2	Variational Information Maximization and Variational Likelihood Training	47
3.3	Information Maximization and Maximum Conditional Likelihood	51

3.3.1	Variational Information Maximization and Exact Conditional Likelihood Training in Feed-Forward Models	53
3.3.2	Variational Information Maximization and Exact Training of Autoencoders	54
3.3.3	Variational Information Maximization and Variational Training of Autoencoders	58
3.4	Summary	62
4	Variational Information Maximization for Linear Dimensionality Reduction	65
4.1	Introduction	66
4.1.1	Optimization of Linsker’s Criterion	67
4.2	Optimization of the Auxiliary Variational Bound	69
4.2.1	Representations	70
4.2.2	Learning Optimal Parameters	71
4.2.3	Learning Optimal Representations in the Augmented $\{y, z\}$ -space	74
4.3	Demonstrations	77
4.3.1	Hand-Written Digits: Comparing the Bounds	77
4.3.2	Hand-Written Digits: Reconstructions	80
4.4	Summary	81
5	Variational Information Maximization for Nonlinear Dimensionality Reduction	83
5.1	Introduction	84
5.2	Nonlinear Channels for Information-Theoretic Clustering	86
5.2.1	Variational Information Maximization for Gaussian Mixtures	86
5.2.2	Information Maximization for Clustering with Encoder Models	90
5.2.3	Information Maximization for Clustering with Kernelized Encoder Models	93
5.3	Nonlinear Gaussian Channels	96
5.3.1	Analytic Properties of IM Solutions	97
5.4	Demonstrations	102
5.4.1	IM for Gaussian Mixture Models	102
5.4.2	Kernelized Information Theoretic Clustering	106
5.5	Summary	113
6	Variational Information Maximization for Learning High-Dimensional Discrete Representations	116
6.1	Introduction	117
6.2	Variational Learning of Population Codes	118
6.2.1	Local Iterative Learning	119
6.2.2	Optimal Encoder Models	121
6.3	Fisher Information and Mutual Information	123
6.3.1	Fisher Approximation of Mutual Information	124
6.3.2	Local Approximation of Mutual Information	126

6.3.3	Constraints on the Encoder Parameters	129
6.4	Variational Lower Bound <i>vs.</i> Fisher and Local Approximations of Mutual Information	131
6.5	Demonstrations	132
6.6	Summary	138
7	Auxiliary Variational Inference and Variational Mutual Informa- tion Maximization	140
7.1	Introduction	141
7.1.1	Existing Variational Approximations	142
7.2	Auxiliary Variational Method	143
7.2.1	Optimizing the Auxiliary Variational Bound	145
7.2.2	Specific Auxiliary Representations	147
7.3	Relation to Variational Mixture Models	150
7.4	Demonstrations	151
7.5	Summary	156
8	Discussion	159
A	Linsker’s Bound: Centering in the Code Space	170
B	Analysis of Variational Information Maximization	172
B.1	Variational Information Maximization and Feed-Forward Models .	172
B.2	Variational Information Maximization and Flat Mixture Models .	174
C	Variational Information Maximization for Isotropic Gaussian Chan- nels	177
C.1	Linear Gaussian Channels: Linear Decoders	178
C.1.1	Nature of Optimal Solutions	179
C.2	Nonlinear Gaussian Channels: Linear Decoders	180
C.2.1	Kernelized Representation	180
C.2.2	Nature of optimal solutions	182
C.3	Nonlinear Gaussian Channels: Nonlinear Decoders and KPCA . .	184
C.3.1	Constraints on the Feature Mappings $p(\mathbf{f} \mathbf{x})$	185
C.3.2	Variational Bounds on $H(\mathbf{f})$	186
C.3.3	Kernelized Representation	187
C.3.4	Nature of Optimal Solutions	187
C.3.5	Optimal Kernel Functions	188
	Bibliography	191

List of Figures

1.1	A schematic representation of a noisy channel. The source vector \mathbf{x} is transformed into a codeword \mathbf{t} (for example, by introducing a known systematic redundancy). The codeword is transmitted over the noisy channel, leading the perturbed representation \mathbf{y} . The receiver decodes \mathbf{y} to produce an estimate $\tilde{\mathbf{x}}$ of the original source vector.	8
1.2	Graphical model of a noisy communication channel. Generally, we assume that we cannot eliminate the channel noise, which is implied in the parameterization of $p(\mathbf{y} \mathbf{t})$	11
2.1	Simple variational decoders. <i>Left</i> : A noisy encoder $p(\mathbf{y} \mathbf{x})$; <i>Right</i> : constrained variational decoders $q^{(1)}(\mathbf{x} \mathbf{y})$, $q^{(2)}(\mathbf{x} \mathbf{y})$, and $q^{(3)}(\mathbf{x} \mathbf{y})$. The shaded nodes correspond to the hidden encoded representations $\{\mathbf{y}\}$. The transparent nodes correspond to the sources $\{\mathbf{x}\}$ and their reconstructions. (Unless mentioned otherwise, we will use the shaded nodes to indicate unobservable variables).	29
2.2	Structured variational decoders. (a) A noisy encoder $p(\mathbf{y} \mathbf{x})$; (b) the corresponding fully-structured exact decoder $p(\mathbf{x} \mathbf{y})$; (c) a sparse variational decoder $q(\mathbf{x} \mathbf{y})$. The shaded nodes correspond to the encodings $\{\mathbf{y}\}$. The sparse variational decoder retains some of the structure of the exact posterior $p(\mathbf{x} \mathbf{y})$	32
2.3	A noisy channel $p(\mathbf{y} \mathbf{x})$ with a structured mixture-type decoder $q(\mathbf{x} \mathbf{y})$. (The states of the reconstructed variables are denoted by $\tilde{\mathbf{x}}$). The role of the auxiliary variables \mathbf{z} is to introduce additional structure into the decoder; they are <i>not</i> transmitted across the channel $p(\mathbf{y} \mathbf{x})$ and do not explicitly constrain $p(\mathbf{x}, \mathbf{y})$. The dashed lines show the mappings to and from the auxiliary space. The auxiliary node is shown by the double circle.	35
3.1	Generative and encoder models. (a): A generative model $\mathcal{M}_L \stackrel{\text{def}}{=} p(\mathbf{y})p(\mathbf{x} \mathbf{y})$ trained by maximizing the likelihood \mathcal{L} ; (b): a model of a noisy channel $\mathcal{M}_I \stackrel{\text{def}}{=} \tilde{p}(\mathbf{x})p(\mathbf{y} \mathbf{x})$ (trained by maximizing the mutual information $I(\mathbf{x}, \mathbf{y})$). The shaded nodes correspond to the hidden variable representations \mathbf{y}	42

3.2	Mutual information $I(\mathbf{x}, \mathbf{y})$, the variational lower bound $\hat{I}(\mathbf{x}, \mathbf{y})$, and the exact log-likelihood score $\mathcal{L}(\mathcal{M}_L)$ in the parameter space $\Theta = \{p(\mathbf{y}), p(\mathbf{x} \mathbf{y})\}$ (a schematic plot). Intuitively, optimization of the tighter bound $\hat{I}(\mathbf{x}, \mathbf{y})$ on the generally intractable mutual information $I(\mathbf{x}, \mathbf{y})$ leads to a better estimate $\Theta_{\hat{I}}$ of the I -optimal parameters Θ_I . Heuristically, the tightness of a bound may be treated as an approximate criterion of optimality.	45
3.3	Generative and encoder models for i.i.d. patterns. <i>Left</i> : generative model \mathcal{M}_L for M i.i.d. training patterns $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}\}$ with the corresponding latent variable representations $\{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(M)}\}$; <i>Right</i> : the corresponding model of the memoryless channel \mathcal{M}_I . The shaded nodes indicate the implied conditionings on the intrinsically deterministic model parameters.	50
3.4	An autoencoder model for M training patterns $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}\}$ and their reconstructions $\{\tilde{\mathbf{x}}^{(1)}, \dots, \tilde{\mathbf{x}}^{(M)}\}$. The shaded nodes indicate the conditionings on the deterministic model parameters.	51
3.5	Variational information maximization and noiseless autoencoders. Maximization of the conditional likelihood in noiseless autoencoders \mathcal{M}_C reduces to maximization of the variational lower bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ in the corresponding noiseless channels \mathcal{M}_I , where the variational decoder is given by $q(\mathbf{x} = \mathbf{s} \mathbf{y}) = p(\tilde{\mathbf{x}} = \mathbf{s} \mathbf{y})$. Thick arrows show deterministic mappings; the dashed arrow corresponds to the variational distribution.	55
3.6	Variational information maximization and stochastic autoencoders. Maximization of the standard lower bound on the conditional likelihood $\tilde{\mathcal{L}}_{\tilde{\mathbf{x}} \mathbf{x}}$ in \mathcal{M}_C reduces to maximization of the generic lower bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ in the corresponding encoder models $\tilde{\mathcal{M}}_I \stackrel{\text{def}}{=} q_{\mathbf{y} \mathbf{x}}(\mathbf{y} \mathbf{x})\tilde{p}(\mathbf{x})$ if $q_{\mathbf{y} \mathbf{x}=\mathbf{s}}(\mathbf{y} \mathbf{x}) = q(\mathbf{y} \mathbf{x} = \mathbf{s}, \tilde{\mathbf{x}} = \mathbf{s})$. Here $q(\mathbf{y} \mathbf{x}, \tilde{\mathbf{x}})$ is the variational posterior of $\tilde{\mathcal{L}}_{\tilde{\mathbf{x}} \mathbf{x}}$. Dashed arrows show the graphical structure of the variational distributions.	60
4.1	Variational information maximization for noisy constrained dimensionality reduction. (a) : <i>Top curves</i> : Average values of the variational auxiliary bounds $\tilde{I}(\mathbf{x}, \mathbf{y})$, obtained by the IM algorithm started at 10 random model initializations (shown for $ z = 2, \dots, 5$); <i>bottom line</i> : the <i>as-if</i> Gaussian $I_G(\mathbf{x}, \mathbf{y})$ bound (computed numerically). The results are shown for the digits data with $ \mathbf{x} = 196$, $ \mathbf{y} = 6$ for $M = 30$ patterns and $T = 30$ iterations of the IM. (b) : Hinton diagram for $\mathbf{W}\mathbf{W}_{pca}^T(\mathbf{W}\mathbf{W}_{pca}^T)^T \in \mathbb{R}^{6 \times 6}$ for $ z = 3$, $T = 30$. For orthonormal weights spanning identical subspaces, we would expect to see the identity matrix.	78

4.2	<i>Histogram bars</i> : marginal variances of the orthonormal noisy projections of the data patterns $\langle (y_i - \langle y_i \rangle)^2 \rangle_{p(y_i x)\tilde{p}(x)}$ along each of $i = 1, \dots, 6$ dimensions of the code space. The code space $\mathcal{R}^{ y }$ is spanned by W_{pca} and W , for Linsker's and the auxiliary variational objectives respectively. <i>Vertical lines</i> : variances of the average distances for $M = 30$ data patterns. The results are shown for $ x = 196$, $ y = 6$; the auxiliary space size $ z = 3$; the number of iterations $T = 30$	79
4.3	Reconstructions of the source patterns from encoded representations. (a) : A subset of the generic patterns used to generate the source vectors; (b) : the corresponding reconstructions from 6 principal components; (c) : the corresponding \tilde{I}_H -optimal reconstructions at $\langle x \rangle_{q(x y,z)} = U_z y$ for the hybrid $\{y, z\}$ representations ($ y = 6$, $ z = 3$).	80
5.1	\mathcal{L} - and \hat{I} - maximization for a three-component Gaussian mixture ($ y = 3$). <i>Left</i> : clustering with the EM algorithm on \mathcal{L} (the squares show the cluster centers); <i>Middle</i> : clustering with the IM on $\hat{I}(x, y)$; <i>Right</i> : Proportion of the testing points assigned to each of $ y = 3$ clusters under $T = 10$ different runs of the optimization procedure (light gray bars indicate the variance). The exact mutual information computed by the EM- and IM-trained models was $I_L(x, y) \approx 0.75$ and $I_I(x, y) \approx 1.00$ respectively; the corresponding log-likelihoods are $\mathcal{L}_L \approx -4.74$ and $\mathcal{L}_I \approx -4.99$	103
5.2	Samples from the Gaussian mixture model \mathcal{M}_L trained by the EM- and IM-algorithms with $ x = 2$, $ y = 3$. The number of training patterns was $M = 250$. <i>Left</i> : \mathcal{M}_L is trained by maximizing the likelihood \mathcal{L} (resulting in $I_L(x, y) \approx 0.654$); <i>Right</i> : \mathcal{M}_L is trained by maximizing the bound $\hat{I}(x, y)$ (resulting in $I_I(x, y) \approx 0.911$). The ellipses show probability contours of the mixture components at $1/2$ Mahalanobis distance from the means μ_j for $j = 1, \dots, y $. The dashed lines indicate the boundary of the convex region used to generate the underlying data.	104
5.3	$\hat{I}(x, y)$ -maximization <i>vs.</i> K-means. <i>Solid lines</i> : eigenvectors of the covariances Σ_j obtained by the IM. <i>Dashed lines</i> : eigenvectors of the sample covariances Σ_j^{KM} associated with each of the clusters obtained by the deterministic k-means algorithm. All the eigenvectors are weighted by the corresponding eigenvalues and centered at the components' means (μ_j and μ_j^{KM} for the IM and k-means respectively). The IM was initialized at μ_j^{KM} and Σ_j^{KM} . The filled circles show the training patterns which are classified differently by the k-means and the IM started at the considered initialization.	105

5.4	Clustering for $ y = 3$. <i>Left</i> : Gaussian mixtures; <i>Middle</i> : K-means; <i>Right</i> : information-maximization for the (RBF-)kernelized encoder. Light, medium, and dark-gray squares show the cluster colors corresponding to deterministic cluster allocations. The color intensity of each training point $\mathbf{x}^{(m)}$ is the average of the pure cluster intensities, weighted by the responsibilities $p(y_j \mathbf{x}^{(m)})$. Nearly indistinguishable dark colors for the majority of patterns under the Gaussian mixture clustering indicate soft cluster assignments (see also Figure 5.5). Note that by applying the kernelized IM algorithm, we obtain nearly deterministic cluster assignments to locally smooth data regions.	107
5.5	Probabilities $p(y_j \mathbf{x}^{(m)})$ for the spiral data, $\mathbf{x} \in \mathbb{R}^2$, $y \in \{y_1, y_2, y_3\}$. The total number of patterns $M = 70$. <i>Left</i> : Gaussian mixtures; <i>Middle</i> : K-means; <i>Right</i> : information-maximization for the kernelized encoder model (RBF kernel with the learned parameter $\beta = 0.825$).	108
5.6	Learning cluster allocations for $ y = 2$. Where appropriate, the stars show the cluster centers. <i>Top left</i> : clustering with the two-component Gaussian mixture trained by the EM algorithm on \mathcal{L} ; <i>Top right</i> : clustering with the k-means; <i>Bottom left</i> : clustering with the encoder model $p(y_j \mathbf{x}) \propto \exp\{-\ \mathbf{x} - \mathbf{w}_j\ ^2/s_j\}$ trained by maximizing mutual information $I(\mathbf{x}, y)$; <i>Bottom right</i> : clustering with the kernelized encoder model $p(y_j \mathbf{x}) \propto \exp\{-\ \phi(\mathbf{x}) - \mathbf{w}_j\ ^2/s_j\}$ trained by maximizing $I(\mathbf{x}, y)$ (the results are shown for the RBF kernel). For the kernelized model, the inverse variance β of the RBF kernel varied from $\beta_0 = 1$ (at the initialization) to $\beta \approx 0.604$ after convergence.	109
5.7	Learning cluster allocations for $ y = 3$. Where appropriate, the stars show the cluster centers. <i>Top left</i> : clustering with the three-component Gaussian mixture trained by the EM algorithm on \mathcal{L} ; <i>Top right</i> : clustering with the k-means algorithm; <i>Bottom left</i> : clustering in the kernelized encoder model (the results are shown for the RBF kernel with the inverse variance fixed at $\beta_0 = 1$); <i>Bottom right</i> : clustering with the kernelized encoder model with the adaptable kernel parameter (the inverse variance of the RBF kernel varied from $\beta_0 = 1$ (at the initialization) to $\beta \approx 0.579$ after convergence). Note that by learning the kernel parameter β we obtain higher values of the mutual information ($I \approx 1.10$ vs $I \approx 1.02$) and more intuitive cluster assignments.	110

5.8	Learning cluster allocations by the feature-space k-means. The results are shown for the fixed RBF kernels with the inverse variance parameter β . <i>Left</i> : $ y = 2$, $\beta_1 = 1$ and $\beta_2 = 0.604$; <i>Right</i> : $ y = 3$, $\beta_1 = 1$ and $\beta_2 = 0.579$. The number of training patterns was $M = 210$. The patterns shown by the double circles \odot are the only ones clustered differently for the two settings of the kernel parameters (β_1 and β_2), assuming identical initializations. The parameters β_2 were set to the converged values of β for the corresponding encoder models (<i>cf</i> Figure 5.6 (<i>bottom left</i>), Figure 5.7 (<i>bottom left</i>)).	111
6.1	Changes in the exact mutual information $I(\mathbf{x}, \mathbf{y})$ for parameters of the encoder $p(\mathbf{y} \mathbf{x})$ obtained by maximizing the variational lower bound $\tilde{I}(\mathbf{x}, \mathbf{y})$, the Fisher information criterion $\tilde{I}_F(\mathbf{x}, \mathbf{y})$, and the local approximation criterion $\tilde{I}_L(\mathbf{x}, \mathbf{y})$. The number of training stimuli $M = 20$. The results are averaged over 30 runs from random parameter initializations. <i>Left</i> : $ \mathbf{x} = 3$, $ \mathbf{y} = 5$ <i>Right</i> : $ \mathbf{x} = 7$, $ \mathbf{y} = 5$	133
6.2	Optimal encoder weights $\mathbf{W} \in \mathbb{R}^{ \mathbf{y} \times \mathbf{x} }$ obtained by maximizing $\tilde{I}(\mathbf{x}, \mathbf{y})$, $\tilde{I}_F(\mathbf{x}, \mathbf{y})$, and $\tilde{I}_L(\mathbf{x}, \mathbf{y})$. Here $ \mathbf{x} = 3$, $ \mathbf{y} = 5$, and $M = 20$. The white squares correspond to the settings $w_{ij} \approx 0.5$. The black squares correspond to the settings $w_{ij} \approx -0.5$. For the illustrated case, $\max_{ij} \{ w_{ij} - 0.5 \} \sim O(10^{-4})$	135
6.3	Exact mutual information as a function of (approximately) discrete-valued weights. <i>Left</i> : typical changes in the exact $I(\mathbf{x}, \mathbf{y})$ under maximization of $\tilde{I}(\mathbf{x}, \mathbf{y})$, $\tilde{I}_F(\mathbf{x}, \mathbf{y})$, and $\tilde{I}_L(\mathbf{x}, \mathbf{y})$ for the auxiliary parameters $\{a_{ij}\}$. <i>Right</i> : values of the exact mutual information computed at $\mathbf{W} = \{-0.5, 0.5\}^{ \mathbf{y} \times \mathbf{x} }$ around the optimal \mathbf{W}_i , \mathbf{W}_F , and \mathbf{W}_L . The x -axis corresponds to the decimal representations of the binary encodings of the discrete weight vector $\mathbf{w} \in \{-0.5, 0.5\}^{ \mathbf{x} \times \mathbf{y} }$ obtained by flattening $\mathbf{W} \in \mathbb{R}^{ \mathbf{y} \times \mathbf{x} }$	136
6.4	Stochastic binary encodings of the continuous input stimuli. (<i>a</i>): a subset of the original visual stimuli. (<i>b</i>): 30 samples from the corresponding encoding distribution $p(\mathbf{y} \mathbf{x})$ for $ \mathbf{y} = 10$ spiking units (black and white dashes correspond to silent and spiking output units). (<i>c</i>): Reconstructions of the original sources from 50 samples of neural spikes. Note that we imposed soft constraints on the variances of firings.	137

7.1	Undirected models and their approximations (a schematic plot). (a) A fully connected pairwise Markov network representing the intractable distribution $p(\mathbf{x})$; (b) a standard mean field approximation $q_{MF}(\mathbf{x})$; (c) a structured mean field approximation $q_{SMF}(\mathbf{x})$; (d) a mixture of mean field models. (All the variables \mathbf{x} are coupled through the mixture label y . The dotted lines serve to indicate that the marginal $q_{MMF}(\mathbf{x})$ expressed from $q_{MMF}(\mathbf{x}, y)$ is in general fully connected). The transparent node y shows the auxiliary variable. The shaded nodes indicate the variables forming the space $\{\mathbf{x}\}$ of the original model $p(\mathbf{x})$	143
7.2	An auxiliary mean field approximation. The target distribution $p(\mathbf{x}, y)$ is approximated by $q(\mathbf{x}, y)$, which is structured in the <i>augmented space</i> . Note that the mapping $p(y \mathbf{x})$ was introduced in a way which does not affect the original fully connected pairwise target distribution $p(\mathbf{x})$. The transparent node y shows the auxiliary variable. The shaded nodes indicate the variables forming the space $\{\mathbf{x}\}$ of the original model $p(\mathbf{x})$	144
7.3	An auxiliary variational framework for approximate inference. (a) Systematic changes in the mean squared error for the estimates of the second-order moments with the growth in the number of mixture states M ; (b) <i>Top</i> : re-weighting coefficients for a set of fixed structured approximators (each $q(\mathbf{x} y)$ is a uniquely structured spanning tree); <i>Bottom</i> : the lower bounds on $\log Z$ obtained by each individual tree. Note that greater mixing coefficients were assigned to the tree-structured distributions which individually had resulted in higher values of lower bounds on $\log Z$	152
7.4	Influence of the structure of the auxiliary conditional $p(y \mathbf{x})$ and the variational $q(\mathbf{x} y)$ distributions on the theoretically achievable improvement of the auxiliary variational lower bound on $\log Z$ over a mixture of simple bounds. The x -axis correspond to the “flip rate” $q(y)$. The y -axis correspond to the bound $\tilde{I}(\mathbf{x}, y)$. The results are computed for $ \mathbf{x} = 20$, $ \mathbf{y} = 40$. The curves correspond to the exact values of the bound $\tilde{I}(\mathbf{x}, y)$ for random structures of $q(\mathbf{x} y)$ and $p(y \mathbf{x})$ which satisfy the specified sparsity constraints ($ \boldsymbol{\pi}_y(x) = 3$, $ \boldsymbol{\pi}_x(y) $ and $ \boldsymbol{\pi}_y(y) $ are shown in the legend box). Note that a choice of a richer auxiliary conditional $p(y \mathbf{x})$ generally leads to higher values of $\tilde{I}(\mathbf{x}, y)$	153

Chapter 1

Introduction

Possibly one of the most important tasks faced by an intelligent organism may be formulated as learning about the environment, and using the knowledge to address challenges offered by the external world. Knowledge about the world is represented internally by means of synaptic structures or neural activations. Generally, it is believed that internal representations are formed in a way which allows them to capture useful information about the external environment. The goal of learning may in this case be formulated as finding informative representations about the regularities occurring in the external world.

Many of machine learning methods address fundamentally similar tasks, where one of the key goals is to find informative representations of the observations about the environment. One obvious application related to modeling of biological systems is learning neural population codes for a set of input stimuli. Some other applications include speech analysis and object tracking, face recognition and automated medical diagnostics, where the extracted representations can be used for making important practical decisions. Unfortunately, the problem of finding internal representations of the environment may become notoriously difficult in the presence of noise. In this thesis we consider a class of machine learning methods making a recourse to information theory, and address the fundamental computational problems of applying the methods for finding informative regularities in stochastic environments.

1.1 Background

In recent years we witnessed a significant transformation of methods aimed at designing intelligent systems. Historically, such methods were based mainly on modeling heuristic expert knowledge rather than the problem domain. While being sometimes efficient for addressing relatively simple problems of local inference, the approaches lacked formal mechanisms for resolving global queries not representable in the basis of specified heuristics and axioms. Moreover, they could not be easily applied in incomplete and uncertain domains, i.e. in the majority of practically interesting and challenging situations (e.g. Jaynes (1996), Jensen (1996)). Late 80s and early 90s were marked by an increased popularity of non-symbolic methods and development of the learning paradigm, where one of the

principal objectives was to define a utility function and optimize its expectation with respect to a (generally, parametric) model of the environment. While definitions of the objective criteria used by the earlier models were not always properly understood, it was soon realized that many of the sensible optimization criteria for *supervised*¹ training of early connectionist models had strong probabilistic foundations (e.g. Neal (1992), Bishop (1995)). The probabilistic view was important, as it helped to shift the heuristics from the level of function definition to the level of model specification, which was arguably easier to understand by the majority of experts. Moreover, it suggested a unifying framework for interpretation and analysis of the regression and classification algorithms.

At the same time, a particular level of attention was paid to the development of *unsupervised*² learning techniques (see e.g. Hinton and Sejnowski (1999)), where the idea was to learn statistical models, or extract statistical regularities of the observations (see e.g. Dayan (2001) and references therein for a high-level introduction). Possibly the largest class of the unsupervised methods is based on estimation of probability densities of the observed data, where the fundamental idea is to estimate parameters of a model which would give rise to the observations. Specifically, *maximum likelihood* methods aim at fitting a constrained distribution $p(\mathbf{x}|\Theta)$ to a set of the empirical observations $\{\mathbf{x}\}$ by estimating the deterministic parameters Θ . One way of introducing constraints on $p(\mathbf{x})$ is to consider a generative *latent variable* model $\mathcal{M}_L = p(\mathbf{x}|\mathbf{y}, \Theta_{x|y})p(\mathbf{y}|\Theta_y)$, where \mathbf{y} defines a vector of latent variables, $p(\mathbf{y}|\Theta_y)$ is the marginal distribution of the latent variables parameterized by Θ_y , and $p(\mathbf{x}|\mathbf{y}, \Theta_{x|y})$ is the conditional distribution parameterized by $\Theta_{x|y}$. The goal of learning by maximizing the likelihood would in this case involve fitting $p(\mathbf{x}|\Theta_y, \Theta_{x|y}) = \int_{\mathbf{y}} p(\mathbf{y}|\Theta_y)p(\mathbf{x}|\mathbf{y}, \Theta_{x|y})d\mathbf{y}$ to the observed data, which would involve maximization of the likelihood for Θ_y and $\Theta_{x|y}$ (see e.g. Everitt (1984), Hertz et al. (1991), Bishop (1995), Hinton and Sejnowski (1999) for references, comprehensive descriptions of common unsupervised learning techniques, and general discussions of density estimation methods).

Another class of unsupervised learning methods is based on the idea of finding *most informative* descriptions \mathbf{y} of the training patterns \mathbf{x} , for example by learning the optimal mapping $p(\mathbf{y}|\mathbf{x}, \Theta_{y|x})$. The coding efficiency is typically quantified by *mutual information* between \mathbf{x} and \mathbf{y} , which is defined as $I(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} H(\mathbf{y}) - H(\mathbf{y}|\mathbf{x})$. Here $H(\mathbf{y}|\mathbf{x}) \stackrel{\text{def}}{=} - \int_{\mathbf{x}} \int_{\mathbf{y}} p(\mathbf{y}|\mathbf{x}, \Theta_{y|x})p(\mathbf{x}) \log p(\mathbf{y}|\mathbf{x}, \Theta_{y|x})d\mathbf{y}d\mathbf{x}$ is the *conditional entropy* and $H(\mathbf{y}) \stackrel{\text{def}}{=} - \int_{\mathbf{y}} p(\mathbf{y}) \log p(\mathbf{y})d\mathbf{y}$ is the *marginal entropy*, which for discrete variables may be thought of as measures of uncertainty (see a more detailed discussion in Section 1.3). Additionally, $p(\mathbf{x})$ defines the distribution of the source patterns, and $p(\mathbf{y}) = \int_{\mathbf{x}} p(\mathbf{y}|\mathbf{x}, \Theta_{y|x})p(\mathbf{x})d\mathbf{x}$ is the marginal distribution of their descriptions. The mutual information may therefore be interpreted as the decrease in uncertainty of the training patterns \mathbf{x} given their representations \mathbf{y} , which is

¹In many cases supervised learning may be thought of as an optimization framework for evaluating parameters of the mapping from the inputs to the explicitly defined target outputs.

²Unsupervised learning defines an optimization framework for the cases when the observations cannot be represented as associative pairs of inputs and the corresponding target outputs (or environmental values, cf reinforcement learning).

one of the fundamental measures of information theory (see e.g. McEliece (1977), Cover and Thomas (1991)). While the idea of maximizing mutual information has been long known in communication and information theory (see e.g. Blahut (1972), Arimoto (1972), McEliece (1977)), it was introduced to the machine learning community relatively recently (Linsker (1988)). Since then it has been used in a variety of contexts, including clustering (e.g. Dhillon and Guan (2003), Jenssen et al. (2003)), population coding (e.g. Brunel and Nadal (1998), Stocks and Mannella (2001)), and feature extraction (Principe et al. (2000), Torkkola and Campbell (2000), Gokcay and Principe (2002)). A variety of other unsupervised learning techniques making a recourse to information theory are conceptually different from Linsker’s *infomax* (Linsker (1988), Linsker (1989a), Linsker (1989b)), and utilize the mutual information in different contexts. For example, the family of “information bottleneck” methods (Tishby et al. (1999), Slonim et al. (2001), Chechik and Tishby (2002), Dhillon et al. (2002), Still and Bialek (2004)) proposes to maximize the amount of information which the compressed codes of the original data contain about some variables of relevance, while minimizing the amount of information between the data and the codes. The semi-supervised learning methods of Szummer and Jaakkola (2002) and Corduneanu and Jaakkola (2003) suggest to use the local mutual information (computed for each small region of the data space) as a local regularizer of a conditionally trained model. A fundamentally different *Imax* objective criterion (see e.g. Becker (1992), Becker and Hinton (1992)) suggests to optimize the information content between the code variables, rather than between the sources and the codes. The *redundancy reduction* methods (Barlow (1989), Atick (1992), Field (1994)) may also be viewed as specific approximations of the exact mutual information between the source patterns and their encoded representations (Nadal and Parga (1994), Nadal et al. (1998)). These are just a few of the machine learning methods to mention which require computations of the mutual information under various modeling assumptions.

We may generally view these learning methods as optimization procedures, which optimize different objective criteria defined for specific probabilistic models. The models may be conveniently described within the probabilistic *graphical modeling* paradigm, which offers a convenient framework for graphical representations of joint probability distributions via local constraints (Pearl (1988), Lauritzen and Spiegelhalter (1988), Cowell et al. (1999)). The objective functions, including the likelihood and the mutual information, may then be expressed by computing specific expectations for the considered probabilistic models (generally, the objectives may be interpreted as functionals of the joint probability distribution). The probabilistic view offers many advantages, including a formal generalization of inference (performed by the sound rules of probability theory), and theoretically motivated definitions of objective criteria for optimization. Moreover, reliance on the formal calculus of uncertainty makes the graphical models suitable for reasoning in uncertain (e.g. noisy or partially hidden) domains.

Unfortunately, this theoretical rigor comes at the expense of computational complexity of learning, which may be prohibitively high for many problems of interest. This motivates necessity of accurate and computationally tractable

approximations of the optimized criteria. The fundamental cause of computational intractability of rigorous learning methods in most of the high-dimensional densely-connected graphical models is the high complexity of computing the expectations over the hidden variables (typically, this complexity would be low only for a limited number of special cases). For example, fitting the model $p(\mathbf{x}, \mathbf{y})$ to a set of training patterns $\{\mathbf{x}\}$ by maximizing the exact likelihood would require computations of the marginal $p(\mathbf{x})$, which involves averaging $p(\mathbf{x}|\mathbf{y})$ over the marginal distribution of the hidden variables $p(\mathbf{y})$. Analogously, the exact evaluation of the mutual information between the sources \mathbf{x} and the codes \mathbf{y} requires computations of the entropy of the hidden variables, which is a negated expectation of $\log p(\mathbf{y})$ over $p(\mathbf{y})$. Clearly, as long as the integrals cannot be computed analytically (which happens, for example, when $p(\mathbf{y})$ is not in the tractable family and the integrands do not decouple in \mathbf{y}), the complexity of computing the objectives is generally exponential in $|\mathbf{y}|$.

In this work we will focus particularly on a discussion of variational lower bounds on the generally intractable mutual information for intrinsically intractable stochastic channels. While variational bounds on the generally intractable likelihoods are well-known and widely applied in a variety of approximate learning problems (see e.g. Jaakkola (1997), Saul and Jordan (1998), Barber and Wiergerinck (1998), Neal and Hinton (1998) and many others (see Jordan et al. (1998) for references and overview)), surprisingly little work in machine learning appears to have been done on developing a general framework for optimizing the mutual information for the cases when it cannot be computed exactly. In fact, most of the current methods consider maximizing the mutual information for low-scale channels (when the computations may be tractable), or consider various numerical approximations which are generally accurate only in the asymptotic limits (see e.g. Brunel and Nadal (1998), Kang and Sompolinsky (2001), Hoch et al. (2003)). Other somewhat heuristic objective criteria (e.g. “*information potentials*”) are sometimes used as surrogate objectives (Torkkola and Campbell (2000), Gokcay and Principe (2002)), but generally little is known about the conditions when such approximations may be accurate for maximizing the exact mutual information (Torkkola (2000)). Other methods suggest to optimize proper bounds, but may only be applied for a specific choice of the variational distribution (e.g. Linsker (1992), Jaakkola and Jordan (1998), Lawrence et al. (1998)), which complicates their extensions to richer families of methods.

In this work we will describe a family of variational *lower bounds* on mutual information $I(\mathbf{x}, \mathbf{y})$, which gives rise to a formal and theoretically justified approach to information maximization in noisy channels. We will outline a simple variational Information Maximization (IM) algorithm, reminiscent of the generalized variational EM algorithm for likelihood training (Neal and Hinton (1998)), and demonstrate that it provides a simple and effective tool for learning encoders and variational decoders in a principled manner. We will show that in the simplest case when the variational distribution is unconstrained, the proposed variational optimization procedure reduces to a generally intractable form of the Arimoto-Blahut (Arimoto (1972), Blahut (1972)) algorithm for maximization of channel capacity. In the case when the variational decoder is constrained to be a linear

Gaussian, a specific application of the bound reduces to Linsker’s *as-if* Gaussian criterion (Linsker (1992)). Then we will discuss a richer family of the *auxiliary variational* lower bounds on the mutual information, which formally generalizes on the simple variational bounds.

In the later chapters, we will investigate the relation of our approach to the variational likelihood training in generative models and conditional likelihood training in stochastic autoencoders. Specifically, we will show that conventional methods of maximizing the variational Jensen’s bounds on the conditional likelihood in stochastic autoencoders may be interpreted as a complicated way of optimizing a simple bound on the mutual information. Additionally, standard approaches to training (noiseless) autoencoders may be shown to optimize the simplest form of the variational lower bound on the mutual information between the sources and the codes for a noiseless encoding mapping. Furthermore, we will discuss applications of the considered framework to dimensionality reduction, clustering, and population coding. We will show that for a specific parameterization of the encoder distribution, our method favorably compares with the common approximation of Brunel and Nadal (1998) and Kang and Sompolinsky (2001). We will also compare the variational method with an alternative approximation of the mutual information inspired by the recent work of Szummer and Jaakkola (2002) and Corduneanu and Jaakkola (2003), and empirically demonstrate that in the considered cases our variational technique is more preferable.

Probably the most important contribution of this work is a simple and general framework for maximizing the mutual information in high-dimensional stochastic channels. We show that it formally generalizes and outperforms some of the common information-maximizing techniques. Curious side-products of the framework are the general relations between the conditional and the information-theoretic learning methods (Chapter 3), a compression method for continuous data (Chapter 4), an information-theoretic clustering method applicable for learning kernel functions (Chapter 5), and a local learning rule for population coding in a set of point-neuron models (Chapter 6).

In the rest of this chapter we will briefly discuss general differences of Maximum Likelihood and Maximum Mutual Information methods, and introduce the basic concepts of information theory. The presentation in Section 1.2 is high-level; we will return to it in more details in Chapter 3.

1.2 Maximum Likelihood and Maximum Mutual Information

One of the principal goals of generative graphical modeling is finding a model which describes how the observed data is generated from hidden variable representations. A principled objective function which is typically optimized during learning in this case is the likelihood of the model’s parameters. The probabilistic graphical formulation (see e.g. Pearl (1988), Jordan (1998), Jordan (2005)) may be used to simplify computations of the likelihood in a principled way, where the optimized objective may be formally obtained from the model’s structure and

parameterization. At the stage of problem specification, this helps to shift the focus of applying expert heuristics from developing *ad hoc* algorithms and *ad hoc* optimization criteria to designing graphical *models* (which may have a semantic meaning and may arguably be easier to specify, analyze, and apply for answering inferential queries).

The choice of the likelihood as an objective function for optimization is attractive for a number of reasons. In particular, many practically useful objective functions used rather heuristically in a variety of applied domains (from engineering to neuroscience) may be shown to correspond to likelihood maximization in specific graphical models (e.g. Bishop (1995)). Moreover, if the joint distribution defined by a chosen model has the same parametric form as the true distribution, then likelihood-based learning may be viewed as obtaining statistical estimators of the unknown parameters from a set of samples given by the training set. Subject to very general regulatory conditions, the resulting maximum-likelihood estimators have attractive statistical properties, such as asymptotic consistency, efficiency, and normality (Fisher (1922), Fisher (1925), Cramer (1946), Fisher (1950)).

In many cases, the distribution of the data expressed from a specified model does not perfectly match the underlying process (for example, due to the restrictiveness of the inductive bias). However, optimization of the likelihood is still justifiable in this case. Specifically, for i.i.d. observations, maximization of the likelihood with respect to the deterministic parameters of the model may be viewed as minimization of the Kullback-Leibler divergence (Kullback and Leibler (1951), Kullback (1959)) between the empirical distribution and the model. Since the KL divergence is theoretically minimal if and only if both distributions are identical, we may expect that a trained model will approximate the empirical distribution, subject to the modeling constraints. As the size of a training dataset increases, a properly chosen model trained by maximizing the likelihood will tend to approximate the underlying data-generating process. An optimized model may then be used for generating new data, inferring latent variable representations of the observations, etc.

It is intuitive that for a latent variable model to provide a good fit to the observations, the latent variables should capture useful statistical properties of the visible variables. By approximating the empirical distribution, generative models trained by maximizing the likelihood tend to find latent variable representations which are useful for generating the training data. This, however, does not necessarily quantify predictability of the hidden variables from the observations. In many practical situations it may be useful not only to be able to generate the data and answer inferential queries about the hidden states, but also to quantify how much information about the hidden variables is contained in a given set of observations. For example, one may be interested in obtaining a model where the observed symptoms are most informative about the hidden diseases, or where the training data points are most predictive of their possible cluster labels. It is known from information theory that a principal measure of informativeness is given by the mutual information (e.g. McEliece (1977), Cover and Thomas (1991)). Therefore, if the goal of modeling is re-defined as finding the most in-

formative representations of the training data (in the entropy loss sense) then maximization of the mutual information (Linsker (1988)) between the hidden and the visible domains is, intuitively, a reasonable learning strategy to consider.

An alternative justification for maximization of the mutual information comes from the field of communication and information theory, which is closely related to graphical modeling and machine learning (e.g. MacKay (2003)). In this context we can construct a probabilistic graphical model for a noisy communication problem where the visible data vectors $\{\mathbf{x}\}$, passed through a noisy channel, give rise to perturbed received representations $\{\mathbf{y}\}$. A classic result in information and communication theory is that the mutual information $I(\mathbf{x}, \mathbf{y})$ between the transmitted and the received vectors gives a lower bound on the channel capacity, which defines the maximum source rate for which reliable communication may theoretically be achieved (Shannon (1948), McEliece (1977)). Then it is clear that maximization of mutual information between the transmitted and the received vectors may be interpreted as maximizing the lower bound on the information transmission over the noisy channel (Arimoto (1972), Blahut (1972), Cover and Thomas (1991)), which is a formally justifiable, but generally computationally intractable procedure.

1.3 Information-Maximization in Noisy Channels

The reliable communication of information over noisy channels is a fundamental problem, ranging from the construction of good error-correcting codes (see e.g. Pretzel (1996), MacKay (2003)) to neural sensory processing (Barlow (1989), Nadal et al. (1998)). A principal goal of information transmission is to communicate the source data over a channel without error. The problem is complicated in the presence of a channel noise, which transforms the data sent through the channels. In many practical cases we have a limited control over the communication system and cannot eliminate the noise completely. We would therefore like to find an efficient way to minimize possible effects of the inherently irreducible channel noise.

One of the key ideas of information theoretic approaches to communication in noisy channels is to reduce possible effects of the noise by encoding the source vectors $\{\mathbf{x}\}$ into codewords $\{\mathbf{t}\}$, which is typically addressed by introducing a known form of a systematic redundancy. At the communication stage, the codewords $\{\mathbf{t}\}$ are transmitted through the noisy channel, which leads to the perturbed representations $\{\mathbf{y}\}$ at the receiver's end. Then the receiver applies a known decoding rule to obtain reconstructions $\{\tilde{\mathbf{x}}\}$ of the original sources (see figure 1.1). Intuitively, we would like to construct the codes in such a way that despite the noise, the received vectors preserve much of the useful information about the original sources; in other words, $\{\mathbf{y}\}$ should be strongly predictive of $\{\mathbf{x}\}$. The focus of coding theory and information theory is mainly on finding optimal encoder and decoder systems and analyzing their performance (McEliece (1977), Cover and Thomas (1991), MacKay (2003)).

Traditionally, information theoretic approaches were applied to data compression and communication problems (Shannon (1948), McEliece (1977)). However,



Figure 1.1: A schematic representation of a noisy channel. The source vector \mathbf{x} is transformed into a codeword \mathbf{t} (for example, by introducing a known systematic redundancy). The codeword is transmitted over the noisy channel, leading to the perturbed representation \mathbf{y} . The receiver decodes \mathbf{y} to produce an estimate $\tilde{\mathbf{x}}$ of the original source vector.

relatively recently they have been used for a wider range of machine learning tasks, from feature extraction (Becker (1992), Becker and Hinton (1992), Fisher and Principe (1998), Torkkola and Campbell (2000)) to clustering (Dhillon et al. (2002), Dhillon and Guan (2003), Jenssen et al. (2003)), neural population coding (Linsker (1989b), Linsker (1992), Brunel and Nadal (1998), Pouget et al. (1998), Zhang and Sejnowski (1999)), audio-video fusion (Fisher et al. (2000)), and independent components analysis (Bell and Sejnowski (1995), Shriki et al. (2002)). The family of *information bottleneck* (Tishby et al. (1999), Friedman et al. (2001)) approaches used a specific information-theoretic definition of the rate distortion function, with a number of applications to feature extraction (Slonim et al. (2001), Slonim (2002), Chechik and Tishby (2002)) and dimensionality reduction (Chechik et al. (2003), Still and Bialek (2004)). Moreover, several techniques were developed for combining likelihood and information-theoretic training, with applications to semi-supervised classification (Szummer and Jaakkola (2002), Cordonanu and Jaakkola (2003)). The successful applications of information theoretic techniques for a wide range of machine learning problems advocates their use as a general method for learning in noisy environments, which may go beyond the conventional applications to coding and communication.

In this section we briefly review maximization of information transfer as a general strategy for training probabilistic graphical models, and point out computational difficulties of the exact formulation, as well as a few common ways to address the issue of intractability. Initially, however, we will introduce the fundamental definitions of information theory, which we will use throughout the subsequent discussions.

1.3.1 Information Content, Entropy, and Mutual Information

Before discussing the information maximization principle in the context of learning in probabilistic graphical models, it is instructive to define the *information content* of a random event, as well as the marginal and conditional *entropy* associated with its distributions. First, we will assume that the random variable \mathbf{X} takes discrete values \mathbf{x} in the range \mathcal{R}_X with probability $p(\mathbf{X} = \mathbf{x})$ (to simplify the notation we will write $p(\mathbf{x})$ to denote both the probability of the discrete event \mathbf{x} or, for continuous variables, the probability density function; also, when it is clear from the context, we will use lower-case letters to indicate both random variables

and their current states). Then we will discuss the corresponding definitions for the case of continuous random variables.

1.3.1.1 Information content

We define the information content (or *degree of surprise*) of an event as a continuous monotonically decreasing functional $I(\mathbf{x}) \stackrel{\text{def}}{=} f\{p(\mathbf{x})\} \in [0, \infty)$, $\mathbf{x} \in \mathcal{R}_X$, $f \in \mathcal{F}$ of the distribution $P(\mathbf{x})$, such that a linear decrease in certainty leads to a generally super-linear increase in informativeness. Another requirement which we impose on I is that the combined information content of a set of independent events $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}$ is additive, i.e.

$$I(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}\}) \stackrel{\text{def}}{=} f\{p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)})\} \equiv \sum_{m=1}^M f\{p(\mathbf{x}^{(m)})\}, \quad (1.1)$$

$$\forall i = 1, \dots, M. \mathbf{x}^{(i)} \in \mathcal{R}_X, f \in \mathcal{F}.$$

Here \mathcal{F} indicates the space of the functionals which satisfy the specified constraints. From (1.1) it is clear that

$$f\{p(\mathbf{x})^M\} = Mf\{p(\mathbf{x})\}. \quad (1.2)$$

By the constraints on f and the continuity principle it is easy to show that (1.2) holds $\forall M \in (-\infty, \infty)$ (see e.g. Bishop (1995)). It is also clear that a family of functionals satisfying (1.2) includes $f\{p(\mathbf{x})\} = -\log_b p(\mathbf{x})$, where the base $b > 1$, thus leading to³

$$I(\mathbf{x}) = -\log_b p(\mathbf{x}), \quad b > 1. \quad (1.3)$$

By convention, the information content is expressed in bits ($b = 2$) or in *nats* ($b = e$); unless noted otherwise, we will usually assume that $b = e$. From (1.3) it is clear that the information content of the unlikely events ($p(\mathbf{x}) \rightarrow 0$) is exponentially large, while the information content of events occurring deterministically ($p(\mathbf{x}) = 1$) equals to zero.

1.3.1.2 Marginal entropy of discrete events

The absolute entropy, which is a fundamental measure of uncertainty of a discrete random variable $\mathbf{x} \sim p(\mathbf{x})$, is defined as the average information content

$$H(\mathbf{x}) \stackrel{\text{def}}{=} -\langle \log p(\mathbf{x}) \rangle_{p(\mathbf{x})}, \quad (1.4)$$

where $\langle f(\mathbf{x}) \rangle_{p(\mathbf{x})} \stackrel{\text{def}}{=} \sum_{\mathbf{x}} f(\mathbf{x})p(\mathbf{x})$ denotes the expectation of $f(\mathbf{x})$ over $p(\mathbf{x})$ (assuming that \mathbf{x} is discrete). If $p(\mathbf{x})$ is a marginal distribution of a more complex model, (1.4) is commonly referred to as the *marginal* entropy. Originally developed in statistical physics for describing macroscopic states of closed systems (see e.g. Landau and Lifshitz (1996) for the statistical mechanics formulation), it has

³For $\mathcal{R}_x \equiv \mathbb{N}$ and the normalized distribution function $p(x) = (1/2)^x$, the information content is expressed in bits as $I(x) = f(p(x)) = xf(1/2) = -\log_2 p(x)$ with $f(p(x)) \equiv -\log_2 p(x)$.

been widely used in information theory (e.g. Shannon (1948), McEliece (1977)) and other applied domains. In this context, it may be formally thought of as a quantification of the average amount of information provided by the observation of a random variable \mathbf{x} , distributed according to $p(\mathbf{x})$.

The entropy defined by (1.4) satisfies a number of useful properties. For example, it has been shown that $H(\mathbf{x})$ lower-bounds an average length of a binary message needed to encode \mathbf{x} (Shannon (1948)). Additionally, it is easy to see that for discrete variables the entropy is non-negative, reaching zero only in the deterministic case. For this case it is also easy to show that $H(\mathbf{x})$ is maximized when $p(\mathbf{x}) = 1/|\mathbf{x}|$ is uniform, leading to $H(\mathbf{x}) \leq \log |\mathbf{x}|$, where $|\mathbf{x}|$ is the number of states. Therefore, informally we may think of $H(\mathbf{x})$ as a measure of “sharpness” of $p(\mathbf{x})$, with the maximum obtained for the flattest (uniform) distribution. A number of other useful properties and interpretations are described in a variety of texts on information theory, e.g. Pinsker (1964), McEliece (1977), Cover and Thomas (1991).

1.3.1.3 Conditional entropy of discrete events

For a pair of random vectors \mathbf{x} , \mathbf{y} , we may analogously define the conditional entropy $H(\mathbf{x}|\mathbf{y})$ (*equivocation* of \mathbf{y} about \mathbf{x}). It is given as a functional of the conditional distribution, averaged over the joint probability of the combined events

$$H(\mathbf{x}|\mathbf{y}) \stackrel{\text{def}}{=} -\langle \log p(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{x},\mathbf{y})} = \langle H(p(\mathbf{x}|\mathbf{y})) \rangle_{p(\mathbf{y})}. \quad (1.5)$$

Here we defined $H(p(\mathbf{x}|\mathbf{y})) \stackrel{\text{def}}{=} -\langle \log p(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{x}|\mathbf{y})}$ to be the entropy of the conditional distribution $p(\mathbf{x}|\mathbf{y})$. For the communication channel shown on Figure 1.1, the conditional entropy (1.5) may be interpreted as the uncertainty of the receiver about the source \mathbf{x} for the fixed received vector \mathbf{y} , averaged over the distribution of the received vectors $p(\mathbf{y})$.

As before, it is possible to interpret the conditional entropy as a quantification of an average sharpness. Note that a large value of $H(p(\mathbf{x}|\mathbf{y}))$, achievable for flat posteriors $p(\mathbf{x}|\mathbf{y})$, would indicate a large uncertainty in decoding a specific codeword \mathbf{y} by using $p(\mathbf{x}|\mathbf{y})$. Analogously, large values of the conditional entropy $H(\mathbf{x}|\mathbf{y})$ would indicate a large uncertainty of decoding with the exact posterior *on average*. Intuitively, the receiver may be interested in channels leading to small values of the conditional entropy $H(\mathbf{x}|\mathbf{y})$, though generally the reconstruction error will depend on the receiver’s decoding strategy.

1.3.1.4 Mutual Information

Finally, the *mutual information* is defined as a change in the average information content of \mathbf{x} due to the knowledge of \mathbf{y} :

$$I(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} H(\mathbf{x}) - H(\mathbf{x}|\mathbf{y}) = H(\mathbf{y}) - H(\mathbf{y}|\mathbf{x}). \quad (1.6)$$

Fundamentally, it may be interpreted as a measure of predictability of \mathbf{y} from \mathbf{x} and, due to the symmetry, of \mathbf{x} from \mathbf{y} . From (1.6) it is easy to obtain the

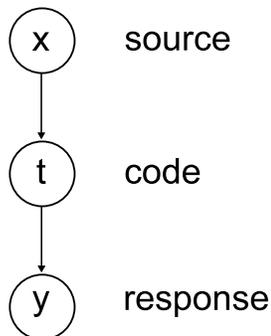


Figure 1.2: Graphical model of a noisy communication channel. Generally, we assume that we cannot eliminate the channel noise, which is implied in the parameterization of $p(\mathbf{y}|\mathbf{t})$.

equivalent formulation

$$I(\mathbf{x}, \mathbf{y}) = KL(p(\mathbf{x}, \mathbf{y}) || p(\mathbf{x})p(\mathbf{y})). \quad (1.7)$$

Clearly, the mutual information is minimized when $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y})$, i.e. when the received vectors \mathbf{y} are independent of the transmitted sources \mathbf{x} . It is easy to see that if both \mathbf{x} and \mathbf{y} are discrete, the mutual information is bounded as $I(\mathbf{x}, \mathbf{y}) \in [0, \log \min\{|\mathbf{x}|, |\mathbf{y}|\}]$, where $|\mathbf{x}|$ and $|\mathbf{y}|$ correspond to the number of states of the input and output variables.

1.3.1.5 Generalization to the continuous case

It may be shown that definitions similar to (1.4) – (1.7) are also applicable for non-discrete random vectors (e.g. Cover and Thomas (1991)). In this case $H(\mathbf{x})$ and $H(\mathbf{x}|\mathbf{y})$ are referred to as *differential entropies*, and the averages in (1.4) and (1.5) are computed by integrating over \mathcal{R}_x and \mathcal{R}_y . It may be shown that in this case the differential entropies diverge to $-\infty$, which makes it difficult to interpret them as a measure of certainty in \mathbf{x} . In fact, the differential entropies are not even invariant under deterministic invertible transformations (McEliece (1977), Bell and Sejnowski (1995)). However, by computing discrete quantizations of the random variables, it may be shown that the mutual information $I(\mathbf{x}, \mathbf{y})$ computed over the continuous range satisfies properties of the mutual information for the discrete case, and therefore has similar properties and a similar fundamentally useful interpretation. Specifically, from definition (1.7) and the property of non-negativity of the Kullback-Leibler divergence, we may see that $I(\mathbf{x}, \mathbf{y}) \geq 0$ for both discrete and continuous quantities.

1.3.1.6 Information maximization in graphical models

Figure 1.2 shows a simple graphical model of a noisy communication channel⁴. In a stochastic context, the quantitative part of the model is specified by the source distribution $p(\mathbf{x})$, the encoding transformation $p(\mathbf{t}|\mathbf{x})$, and a generally noisy non-invertible channel mapping $p(\mathbf{y}|\mathbf{t})$. The key idea of information maximization is to choose a mapping from source variables (inputs) \mathbf{x} to response variables (outputs) \mathbf{y} such that the outputs contain as much information as possible about which of the inputs was transmitted. The principal measure of information transfer in this context is the mutual information $I(\mathbf{x}, \mathbf{y})$, defined by expression (1.6). The general aim will be to set any adjustable parameters of $p(\mathbf{y}|\mathbf{x}) = \langle p(\mathbf{y}|\mathbf{t}) \rangle_{p(\mathbf{t}|\mathbf{x})}$ in order to maximize $I(\mathbf{x}, \mathbf{y})$. Generally, we will presume that we are limited in the way we can adjust $p(\mathbf{y}|\mathbf{t})$; specifically, we may assume that we cannot explicitly modify the channel noise. For example, for deterministic mappings $p(\mathbf{t}|\mathbf{x}) \sim \delta(\mathbf{t} - \mathbf{t}(\mathbf{x}))$ and Gaussian channels $p(\mathbf{y}|\mathbf{t}) \sim \mathcal{N}_y(\mathbf{t}, \sigma_e^2 \mathbf{I})$, we easily get $p(\mathbf{y}|\mathbf{x}) \sim \mathcal{N}_y(\mathbf{t}(\mathbf{x}), \sigma_e^2 \mathbf{I})$. For the fixed noise variance σ_e^2 , our purpose in this specific case would be to modify the encoder mapping $\mathbf{t}(\mathbf{x})$. In a more general context, we will be trying to learn the modifiable parameters of $p(\mathbf{y}|\mathbf{x})$ to maximize the mutual information. With a slight abuse of terminology, we will refer to $p(\mathbf{y}|\mathbf{x})$ as the encoder distribution.

One specific case of interest we consider in this chapter is when the source is given by the empirical distribution

$$\tilde{p}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \delta(\mathbf{x} - \mathbf{x}^{(m)}), \quad (1.8)$$

for a finite set of unique training patterns $\{\mathbf{x}^{(m)} | m = 1, \dots, M\}$ (generalization to the repeating patterns is straight-forward; we omit it here for clarity). Unless stated otherwise, we will also assume that the channel is *memoryless*, i.e. each perturbed pattern $\mathbf{y}^{(i)}$ depends only on the corresponding source vector $\mathbf{x}^{(i)}$. In this case, the theoretically optimal decoder which maximizes the mutual information is given by Bayes rule (e.g. Papoulis (1984)) as

$$p(\mathbf{x}|\mathbf{y}) = \sum_{k=1}^M \delta(\mathbf{x} - \mathbf{x}^{(k)}) \omega_k(\mathbf{y}), \quad \omega_k(\mathbf{y}) \stackrel{\text{def}}{=} \frac{p(\mathbf{y}|\mathbf{x}^{(k)})}{\sum_{m=1}^M p(\mathbf{y}|\mathbf{x}^{(m)})}, \quad (1.9)$$

where we assumed that the adjustable parameters of the encoder $p(\mathbf{y}|\mathbf{x})$ in the definition of $\omega(\mathbf{y})$ are fixed at their optimal values. Clearly, (1.9) is a mixture of Dirac delta functions, which reduces to a single delta peak only for the case when stochastic images \mathcal{I} of distinct source vectors are non-intersecting under the encoder $p(\mathbf{y}|\mathbf{x})$, i.e. $\forall i, j \in \{1, \dots, M\}. i \neq j. \mathcal{I}(\mathbf{x}_i) \cap \mathcal{I}(\mathbf{x}_j) = \emptyset$. Formally, this is never the case if $\forall \mathbf{y}. \exists \mathbf{x}. p(\mathbf{y}|\mathbf{x}) \neq 0$, which leads to a mixture form of the Bayesian decoder $p(\mathbf{x}|\mathbf{y})$.

⁴Several equivalent graphical representations are possible. The model shown on Figure 1.2 is chosen in order to illustrate the similarity with the noisy communication channel (Figure 1.1). Here we presumed that $p(\mathbf{y}|\mathbf{t}) = \langle p(\mathbf{y}|\mathbf{t}, \mathbf{e}) \rangle_{p(\mathbf{e})}$, where $\mathbf{e} \sim p(\mathbf{e})$ is the implied random noise.

1.4 Computational Problems of the Exact Formulation

Unfortunately, the choice of the theoretically optimal decoder (1.9) may complicate evaluation of the true mutual information $I(\mathbf{x}, \mathbf{y})$, as well as learning of the optimal encoder $p(\mathbf{y}|\mathbf{x})$. The reason for it is that in the expression of the mutual information (1.6), the complexity of computing the conditional entropy $H(\mathbf{x}|\mathbf{y})$ for the mixture distribution $p(\mathbf{x}|\mathbf{y})$ is in general exponential in $|\mathbf{y}|$. Indeed, from (1.6) and (1.9) we get

$$I(\mathbf{x}, \mathbf{y}) = \left\langle \log \sum_{m=1}^M p(\mathbf{y}|\mathbf{x}^{(m)}) \right\rangle_{p(\mathbf{y})} - \sum_{m=1}^M \langle \log p(\mathbf{y}|\mathbf{x}^{(m)}) \rangle_{p(\mathbf{y}|\mathbf{x}^{(m)})}, \quad (1.10)$$

which involves integration over the \mathbf{y} variables. For a model shown on Figure 1.2, evaluation of the second term in (1.10) is computationally tractable as long as $p(\mathbf{y}|\mathbf{x})$ has a simple structure (for example, if it is factorized in \mathbf{y}). However, evaluation of the first term in (1.10) (which is, in fact, the marginal entropy of the outputs $H(\mathbf{y})$) is in general problematic, since the marginal $p(\mathbf{y}) \propto \sum_{m=1}^M p(\mathbf{y}|\mathbf{x}^{(m)})$, being a mixture, is coupled in \mathbf{y} . Therefore, computation of the mutual information generally involves integration over an exponential number of states of the output vectors. Moreover, it is important to note that the representational complexity of storing the optimal decoders is also in general exponential in $|\mathbf{y}|$, which motivates a choice of parametric or structural constraints on decoder's specification.

1.4.1 Tractable Special Cases

The problem of optimizing (1.6) has been well explored for a number of tractable special cases. Specifically, it has been looked at in the context of maximizing the channel capacity for systems with very low-dimensional output spaces (Arimoto (1972), Linsker (1989b)). Indeed, for low-dimensional discrete outputs the effective number of states of \mathbf{y} may be reasonably small, so that the summations in (1.10) may be computed exactly. For the continuous case (e.g. for $\mathbf{y} \in \mathbb{R}^{|\mathbf{y}|}$), it may be impossible to express the mutual information analytically. However, due to the fact that \mathbf{y} is low-dimensional, the integrals in (1.10) may be efficiently approximated numerically (for example, by applying the mean-value theorem or using any of the known methods of numerical integration, see e.g. Korn and Korn (1968), Press et al. (1992), Dahlquist (2003)).

Additionally, computations may be tractable in the case when the encoder $p(\mathbf{y}|\mathbf{x}) \sim \delta(\mathbf{y} - \mathbf{g}(\mathbf{x}; \Theta))$ is deterministic and invertible (Bell and Sejnowski (1994), Nadal and Parga (1994)), so that the mutual information is easily expressed as

$$I(\mathbf{y}, \mathbf{x}) = c - \left\langle \log \frac{p(\mathbf{x})}{|\mathbf{J}_{\mathbf{x}}|} \right\rangle_{p(\mathbf{x})} = c + H(\mathbf{x}) + \langle \log |\mathbf{J}_{\mathbf{x}}| \rangle_{p(\mathbf{x})}. \quad (1.11)$$

Here c is a constant which does not affect the optimization surface for $p(\mathbf{y}|\mathbf{x})$, $\mathbf{J}_{\mathbf{x}} = \{\partial g_i(\mathbf{x})/\partial x_j\}$ is the Jacobian of the invertible mapping, and the average is

computed over the source vectors \mathbf{x} . It is easy to see that if $p(\mathbf{x})$ is the empirical distribution (1.8), the objective (1.11) is tractable. Moreover, for specific choices of the invertible mapping $g(\mathbf{x}; \Theta) : \mathbf{x} \rightarrow \mathbf{y}$, optimization of (1.11) with respect to the encoder parameters Θ leads to the *infomax* (Linsker, 1989a) formulation of ICA (Bell and Sejnowski, 1995).

Other tractable cases include linear Gaussian channels with the Gaussian distribution of the sources (Linsker, 1989a), or in fact any model where $p(\mathbf{x}, \mathbf{y})$ is jointly Gaussian. For this case, maximization of the mutual information reduces to learning the structured covariance matrix of the Gaussian decoder $p(\mathbf{x}|\mathbf{y})$. Generally, for all of the discussed special cases, it is tractable to compute $I(\mathbf{x}, \mathbf{y})$ and its gradients exactly. In order to find the optimal encoder, one may use iterative algorithms for maximizing the channel capacity (Arimoto (1972), Blahut (1972)), apply numerical optimization procedures (e.g. Bishop (1995), Luenberger (1998)), or in some cases find the optimal solutions analytically (Linsker (1989a), Linsker (1989b)).

1.5 Common Approximations of Mutual Information

We will now briefly describe several common methods which may be used to address the computational intractability of the exact formulation of the information-maximization problem. We focus specifically on Linsker’s *as-if* Gaussian approximation, as it will prove to be related to the variational approach to information maximization which we will introduce in Chapter 2. See Section 6.3.1 for a detailed discussion of Brunel and Nadal’s Fisher approximation criterion (Brunel and Nadal (1998)).

1.5.1 As-if Gaussian Approximation of $I(\mathbf{x}, \mathbf{y})$

In Section 1.4.1 we mentioned several choices of the encoding distribution leading to the exact computations of the mutual information. For more general encoding distributions, it may be necessary to consider approximations of $I(\mathbf{x}, \mathbf{y})$. Here we will discuss the method introduced by (Linsker (1992)), which utilizes the simple *as-if Gaussian* approximation of the mixture distribution. The key idea here is to approximate $I(\mathbf{x}, \mathbf{y})$ by assuming that the joint distribution $p(\mathbf{x})p(\mathbf{y}|\mathbf{x}) \approx p_G(\mathbf{x}, \mathbf{y}) \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is a Gaussian, independently of the exact form of $p(\mathbf{x}, \mathbf{y})$. Without loss of generality, we will presume that the data is centered, i.e. $\boldsymbol{\mu} = \mathbf{0}$.

Note that the conditional entropy in (1.6) may in this case be approximated by

$$H(\mathbf{x}|\mathbf{y}) \approx H_G(\mathbf{x}|\mathbf{y}) \stackrel{\text{def}}{=} -\langle \log p_G(\mathbf{x}|\mathbf{y}) \rangle_{p_G(\mathbf{x}, \mathbf{y})} = (1/2) \log(2\pi e)^{|\mathbf{x}|} |\boldsymbol{\Sigma}_{\mathbf{x}|\mathbf{y}}|, \quad (1.12)$$

where $\Sigma_{x|y}$ is the covariance of the decoder $p_G(\mathbf{x}|\mathbf{y})$ expressed from the joint Gaussian $p_G(\mathbf{x}, \mathbf{y})$. If the joint covariance is partitioned as

$$\Sigma \stackrel{\text{def}}{=} \langle [\mathbf{x} \ \mathbf{y}][\mathbf{x} \ \mathbf{y}]^T \rangle_{p(\mathbf{x}, \mathbf{y})} - \langle [\mathbf{x} \ \mathbf{y}] \rangle_{p(\mathbf{x}, \mathbf{y})} \langle [\mathbf{x} \ \mathbf{y}]^T \rangle_{p(\mathbf{x}, \mathbf{y})} \quad (1.13)$$

$$\stackrel{\text{def}}{=} \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{pmatrix} \in \mathbb{R}^{(|x|+|y|) \times (|x|+|y|)}, \quad (1.14)$$

we may express the conditional covariance as

$$\Sigma_{x|y} = \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx} \quad (1.15)$$

(see e.g. von Mises (1964)). Then a substitution into the expression for mutual information (1.6) leads to the objective function

$$2I_G(\mathbf{x}, \mathbf{y}) = \log |\Sigma_{xx}| - \log |\Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx}|. \quad (1.16)$$

For the assumed case of $\boldsymbol{\mu} = 0$, the objective (1.16) reduces to

$$2I_G(\mathbf{x}, \mathbf{y}) = \text{const} - \log |\langle \mathbf{x}\mathbf{x}^T \rangle - \langle \mathbf{x}\mathbf{y}^T \rangle \langle \mathbf{y}\mathbf{y}^T \rangle^{-1} \langle \mathbf{y}\mathbf{x}^T \rangle|, \quad (1.17)$$

where the averages are computed over the source and the channel distributions $p(\mathbf{x})$, $p(\mathbf{y}|\mathbf{x})$. Note that if $p(\mathbf{x})$ is the empirical distribution, the as-if Gaussian approximation (1.17) of the mutual information $I(\mathbf{x}, \mathbf{y})$ is a function of the encoder $p(\mathbf{y}|\mathbf{x})$ alone; optimization of (1.17) corresponds to learning the encoder's parameters.

It turns out that the as-if Gaussian objective $I_G(\mathbf{x}, \mathbf{y})$ corresponds to a proper variational lower bound on the true mutual information $I(\mathbf{x}, \mathbf{y})$ under the assumption of a Gaussian decoder (see Section 2.1). In other words, Linsker's approximation (Linsker (1992)) may be seen as a special case of a much more general variational procedure.

We will also show that in some simple cases, the encoder parameters obtained by maximizing (1.16) may be expressed analytically. Specifically, for linear Gaussian channels with the isotropic noise $p(\mathbf{y}|\mathbf{x}) \sim \mathcal{N}_y(\mathbf{W}\mathbf{x}, \sigma^2 \mathbf{I})$, where $\mathbf{W} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$ and $|\mathbf{y}| < |\mathbf{x}|$, the right singular vectors of the optimal weights correspond to rotations of the principle eigenvectors of the sample covariance $\mathbf{S} \stackrel{\text{def}}{=} \langle \mathbf{x}\mathbf{x}^T \rangle$ (see the discussion in Section 4.1.1). Equivalently, the result may be obtained in the noiseless limit of undercomplete linear projections $\mathbf{y} = \mathbf{W}\mathbf{x}$.

1.5.2 Other Approximations of $I(\mathbf{x}, \mathbf{y})$

Note that a general way to evaluate intractable averages over the hidden variables is by drawing independent samples from appropriate distributions, and using Monte-Carlo estimations of the objective criteria (see e.g. Neal (1993), Gamerman (1997)). The method is conceptually attractive, as conceptually similar techniques can be used to approximate arbitrary averages over the latent variables. However, in practice sampling techniques have several important drawbacks. First of all, sampling independent points from the equilibrium distribution may be extremely time-consuming, since the system may get trapped in a local

mode of the distribution (e.g. Hertz et al. (1991), Neal (1993)). This problem is partially addressed by a number of auxiliary sampling techniques, such as hybrid Monte Carlo for continuous spaces (e.g. Neal (1993)), or the Swendsen-Wang (Swendsen and Wang, 1987) and partial decoupling (Higdon (1998), Morris (1999)) algorithms for discrete variable spaces, which may help to improve convergence time and reduce the time gap between independent samples (see e.g. MacKay (1998) for an introductory discussion). A somewhat more fundamental problem is assessing convergence to the equilibrium (and therefore the quality of the resulting approximations). In the future, we may potentially consider applications of sampling methods in the information-theoretic context; however, due to the general difficulties of assessing the quality of the resulting estimates, we do not consider them in the suggested work. In what follows, we will only consider analytical approximations of $I(\mathbf{x}, \mathbf{y})$, and compare our results against other analytical approximations.

One of the relatively recent analytical approximations of mutual information is based on the idea of interpreting different dimensions $\{y_i\}$ of the codewords $\mathbf{y} \in \mathbb{R}^{|\mathbf{y}|}$ as independent samples from a distribution parameterized by the source vectors \mathbf{x} (see Brunel and Nadal (1998)). Specifically, if the encoding distribution $p(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^{|\mathbf{y}|} p_i(y_i|\mathbf{x})$ and $p_i(y_i = s|\mathbf{x}) \equiv p_j(y_j = s|\mathbf{x})$ for all $i, j = 1, \dots, |\mathbf{y}|$, then one may make a recourse to the results of statistical parameter estimation theory (e.g. Cramer (1946)) and approximate a lower bound on mutual information $I(\mathbf{x}, \mathbf{y})$ by applying the data-processing and the Cramer-Rao (Cover and Thomas (1991)) inequalities. While a direct application of the inequalities leads to a computationally intractable bound on $I(\mathbf{x}, \mathbf{y})$, the bound may in some cases be efficiently approximated by considering further numerical relaxations (Brunel and Nadal (1998)). Generally, the suggested approximations of the lower bound on $I(\mathbf{x}, \mathbf{y})$ are accurate under the asymptotic assumption of infinitely large code spaces ($|\mathbf{y}| \rightarrow \infty$), but they may lead to unstable behavior of the resulting learning algorithms for $|\mathbf{x}| < |\mathbf{y}|$. The identical results were also obtained by Kang and Sompolinsky (2001), who considered different numerical approximations applied in the same asymptotic limit $|\mathbf{y}| \rightarrow \infty$.

Another approximation of $I(\mathbf{x}, \mathbf{y})$ may be obtained by a simple re-formulation of the local approximations of Szummer and Jaakkola (2002) and Corduneanu and Jaakkola (2003), who propose to approximate the mutual information numerically for each infinitely small symmetric region in the source space. By applying the chain rule on mutual information, it is straight-forward to show that their criterion may be used to approximate a lower bound on $I(\mathbf{x}, \mathbf{y})$ (see discussion in Section 6.3.2). A number of other somewhat heuristic objective criteria is often used to approximate $I(\mathbf{x}, \mathbf{y})$ (see e.g. Principe et al. (2000), Torkkola and Campbell (2000), Gokcay and Principe (2002)), though the resulting approximations may be rather difficult to justify for maximizing the mutual information (Torkkola (2000)). Other work focuses specifically on noiseless overcomplete communication channels (where $|\mathbf{y}| > |\mathbf{x}|$ and $p(\mathbf{y}|\mathbf{x}) \sim \delta(\mathbf{y} - \mathbf{f}(\mathbf{x}))$ (see Shriki et al. (2002))) and cannot be easily applied in the case of stochastic encoding mappings.

1.6 Thesis Overview

In the following chapters we will discuss a principled and computationally tractable approach to maximization of the information transfer in large-scale stochastic models. Particularly, we will be interested in the family of variational approaches, where the goal may informally be formulated as transforming optimization of a generally intractable functional to optimization of its tractable bound. The tractability of the variational methods is typically achieved by decoupling the degrees of freedom of the original functional at the cost of introducing additional variational parameters. Variational methods have been long described in extensive mathematics and physics literature (see e.g. Gelfand and Fomin (1963), Ewig (1985), Fox (1987), Riley et al. (2002)), and applied to graphical modeling in the context of bounding the partition functions (see e.g. Jaakkola (1997), Jordan et al. (1998), Wainwright et al. (2002), Wainwright (2002)). We will be mainly interested in deriving tractable bounds on the mutual information.

In Chapter 2 we will discuss a simple and general variational approach to maximizing a proper (generic) lower bound on $I(\mathbf{x}, \mathbf{y})$ in a noisy channel, and show that in the special case when the variational distribution is unconstrained, the method gives rise to the family of Blahut-Arimoto type algorithms (Arimoto (1972), Blahut (1972)). Generally, however, constraining the variational distributions is important in order to ensure that the method remains computationally tractable. We will also show that optimization of Linsker’s *as-if* Gaussian objective criterion (Linsker (1992)) corresponds to a specific way of optimizing the variational lower bound on $I(\mathbf{x}, \mathbf{y})$ for a specific choice of the variational decoder distribution. Finally, we will introduce an *auxiliary variational* approach to approximate information maximization, which formally generalizes on the simple generic bound without altering properties of the original channel.

In Chapter 3 we will explore general relations of the variational information-maximizing algorithm to maximum likelihood learning in generative models and conditional likelihood learning in chains. We will outline general differences between encoder models of communication channels and generative models. We will also show that the likelihood of a generative model may be viewed as a lower bound on $I(\mathbf{x}, \mathbf{y})$ for the corresponding model of the stochastic channel, where the encoding distribution is the exact posterior of the generative model. A practical side-effect of this study is an information-theoretic objective for training generative models. We will also demonstrate a close relation between optimizing the generic variational bound on $I(\mathbf{x}, \mathbf{y})$ and conditional likelihood training in stochastic autoencoders. Additionally, we will show that common approaches to training noiseless autoencoders maximize proper lower bounds on $I(\mathbf{x}, \mathbf{y})$ (under the assumption of noiseless encoding mappings).

In Chapter 4 we will consider an application of the variational information maximizing approach to noisy constrained dimensionality reduction. We will prove a simple result that for constrained linear Gaussian channels, optimization of Linsker’s bound cannot improve on the PCA projections. On the other hand, the richer family of the auxiliary variational lower bounds on $I(\mathbf{x}, \mathbf{y})$ leads to significant improvements over Linsker’s (PCA) bounds. This result is potentially

interesting from the communication-theoretic perspective, as it demonstrates a simple and computationally efficient way to produce tighter bounds on the capacity of a communication channel without altering its properties (e.g. without communicating more data across the channels). Additionally, we will discuss a simple information-theoretic approach to constrained dimensionality reduction for hybrid channels $\mathbf{x} \rightarrow \{\mathbf{y}, z\}$ (where $\mathbf{x} \in \mathbb{R}^{|\mathbf{x}|}$, $\mathbf{y} \in \mathbb{R}^{|\mathbf{y}|}$, and $z \in \{1, \dots, |z|\}$), which may significantly improve reconstructions of the sources $\{\mathbf{x}\}$ from their lower-dimensional representations $\{\mathbf{y}\}$ at a small increase in the transmission cost (given by $|z|$). We will also point out a curious link between maximizing $I(\mathbf{x}, \{\mathbf{y}, z\})$ in a hybrid channel and maximizing the likelihood for a mixture of constrained Factor Analysis-type models with the uniform (rather than the spherical) distribution of the factors.

Chapter 5 demonstrates applications of the information-maximizing framework for the case of nonlinear encoder models. Specifically, we will discuss several ways of applying the framework to information-theoretic clustering. We will empirically demonstrate that the resulting information-theoretic clustering approaches favorably compare with the conventional clustering techniques. Moreover, we will show that the information-maximizing framework may be used to learn kernel functions, which may indeed be of a practical benefit for visualizing the underlying structure of the data. Moreover, we will review some of the theoretical properties of the variational information-maximizing framework for nonlinear Gaussian encoding distributions; for example, we will show that the Gaussian Process Latent Variable Model (Lawrence (2003)) arises as a special case of our information-theoretic formulation.

Chapter 6 explores applicability of the variational information-maximizing framework in the context of learning high-dimensional binary representations of continuous source patterns. We will consider the situation when the binary encodings are conditionally independent, and compare our variational approach with Brunel and Nadal’s Fisher approximation of mutual information (Brunel and Nadal (1998)). Moreover, we will use the results of Szummer and Jaakkola (2002) and Corduneanu and Jaakkola (2003) to derive another straight-forward approximation of $I(\mathbf{x}, \mathbf{y})$. Our empirical results indicate that for the considered case, the variational approach is most preferable, while both the *local* approximation of $I(\mathbf{x}, \mathbf{y})$ (based on the work of Corduneanu and Jaakkola (2003)) and the generic variational approach significantly outperform the common Fisher approximation for undercomplete projections (the compression paradigm). Additionally, we will demonstrate that for a considered encoding distribution it is possible to derive a *local* learning rule, which may potentially be attractive from the neuro-biological perspective. This extends the results of Linsker (1997), who has previously showed the existence of local approximations of nonlinear optimization procedures on $I(\mathbf{x}, \mathbf{y})$ for invertible encoding mappings $\mathbf{x} \mapsto \mathbf{y}$ (*cf* Nadal and Parga (1994), Bell and Sejnowski (1995)).

In Chapter 7 we will change the perspective and consider a seemingly unrelated problem of lower-bounding the normalizing constant (log partition function) of a probability distribution. Specifically, we will introduce an auxiliary variable extension of any structured mean field theory, which can be useful in the context

of approximate probabilistic inference. While the method described there is of a potential interest as a general approach to approximate inference, it demonstrates a curious link to our variational information-maximizing framework. In particular, we will show that the improvement of the proposed bound on $\log Z$ over a convex combination of simpler bounds given by the standard theories is defined by a specific form of the generic lower bound on mutual information. The variational information-maximizing framework may therefore be seen as addressing an integral subgoal of variational inference. We will also show that the existing variational mixture methods (see Jaakkola and Jordan (1998), Lawrence et al. (1998)) may be viewed as specific numerically difficult instances of the auxiliary variational approach.

Finally, Chapter 8 summarizes the results and outlines directions for extending the presented work.

Chapter 2

Variational Information Maximization

Here we describe a family of variational lower bounds on mutual information $I(\mathbf{x}, \mathbf{y})$, which gives rise to a formal and theoretically justified approach to information maximization in noisy channels. In Section 2.1 we describe the variational Information Maximization (IM) algorithm, reminiscent of the generalized variational EM algorithm for likelihood training, and demonstrate that it provides a simple and effective tool for learning encoders and variational decoders in a principled manner. We show that in the simplest case when the simplest form of the bound is used and the decoder is unconstrained, the proposed variational optimization procedure reduces to a generally intractable form of the Arimoto-Blahut (Arimoto (1972), Blahut (1972)) algorithm for maximization of channel capacity. However, by choosing appropriate parametric constraints on the encoder-decoder pair, we may avoid intractabilities in a principled manner. The resulting algorithm would be optimizing a simple *generic* lower bound on mutual information, subject to the specific constraints on the variational decoder.

Then we outline general properties of the bound. In Section 2.2.1 we describe effects of choosing specific sparse decoder structures on the generic lower bound on $I(\mathbf{x}, \mathbf{y})$. Not surprisingly, we show that by considering richer decoder structures, we may indeed obtain tighter bounds on the intractable measure of information content. We also discuss a simple relation of the variational bound to the common *as-if Gaussian* approximation of the mutual information (Linsker, 1992). As we show in Section 2.2.2, independently of the choice of the encoder distribution, Linsker’s approach corresponds to a special case of our formulation, where the variational decoder is constrained to be a linear Gaussian.

Finally, in Section 2.3 we introduce a richer family of *auxiliary variational* lower bounds, which generalize on the simpler generic bounds on $I(\mathbf{x}, \mathbf{y})$. The key idea there is to introduce additional variables, which may be used for capturing useful features of the source patterns, and for introducing global dependencies to the decoded sources. Importantly, we show that the projections to the auxiliary space may be defined in a way which does not alter the original channel. By constraining the mappings to the auxiliary space to be in a tractable family and by imposing appropriate constraints on the variational decoders, the resulting bound

is a formal tractable generalization of the simpler generic approach described in Section 2.1.

2.1 Generic Lower Bound on Mutual Information

Since the exact evaluation and optimization of the information transfer is in general computationally intractable, our central aim here will be to maximize a lower bound on the mutual information in a tractable way. Using the formulation $I(\mathbf{x}, \mathbf{y}) = H(\mathbf{x}) - H(\mathbf{x}|\mathbf{y})$, we are interested in optimising $I(\mathbf{x}, \mathbf{y})$ with respect to the encoder $p(\mathbf{y}|\mathbf{x})$. In many cases which we consider $p(\mathbf{x})$ is simply the empirical distribution (1.8). In principle, this may be generalized to any case where averaging over the sources $p(\mathbf{x})$ is tractable. Nevertheless, since the distribution of the source patterns is not a function of the channel parameters, it has no effect on the optimization surface for $p(\mathbf{y}|\mathbf{x})$. Thus, to express the objective function we need to bound the intractable entropic term $H(\mathbf{x}|\mathbf{y})$ suitably.

2.1.1 Definition of a Simple Lower Bound on Mutual Information

In order to derive a simple lower bound on mutual information, we consider the Kullback-Leibler divergence $KL(p(\mathbf{x}|\mathbf{y})||q(\mathbf{x}|\mathbf{y}))$ between the posterior $p(\mathbf{x}|\mathbf{y})$ and its variational approximation $q(\mathbf{x}|\mathbf{y})$. Non-negativity of the divergence (e.g. Cover and Thomas (1991)) implies

$$\langle \log p(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{x}|\mathbf{y})} - \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{x}|\mathbf{y})} \geq 0 \Rightarrow \underbrace{\langle \log p(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{x}|\mathbf{y})p(\mathbf{y})}}_{-H(\mathbf{x}|\mathbf{y})} \geq \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{x}|\mathbf{y})p(\mathbf{y})}. \quad (2.1)$$

This leads to

$$I(\mathbf{x}, \mathbf{y}) \geq H(\mathbf{x}) + \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{x}, \mathbf{y})} \stackrel{\text{def}}{=} \tilde{I}(\mathbf{x}, \mathbf{y}), \quad (2.2)$$

where $q(\mathbf{x}|\mathbf{y})$ is an arbitrary distribution, which saturates the bound for $q(\mathbf{x}|\mathbf{y}) \equiv p(\mathbf{x}|\mathbf{y})$. This agrees with the intuition that for a given channel $p(\mathbf{y}|\mathbf{x})$, the optimal decoder should correspond to the Bayesian posterior $p(\mathbf{x}|\mathbf{y})$, though for this case the computation of the mutual information (and its derivatives) may be intractable.

The bound (2.2) explicitly includes both the encoder $p(\mathbf{y}|\mathbf{x})$ and decoder $q(\mathbf{x}|\mathbf{y})$ and has a form similar to the criterion optimized by the Blahut-Arimoto algorithms for channel capacity (Blahut (1972), Arimoto (1972)). However, by analogy with a related variational extension (Neal and Hinton (1998)) of the classic expectation maximization algorithm (Dempster et al. (1977)), it is practical to extend the conventional formulation and *constrain* the decoder $q(\mathbf{x}|\mathbf{y})$ to lie in a tractable family (Barber and Agakov (2003)). These constraints are important, as they help to avoid computational intractability of the theoretically optimal, but practically infeasible decoders derived in the unconstrained formulation (see e.g. Cover and Thomas (1991)). The bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ is then optimized for both the

encoder $p(\mathbf{y}|\mathbf{x})$ and the variational distribution $q(\mathbf{x}|\mathbf{y})$ (which is generally different from the posterior $p(\mathbf{x}|\mathbf{y})$), subject to the imposed distribution and tractability constraints. The variational distribution $q(\mathbf{x}|\mathbf{y})$, which we use in the specification of the objective function, may be simply thought of as an auxiliary entity to facilitate computation and optimization of the otherwise intractable mutual information $I(\mathbf{x}, \mathbf{y})$. Importantly, our principal goal of maximizing the bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ is, as in the exact formulation, to learn the optimal adjustable parameters of the *encoder* $p(\mathbf{y}|\mathbf{x})$.

Other (for example, mean-field type) relaxations of the mutual information may potentially be considered (see e.g. Jaakkola and Jordan (1998)). However, our current experience suggests that for certain choices of the decoder $q(\mathbf{x}|\mathbf{y})$, the variational bound (2.2) considered above is particularly computationally convenient. We will now outline a procedure for optimizing the bound and discuss some of its fundamental properties.

2.1.2 The Variational IM Algorithm

Let \mathcal{P} and \mathcal{Q} denote families of the encoders $p(\mathbf{y}|\mathbf{x})$ and variational decoders $q(\mathbf{x}|\mathbf{y})$ respectively. To assure applicability of the learning for the case of large scale models, the distribution families must be chosen in such a way that the bound (2.2) is tractable. This may be achieved, for example, by imposing appropriate parametric or structural constraints on \mathcal{Q} in such a way that the averaging over the codewords \mathbf{y} is practically feasible. In our approximate approach to information maximization, we follow the standard iterative variational procedure and maximize $\tilde{I}(\mathbf{x}, \mathbf{y})$ with respect to $p(\mathbf{y}|\mathbf{x}) \in \mathcal{P}$ and $q(\mathbf{x}|\mathbf{y}) \in \mathcal{Q}$. A simple recursive optimization procedure performing information maximization (**IM algorithm**) is then given as follows:

1. For a fixed $q_{X|Y}^{(t)} \in \mathcal{Q}$, find $p_{Y|X}^{(t+1)} = \arg \max_{p_{Y|X} \in \mathcal{P}} \tilde{I}(\mathbf{x}, \mathbf{y}; p_{Y|X}, q_{X|Y}^{(t)})$ [**M-step**];
2. For a fixed $p_{X|Y}^{(t+1)} \in \mathcal{P}$, find $q_{X|Y}^{(t+1)} = \arg \max_{q_{X|Y} \in \mathcal{Q}} \tilde{I}(\mathbf{x}, \mathbf{y}; p_{Y|X}^{(t+1)}, q_{X|Y})$ [**I-step**];
3. Iterate **M-** and **I-** steps until a convergence criterion is met

where we used t for the iteration number and $p_{Y|X}$, $q_{X|Y}$ for the conditional encoder $p(\mathbf{y}|\mathbf{x})$ and decoder $q(\mathbf{x}|\mathbf{y})$ respectively, and implied the distribution constraints on $p_{Y|X}$ and $q_{X|Y}$.

Proposition 2.1. *The variational IM algorithm is guaranteed to maximize or leave unchanged a lower bound on the mutual information.*

Proof. The proof is straight-forward and follows from (2.2) and the algorithm specification. Let t be the iteration number. From the M-step, it is clear that

$$\langle \log q^{(t)}(\mathbf{x}|\mathbf{y}) \rangle_{p^{(t+1)}(\mathbf{y}|\mathbf{x})p(\mathbf{x})} \geq \langle \log q^{(t)}(\mathbf{x}|\mathbf{y}) \rangle_{p^{(t)}(\mathbf{y}|\mathbf{x})p(\mathbf{x})}. \quad (2.3)$$

Analogously, the I-step leads to

$$\langle \log q^{(t+1)}(\mathbf{x}|\mathbf{y}) \rangle_{p^{(t+1)}(\mathbf{y}|\mathbf{x})p(\mathbf{x})} \geq \langle \log q^{(t)}(\mathbf{x}|\mathbf{y}) \rangle_{p^{(t+1)}(\mathbf{y}|\mathbf{x})p(\mathbf{x})}. \quad (2.4)$$

Then from the transitivity we get

$$\tilde{I}(\mathbf{x}, \mathbf{y} | p_{Y|X}^{(t+1)}, q_{X|Y}^{(t+1)}) \geq \tilde{I}(\mathbf{x}, \mathbf{y} | p_{Y|X}^{(t)}, q_{X|Y}^{(t)}), \quad (2.5)$$

where $\tilde{I}(\mathbf{x}, \mathbf{y})$ is defined as in (2.2). Therefore, the IM algorithm is guaranteed to maximize (or leave unchanged) the *lower bound* on the true mutual information given by $\tilde{I}(\mathbf{x}, \mathbf{y})$. \square

In general, applicability of proposition 2.1 depends on the specifics of optimization methods used during the *I*- and *M*- steps. For example, if optimization is performed with respect to parameters of $p(\mathbf{y}|\mathbf{x})$ and $q(\mathbf{x}|\mathbf{y})$ by a numerical ascent in the parameter space, one may encounter the usual problem of non-monotonic convergence (e.g. for inappropriate learning rates, etc.) Practically, the proposition holds if it is possible to express the IM steps in terms of closed-form fixed-point updates (i.e. at each iteration the optima are expressed analytically), or if the steps apply numerical optimization procedures which always lead to monotonic changes in the objective functions (Bishop (1995), Galeev and Tihomirov (2000)).

2.1.2.1 Relation to the Blahut-Arimoto algorithm for maximizing channel capacity

In the cases when the exact mutual information is tractable to compute, one may apply any of the known numerical optimization techniques (e.g. Galeev and Tihomirov (2000)) to optimize $I(\mathbf{x}, \mathbf{y})$ directly. Alternatively, one may optimize $I(\mathbf{x}, \mathbf{y})$ by applying an iterative procedure, known as the Blahut-Arimoto algorithm for maximizing the channel capacity (Arimoto (1972), Blahut (1972)). For a fixed distribution of the source variables $p_X(\mathbf{x})$, the algorithm is given as

1. Find $p_{Y|X}^{(t+1)} = \arg \max_{p_{Y|X} \in \mathcal{P}} \langle \log p_{X|Y}^{(t)}(\mathbf{x}|\mathbf{y}) \rangle_{p_{Y|X}(\mathbf{y}|\mathbf{x})p_X(\mathbf{x})}$;
2. Iterate until convergence.

Here it is assumed that $p_{X|Y}^{(t)}(\mathbf{x}|\mathbf{y}) \propto p_{Y|X}^{(t)}(\mathbf{y}|\mathbf{x})p_X(\mathbf{x})$ is the exact Bayes-optimal decoder expressed from the channel distribution at the t^{th} iteration of the training algorithm.

It is easy to see that the Blahut-Arimoto algorithm corresponds to a special case of the variational IM algorithm described above. Indeed, if the decoder is unconstrained, i.e. if \mathcal{Q} defines a family of *all* possible conditional probability density functions $q(\mathbf{x}|\mathbf{y})$ for $\mathbf{x} \in \mathcal{R}_x$ and $\mathbf{y} \in \mathcal{R}_y$, then the information-maximization algorithm reduces to a form of the Blahut-Arimoto algorithm for learning channel distributions. The proof is straight-forward, and follows from the convexity of $\tilde{I}(\mathbf{x}, \mathbf{y})$ in $q(\mathbf{x}|\mathbf{y})$ and the fact that the bound is saturated for $q(\mathbf{x}|\mathbf{y}) \equiv p(\mathbf{x}|\mathbf{y})$

(see expression (2.2)). Clearly, a reduction of the IM algorithm to the Blahut-Arimoto procedure also occurs when $q(\mathbf{x}|\mathbf{y})$ is constrained to be identical to the true posterior $p(\mathbf{x}|\mathbf{y})$ at each iteration of learning, i.e.

$$q^{(t)}(\mathbf{x}|\mathbf{y}) \equiv p^{(t)}(\mathbf{x}|\mathbf{y}). \quad (2.6)$$

In both of the described special cases, the iterative optimization procedure is guaranteed to maximize or leave unchanged the *exact* value of the mutual information. Indeed, a straight-forward substitution of (2.6) into (2.3) and (2.4) implies

$$\langle \log p^{(t+1)}(\mathbf{x}|\mathbf{y}) \rangle_{p^{(t+1)}(\mathbf{y}|\mathbf{x})p(\mathbf{x})} \geq \langle \log p^{(t)}(\mathbf{x}|\mathbf{y}) \rangle_{p^{(t)}(\mathbf{y}|\mathbf{x})p(\mathbf{x})}, \quad (2.7)$$

i.e. $I^{(t+1)}(\mathbf{x}, \mathbf{y}) \geq I^{(t)}(\mathbf{x}, \mathbf{y})$. An alternative proof shows that the Blahut-Arimoto algorithm does not decrease the true mutual information, and is given here for completeness.

Proposition 2.2. *The Blahut-Arimoto algorithm is guaranteed to maximize or leave unchanged the exact value of mutual information.*

Proof. The iteration step of the algorithm implies

$$\langle \log p^{(t)}(\mathbf{x}|\mathbf{y}) \rangle_{p^{(t+1)}(\mathbf{y}|\mathbf{x})p(\mathbf{x})} \geq \langle \log p^{(t)}(\mathbf{x}|\mathbf{y}) \rangle_{p^{(t)}(\mathbf{y}|\mathbf{x})p(\mathbf{x})}. \quad (2.8)$$

Then, from non-negativity of $\langle KL(p^{(t+1)}(\mathbf{x}|\mathbf{y})||p^{(t)}(\mathbf{x}|\mathbf{y})) \rangle_{p(\mathbf{x})}$, we get

$$\langle \log p^{(t+1)}(\mathbf{x}|\mathbf{y}) \rangle_{p^{(t+1)}(\mathbf{y}|\mathbf{x})p(\mathbf{x})} \geq \langle \log p^{(t)}(\mathbf{x}|\mathbf{y}) \rangle_{p^{(t+1)}(\mathbf{y}|\mathbf{x})p(\mathbf{x})}. \quad (2.9)$$

As before, for a fixed $H(\mathbf{x})$, a combination of (2.8) with (2.9) leads to $I(\mathbf{x}, \mathbf{y}|p_{Y|X}^{(t+1)}) \geq I(\mathbf{x}, \mathbf{y}|p_{Y|X}^{(t)})$. \square

It is known empirically that for convex objective functions, iterative optimizers may be slow to converge (see e.g. the empirical results (Minka (2003)) on convergence of the *iterative scaling* algorithms (Darroch and Ratcliff (1972), Berger et al. (1996), Collins et al. (2002)) for log-linear models). If $I(\mathbf{x}, \mathbf{y})$ is convex in the encoder's parameters, we may hypothesize that alternative optimization techniques may outperform the Blahut-Arimoto algorithm in terms of the convergence speed; moreover, computational tractability of the Blahut-Arimoto algorithm would typically imply tractability of alternative optimization methods. Generally, it is clear that in tractable channels one should have a flexibility in choosing a specific optimization procedure.

Importantly, we emphasize once again that for many interesting channels, the Blahut-Arimoto algorithm (and other numerical procedures optimizing $I(\mathbf{x}, \mathbf{y})$ directly) could be computationally difficult to apply, as it would require optimization of the generally intractable entropy of a mixture in the optimization steps. The proposed IM algorithm (and its numerical variations) address the problems of computational intractability by introducing an additional functional parameter – the variational decoder $q(\mathbf{x}|\mathbf{y})$, which is constrained to lie in a tractable family. Then our goal is to maximize a proper lower bound on the intractable objective $I(\mathbf{x}, \mathbf{y})$ with respect to the variational decoder and the adjustable parameters of the channel distribution, subject to the imposed constraints.

2.1.3 Reconstruction of the Source Patterns

As discussed in Section 1.3, by maximizing the exact mutual information for $p(\mathbf{y}|\mathbf{x})$, we try to establish an optimal way of *encoding* the sources for transmitting them across the channel. Our goal remains unchanged in the approximate variational formulation (though the optimization surface for $p(\mathbf{y}|\mathbf{x})$ will now be affected by a specific choice of the variational distribution). Importantly, the bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ does not explicitly specify which decoder should be used at the receiver's end. However, just as in variational approaches to likelihood maximization, it provides an option of reconstructing the source patterns with the learned variational posterior $q(\mathbf{x}|\mathbf{y})$. Here we briefly discuss other possible ways to perform the reconstructions.

2.1.3.1 Reconstructions using the exact posterior

Once the optimal encoder is learned, the receiver may use it to compute the Bayes-optimal decoder $p(\mathbf{x}|\mathbf{y})$, and apply it to reconstructing the original sources $\{\mathbf{x}\}$ for the received representations $\{\mathbf{y}\}$. From (2.2) it is easy to see that using any other decoder (for a fixed learned channel) would indeed be \tilde{I} -suboptimal. Unfortunately, using the exact posterior $p(\mathbf{x}|\mathbf{y})$ for decoding may not be plausible for many practical cases. For example, if $p(\mathbf{x})$ is the empirical distribution, computation of $p(\mathbf{x}|\mathbf{y})$ makes a recourse to the training data $\{\mathbf{x}^{(i)}|i = 1, \dots, M\}$ (see (1.9)). In practice, this means that in order to perform the reconstructions, the receiver needs to have an access to the empirical distribution $p(\mathbf{x})$ used during training, which requires storing the entire training set. Moreover, even if the training set is available at the receiver's end of the channel, using $p(\mathbf{x}|\mathbf{y})$ for decoding would constrain the reconstructions to lie in the training set, which in some cases may be restrictive (effectively, this case would correspond to extracting training patterns $\{\mathbf{x}^{(m)}|m = 1, \dots, M\}$ from their noisy *encoded* representations $\{\mathbf{y}\}$).

Finally, for many interesting practical problems, computation of $p(\mathbf{x}|\mathbf{y})$ may be intractable *despite* the tractability of averaging in (2.2). One practical case when this may happen is the problem of syndrome decoding in binary symmetric channels (see e.g. McEliece (1977)). Note that if all the source vectors are binary and equiprobable, i.e. $\forall \mathbf{x} \in \{0, 1\}^{|\mathbf{x}|}$. $p(\mathbf{x}) = 1/2^{|\mathbf{x}|}$, computation of the exact posterior

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{x})}{\sum_{m=1}^{2^{|\mathbf{x}|}} p(\mathbf{y}|\mathbf{x}^{(m)})} \quad (2.10)$$

involves a summation over an exponential number of states of the source variables $\mathbf{x} \in \{0, 1\}^{|\mathbf{x}|}$, which in most cases of interest is prohibitively expensive. Depending on the structure of the graph, the exact posteriors may in some cases be accurately approximated (Gallager (1963), MacKay (1999a)) by iterative algorithms (e.g. Gallager (1963), Pearl (1988), Murphy et al. (1999)). However, for bipartite graphs, the accuracy of approximations will typically drop with an increase in the graph's connectivity (Burshtein and Miller (2002)), and convergence will generally be a problem.

2.1.3.2 Reconstructions using the approximate posterior

It is clear that in the cases when it is impossible to compute the Bayes-optimal decoder $p(\mathbf{x}|\mathbf{y})$, one needs to consider using alternative decoding schemes. For example, one may need to approximate the exact posteriors $p(\mathbf{x}|\mathbf{y})$, e.g. by using (structured) mean field models (Jaakkola (1997), Barber and Wiergerinck (1998), Jordan et al. (1998), Saad and Opper (2001)), or applying iterative algorithms for approximating the posteriors (Pearl (1988), Yedidia et al. (2000a), Yedidia et al. (2000b), Weiss and Freeman (2001)). Generally, these methods would require approximations of a new posterior for each new received vector \mathbf{y} , which may be plausible, but expensive to compute in practice. In the case of applying iterative message-passing algorithms, it may be difficult to ensure the convergence (see Ihler et al. (2005), Mooij and Kappen (2005) for the discussion of sufficient conditions), while approximations yielded by the provenly convergent algorithms (Heskes (2002), Yuille (2002)) may not necessarily be accurate estimates of the exact posteriors. On the strong side, in many practical situations when the propagation algorithms do converge, they often lead to accurate approximations of the posteriors (MacKay and Neal (1999), Yedidia et al. (2000a)), as unlike most of the variational approximations they do not impose explicit structural or parametric constraints on the posteriors.

An attractive feature of the bound (2.2) is availability of the variational decoder $q(\mathbf{x}|\mathbf{y})$, which is learned along with the optimal parameters of the encoder $p(\mathbf{y}|\mathbf{x})$. A quick and simple way to reconstruct the transmitted source vectors would be to use $q(\mathbf{x}|\mathbf{y})$ with the optimized variational parameters. This application of the variational distribution $q(\mathbf{x}|\mathbf{y})$ is analogous to using the variational posterior for inference in generative models after training with the variational EM (Neal and Hinton (1998)) algorithm. The principal advantage of variational decoding is the simplicity of the resulting inference procedure. Clearly, subject to the constraints on the variational decoder, the resulting variational reconstructions are optimal, since optimally the objective function (2.2) is maximized when $q(\mathbf{x}|\mathbf{y}) \equiv p(\mathbf{x}|\mathbf{y})$. Moreover, since the bound (2.2) is based on the KL divergence between the true and the approximating posteriors, optimization for the variational decoder is equivalent to a moment matching approximation of $p(\mathbf{x}|\mathbf{y})$ by $q(\mathbf{x}|\mathbf{y})$, averaged over $p(\mathbf{y})$. This fact may potentially be beneficial in terms of decoding, since the more successful decoding algorithms tend to approximate the mean of the posterior $p(\mathbf{x}|\mathbf{y})$ (Saad and Opper, 2001), whilst standard mode matching approaches (such as mean-field theory) typically get trapped in the one of many sub-optimal modes. We stress, however, that using the variational posterior for inference is just one of the possible ways to reconstruct the sources.

2.1.4 Posterior Approximations

There is an interesting relationship between maximizing the bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ in a noisy channel, and computing an optimal estimate of an intractable posterior in a generative graphical model with observations \mathbf{y} and hidden variables \mathbf{x} .

One of the goals of inference in graphical models is to approximate moments of the generally intractable model-specific posterior $p(\mathbf{x}|\mathbf{y})$, where \mathbf{x} is a vector of

hidden variables, and \mathbf{y} is a vector of observations. In general, this computation is intractable, and approximations are required. A standard mean field approach approximates the posterior marginal by minimizing the KL divergence:

$$KL(q(\mathbf{x}|\mathbf{y})||p(\mathbf{x}|\mathbf{y})) = \sum_{\mathbf{x}} \{q(\mathbf{x}|\mathbf{y}) \log q(\mathbf{x}|\mathbf{y}) - q(\mathbf{x}|\mathbf{y}) \log p(\mathbf{x}|\mathbf{y})\} \quad (2.11)$$

where $q(\mathbf{x}|\mathbf{y}) = \prod_i q(x_i|\mathbf{y})$. In this case, the KL divergence and its functional gradients with respect to $q(x_i|\mathbf{y})$ are usually tractably computable (up to a neglectable prefactor). Assuming that $q(\mathbf{x}_{\setminus i}|\mathbf{y}) \stackrel{\text{def}}{=} \prod_{j \neq i} q(x_j|\mathbf{y})$ is fixed, we may optimize $KL(q(\mathbf{x}|\mathbf{y})||p(\mathbf{x}|\mathbf{y}))$ for $q(x_i|\mathbf{y})$, and express the variational posterior marginal analytically (from the convexity of the KL). Typically, this results in $q(x_i|\mathbf{y})$ approximating any one of a very large number of local modes of the model-specific posterior $p(x_i|\mathbf{y})$. Clearly, if the goal is to approximate moments of the exact posterior, using mean field decoding which optimizes (2.11) is generally suboptimal.

Alternatively, we may consider

$$\begin{aligned} KL(p(\mathbf{x}|\mathbf{y})||q(\mathbf{x}|\mathbf{y})) &= \sum_{\mathbf{x}} (p(\mathbf{x}|\mathbf{y}) \log p(\mathbf{x}|\mathbf{y}) - p(\mathbf{x}|\mathbf{y}) \log q(\mathbf{x}|\mathbf{y})) \\ &= - \sum_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}) \log q(\mathbf{x}|\mathbf{y}) + c, \end{aligned} \quad (2.12)$$

where c is an irrelevant constant (the inverse entropy of $p(\mathbf{x}|\mathbf{y})$, which is not a function of $q(\mathbf{x}|\mathbf{y})$). This is the correct KL divergence in the sense that, optimally, $q(x_i|\mathbf{y}) = p(x_i|\mathbf{y})$, i.e. the posterior mean is correctly calculated. The difficulty lies in performing averages with respect to the exact posterior $p(\mathbf{x}|\mathbf{y})$, which are generally intractable (since the posterior is fully coupled in \mathbf{x}). However, it may be possible to minimize the *average* divergence $\langle KL(p(\mathbf{x}|\mathbf{y})||q(\mathbf{x}|\mathbf{y})) \rangle_{p(\mathbf{y})}$, which leads to maximization of the generic bound (2.2) on the mutual information

$$\sum_{\mathbf{y}} \sum_{\mathbf{x}} p(\mathbf{y}) p(\mathbf{x}|\mathbf{y}) \log q(\mathbf{x}|\mathbf{y}) = \sum_{\mathbf{y}} \sum_{\mathbf{x}} p(\mathbf{x}) p(\mathbf{y}|\mathbf{x}) \log q(\mathbf{x}|\mathbf{y}), \quad (2.13)$$

where we ignored the irrelevant entropy of the source vectors. While for any given \mathbf{y} , the best posterior mean estimate may be difficult to compute, we may apply the bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ to calculate the best posterior mean estimate *on average*.

2.2 Tractable Choices of Variational Decoders

In Section 2.1 we discussed a simple variational approach to maximization of mutual information in noisy channels, where the problem of optimizing the computationally intractable objective $I(\mathbf{x}, \mathbf{y})$ was transformed to optimization of a proper lower bound on the objective criterion. It is intuitive that the tightness and the tractability of the described lower bound (2.2) depends on a specific choice of the variational distribution $q(\mathbf{x}|\mathbf{y})$. Here we discuss several of such choices, namely structured and Gaussian variational posteriors.

First, we will show that by considering structured variational posteriors, we may indeed improve on simple factorized variational approximations. Then we

re-formulate the bound (2.2) for the special case when the variational decoder is a linear Gaussian, and show that a specific way of optimizing the resulting objective reduces to maximization of Linsker’s *as-if Gaussian* criterion (1.16).

2.2.1 Structured Decoders

If the codes \mathbf{y} are predictive of the sources \mathbf{x} (which is something that we hope to achieve when we maximize the mutual information for the encoder parameters), then $p(\mathbf{x}|\mathbf{y})$ will typically be sharply peaked around its mode. This motivates a choice of simple (for example, uni-modal) approximations $q(\mathbf{x}|\mathbf{y})$ to the posterior, which in practice may significantly reduce computational complexity of optimization. In general, however, it is intuitive that the tightness of the bound may depend on the choice of the decoder $q(\mathbf{x}|\mathbf{y})$. For computational and representational purposes it is often convenient to assume that the approximate posterior $q(\mathbf{x}|\mathbf{y})$ in the variational bound (2.2) is factorized¹ in \mathbf{x} . Here we briefly describe effects which such a choice of the decoder structure may have on the bound. We also point out natural extensions of (2.2) to chain-type decoders, which overcome some of the restrictions of the factorized assumption.

2.2.1.1 Effects of the Decoder Structure

First of all, to demonstrate an influence which a choice of the decoder structure may have on the bound (2.2), we consider a simple case of factorized, mean-field type decoders for a model with $|\mathbf{x}|$ input and $|\mathbf{y}|$ output variables. For illustration purposes, we considered three decoders with different structures shown on Figure 2.1, so that

$$q^{(1)}(\mathbf{x}|\mathbf{y}) = \prod_{i=1}^{|\mathbf{x}|} q(x_i|\mathbf{y}), \quad q^{(2)}(\mathbf{x}|\mathbf{y}) = \prod_{i=1}^{|\mathbf{x}|} q(x_i|y_i), \quad q^{(3)}(\mathbf{x}|\mathbf{y}) = \prod_{i=1}^{|\mathbf{x}|} q(x_i), \quad (2.14)$$

where the products are computed over all the dimensions of \mathbf{x} . Note that in all these cases, different dimensions of the reconstructed vectors are conditionally independent given the codes. Arguably, this makes the considered decoders intrinsically less powerful than the theoretically optimal decoder defined by the true posterior $p(\mathbf{x}|\mathbf{y})$, where the source variables are generally conditionally dependent. However, it is easy to see that the considered structural constraints on $q(\mathbf{x}|\mathbf{y})$ may significantly simplify the resulting optimization procedures.

By computing the functional derivatives of the bound (2.2) with respect to the specified decoders (2.14) subject to the normalizing constraints, it is easy to show that for each of the three cases, the optimal variational distributions are defined as

$$q^{(1)}(x_i|\mathbf{y}) = p(x_i|\mathbf{y}), \quad q^{(2)}(x_i|\mathbf{y}) = q^{(2)}(x_i|y_i) = p(x_i|y_i), \quad q^{(3)}(x_i|\mathbf{y}) = q^{(3)}(x_i) = p(x_i). \quad (2.15)$$

¹The assumption becomes particularly important when the set of possible vectors generated by $p(\mathbf{x})$ is exponentially large; for example, when $\mathbf{x} \in \{0,1\}^{|\mathbf{x}|}$ and $p(\mathbf{x})$ is uniform.

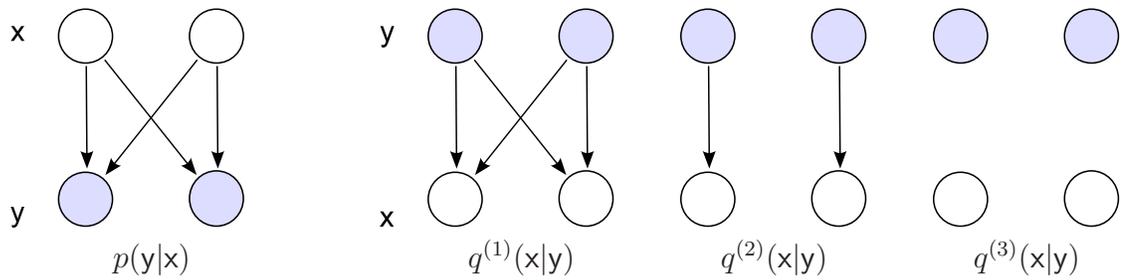


Figure 2.1: Simple variational decoders. *Left*: A noisy encoder $p(y|x)$; *Right*: constrained variational decoders $q^{(1)}(x|y)$, $q^{(2)}(x|y)$, and $q^{(3)}(x|y)$. The shaded nodes correspond to the hidden encoded representations $\{y\}$. The transparent nodes correspond to the sources $\{x\}$ and their reconstructions. (Unless mentioned otherwise, we will use the shaded nodes to indicate unobservable variables).

Then it is clear that for the variational decoder $q^{(1)}(x|y)$, the optimal achievable bound (2.2) on $I(x, y)$ is given by

$$\tilde{I}^{(1)}(\mathbf{x}, \mathbf{y}) = H(\mathbf{x}) - \sum_{i=1}^{|\mathbf{x}|} H(x_i|y), \quad (2.16)$$

where $H(x_i|y) \stackrel{\text{def}}{=} -\langle \log p(x_i|y) \rangle_{p(y|x)p(x)}$. Analogously, for $q^{(2)}(x_i|y)$ and $q^{(3)}(x_i|y)$ we obtain

$$\tilde{I}^{(2)}(\mathbf{x}, \mathbf{y}) = H(\mathbf{x}) - \sum_{i=1}^{|\mathbf{x}|} H(x_i|y_i), \quad \tilde{I}^{(3)}(\mathbf{x}, \mathbf{y}) = H(\mathbf{x}) - \sum_{i=1}^{|\mathbf{x}|} H(x_i), \quad (2.17)$$

with the similar definitions of the conditional entropies. It is easy to see that since the conditioning decreases the entropy (see e.g. Cover and Thomas (1991)), and all the variables for the three models are defined on identical domains, we get $H(x_i|y) \leq H(x_i|y_i) \leq H(x_i)$. This leads to the intuitive relation

$$I(\mathbf{x}, \mathbf{y}) \geq \tilde{I}^{(1)}(\mathbf{x}, \mathbf{y}) \geq \tilde{I}^{(2)}(\mathbf{x}, \mathbf{y}) \geq \tilde{I}^{(3)}(\mathbf{x}, \mathbf{y}), \quad (2.18)$$

which demonstrates that simplifications of the decoder structure (see Figure 2.1) result in relaxations of the theoretically achievable lower bounds.

In the case that the codes \mathbf{y} contain little information about the sources \mathbf{x} (for example, if the channel noise is high or dimensionality of the codewords $|\mathbf{y}|$ is low), factorized uni-modal approximations of the posterior $p(\mathbf{x}|\mathbf{y})$ may be quite restrictive. In this case we may expect that a decrease in the theoretically optimal bound due to a sub-optimal choice of the decoder structure may be quite significant.

Example: To demonstrate the effects numerically, we have considered a simple binary model with $|\mathbf{x}| = 2$ source and $|\mathbf{y}| = 2$ output variables. The encoder mapping was defined in such a way that both units $y_1, y_2 \in \{0, 1\}$ were solving

a noisy XOR problem with the noise levels $\epsilon_1 = 0.2$ and $\epsilon_2 = 0.3$. The channel probabilities for this case were defined as $p(y_i = 1|x_1 = x_2) = \epsilon_i$ and $p(y_i = 1|x_1 \neq x_2) = 1 - \epsilon_i$, where the conditioning indicates whether the source units $x_1, x_2 \in \{0, 1\}$ are set to the same binary states. The source variable marginals were defined as $p(x_1 = 1) = 0.2$ and $p(x_2 = 1) = 0.5$. It turns out that in this case the exact value of the mutual information is given by $I \approx 0.3548$, with the bounds $I^{(1)} \approx 0.1161$, $I^{(2)} \approx 0.0630$, and $I^{(3)} = 0$.

The result that a richer structure of a variational decoder $q(\mathbf{x}|\mathbf{y})$ leads to an improvement of the lower bound on $I(\mathbf{x}, \mathbf{y})$ is in full agreement with our expectations and with (2.18). It is also related to previous results in approximate inference and variational likelihood learning, which show that compared to the fully factorized approximations, structured mean field models (e.g. Barber and Wiergerinck (1998), Saad and Opper (2001)) often lead to tighter bounds on the log partition function. We will now consider a simple choice of the decoder structure which helps to go beyond the simple factorized approximations.

2.2.1.2 Chain Decoders

A possible limitation of the considered variational decoders $q(\mathbf{x}|\mathbf{y})$ is the factorized assumption $q(\mathbf{x}|\mathbf{y}) = \prod_{i=1}^{|\mathbf{x}|} q(x_i|\mathbf{y})$, which may lead to inaccurate approximations of the exact posterior. Specifically, for the factorized channel $p(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^{|\mathbf{y}|} p(y_i|\mathbf{x})$, the posterior $p(\mathbf{x}|\mathbf{y})$ is fully coupled in \mathbf{x} (see Figures 2.2 (a), (b)). While in some cases the factorized constraints may be plausible and indeed physically and computationally justifiable, in general they may be too restrictive. Intuitively, the bound on $I(\mathbf{x}, \mathbf{y})$ could be made tighter by retaining some of the structure in the specification of the variational distribution $q(\mathbf{x}|\mathbf{y})$. We will now show how this may be achieved in the context of information maximization.

A straight-forward application of the chain rule for probabilities of multivariate distributions

$$p(\mathbf{x}) = p(x_1) \prod_{i=2}^{|\mathbf{x}|} p(x_i|x_1, \dots, x_{i-1}) \quad (2.19)$$

leads to the corresponding rule for entropies

$$\begin{aligned} H(\mathbf{x}) &= -\langle \log p(x_1) \rangle_{p(\mathbf{x})} - \sum_{i=2}^{|\mathbf{x}|} \langle \log p(x_i|x_1, \dots, x_{i-1}) \rangle_{p(\mathbf{x})} \\ &= H(x_1) + \sum_{i=2}^{|\mathbf{x}|} H(x_i|x_1, \dots, x_{i-1}). \end{aligned} \quad (2.20)$$

Analogously, by conditioning on \mathbf{y} and averaging over $p(\mathbf{y})$, we may obtain a simple expression for the conditional entropy $H(\mathbf{x}|\mathbf{y})$:

$$H(\mathbf{x}|\mathbf{y}) = H(x_1|\mathbf{y}) + \sum_{i=2}^{|\mathbf{x}|} H(x_i|x_1, \dots, x_{i-1}, \mathbf{y}). \quad (2.21)$$

A straight-forward application of the definition of the mutual information (1.6) leads to

$$I(\mathbf{x}, \mathbf{y}) = I(x_1, \mathbf{y}) + \sum_{i=2}^{|\mathbf{x}|} I(x_i, \mathbf{y} | x_1, \dots, x_{i-1}), \quad (2.22)$$

which is just a chain rule on the mutual information (see e.g. Cover and Thomas (1991)). Equivalently, (2.22) may be expressed as

$$I(\mathbf{x}, \mathbf{y}) = H(\mathbf{x}) + \langle \log p(x_1 | \mathbf{y}) \rangle_{p(x_1, \mathbf{y})} + \sum_{i=2}^{|\mathbf{x}|} \langle \log p(x_i | x_1, \dots, x_{i-1}, \mathbf{y}) \rangle_{p(x_1, \dots, x_{i-1}, \mathbf{y})}. \quad (2.23)$$

Clearly, the marginal entropy $H(\mathbf{x})$ in (2.23) is independent of the encoder distribution $p(\mathbf{y} | \mathbf{x})$ and therefore has no effect on the optimization surface. However, the complexity of evaluating each average of $\log p(x_i | x_1, \dots, x_{i-1}, \mathbf{y})$ over the joint distribution $p(\mathbf{x}, \mathbf{y})$ is in general exponential in the number of parents of each variable x_i . One obvious way to bound the mutual information in a tractable way, while retaining some of the structure of $p(\mathbf{x} | \mathbf{y})$ in (2.23), would be to limit the number of parental connections in the variational approximation of $p(x_i | x_1, \dots, x_{i-1}, \mathbf{y})$. If $\boldsymbol{\pi}^x(x_i) \subseteq \{x_i | i = 1, \dots, |\mathbf{x}|\}$, $\boldsymbol{\pi}^y(x_i) \subseteq \{y_i | i = 1, \dots, |\mathbf{y}|\}$ are the x - and y -parents of the i^{th} reconstructed variable x_i , we may define the variational distribution to satisfy

$$q(x_i | \{x\} \setminus x_i, \{y\}) = q(x_i | \boldsymbol{\pi}^x(x_i), \boldsymbol{\pi}^y(x_i)). \quad (2.24)$$

It is clear that the definition (2.24) formally generalizes (2.15). Moreover, in general the resulting conditional $q(\mathbf{x} | \mathbf{y})$ is not factorized in x_i (see Figure 2.2).

It is easy to see that the choice of (2.24) gives rise to the bound

$$I(\mathbf{x}, \mathbf{y}) \geq \tilde{I}(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} H(\mathbf{x}) + \sum_i \langle \log q(x_i | \boldsymbol{\pi}^x(x_i), \boldsymbol{\pi}^y(x_i)) \rangle_{p(\mathbf{x}_i, \boldsymbol{\pi}^x(x_i), \boldsymbol{\pi}^y(x_i))}, \quad (2.25)$$

Clearly, if the number of parents $|\boldsymbol{\pi}^x(x_i)| + |\boldsymbol{\pi}^y(x_i)|$ is small, the summations in (2.25) may be performed exactly (for discrete domains) or efficiently approximated numerically (e.g. by using standard techniques of numerical integration for continuous domains). A simple structured decoder can be characterized by a sparse mapping from \mathbf{y} to \mathbf{x} and a chain in the x 's (see Figure 2.2), though any choice of parents satisfying the tractability constraints is possible. In general, an extra care must be taken to ensure that $q(\mathbf{x} | \mathbf{y})$ remains a proper distribution (i.e. its graph does not have directed cycles). Clearly, some of the possible choices for the variational decoder $q(\mathbf{x} | \mathbf{y})$ are trees or higher-order Markov chains, with an additional constraint on $|\boldsymbol{\pi}^y(x_i)|$.

By analogy with the simple factorized example (2.15) described above, by computing functional derivatives of the bound (2.25) with respect to $q(x_i | \boldsymbol{\pi}^x(x_i), \boldsymbol{\pi}^y(x_i))$ subject to the normalization constraints, we can easily find that the optimal settings would give rise to $q(x_i | \boldsymbol{\pi}^x(x_i), \boldsymbol{\pi}^y(x_i)) = p(x_i | \boldsymbol{\pi}^x(x_i), \boldsymbol{\pi}^y(x_i))$. This would lead to the optimal theoretically achievable lower bound

$$I(\mathbf{x}, \mathbf{y}) \geq H(\mathbf{x}) - \sum_i H(x_i | \boldsymbol{\pi}^x(x_i), \boldsymbol{\pi}^y(x_i)). \quad (2.26)$$

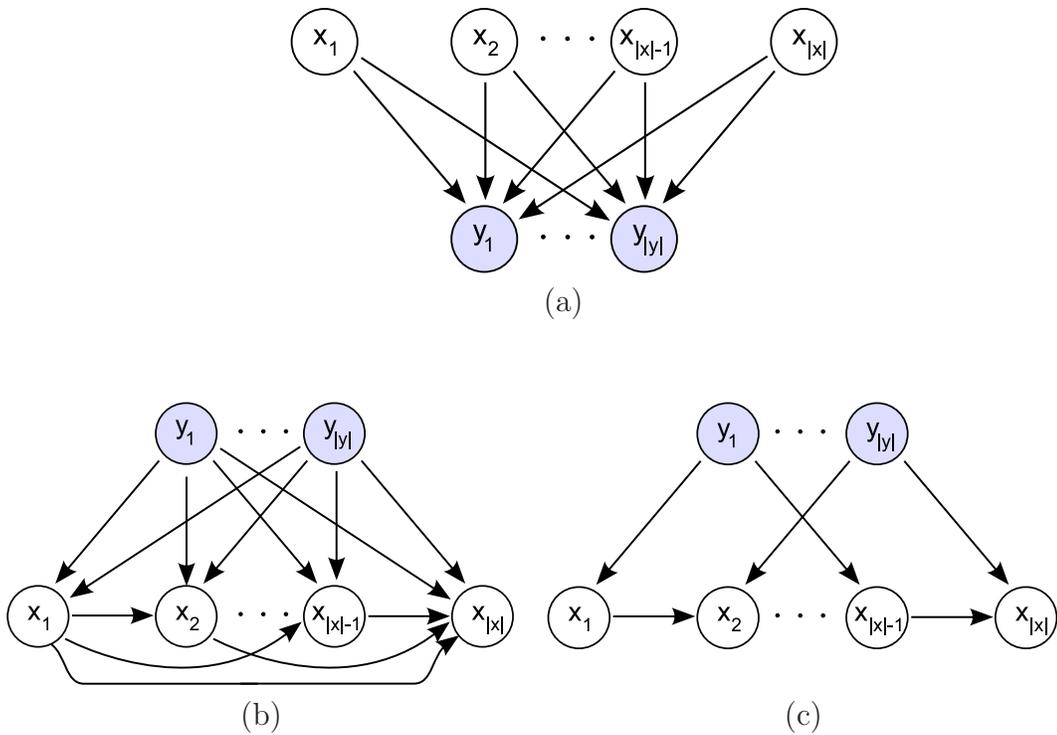


Figure 2.2: Structured variational decoders. (a) A noisy encoder $p(y|x)$; (b) the corresponding fully-structured exact decoder $p(x|y)$; (c) a sparse variational decoder $q(x|y)$. The shaded nodes correspond to the encodings $\{y\}$. The sparse variational decoder retains some of the structure of the exact posterior $p(x|y)$.

(compare with (2.23)). Again, from the non-negativity of the KL divergence it is easy to see that unless the variables in $\pi^x(x_i)$, $\pi^y(x_i)$ are independent, the conditioning on additional variables always decreases the entropy, independently of the specifics of channel parameterization. Therefore, by retaining more structure in the variational decoders, we are guaranteed to improve on the theoretically achievable lower bounds on $I(\mathbf{x}, \mathbf{y})$.

Finally, note that even in the cases when the models are defined over discrete spaces with finite alphabets of the code and source variables, it may be difficult to compute the entropic terms in (2.26); for example, this may happen when averaging over $p(\mathbf{x})$ is expensive. However, by imposing sparsity constraints on the *encoding* distribution $p(y|x)$, we may be able to perform the computations exactly (see the discussion in Section 7.2.2). Obviously, for continuous variable spaces it is necessary to impose additional parametric constraints on $p(y|x)$ and $q(x|y)$.

2.2.2 Gaussian Decoders and the Link to Linsker’s Criterion

Here we consider a special case of the bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ for linear Gaussian variational decoders. We also show that optimization of Linsker’s *as-if* Gaussian objective (Linsker (1992)), approximating the exact mutual information $I(\mathbf{x}, \mathbf{y})$,

corresponds to a special way of optimizing the variational lower bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ independently of the choice of the encoder distribution $p(\mathbf{y}|\mathbf{x})$.

Let us define the variational decoder to be a Gaussian $q(\mathbf{x}|\mathbf{y}) \sim \mathcal{N}_{\mathbf{x}}(\mathbf{U}\mathbf{y}, \mathbf{\Sigma})$, where $\mathbf{U} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{y}|}$. Clearly, in this case the variational lower bound (2.2) is expressed as

$$\begin{aligned} \tilde{I}(\mathbf{x}, \mathbf{y}) &= \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})p(\mathbf{x})} + H(\mathbf{x}) \\ &= -\frac{1}{2} \langle \text{tr} \{ \mathbf{\Sigma}^{-1}(\mathbf{x} - \mathbf{U}\mathbf{y})(\mathbf{x} - \mathbf{U}\mathbf{y})^T \} \rangle_{p(\mathbf{y}|\mathbf{x})p(\mathbf{x})} - \frac{1}{2} \log |\mathbf{\Sigma}| + c \end{aligned} \quad (2.27)$$

where c is an irrelevant constant. One way to optimize (2.27) would be by applying the iterative IM algorithm for the encoder $p(\mathbf{y}|\mathbf{x})$ and parameters of the decoder $\mathbf{U}, \mathbf{\Sigma}$. Alternatively, at the first iteration of the algorithm we may utilize Gaussianity of the decoder $q(\mathbf{x}|\mathbf{y})$ to find analytical expressions of the parameters \mathbf{U} and $\mathbf{\Sigma}$, expressing them as functions of the fixed encoder $p(\mathbf{y}|\mathbf{x})$. By substituting the results into (2.27), we may re-define the bound as a function of the encoder alone. Then the new objective function may be optimized numerically with respect to encoder parameters. It turns out that this approach to maximizing the lower bound (2.27) leads to optimization of Linsker's *as-if* Gaussian criterion (1.16).

Proposition 2.3. *For any channel $p(\mathbf{y}|\mathbf{x})$, optimization of Linsker's as-if Gaussian criterion $I_G(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} -\log |\langle \mathbf{x}\mathbf{x}^T \rangle - \langle \mathbf{x}\mathbf{y}^T \rangle \langle \mathbf{y}\mathbf{y}^T \rangle^{-1} \langle \mathbf{y}\mathbf{x}^T \rangle|$ for the encoder parameters corresponds to a specific way of optimizing the variational lower bound on the mutual information $\tilde{I}(\mathbf{x}, \mathbf{y}) = \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})p(\mathbf{x})}$ with linear Gaussian decoders² $q(\mathbf{x}|\mathbf{y}) \sim \mathcal{N}_{\mathbf{x}}(\mathbf{U}\mathbf{y}, \mathbf{\Sigma})$.*

Proof. Optimizing (2.27) for the decoder covariance $\mathbf{\Sigma}$ (assumed to be non-singular), we obtain the extremum condition

$$\mathbf{\Sigma}^{-1} \langle (\mathbf{x} - \mathbf{U}\mathbf{y})(\mathbf{x} - \mathbf{U}\mathbf{y})^T \rangle_{p(\mathbf{x})p(\mathbf{y}|\mathbf{x})} \mathbf{\Sigma}^{-1} = \mathbf{\Sigma}^{-1}. \quad (2.28)$$

Clearly, this results in the optimal decoder's covariance

$$\mathbf{\Sigma} = \langle (\mathbf{x} - \mathbf{U}\mathbf{y})(\mathbf{x} - \mathbf{U}\mathbf{y})^T \rangle_{p(\mathbf{x})p(\mathbf{y}|\mathbf{x})}. \quad (2.29)$$

It is easy to verify that (2.29) is indeed a maximum. Substitution into (2.27) leads to

$$\tilde{I}(\mathbf{x}, \mathbf{y}) \propto -\log \left| \langle (\mathbf{x} - \mathbf{U}\mathbf{y})(\mathbf{x} - \mathbf{U}\mathbf{y})^T \rangle_{p(\mathbf{x}, \mathbf{y})} \right|, \quad (2.30)$$

where we ignored irrelevant additive and multiplicative constants. Computing the derivatives for the decoder weights \mathbf{U} , we obtain

$$\mathbf{\Sigma}^{-1} \langle \mathbf{x}\mathbf{y}^T \rangle = \mathbf{\Sigma}^{-1} \mathbf{U} \langle \mathbf{y}\mathbf{y}^T \rangle. \quad (2.31)$$

Assuming that $\langle \mathbf{y}\mathbf{y}^T \rangle \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{y}|}$ is non-singular (e.g. this is the case when the number of linearly independent codes exceeds their dimensionality), we may express

²Here we will focus on the discussion of the centered data and centered codes, i.e. $\langle \mathbf{x} \rangle = 0$, $\langle \mathbf{y} \rangle = 0$. The more general case may be obtained analogously by considering variational decoders of the form $q(\mathbf{x}|\mathbf{y}) \sim \mathcal{N}_{\mathbf{x}}(\mathbf{U}(\mathbf{y} - \langle \mathbf{y} \rangle), \mathbf{\Sigma})$ (see a brief discussion in Appendix A).

the optimal weights as $\mathbf{U} = \langle \mathbf{x}\mathbf{y}^T \rangle \langle \mathbf{y}\mathbf{y}^T \rangle^{-1} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{y}|}$. By substituting the optimal decoder weights into (2.30), we obtain the objective function

$$\tilde{I}(\mathbf{x}, \mathbf{y}) = -\log |\langle \mathbf{x}\mathbf{x}^T \rangle - \langle \mathbf{x}\mathbf{y}^T \rangle \langle \mathbf{y}\mathbf{y}^T \rangle^{-1} \langle \mathbf{y}\mathbf{x}^T \rangle|, \quad (2.32)$$

which, up to irrelevant constants, is exactly the Linsker’s Gaussian criterion (1.16). \square

In general the described modification of the objective criterion may complicate the analysis of optimal solutions, as (2.32) may be highly nonlinear in the encoder parameters. Moreover, by expressing the bound as a function of the encoder alone, we may change the optimization surface (for example, if the bound is not a convex function of encoder parameters), which may affect the obtained optima. Nevertheless, by an argument similar to Proposition 2.1, it is easy to see that the described optimization procedure is guaranteed not to weaken $\tilde{I}(\mathbf{x}, \mathbf{y})$, though the obtained solutions will depend on the specifics of the optimization strategies used.

Finally, we stress once again that optimization of the *as-if* Gaussian approximation of the mutual information (Linsker (1992)) corresponds to one specific way of optimizing the variational lower bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ under the specific assumption that the variational decoder is a linear Gaussian. It is clear that the variational formulation is a powerful generalization of Linsker’s approach. First, the generic lower bound (2.2) gives us a flexibility in choosing the family of decoder distributions \mathcal{Q} . Secondly, for any decoder in the family $q(\mathbf{x}|\mathbf{y}, \Theta) \in \mathcal{Q}$, we are free to choose a specific optimization procedure for the encoder and decoder parameters. Optimization of Linsker’s bound corresponds to one of such choices for the family of linear Gaussian decoders. As we show in the following sections, by considering a richer family of variational decoders (and in fact a richer family of bounds on the mutual information) we may significantly improve on simple approximations.

2.3 An Auxiliary Variational Lower Bound on Mutual Information

In section 2.2.1 we showed that by retaining in the variational decoder $q(\mathbf{x}|\mathbf{y})$ some of the local sub-structure of the fully-coupled exact posterior $p(\mathbf{x}|\mathbf{y})$, we could indeed obtain tighter bounds on the mutual information. We also noted that under the specific structural constraints, the computations should indeed remain tractable, thus resulting in simple ways of improving on factorized mean field bounds on the mutual information. Here we discuss a formal generalization of structured bounds on $I(\mathbf{x}, \mathbf{y})$ for mixture-of-experts type decoders.

2.3.1 Representations

A possible way to further improve on the structured bounds (2.25) is to increase the representational power of structured variational decoders $q(\mathbf{x}|\mathbf{y})$ by capturing *global* dependencies between the reconstructed variables. One way to achieve it is

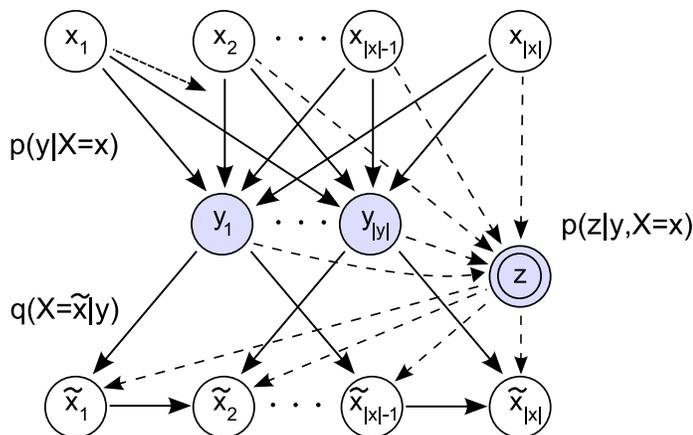


Figure 2.3: A noisy channel $p(y|x)$ with a structured mixture-type decoder $q(x|y)$. (The states of the reconstructed variables are denoted by \tilde{x}). The role of the auxiliary variables z is to introduce additional structure into the decoder; they are *not* transmitted across the channel $p(y|x)$ and do not explicitly constrain $p(x, y)$. The dashed lines show the mappings to and from the auxiliary space. The auxiliary node is shown by the double circle.

to consider multi-modal decoders $q(x|y) = \langle q(x|y, z) \rangle_{q(z|y)}$, where the introduced *auxiliary variables* z are effectively the unknown mixture labels. Effectively, this choice of the variational decoder corresponds to a form of mixture-of-experts-type models (Jacobs et al. (1991), Neuneier et al. (1994), Bishop (1994), Bishop (1995)) for modeling of the conditional distributions. Clearly, the fully-coupled structure of the resulting variational distribution $q(x|y)$ will qualitatively agree with the dependency structure of the exact posterior, as different dimensions of the reconstructed vectors x would be coupled through the auxiliary variables z . Moreover, for any interesting choice of the auxiliary space $\{z\}$, the decoder $q(x|y)$ will typically be multi-modal, which would agree with the common mixture-type structure of $p(x|y)$ (see (1.8)). We may therefore intuitively hope that this choice of the variational posterior will generally result in tighter bounds on $I(x, y)$.

Despite the apparent merits of the mixture-type variational decoder $q(x|y)$, its possible disadvantage relates to the fact that specifying the conditional mixing coefficients $q(z|y)$ in a principled manner may be rather difficult. Moreover, if the auxiliary variables z are conditionally independent from the *original* source patterns given the codes, any noise in the stochastic encodings y will affect determining of the mixing states. Intuitively, this may have an overwhelming negative effect on decoding, which may cause relaxations of the bound on $I(x, y)$. We may therefore wish to reduce the effects which the noise in the encodings has on the specification of the decoder $q(x|y)$.

One way to address this matter is by introducing an additional mapping to the auxiliary variable space $p(z|x, y)$, which may be thought of as an additional variational parameter (see Figure 2.3). Indeed, even when the channel is noisy, the conditional dependence of the auxiliary variables z on the unperturbed source patterns could result in an accurate detection of the states of the auxiliary vari-

ables. Note that the *auxiliary conditional* distribution $p(\mathbf{z}|\mathbf{x}, \mathbf{y})$ is defined in a way that does not affect the original noisy channel $p(\mathbf{y}|\mathbf{x})$, as the channel would remain a marginal of the joint distribution of the original sources, codes, and auxiliary variables

$$p(\mathbf{x}, \mathbf{y}, \mathbf{z}) = p(\mathbf{x})p(\mathbf{y}|\mathbf{x})p(\mathbf{z}|\mathbf{x}, \mathbf{y}). \quad (2.33)$$

The role of the auxiliary variables \mathbf{z} in this context would be to capture global features of the transmitted sources, and use these features for choosing optimal experts in the mixture-of-experts type³ decoder. Importantly, the auxiliary variables \mathbf{z} are *not* transmitted across the channel. Their purpose here is to define a richer family of bounds on $I(\mathbf{x}, \mathbf{y})$ which would generalize over the simpler bounds (2.16)–(2.17), (2.25) with factorial or structured variational decoders.

It is easy to see that for an unconstrained choice of the variational parameters, a straight-forward application of the generic bound (2.2) for the discussed formulation of the mixture-type decoder may result in generally intractable integrals over \mathbf{x} and \mathbf{y} . A simple way to handle the intractability is to consider further variational relaxations of the lower bound (2.2) by constraining all the variational distributions to be tractable. We will now show that this indeed leads to a tractable generalization over (2.2).

2.3.2 An Auxiliary Variational Lower Bound on $I(\mathbf{x}, \mathbf{y})$

An *auxiliary variational* lower bound on $I(\mathbf{x}, \mathbf{y})$ may be obtained by considering the general properties of the mutual information. It may also be shown that the resulting bound corresponds to a tractable variational relaxation of the generic form (2.2) for the mixture-of-experts type decoders.

Let $p(\mathbf{x}, \mathbf{y}, \mathbf{z})$ define a general joint distribution over the sources \mathbf{x} , the encodings \mathbf{y} , and the auxiliary variables (features) \mathbf{z} , where we parameterize the original channel $p(\mathbf{y}|\mathbf{x})$ and the auxiliary variational conditional $p(\mathbf{z}|\mathbf{x}, \mathbf{y})$, and assume that $p(\mathbf{x})$ is known and fixed (see expression (2.33) and Figure 2.3). From the chain rule for mutual information (2.22), we may express $I(\mathbf{y}, \mathbf{x})$ as

$$I(\mathbf{y}, \mathbf{x}) = I(\{\mathbf{z}, \mathbf{y}\}, \mathbf{x}) - I(\mathbf{x}, \mathbf{z}|\mathbf{y}), \quad (2.34)$$

where

$$I(\{\mathbf{z}, \mathbf{y}\}, \mathbf{x}) \stackrel{\text{def}}{=} H(\mathbf{x}) - H(\mathbf{x}|\mathbf{z}, \mathbf{y}) \quad (2.35)$$

is the amount of information that the features \mathbf{z} and codes \mathbf{y} jointly contain about the sources, and

$$I(\mathbf{x}, \mathbf{z}|\mathbf{y}) \stackrel{\text{def}}{=} H(\mathbf{z}|\mathbf{y}) - H(\mathbf{z}|\mathbf{x}, \mathbf{y}) \quad (2.36)$$

is the conditional mutual information. Substituting (2.35) and (2.36) into (2.34), we obtain a general expression of the mutual information $I(\mathbf{x}, \mathbf{y})$ as a function of conditional entropies of the sources, codes, and auxiliary variables

$$I(\mathbf{y}, \mathbf{x}) = H(\mathbf{x}) + H(\mathbf{z}|\mathbf{x}, \mathbf{y}) - H(\mathbf{x}|\mathbf{y}, \mathbf{z}) - H(\mathbf{z}|\mathbf{y}). \quad (2.37)$$

³One may immediately notice that the optimal *unconstrained* variational distribution of the mixing labels $q(\mathbf{z}|\mathbf{y})$ of the decoder $q(\mathbf{x}|\mathbf{y}) = \langle q(\mathbf{x}|\mathbf{y}, \mathbf{z}) \rangle_{q(\mathbf{z}|\mathbf{y})}$ should itself be a mixture-of-experts model, defined as $q(\mathbf{z}|\mathbf{y}) = \langle p(\mathbf{z}|\mathbf{x}, \mathbf{y}) \rangle_{p(\mathbf{x}|\mathbf{y})}$. This is confirmed at a later stage (see the objective (2.38)).

Clearly, $H(\mathbf{x})$ in (2.37) is not a function of the channel parameters, and it may be safely ignored for the purpose of learning the optimal encoder $p(\mathbf{x}|\mathbf{y})$. The entropic term $H(\mathbf{z}|\mathbf{x}, \mathbf{y}) = -\langle \log p(\mathbf{z}|\mathbf{x}, \mathbf{y}) \rangle_{p(\mathbf{z}, \mathbf{x}, \mathbf{y})}$ is a functional of the auxiliary conditional mapping to the feature space $p(\mathbf{z}|\mathbf{x}, \mathbf{y})$, which may be chosen to lie in a tractable family (or defined to be a deterministic function of \mathbf{x} and \mathbf{y}). Computations of the remaining terms in (2.37) are in general problematic, as both posteriors $p(\mathbf{x}|\mathbf{y}, \mathbf{z})$ and $p(\mathbf{z}|\mathbf{y})$ are effectively mixture distributions (coupled in \mathbf{x} and \mathbf{z} respectively). However, by analogy with (2.2) we may bound both terms, which would result in

$$I(\mathbf{x}, \mathbf{y}) \geq H(\mathbf{x}) + H(\mathbf{z}|\mathbf{x}, \mathbf{y}) + \langle \log q(\mathbf{x}|\mathbf{y}, \mathbf{z}) \rangle_{p(\mathbf{x}, \mathbf{y}, \mathbf{z})} + \langle \log q(\mathbf{z}|\mathbf{y}) \rangle_{p(\mathbf{y}, \mathbf{z})}. \quad (2.38)$$

This may be further recognized as a variational relaxation of the generic criterion (2.2) for a mixture-of-experts decoder $q(\mathbf{x}|\mathbf{y}) = \langle q(\mathbf{x}|\mathbf{y}, \mathbf{z}) \rangle_{q(\mathbf{z}|\mathbf{y})}$, where the auxiliary variational conditional $p(\mathbf{z}|\mathbf{x}, \mathbf{y})$ plays the role of the additional variational parameter.

Again, to ensure that the averages in (2.38) are tractable, we need to constrain the variational decoders $q(\mathbf{x}|\mathbf{y}, \mathbf{z})$ and $q(\mathbf{z}|\mathbf{y})$ (in addition to constraining the auxiliary mapping $p(\mathbf{z}|\mathbf{x}, \mathbf{y})$). The objective (2.38) needs to be optimized for the channel encoder, variational decoder, and the auxiliary conditional distributions, subject to the imposed constraints. One way to perform the optimization is by considering a straight-forward extension of the variational IM algorithm to include the auxiliary mappings. By analogy with proposition 2.1, it is easy to show that the resulting iterative optimization procedure maximizes or leaves unchanged the bound on the mutual information $I(\mathbf{x}, \mathbf{y})$. Note that the role of the auxiliary vectors \mathbf{z} in this context is to capture global dependencies in the reconstructed variables; the auxiliary variables are *not* transmitted across the noisy channel, which for our case is defined by $p(\mathbf{y}|\mathbf{x})$.

It is important to note that the objective (2.38) is a formal generalization of the generic lower bound on the mutual information. Indeed, it is clear that if \mathbf{z} is a vector of deterministic variables taking a single state, the bound (2.38) reduces to $H(\mathbf{x}) + \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{x}, \mathbf{y})}$ (*cf* expression (2.2)). Simpler special cases of the objective may be obtained by imposing further constraints on $p(\mathbf{z}|\mathbf{x}, \mathbf{y})$ and $q(\mathbf{z}|\mathbf{y})$. For example, by constraining the variational distribution of the coefficients to be $q(\mathbf{z}|\mathbf{y}) \equiv q(\mathbf{z})$, we may obtain a mixture bound on $I(\mathbf{x}, \mathbf{y})$ as a special case. By further considering a constrained auxiliary mapping $p(\mathbf{z}|\mathbf{x}, \mathbf{y}) \equiv p(\mathbf{z}|\mathbf{y})$, we may transform (2.38) to the Jensen's relaxation of (2.2) for mixture decoders, etc.

In a more general case, the resulting decoder $q(\mathbf{x}|\mathbf{y}) = \langle q(\mathbf{x}|\mathbf{y}, \mathbf{z}) \rangle_{q(\mathbf{z}|\mathbf{y})}$ is a generalization of the mixture-of-experts model, in the sense that the variational prior over the mixing coefficient $q(\mathbf{z}|\mathbf{y})$ is itself a mixture. Indeed, from (2.38) it is clear that the local unconstrained optimization for the mixture prior leads to $q(\mathbf{z}|\mathbf{y}) = p(\mathbf{z}|\mathbf{y})$, as expressed from the channel $p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$ and the auxiliary variational conditional $p(\mathbf{z}|\mathbf{x}, \mathbf{y})$. As the exact decoder $p(\mathbf{x}|\mathbf{y})$ will typically have a mixture form with the data-dependent coefficients (see expression (1.9)), we may informally hope that a choice of a multi-modal variational decoder will generally lead to tighter bounds on $I(\mathbf{x}, \mathbf{y})$.

One may further extend (2.38) by considering models with a richer hierarchy, though an extra care should be taken to ensure tractability. For example, in order for the bound (2.38) to be computationally tractable, one may need to impose additional structural constraints on the conditionals, analogous to the ones discussed in section 2.2.1. For example, if $\pi^x(z_i) \subseteq \{\mathbf{x}\}$, $\pi^y(z_i) \subseteq \{\mathbf{y}\}$ are lower-dimensional x - and y - parents of z_i , one may consider a sparse factorized form of the auxiliary conditional

$$p(\mathbf{z}|\mathbf{y}, \mathbf{x}) = \prod_i p(z_i|\pi^x(z_i), \pi^y(z_i)). \quad (2.39)$$

In this case, the entropic term $-H(z_i|\mathbf{x}, \mathbf{y})$ may be exactly expressed as

$$\begin{aligned} \langle \log p(z_i|\pi^x(z_i), \pi^y(z_i)) \rangle_{p(z_i, \pi^x(z_i), \pi^y(z_i))} &= \sum_{\eta^x(z_i)} \prod_{k \in \eta^x(z_i)} p(x_k) \sum_{\pi^y(z_i)} \prod_{j \in \pi^y(z_i)} p(y_j|\pi^x(y_j)) \times \\ &\quad \sum_{z_i} p(z_i|\pi^x(z_i), \pi^y(z_i)) \log p(z_i|\pi^x(z_i), \pi^y(z_i)), \end{aligned} \quad (2.40)$$

where $\eta^x(z_i) \stackrel{\text{def}}{=} \pi^x(\pi^y(z_i)) \cup \pi^x(z_i)$ and the definition of the parents π was extended for sets of variables in the obvious manner. As before, tractability of this computation may be ensured by choosing an appropriate parameterization (or imposing appropriate sparse structural constraints on discrete models). Analogously, we may choose a tractable structure for the remaining variational parameters $q(\mathbf{z}|\mathbf{y})$ and $q(\mathbf{x}|\mathbf{y}, \mathbf{z})$.

2.4 Summary

We described a general theoretically justified approach to information maximization in noisy channels, which extends the family of Blahut-Arimoto type algorithms (Arimoto (1972), Blahut (1972)) by considering a constrained tractable family of decoder distributions. We also showed that the constraints are fundamentally important for avoiding intractability of computing and optimizing the mutual information in the presence of noise. By constraining the approximate posterior, we effectively re-define the optimized objective criterion to be a proper variational lower bound on the mutual information. The suggested iterative optimization procedure, the **IM** algorithm, is guaranteed not to weaken the bound, and has a form reminiscent of the variational EM algorithm for approximate likelihood maximization.

Whilst the generic formulation of the objective criterion is straightforward, it appears to have attracted little previous attention as a practical tool for maximizing mutual information. The suggested formulation is conceptually simple and general; moreover, it implies optimization of a proper *lower bound*, rather than an approximation of the true mutual information. The generality of the approach allows a flexibility in the choice of variational decoders or specific optimization procedures, which suggests that the method may naturally generalize

other techniques for approximate information maximization. Indeed, we showed that optimization of Linsker’s *as-if* Gaussian objective criterion (Linsker (1992)) corresponds to a specific way of optimizing the variational lower bound on mutual information for a specific choice of linear Gaussian decoders.

Furthermore, we considered extensions of the simple bound on $I(\mathbf{x}, \mathbf{y})$ to the case of structured decoders, and showed that by retaining local dependencies, we could indeed tighten the theoretically achievable lower bounds on $I(\mathbf{x}, \mathbf{y})$.

Finally, we introduced a new *auxiliary variational* approach to information maximization. The key idea of the suggested auxiliary method is to introduce additional variables, which may be used for capturing useful features of the source patterns, and for introducing global dependencies to the decoded sources. Importantly, the variables and the projections to the auxiliary space are defined in a way which does not impose explicit constraints on the original channel. It is intuitive that these richer representations help to retain some of the structure (and, for the auxiliary formulation, multimodality) of the exact decoder. Moreover, for any constrained choice of the variational distributions $q(\mathbf{x}|\mathbf{y}, \mathbf{z})$, they include simpler generic bounds as special cases. We will subsequently confirm that this more general family of bounds may indeed significantly improve on simpler approaches. However, first we will discuss general relations between optimizing the mutual information in encoder models of noisy channels, optimizing the likelihood in generative latent variable models, and optimizing the conditional likelihood in stochastic autoencoders.

Chapter 3

Likelihood, Conditional Likelihood, and Information Maximization

3.1 Introduction

The problem of finding informative lower-dimensional representations of the observed data may be addressed in a number of different ways. One way to model a relationship between hidden and visible variables (codes and sources) \mathbf{y} and \mathbf{x} is by considering a generative model, and training it by maximizing the likelihood¹ \mathcal{L} . By trying to fit the empirical distribution (subject to the modeling constraints), generative models tend to find latent representations which could be useful for generating the observed patterns. The lower-dimensional latent-space representations may then be obtained by applying Bayes rule. A classic alternative to generative models for finding informative codes is given by the self-extracting autoencoder-type networks, where the sources \mathbf{x} are passed through the bottleneck of hidden variable representations \mathbf{y} , and extracted at the decoding end to yield the reconstructions $\tilde{\mathbf{x}}$. A standard approach is to train such models by minimizing the reconstruction error, which under certain modeling assumptions may be seen as conditional likelihood learning (Bishop (1995)).

In a sense, generative latent variable models $\mathbf{y} \rightarrow \mathbf{x}$ and autoencoders $\mathbf{x} \rightarrow \mathbf{y} \rightarrow \tilde{\mathbf{x}}$ may be viewed as useful frameworks for finding informative hidden variable representations of the data. However, an arguably more principled measure of informativeness is given by the mutual information $I(\mathbf{x}, \mathbf{y})$, defined as an average reduction in uncertainty of one variable given the other in the stochastic encoder model $\mathbf{x} \rightarrow \mathbf{y}$. The fundamental Fano's result (Fano (1961)) states that independently of the specific parameterization of a (discrete) encoder distribution, the probability of incorrect reconstruction of the sources \mathbf{x} from the codes \mathbf{y} is bounded by the negative mutual information as

$$p_e(\tilde{\mathbf{x}} \neq \mathbf{x}) \geq (H(\mathbf{x}) - I(\mathbf{x}, \mathbf{y}) - 1) / \log M \in [0, 1], \quad (3.1)$$

where M is the number of possible distinct reconstructions (size of the input alphabet). In other words, in discrete channels the mutual information $I(\mathbf{x}, \mathbf{y})$

¹As a monotonic transformation of an objective function does not affect the optimization surface, we will sometimes use the term *likelihood* to refer to the scaled log-likelihood.

is lower-bounded by the probability of the correct reconstruction². Intuitively, this result suggests a relation between information maximization and conditional likelihood training in *stochastic* autoencoders, where the projection into the code space corresponds to the noisy encoding distribution of a channel model. It is interesting to try to understand how all these seemingly different measures relate to each other in completely general settings.

While in some specific cases likelihood and mutual information may indeed have identical optima (see e.g. Cardoso (1997) for a discussion of the noiseless squared ICA case), the optimization surfaces which they define seem to be very different in general. One may understand conceptual differences between the optima provided by likelihood and mutual information maximization for the case of mixture models, where y is a hidden mixture index and \mathbf{x} is a visible pattern. It is well known that the likelihood may be trivially maximized by fitting a component to a local segment of the data (as small as a single pattern), provided that the other patterns are explained by the remaining components. One may think of this as an artifact of the definition of the likelihood for under-constrained models (where the theoretically optimal *unconstrained* model is simply the empirical distribution). Intuitively, such solutions should be suboptimal in the information theoretic sense, as they would not be useful for reconstructing most of the source patterns. Analytically, they would result in a decrease in the marginal entropy of the codes $H(y)$, leading to a general reduction in $I(\mathbf{x}, y)$. On the other hand, maximization of mutual information in unconstrained models may result in different forms of trivial optima, which may be characterized by noiseless projections of the training patterns to maximally uniform encodings.

In principle, for any specific parameterization of the generative and encoder (recognition) models, the solutions provided by the likelihood and mutual information maximization may be related by comparing the specific extrema conditions for the parameters. Most of the previous work on relating maximum-likelihood learning in generative models and autoencoders to the information-theoretic learning in encoder models of noisy channels focused primarily on specific invertible mappings (Pearlmutter and Parra (1996), Cardoso (1997), MacKay (1999b)), or constrained linear models (Cottrell et al. (1987), Baldi and Hornik (1989), Oja (1989)). Other more recent work addressed the problem of finding links between a different, but somewhat related family of *Information Bottleneck* approaches and the maximum likelihood methods for a specific tractable mixture model (Slonim and Weiss (2002)). Our goal here is to discuss a relation between maximum-likelihood and information-maximization methods for both exact and approximate cases, independently of the specific parameterizations (see Section 3.2). Then in Section 3.3 we will compare mutual information maximization with the conditional likelihood training of feed-forward models.

²Note that originally Fano's inequality was derived under the assumption of discrete source variables \mathbf{x} and deterministic reconstructions from the encoded representations, i.e. $H(\tilde{\mathbf{x}}|y) = 0$, see e.g. Fano (1961), Cover and Thomas (1991).

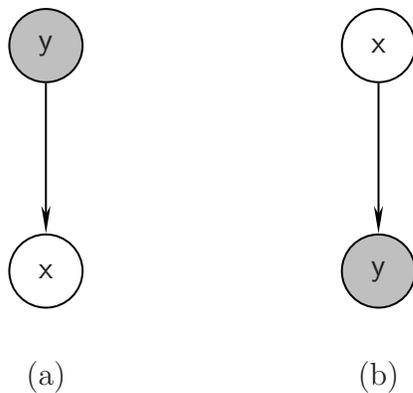


Figure 3.1: Generative and encoder models. (a): A generative model $\mathcal{M}_L \stackrel{\text{def}}{=} p(y)p(x|y)$ trained by maximizing the likelihood \mathcal{L} ; (b): a model of a noisy channel $\mathcal{M}_I \stackrel{\text{def}}{=} \tilde{p}(x)p(y|x)$ (trained by maximizing the mutual information $I(x, y)$). The shaded nodes correspond to the hidden variable representations y .

3.2 Information Maximization and Maximum Likelihood

There are fundamental conceptual differences in how we parameterize and train generative and recognition models. In order to define a generative latent variable model \mathcal{M}_L , we need to specify the prior on the latent codes $p(y|\Theta_y)$, and the conditional distribution of the observations given the codes $p(x|y, \Theta_{x|y})$ (see Figure 3.1 (a)). Throughout the discussion, we will assume that $p(\Theta_y) \sim \delta(\Theta_y - \Theta_y^*)$ and $p(\Theta_{x|y}) \sim \delta(\Theta_{x|y} - \Theta_{x|y}^*)$ (for some optimal settings $\Theta_y^*, \Theta_{x|y}^*$), and drop the conditioning on the parameters $\Theta_y, \Theta_{x|y}$. The distribution of the data under the model is given by the marginal $p(x) = \langle p(x|y) \rangle_{p(y)}$. As discussed earlier, we may train such models by maximizing the log-probability of generating a fixed set of observations $\mathcal{L} \stackrel{\text{def}}{=} \log p(\{x^{(1)}, \dots, x^{(M)}\})$, subject to the constraints on $p(y)$ and $p(x|y)$. Obviously, for independent identically distributed patterns, maximum likelihood training corresponds to minimization of the Kullback-Leibler divergence between the empirical distribution and the model of the observations $p(x)$.

In contrast, to specify a recognition model, we need to define the distribution of the sources $\tilde{p}(x)$ and the generally noisy encoder $\tilde{p}(y|x)$ (see Figure 3.1 (b)). In a vast number of practical applications of information-theoretic training of such models, it is presumed that $\tilde{p}(x)$ is the empirical distribution (see e.g. Fano (1961), McEliece (1977), Linsker (1989a), Nadal and Parga (1994), Bell and Sejnowski (1995), Linsker (1997), Principe et al. (2000), Torkkola and Campbell (2000), Torkkola (2001), Szummer and Jaakkola (2002)). Here we will focus primarily on the discussion of this case. In principle, we may parameterize the encoder $\tilde{p}(y|x)$ arbitrarily, as it is a part of the model's specification. However, for the purpose of comparing the maximum-likelihood training in generative models with the mutual information maximization in noisy channels, we will set $\tilde{p}(y|x)$ to be

the exact posterior of \mathcal{M}_L , i.e. $\tilde{p}(y|x) \propto p(x|y)p(y)$. (For simplicity, we will assume that we can compute and average over the posterior $p(y|x)$, i.e. $p(x, y)$ is in a tractable family – we will discuss a variational extension of this case at a later stage). Effectively, this setting indicates equivalence of the inference under \mathcal{M}_I and \mathcal{M}_L for identical parameter settings; furthermore, both the likelihood³ \mathcal{L} and the mutual information $I(x, y)$ will be functionals of the same distributions, $p(y)$ and $p(x|y)$.

To summarize, the distributions defined by the generative and encoder models for the considered case would be given by

$$\mathcal{M}_L \stackrel{\text{def}}{=} p(y)p(x|y), \quad \mathcal{M}_I \stackrel{\text{def}}{=} \tilde{p}(x)p(y|x), \quad \tilde{p}(x) \propto \sum_{i=1}^M \delta(x - x^{(i)}) \quad (3.2)$$

respectively. For these otherwise unconstrained settings, we will be interested in finding a simple relation between maximization of the mutual information $I(x, y)$ in the channel \mathcal{M}_I and maximization of the log-likelihood \mathcal{L} in the corresponding latent variable model \mathcal{M}_L . In the rest of this section, by referring to a specific objective function (I or \mathcal{L}) we will also implicitly refer to the corresponding model.

There are obvious practical implications of relating these two modeling methods. Indeed, let us assume that we are given parametric constraints on $p(y)$ and $p(x|y)$, and a set of training patterns $\{x\}$. If our goal is to find the most informative latent variable representations $\{y\}$ of the training set $\{x\}$, an obvious approach would be to maximize the mutual information in \mathcal{M}_I . However, as discussed in Section 1.3, the exact optimization would typically be intractable, and approximations or tractable lower bounds on $I(x, y)$ would need to be considered. Such tractable bounds may prove to be weak, which may nullify the initial advantage of trying to optimize a proper measure of information in the first place. In contrast, it is usually significantly easier to maximize the likelihood in the corresponding generative model, and classic probabilistic methods of maximum-likelihood training may eventually prove to be better (in terms of the retained information content) than optimizing any tractable relaxation of $I(x, y)$. We will now show that this is generally *not* the case (at least, not in terms of relaxations of the exact mutual information), and we may indeed consider optimizing a special tractable form of the variational lower bound on $I(x, y)$ in order to maximize the amount of information which $\{y\}$ contain about $\{x\}$.

3.2.1 Variational Information Maximization and Exact Likelihood Training

Let us assume that the data is i.i.d., and the models are defined as in expression (3.2). Then it is easy to show that the exact likelihood, expressed from the

³We extend Fisher’s view of the likelihood (Fisher (1925)) and treat it as a functional of the model $p(x)$ under the specified modeling constraints. In other words, for us it would make perfect sense to integrate the likelihood over distributions of random variables. To indicate the dependence on the functional parameters, we will sometimes denote the likelihood as $\mathcal{L}(p(x))$.

generative model \mathcal{M}_L , is a lower bound on the exact mutual information in the corresponding memoryless communication channel \mathcal{M}_I (up to irrelevant constants which have no effects on the optimization). It is also straight-forward to see that for any tractable choice of \mathcal{M}_L , it is possible to define a tractable lower bound on $I(\mathbf{x}, \mathbf{y})$, which is *at least* as tight as the likelihood. Intuitively, this means that by maximizing the likelihood in generative models, we indeed maximize a proper lower bound on the mutual information in the corresponding channel, but we do it suboptimally in terms of the retained information content. This also suggests a tractable information-theoretic alternative to maximum-likelihood training.

Proposition 3.1. *For i.i.d. patterns $\{\mathbf{x}\}$, maximum likelihood training in the generative model \mathcal{M}_L corresponds to maximization of a **lower bound** on the mutual information in \mathcal{M}_I . Up to irrelevant constants, this bound is weaker or as tight as $\hat{I}(\mathbf{x}, \mathbf{y}) = \langle \log p(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})}$.*

Proof. The results follow immediately from applying straight-forward transformations of the likelihood. Let $\{\mathbf{x}\} \stackrel{\text{def}}{=} \{\mathbf{x}^{(i)} | i = 1, \dots, M\}$ denote a set of training patterns. For i.i.d. data, the average log-likelihood is given by

$$\begin{aligned} \mathcal{L} \stackrel{\text{def}}{=} \log p(\{\mathbf{x}\})/M &= \langle \log p(\mathbf{x}) \rangle_{\tilde{p}(\mathbf{x})} \\ &= \langle \log p(\mathbf{y}) + \log p(\mathbf{x}|\mathbf{y}) - \log p(\mathbf{y}|\mathbf{x}) \rangle_{\tilde{p}(\mathbf{x})}, \end{aligned} \quad (3.3)$$

where $\tilde{p}(\mathbf{x})$ is the empirical distribution, and \mathbf{y} is a latent variable in the joint distribution $p(\mathbf{x}, \mathbf{y})$. By construction, $p(\mathbf{y})$ and $p(\mathbf{x}|\mathbf{y})$ are parts of the model's specification which need to be learned, and the exact posterior $p(\mathbf{y}|\mathbf{x})$ is simply given by Bayes rule. By averaging both parts of (3.3) over $p(\mathbf{y}|\mathbf{x})$, we obtain

$$\mathcal{L} = \langle \log p(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})} - \langle KL(p(\mathbf{y}|\mathbf{x}) \| p(\mathbf{y})) \rangle_{\tilde{p}(\mathbf{x})}. \quad (3.4)$$

Note that the non-negative Kullback-Leibler divergence term in (3.4) would have reduced to the mutual information (expressed from \mathcal{M}_L), had the averages over the empirical distribution $\tilde{p}(\mathbf{x})$ been replaced by the averages over the true model-based marginal $p(\mathbf{x})$. (Indeed, by averaging (3.4) over $p(\mathbf{x})$, we would recover the well-known mean of average log-likelihood for i.i.d. data $\langle \mathcal{L}(\mathbf{x}) \rangle_{p(\mathbf{x})} = -H(\mathbf{x})$).

We will now express $I(\mathbf{x}, \mathbf{y})$ in the corresponding model \mathcal{M}_I of a noisy memoryless channel. It is easy to see that for the considered specification (3.2), the exact mutual information in \mathcal{M}_I is given as

$$I(\mathbf{x}, \mathbf{y}) = H_{\tilde{p}}(\mathbf{x}) - H_{\tilde{p}}(\mathbf{x}|\mathbf{y}) = H_{\tilde{p}}(\mathbf{x}) + \langle \log \tilde{p}(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})}, \quad (3.5)$$

where $H_{\tilde{p}}(\mathbf{x}) \stackrel{\text{def}}{=} -\langle \log \tilde{p}(\mathbf{x}) \rangle_{\tilde{p}(\mathbf{x})}$, $H_{\tilde{p}}(\mathbf{x}|\mathbf{y}) \stackrel{\text{def}}{=} -\langle \log \tilde{p}(\mathbf{x}|\mathbf{y}) \rangle_{\tilde{p}(\mathbf{x}, \mathbf{y})}$, and $\tilde{p}(\mathbf{x}|\mathbf{y})$ is the Bayes-optimal decoder for the training patterns

$$\tilde{p}(\mathbf{x}|\mathbf{y}) \propto \tilde{p}(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x}) = p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x}), \quad (3.6)$$

(as $\tilde{p}(\mathbf{y}|\mathbf{x}) \equiv p(\mathbf{y}|\mathbf{x})$ by construction). By definition, $\tilde{p}(\mathbf{x})$ has a mixture form, which leads to a mixture form for the decoder $\tilde{p}(\mathbf{x}|\mathbf{y})$. Generally, this complicates evaluations of the conditional entropy $H_{\tilde{p}}(\mathbf{x}|\mathbf{y})$ in (3.5), which may be intractable

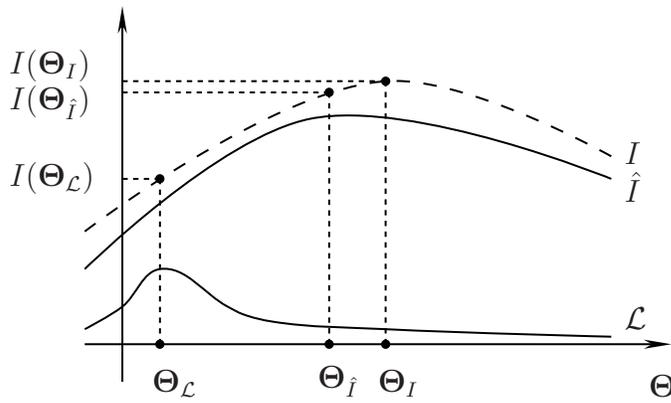


Figure 3.2: Mutual information $I(x, y)$, the variational lower bound $\hat{I}(x, y)$, and the exact log-likelihood score $\mathcal{L}(\mathcal{M}_L)$ in the parameter space $\Theta = \{p(y), p(x|y)\}$ (a schematic plot). Intuitively, optimization of the tighter bound $\hat{I}(x, y)$ on the generally intractable mutual information $I(x, y)$ leads to a better estimate $\Theta_{\hat{I}}$ of the I -optimal parameters Θ_I . Heuristically, the tightness of a bound may be treated as an approximate criterion of optimality.

even for the case of a tractable channel distribution $p(y|x)$. By applying the variational lower bound on the mutual information (2.2), we may bound (3.5) as

$$I(x, y) \geq H_{\tilde{p}}(\mathbf{x}) + \langle \log p(\mathbf{x}|y) \rangle_{p(y|\mathbf{x})\tilde{p}(\mathbf{x})} \quad (3.7)$$

$$\geq H_{\tilde{p}}(\mathbf{x}) + \langle \log p(\mathbf{x}|y) \rangle_{p(y|\mathbf{x})\tilde{p}(\mathbf{x})} - \langle KL(p(y|\mathbf{x})||p(y)) \rangle_{\tilde{p}(\mathbf{x})} \quad (3.8)$$

$$= H_{\tilde{p}}(\mathbf{x}) + \mathcal{L} \Rightarrow \quad (3.9)$$

$$I(x, y) \geq \hat{I}(x, y) \geq H_{\tilde{p}}(\mathbf{x}) + \mathcal{L}, \quad (3.10)$$

where the tighter lower bound on the mutual information is given by

$$\hat{I}(x, y) \stackrel{\text{def}}{=} H_{\tilde{p}}(\mathbf{x}) + \langle \log p(\mathbf{x}|y) \rangle_{p(y|\mathbf{x})\tilde{p}(\mathbf{x})}. \quad (3.11)$$

Note that (3.11) is just the generic lower bound on the mutual information (2.2) in \mathcal{M}_I , with the decoder given by the generative conditional $p(\mathbf{x}|y)$. Additionally, since $\tilde{p}(\mathbf{x})$ is the empirical distribution, the entropic term $H_{\tilde{p}}(\mathbf{x})$ in (3.10) and (3.11) has no effect on the optimization surface for the encoder $p(y|\mathbf{x})$. Thus, up to irrelevant constants, the exact log-likelihood score \mathcal{L} defines a proper lower bound on $I(x, y)$ in the recognition model \mathcal{M}_I . This bound is weaker than the variational relaxation given by $\hat{I}(x, y)$. \square

Proposition 3.1 shows that by maximizing the likelihood in \mathcal{M}_L we indeed optimize a proper lower bound on the generally intractable mutual information between the source patterns $\{\mathbf{x}^{(i)}|i = 1, \dots, M\}$ and their encodings $\{\mathbf{y}^{(i)}|i = 1, \dots, M\}$ generated by the exact posterior $p(y|\mathbf{x})$. It also suggests that the likelihood bound on $I(x, y)$ may in general be weak, which follows from the non-negativity of the Kullback-Leibler term in (3.8). Intuitively, the weakness of the bound may cause significant differences between solutions obtained by maximizing the likelihood and the mutual information (see the schematic plot of Figure

3.2). More importantly, proposition 3.1 suggests that there exists a proper variational lower bound on the mutual information, which is (a): tractable (whenever $p(\mathbf{x}, \mathbf{y})$ is in a tractable family); (b): tighter than the likelihood. This bound may be used for information-theoretic training of *both* the encoder and the generative models specified by expression (3.2). Specifically, for training the generative model \mathcal{M}_L , the idea would be to optimize (3.11) alone, which defines a tighter bound on the exact mutual information. We may consider optimizing the same bound⁴ for the encoder model \mathcal{M}_I . In fact, it is irrelevant whether we choose to optimize (3.11) in \mathcal{M}_I or in \mathcal{M}_L , as in both cases we would be looking at the same objective functional, optimized for the same set of functional parameters.

Note that the non-constant term in the variational lower bound $\hat{I}(\mathbf{x}, \mathbf{y})$ may be thought of as an empirical estimate of the conditional entropy $H(\mathbf{x}|\mathbf{y})$, as expressed from \mathcal{M}_L . Moreover, from (3.10) and (3.11) we get

$$\hat{I}(\mathbf{x}, \mathbf{y}) = \mathcal{L} + \langle KL(p(\mathbf{y}|\mathbf{x})\|p(\mathbf{y})) \rangle_{\tilde{p}(\mathbf{x})} + H_{\tilde{p}}(\mathbf{x}), \quad (3.12)$$

$$= \langle KL(p(\mathbf{y}|\mathbf{x})\|p(\mathbf{y})) \rangle_{\tilde{p}(\mathbf{x})} - KL(\tilde{p}(\mathbf{x})\|p(\mathbf{x})) \quad (3.13)$$

where the Kullback-Leibler term explicitly favours large *average* deviations of the posteriors $p(\mathbf{y}|\mathbf{x})$ from the priors $p(\mathbf{y})$ for a fixed set of training patterns $\{\mathbf{x}\}$. (Clearly, the term is *minimized* when the encodings \mathbf{y} are independent of the sources \mathbf{x}). Note that \hat{I} - and \mathcal{L} -learning are equivalent when the Kullback-Leibler term in (3.12) is independent of the optimized parameters.

Informally, we can see that optimization of $\hat{I}(\mathbf{x}, \mathbf{y})$ corresponds to maximization of the exact log-likelihood, with a certain bias towards deterministic and spread-out encoded representations. Then by expanding (3.12), we get

$$\hat{I}(\mathbf{x}, \mathbf{y}) = \mathcal{L} - \langle \log p(\mathbf{y}) \rangle_{\tilde{p}(\mathbf{y})} - \sum_{i=1}^M H(p(\mathbf{y}|\mathbf{x}^{(i)}))/M + H_{\tilde{p}}(\mathbf{x}), \quad (3.14)$$

where $\tilde{p}(\mathbf{y}) \stackrel{\text{def}}{=} \langle p(\mathbf{y}|\mathbf{x}) \rangle_{\tilde{p}(\mathbf{x})}$ is just the empirical average of the posterior (corresponding to the marginal distribution of the codes in the recognition model \mathcal{M}_I). The last term in (3.14) is an empirical average of the entropies of $p(\mathbf{y}|\mathbf{x}^{(i)})$, which favour deterministic encodings (i.e. sharp posteriors in generative models). To show that the objective $\hat{I}(\mathbf{x}, \mathbf{y})$ indeed favours spread-out representations in $\{\mathbf{y}\}$, we may once again make use of the non-negativity of the KL-divergence

$$\langle \log \tilde{p}(\mathbf{y}) \rangle_{\tilde{p}(\mathbf{y})} - \langle \log p(\mathbf{y}) \rangle_{\tilde{p}(\mathbf{y})} \geq 0 \Rightarrow -\langle \log p(\mathbf{y}) \rangle_{\tilde{p}(\mathbf{y})} \geq H(\tilde{p}(\mathbf{y})). \quad (3.15)$$

Substituting (3.15) into the objective (3.14), we obtain a relaxation of $\hat{I}(\mathbf{x}, \mathbf{y})$, given by

$$\hat{I}(\mathbf{x}, \mathbf{y}) \geq \mathcal{L} + H(\tilde{p}(\mathbf{y})) - \sum_{i=1}^M H(p(\mathbf{y}|\mathbf{x}^{(i)}))/M + H(\tilde{p}(\mathbf{x})). \quad (3.16)$$

⁴Despite the fact that any tractable choice of the variational decoder $q(\mathbf{x}|\mathbf{y})$ could potentially be used for training $\mathcal{M}_I = \tilde{p}(\mathbf{x})p(\mathbf{y}|\mathbf{x})$ (see expression (2.2)), not every such decoder may lead to an improvement over the bound provided by the likelihood score. From proposition 3.1 we see that one choice of the decoder when the improvement does happen is to set $q(\mathbf{x}|\mathbf{y}) = p(\mathbf{x}|\mathbf{y})$.

Clearly, the entropy of $\tilde{p}(\mathbf{y}) = \sum_{i=1}^M p(\mathbf{y}|\mathbf{x}^{(i)})/M$ in (3.16) would be maximized for uniform representations of the training patterns in the code space $\{\mathbf{y}\}$. Since a proper lower bound on the objective criterion $\hat{I}(\mathbf{x}, \mathbf{y})$ biases the solutions towards spread-out codes, one may informally hope that the objective $\hat{I}(\mathbf{x}, \mathbf{y})$ will establish a similar behaviour. As we mentioned earlier, this bias towards a uniform distribution of the latent variables may be useful (for example, for avoiding learning non-informative representations due to singularities and degeneracies of maximum-likelihood solutions).

Example: It is easy to see that there are interesting cases when optimization of the exact mutual information $I(\mathbf{x}, \mathbf{y})$ in the encoder model \mathcal{M}_I is intractable *despite* the tractability of both $p(\mathbf{y}|\mathbf{x})$ and $p(\mathbf{x}|\mathbf{y})$. For a quick illustration of this, let \mathcal{M}_L define the factor analysis model (e.g. Bartholomew (1987), Ghahramani and Hinton (1996)), given as $\mathbf{x} = \mathbf{W}\mathbf{y} + \mathbf{e}$; $\mathbf{y} \sim \mathcal{N}(0, \mathbf{1})$, $p(\mathbf{e}) \sim \mathcal{N}(0, \mathbf{\Psi})$. Clearly, for this model we get $p(\mathbf{x}|\mathbf{y}) \sim \mathcal{N}(\mathbf{W}\mathbf{y}, \mathbf{\Psi})$, $p(\mathbf{x}) \sim \mathcal{N}(0, \mathbf{W}\mathbf{W}^T + \mathbf{\Psi})$, and $p(\mathbf{y}|\mathbf{x}) \sim \mathcal{N}((\mathbf{I} + \mathbf{W}^T\mathbf{\Psi}^{-1}\mathbf{W})^{-1}\mathbf{W}^T\mathbf{\Psi}^{-1}\mathbf{x}, (\mathbf{I} + \mathbf{W}^T\mathbf{\Psi}^{-1}\mathbf{W})^{-1})$. Note that for the corresponding encoder model \mathcal{M}_I , it is intractable to optimize the exact mutual information (3.5), as this would require evaluation of the entropy of a mixture of Gaussians. On the other hand, it is computationally tractable to optimize the variational lower bound $\hat{I}(\mathbf{x}, \mathbf{y})$ on the mutual information, as this would only require computing Gaussian and empirical averages of the quadratic terms $\langle \mathbf{y}\mathbf{x}^T \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})}$, $\langle \mathbf{y}\mathbf{y}^T \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})}$. By substituting $p(\mathbf{x}|\mathbf{y})$ and $p(\mathbf{y}|\mathbf{x})$ into the expressions for \mathcal{L} and $\hat{I}(\mathbf{x}, \mathbf{y})$, we obtain generally different objective criteria. In the special case of $\mathbf{\Psi} = \sigma^2\mathbf{I}$, the relation between the \hat{I} - and \mathcal{L} -optimal weights may be easily established analytically; in fact, for this special case both objectives would give rise to PCA solutions for \mathbf{W} .

Another example, illustrating differences between likelihood and information-theoretic training, is discussed in Section 5.2.1. There we compare \mathcal{L} - and \hat{I} -maximization for training Gaussian mixture models.

3.2.2 Variational Information Maximization and Variational Likelihood Training

We have shown that optimization of the exact likelihood in the generative model $\mathcal{M}_L = p(\mathbf{y})p(\mathbf{x}|\mathbf{y})$ may be motivated in the variational information-theoretic sense (though the likelihood bounds on $I(\mathbf{x}, \mathbf{y})$ may potentially be weak). We have also discussed a simple re-definition of the objective criterion which would lead to a tighter variational lower bound on the mutual information in the corresponding encoder model \mathcal{M}_I . These theoretical results are intuitive, but limited to the cases when it is possible to compute and average over the posterior $p(\mathbf{y}|\mathbf{x})$. However, for more general latent variable models, it may be intractable to maximize the likelihood \mathcal{L} or the bound $\hat{I}(\mathbf{x}, \mathbf{y})$ directly.

The usual way to handle the intractability of exact maximum likelihood learning in the generative model \mathcal{M}_L is to consider the standard variational lower

bound on the likelihood

$$\begin{aligned}\mathcal{L} &= \left\langle \log \int_{\tilde{p}(\mathbf{x})} p(\mathbf{x}, \mathbf{y}) d\mathbf{y} \right\rangle_{\tilde{p}(\mathbf{x})} \\ &= \left\langle \log \int_{\tilde{p}(\mathbf{x})} p(\mathbf{x}, \mathbf{y}) \frac{q(\mathbf{y}|\mathbf{x})}{q(\mathbf{y}|\mathbf{x})} d\mathbf{y} \right\rangle_{\tilde{p}(\mathbf{x})} \geq \tilde{\mathcal{L}},\end{aligned}$$

where we have defined

$$\begin{aligned}\tilde{\mathcal{L}} &= \langle \log p(\mathbf{x}, \mathbf{y}) \rangle_{q(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})} - \langle \log q(\mathbf{y}|\mathbf{x}) \rangle_{q(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})} \\ &= \langle \log p(\mathbf{x}|\mathbf{y}) \rangle_{q(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})} - \langle KL(q(\mathbf{y}|\mathbf{x})||p(\mathbf{y})) \rangle_{\tilde{p}(\mathbf{x})}.\end{aligned}\tag{3.17}$$

Here $\tilde{p}(\mathbf{x})$ is the empirical distribution, and $q(\mathbf{y}|\mathbf{x})$ is an arbitrary variational distribution approximating the true posterior $p(\mathbf{y}|\mathbf{x})$ (see Zemel and Hinton (1994), Zemel and Hinton (1995), Dayan et al. (1995) for the related *Helmholtz machine* formulation, and e.g. Saul et al. (1996), Jaakkola (1997) for the general statement of the variational problem). The standard variational extensions of the expectation-maximizing algorithm (see e.g. Neal and Hinton (1998)) train the generative model \mathcal{M}_L by iteratively optimizing (3.17) with respect to the model parameters $p(\mathbf{y})$, $p(\mathbf{x}|\mathbf{y})$, and the variational posterior $q(\mathbf{y}|\mathbf{x})$. In order to simplify computations of the bound (3.17), the variational distribution $q(\mathbf{y}|\mathbf{x})$ is constrained to ensure tractability of computing the average energy $\langle \log p(\mathbf{x}, \mathbf{y}) \rangle_{q(\mathbf{y}|\mathbf{x})}$, which for directed latent variable models implies tractability of computing $\langle \log p(\mathbf{x}|\mathbf{y}) \rangle_{q(\mathbf{y}|\mathbf{x})}$ and $\langle \log p(\mathbf{y}) \rangle_{q(\mathbf{y}|\mathbf{x})}$.

In order to compare variational approaches to maximization of the likelihood and mutual information, we will define a recognition model

$$\tilde{\mathcal{M}}_I \stackrel{\text{def}}{=} q(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x}),\tag{3.18}$$

which ensures equivalence of the approximate *variational* inference⁵ in \mathcal{M}_L and $\tilde{\mathcal{M}}_I$. Then by analogy with proposition 3.1, it is easy to see that the lower bound $\tilde{\mathcal{L}}$ on the likelihood in \mathcal{M}_L is in fact a proper lower bound on the mutual information in $\tilde{\mathcal{M}}_I$. Moreover, we can easily find a tractable lower bound on $I(\mathbf{x}, \mathbf{y})$ which is tighter than $\tilde{\mathcal{L}}$.

Proposition 3.2. *For i.i.d. patterns $\{\mathbf{x}\}$, maximization of the standard variational lower bound on the likelihood in the generative model \mathcal{M}_L corresponds to maximization of a lower bound on the mutual information in $\tilde{\mathcal{M}}_I$. Up to irrelevant constants, this bound is weaker or as tight as $\hat{I}_q(\mathbf{x}, \mathbf{y}) = \langle \log p(\mathbf{x}|\mathbf{y}) \rangle_{q(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})}$, where $q(\mathbf{y}|\mathbf{x})$ is the approximate posterior of the generative model.*

Proof. By definition, the exact value of mutual information $I(\mathbf{x}, \mathbf{y})$ for model $\tilde{\mathcal{M}}_I$ is given by

$$I(\mathbf{x}, \mathbf{y}) = H_{\tilde{p}}(\mathbf{x}) + \langle \log \tilde{p}(\mathbf{x}|\mathbf{y}) \rangle_{q(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})},\tag{3.19}$$

⁵Since computing the exact posterior $p(\mathbf{y}|\mathbf{x})$ in \mathcal{M}_L will now be intractable, we will need approximations to perform the inference. The choice of the variational posterior $q(\mathbf{y}|\mathbf{x})$ as a model for the encoder ensures tractability and helps to establish a relationship between variational likelihood and mutual information maximization under the assumption of equivalence of approximate inference.

where $\tilde{p}(\mathbf{x}|\mathbf{y}) \propto \tilde{p}(\mathbf{x})q(\mathbf{y}|\mathbf{x})$ is the exact posterior expressed from the encoder model. Then, by analogy with the proof of proposition 3.1, we get

$$I(\mathbf{x}, \mathbf{y}) \geq H_{\tilde{p}}(\mathbf{x}) + \langle \log p(\mathbf{x}|\mathbf{y}) \rangle_{q(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})} - \langle KL(q(\mathbf{y}|\mathbf{x})||p(\mathbf{y})) \rangle_{\tilde{p}(\mathbf{x})} \Rightarrow \quad (3.20)$$

$$I(\mathbf{x}, \mathbf{y}) \geq \hat{I}_q(\mathbf{x}, \mathbf{y}) \geq H_{\tilde{p}}(\mathbf{x}) + \tilde{\mathcal{L}}, \quad (3.21)$$

where $\tilde{\mathcal{L}}$ is defined as in (3.17), and the tighter lower bound on $I(\mathbf{x}, \mathbf{y})$ is defined as

$$\hat{I}_q(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} H_{\tilde{p}}(\mathbf{x}) + \langle \log p(\mathbf{x}|\mathbf{y}) \rangle_{q(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})}, \quad (3.22)$$

which is tractable by construction. \square

Clearly, the relation (3.22) between variational approaches to mutual information and likelihood maximization in $\tilde{\mathcal{M}}_I$ and \mathcal{M}_L is analogous to the similar relation (3.10) for the exact likelihood computations. Effectively, proposition 3.2 states that standard variational approaches to likelihood maximization may indeed be seen as a way to optimize a specific variational lower bound on $I(\mathbf{x}, \mathbf{y})$ in the corresponding encoder model $\tilde{\mathcal{M}}_I$. Again, a tighter (yet tractable) bound would be given by $\hat{I}_q(\mathbf{x}, \mathbf{y})$, which is exactly the generic variational lower bound on the mutual information in $\tilde{\mathcal{M}}_I$ for the specific parameterization of the decoder distribution. Moreover, from (3.20) it is clear that maximization of the lower bound $\hat{I}_q(\mathbf{x}, \mathbf{y})$ on the mutual information and the variational Jensen’s bound $\tilde{\mathcal{L}}$ on the likelihood lead to the same fixed points when the Kullback-Leibler divergence in (3.20) is a constant (also see discussion in Appendix B.2).

It is interesting to note that by optimizing the bound on the likelihood $\tilde{\mathcal{L}}$ for the latent space prior $p(\mathbf{y})$ and substituting the optimal functional parameters back into (3.20), we may transform (3.20) to

$$I(\mathbf{x}, \mathbf{y}) \geq \underbrace{H_{\tilde{p}}(\mathbf{x}) + \langle \log p(\mathbf{x}|\mathbf{y}) \rangle_{q(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})}}_{\hat{I}_q(\mathbf{x}, \mathbf{y})} - \langle KL(q(\mathbf{y}|\mathbf{x})||\langle q(\mathbf{y}|\mathbf{x}) \rangle_{\tilde{p}(\mathbf{x})}) \rangle_{\tilde{p}(\mathbf{x})} \quad (3.23)$$

$$\geq H_{\tilde{p}}(\mathbf{x}) + \tilde{\mathcal{L}}. \quad (3.24)$$

(Here the second inequality (3.24) reduces to an identity if and only if the current settings of $p(\mathbf{y})$ in the generative model \mathcal{M}_L are indeed $\tilde{\mathcal{L}}$ -optimal). Note that the Kullback-Leibler term in (3.23) is generally nonzero. This results in a relaxation of the variational bound on $I(\mathbf{x}, \mathbf{y})$, i.e. apart from simple special cases⁶, the bound $H_{\tilde{p}}(\mathbf{x}) + \tilde{\mathcal{L}}$ will be strictly weaker than $\hat{I}_q(\mathbf{x}, \mathbf{y})$.

Informally, from (3.20) we can see that optimization of the bound on the likelihood in \mathcal{M}_L could reduce to optimizing the generic lower bound $\hat{I}_q(\mathbf{x}, \mathbf{y})$ on the mutual information in $\tilde{\mathcal{M}}_I$, if the marginal distribution $p(\mathbf{y}^{(i)})$ of each latent vectors $\mathbf{y}^{(i)}$ was an independent free parameter. Indeed, the variational bound on

⁶It is easy to see that $\hat{I}_q(\mathbf{x}, \mathbf{y}) = H_{\tilde{p}}(\mathbf{x}) + \tilde{\mathcal{L}}$ may be achieved when $q(\mathbf{y}|\mathbf{x}) = q(\mathbf{y})$, i.e. both the likelihood and the information-theoretic bound use an extremely weak definition of the encoding distribution. Similarly, both bounds are identical if $M = 1$, i.e. there is a single pattern to encode. A brief discussion of other special cases when both bounds lead to identical optima is given in Appendix B.2.

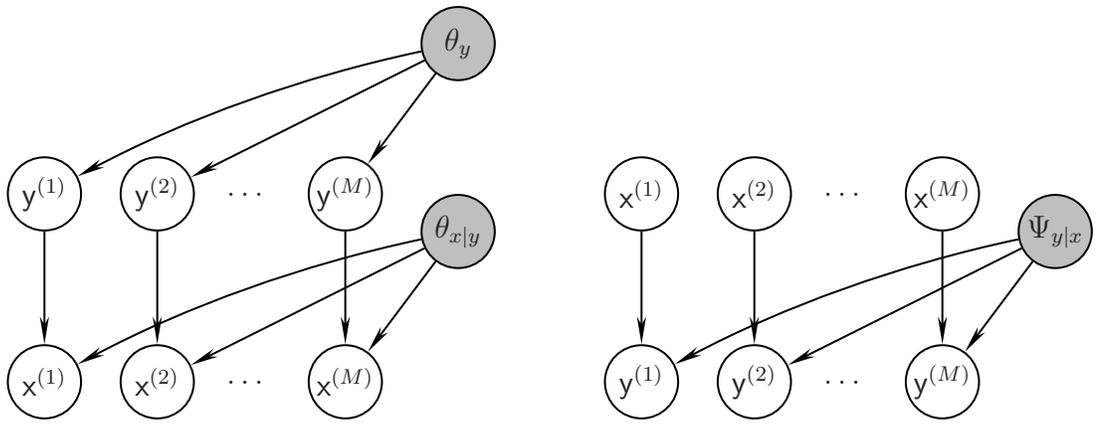


Figure 3.3: Generative and encoder models for i.i.d. patterns. *Left:* generative model \mathcal{M}_L for M i.i.d. training patterns $\{x^{(1)}, \dots, x^{(M)}\}$ with the corresponding latent variable representations $\{y^{(1)}, \dots, y^{(M)}\}$; *Right:* the corresponding model of the memoryless channel \mathcal{M}_I . The shaded nodes indicate the implied conditionings on the intrinsically deterministic model parameters.

the likelihood (3.17) could in this case be expressed as

$$\mathcal{L} \propto \sum_{i=1}^M \langle \log p(x^{(i)}|y) \rangle_{q(y|x^{(i)})} - \sum_{i=1}^M KL(q(y|x^{(i)}) || p(y^{(i)})). \quad (3.25)$$

If for all training patterns $i = 1, \dots, M$ we were free to change each latent-space marginal individually, we could optimally set it to the variational posterior as

$$p(y^{(i)} = y) \stackrel{\text{def}}{=} p(y^{(i)} = y | \Theta_y^{(i)}) = q(y|x^{(i)}, \Psi_{y|x}) \stackrel{\text{def}}{=} q(y|x^{(i)}), \quad (3.26)$$

which would lead to $\tilde{\mathcal{L}} = \hat{I}_q(x, y)$ at the optimum of (3.25). It is clear that (3.26) is exactly the definition of the latent-space marginals $\tilde{p}(y^{(i)})$ in the encoder model $\tilde{\mathcal{M}}_I$ (see Figure 3.3 (right)). It is also clear that (3.26) does not generally apply for the generative model \mathcal{M}_L and i.i.d observations, where

$$\forall \Theta_y. \forall i, j \in \{1, \dots, M\}. p(y^{(i)} = y | \Theta_y) = p(y^{(j)} = y | \Theta_y) \quad (3.27)$$

(see Figure 3.3 (left)).

Generally, we can say that learning by optimizing the exact likelihood (or standard variational lower bounds on the likelihood) in generative models could be interpreted as optimization of a specific lower bound on the generally intractable mutual information in the corresponding encoder models (3.2), (3.18). However, the optimization surfaces and the obtained solutions will generally be different from those given by the variational approach to information-maximization in the corresponding models of memoryless channels. This difference in the objectives is quantified as a summation of generally non-constant KL divergences between the posteriors $p(y|x^{(i)})$ and their empirical average $\langle p(y|x) \rangle_{\tilde{p}(x)}$, as expressed from \mathcal{M}_L for i.i.d. data. Arguably, one of the principal causes of the difference is the

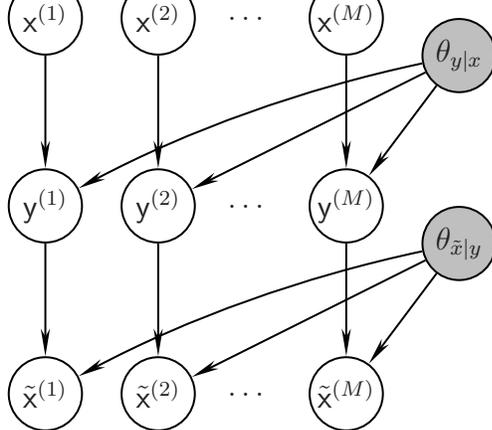


Figure 3.4: An autoencoder model for M training patterns $\{x^{(1)}, \dots, x^{(M)}\}$ and their reconstructions $\{\tilde{x}^{(1)}, \dots, \tilde{x}^{(M)}\}$. The shaded nodes indicate the conditionings on the deterministic model parameters.

identical distribution of the latent vectors in \mathcal{M}_L , which is implied by the specific form of the prior on the latent vectors. We may further notice a closer correspondence between optimizing the mutual information in encoder models (Figure 3.3 (right)) and the conditional likelihood training in autoencoders (see Figure 3.4), which are characterized by generally different marginal distributions of the encodings, and the model specification more similar to the encoder paradigm. Presuming that our intuition here is correct, we may perceive the existence of optimization procedures where the conditional training of autoencoders reduces to the variational information maximization.

3.3 Information Maximization and Maximum Conditional Likelihood

A classic method of obtaining informative lower-dimensional representations of higher-dimensional source vectors is to consider an autoencoder with the parameterized mappings from the source vectors to the codes ($x \rightarrow y$), and from the codes to the reconstructions of the source vectors ($y \rightarrow \tilde{x}$) (Hinton (1989), Baldi and Hornik (1989), Bishop (1995)). The goal is to learn parameters of encoder and decoder mappings which would lead to accurate reconstructions of the sources from their encoded representations. Most of the conventional approaches to training such models (e.g. minimization of the mean squared reconstruction error, minimization of the cross-entropy for multinomial vectors, etc.) are equivalent to maximizing the conditional likelihood $p(\{\tilde{x}\}|\{x\})$ in generalized chains $x \rightarrow y \rightarrow \tilde{x}$ under specific parametric constraints on the encoder and decoder distributions. Several other methods for optimizing parameters of such networks may be shown to correspond to variational likelihood learning in mixture models (Zemel (1993), Zemel and Hinton (1994)), while in some other cases the optimized objectives do not seem to have a clear probabilistic interpretation (Doi and Lewicki, 2004).

Here we will focus primarily on the discussion of conventional conditional training of feed-forward models, and the relations which the resulting objective functions may have to variational bounds on the mutual information.

While it is commonly presumed that the encoding part of the autoencoder is noiseless (i.e. $p(\mathbf{y}|\mathbf{x}) \sim \delta$), probabilistic extensions may easily be considered. By analogy with the marginal likelihood training of arbitrarily structured graphical models, a proper probabilistic method of training stochastic autoencoders would involve integration of the hidden encodings and maximization of the likelihood with respect to deterministic parameters of the model subject to specific constraints. Typically, the goal of learning would be to find the optimal encoder $p(\mathbf{y}|\mathbf{x})$ and decoder $p(\tilde{\mathbf{x}}|\mathbf{y})$ distributions (which effectively reduces to maximizing the conditional likelihood). Clearly, this formulation is somewhat similar to variational approaches to mutual information maximization, where the objective function is optimized with respect to the encoder and variational decoder. Moreover, as we have mentioned earlier, we may perceive another link between mutual information and conditional likelihood maximization from Fano's inequality (3.1), which relates the mutual information with an *upper bound* on the probability of correct reconstructions. Our goal here is to explore these apparent similarities and establish a possible relation between variational information maximization and conditional likelihood training in generalized chains independently of specific parameterizations.

We will begin the discussion by outlining a general relation between conditional likelihood and conditional mutual information maximization in general feed-forward models $\mathbf{x} \rightarrow \mathbf{y} \rightarrow \tilde{\mathbf{x}}$. Clearly, this view extends the reconstruction paradigm of autoencoders, since for general chains the sources \mathbf{x} and the outputs $\tilde{\mathbf{x}}$ may generally lie in different domains. By analogy with the results of section 3.2, we will show that for i.i.d. patterns, the exact conditional likelihood for the considered chains defines a proper lower bound on the *conditional* mutual information in memoryless channels. Effectively, this formulation corresponds to maximizing the amount of information which the encodings contain about the *outputs* for the given set of source patterns $I(\{\tilde{\mathbf{x}}\}, \{\mathbf{y}\}|\{\mathbf{x}\})$.

Then we will consider a special case of stochastic autoencoders, where the empirical distribution is constrained so that $\tilde{p}(\tilde{\mathbf{x}}|\mathbf{x}) \sim \delta(\mathbf{x} - \tilde{\mathbf{x}})$, i.e. the decoded patterns $\tilde{\mathbf{x}}$ are the exact uncorrupted copies of the sources \mathbf{x} . For this specific case, it is easy to establish a simple relation between conditional likelihood training in stochastic autoencoders and maximization of the mutual information in the channel model $\mathbf{x} \rightarrow \mathbf{y}$. Specifically, we will show that the exact mutual information $I(\mathbf{x}, \mathbf{y})$ is a proper *lower bound* on the probability of correct reconstructions in stochastic autoencoders where the reconstructing distribution is given by the exact posterior. Moreover, we will show that if the encoding mappings of such autoencoders are noiseless, maximum conditional likelihood training is equivalent to maximizing the exact mutual information in the corresponding noiseless channel.

Finally, we will consider the situation when the conditional likelihood cannot be computed exactly, and show that the standard variational approaches to maximizing the conditional likelihood in stochastic autoencoders reduce to a specific

instance of the generic IM algorithm.

3.3.1 Variational Information Maximization and Exact Conditional Likelihood Training in Feed-Forward Models

To establish a relationship between conditional likelihood and mutual information training, we will now consider a simple feed-forward model, which defines the joint distribution over the latent variables \mathbf{y} and outputs $\tilde{\mathbf{x}}$ conditionally on the sources \mathbf{x} , i.e.

$$\mathcal{M}_C \stackrel{\text{def}}{=} p(\mathbf{y}|\mathbf{x})p(\tilde{\mathbf{x}}|\mathbf{y}). \quad (3.28)$$

The functional parameters $p(\mathbf{y}|\mathbf{x})$ and $p(\tilde{\mathbf{x}}|\mathbf{y})$ are determined by maximizing the conditional likelihood

$$\mathcal{L}_{\tilde{\mathbf{x}}|\mathbf{x}} = \log p(\{\tilde{\mathbf{x}}\}|\{\mathbf{x}\})/M = \langle \log p(\tilde{\mathbf{x}}|\mathbf{x}) \rangle_{\tilde{p}(\mathbf{x},\tilde{\mathbf{x}})}, \quad (3.29)$$

where $\tilde{p}(\tilde{\mathbf{x}}, \mathbf{x})$ is the empirical distribution of the input and output pairs. In the general formulation, we will not impose specific constraints on the exact form of the empirical distribution, so the ranges $\mathcal{R}_{\mathbf{x}}$, $\mathcal{R}_{\tilde{\mathbf{x}}}$ of the source and output variables \mathbf{x} , $\tilde{\mathbf{x}}$ may in general be different. Once the model is trained, the states of the latent variables \mathbf{y} corresponding to any given source-output pair are inferred from the exact posterior $p(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}}) \propto p(\tilde{\mathbf{x}}|\mathbf{y})p(\mathbf{y}|\mathbf{x})$.

We will now define a recognition model

$$\mathcal{M}_{IC} \stackrel{\text{def}}{=} \tilde{p}(\mathbf{x}, \tilde{\mathbf{x}})p(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}}), \quad (3.30)$$

where $p(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}})$ is the exact posterior of the feed-forward model \mathcal{M}_C . Clearly, the encoder model's specification (3.30) leads to equivalence of the exact inference in \mathcal{M}_C and \mathcal{M}_{IC} for all the visible source-output pairs $\{\mathbf{x}, \tilde{\mathbf{x}}\}$. While this formulation might not have an easily interpretable link to communication theory⁷, it does have a direct relation to conditional likelihood training in feed-forward models, which generalize the reconstruction paradigm of stochastic autoencoders. Indeed, it is easy to show that the exact conditional likelihood in \mathcal{M}_C is in fact a lower bound on the conditional mutual information $I(\tilde{\mathbf{x}}, \mathbf{y}|\mathbf{x})$ in \mathcal{M}_{IC} .

Proposition 3.3. *For i.i.d. patterns $\{\mathbf{x}, \tilde{\mathbf{x}}\}$, conditional likelihood learning in the feed-forward model \mathcal{M}_C corresponds to maximization of a **lower bound** on the conditional mutual information $I(\tilde{\mathbf{x}}, \mathbf{y}|\mathbf{x})$ in \mathcal{M}_{IC} . Up to irrelevant constants, this bound is weaker or as tight as $\hat{I}_C(\tilde{\mathbf{x}}, \mathbf{y}|\mathbf{x}) \stackrel{\text{def}}{=}} \langle \log p(\tilde{\mathbf{x}}|\mathbf{x}, \mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}})\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}})}$.*

We prove and discuss this result in Appendix B.1. Note, however, that in the special case when the chain model \mathcal{M}_C and the empirical distribution $\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}})$ define an autoencoder, the bound $\hat{I}_C(\tilde{\mathbf{x}}, \mathbf{y}|\mathbf{x})$ gives a *model-specific* upper bound on the probability of correct reconstructions

$$\mathcal{L}_{\tilde{\mathbf{x}}|\mathbf{x}} \leq \hat{I}_C(\tilde{\mathbf{x}}, \mathbf{y}|\mathbf{x}) + H_{\tilde{p}}(\tilde{\mathbf{x}}|\mathbf{x}) \quad (3.31)$$

⁷One may potentially view \mathcal{M}_{IC} as a model of the *reverse transmission path* for half-duplex channels (Glover and Grant (2003)), where the received vectors are sent back across a noisy channel whose specific properties may in general be affected by the original sources.

(cf Fano's inequality (3.1)).

We will now use the general result of proposition 3.3 to show that for several specific autoencoders, the conditional likelihood training in \mathcal{M}_C has a strong correspondence to mutual information maximization in *unconditional* encoder models $\mathcal{M}_I = \tilde{p}(\mathbf{x})p(\mathbf{y}|\mathbf{x})$, which has a stronger relation to the variational IM formulation (see Section 2.1.2).

3.3.2 Variational Information Maximization and Exact Training of Autoencoders

If $\mathcal{R}_x \equiv \mathcal{R}_{\tilde{x}}$ and the outputs $\{\tilde{x}\}$ are the exact uncorrupted copies of the sources $\{\mathbf{x}\}$ for all the training patterns $(\mathbf{x}^{(i)}, \tilde{x}^{(i)}) \in \mathcal{X}_C$, $i \in \{1, \dots, M\}$, then the feed-forward model \mathcal{M}_C reduces to a stochastic autoencoder. We will now outline a few simple relations between mutual information maximization in simple channel models \mathcal{M}_I and conditional likelihood training in autoencoders⁸ \mathcal{M}_C for several specific constraints on the encoder and decoder distributions. However, first we will prove a simple lemma which will be extensively used throughout the discussions in the rest of this chapter. Effectively, it will help us to reduce computations of average quantities in chains $\mathbf{x} \rightarrow \mathbf{y} \rightarrow \tilde{x}$ to computing specific integrals in the encoder models $\mathbf{x} \rightarrow \mathbf{y}$.

Lemma 3.1. *Let the outputs \tilde{x} be the exact copies of the sources \mathbf{x} , i.e. $\mathcal{R}_x \equiv \mathcal{R}_{\tilde{x}}$ and $\tilde{p}(\tilde{x}|\mathbf{x}) \sim \delta(\mathbf{x} - \tilde{x})$. Additionally, let $\tilde{p}(\mathbf{x}) = \sum_{i=1}^M \delta(\mathbf{x} - \mathbf{x}^{(i)})/M$ be the empirical distribution of the sources, and $q_{\mathbf{x}|\mathbf{y}}(\mathbf{x} = \mathbf{s}|\mathbf{y}) = q_{\tilde{x}|\mathbf{y}}(\tilde{x} = \mathbf{s}|\mathbf{y})$ for all patterns $\mathbf{s} \in \mathcal{R}_x$ and encodings $\mathbf{y} \in \mathcal{R}_y$. Then $\langle \mathcal{F}\{q_{\tilde{x}|\mathbf{y}}(\tilde{x}|\mathbf{y})\} \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x}, \tilde{x})} = \langle \mathcal{F}\{q_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y})\} \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})}$, where \mathcal{F} is an arbitrary \mathbf{y} -integrable functional of the conditionals $q(\tilde{x}|\mathbf{y})$ and $q(\mathbf{x}|\mathbf{y})$.*

Proof. By the direct substitution, we obtain

$$\begin{aligned} \langle \mathcal{F}\{q_{\tilde{x}|\mathbf{y}}(\tilde{x}|\mathbf{y})\} \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x}, \tilde{x})} &= \frac{1}{M} \sum_{i=1}^M \langle \mathcal{F}\{q_{\tilde{x}|\mathbf{y}}(\tilde{x}|\mathbf{y})\} \rangle_{p(\mathbf{y}|\mathbf{x}^{(i)})\tilde{p}(\tilde{x}|\mathbf{x}^{(i)})} \\ &= \frac{1}{M} \sum_{i=1}^M \int_{\mathbf{y}} p(\mathbf{y}|\mathbf{x}^{(i)}) \int_{\tilde{x}} \mathcal{F}\{q_{\tilde{x}|\mathbf{y}}(\tilde{x}|\mathbf{y})\} \delta(\tilde{x} - \mathbf{x}^{(i)}) d\tilde{x} d\mathbf{y} \\ &= \frac{1}{M} \sum_{i=1}^M \int_{\mathbf{y}} p(\mathbf{y}|\mathbf{x}^{(i)}) \mathcal{F}\{q_{\mathbf{x}|\mathbf{y}}(\mathbf{x}^{(i)}|\mathbf{y})\} d\mathbf{y} \end{aligned} \quad (3.32)$$

$$= \langle \mathcal{F}\{q_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y})\} \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})}, \quad (3.33)$$

where the identity (3.32) follows from the condition that $q_{\mathbf{x}|\mathbf{y}}(\mathbf{x} = \mathbf{s}|\mathbf{y}) = q_{\tilde{x}|\mathbf{y}}(\tilde{x} = \mathbf{s}|\mathbf{y})$. $\forall \mathbf{s} \in \mathcal{R}_x, \forall \mathbf{y} \in \mathcal{R}_y$. \square

⁸As one-layer autoencoders will have the same graphical structure and parameterization as the general chain $\mathbf{x} \rightarrow \mathbf{y} \rightarrow \tilde{x}$, we will refer to them as \mathcal{M}_C (see expression (3.28)). The only things different from the more general feed-forward models are the constraints on the empirical distribution $\tilde{p}(\mathbf{x}, \tilde{x})$ and the ranges of the source-output variables, which do not form a part of the graphical specification.

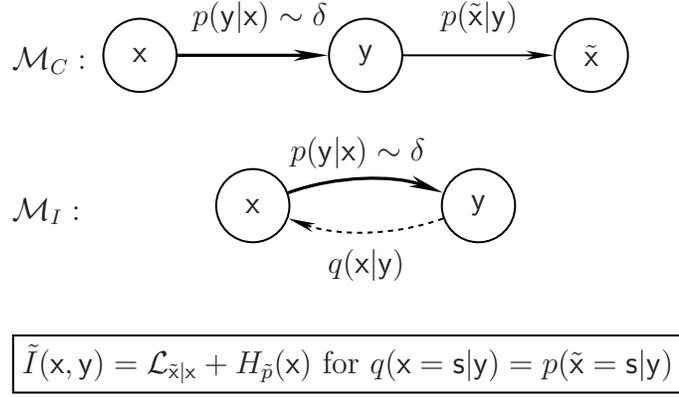


Figure 3.5: Variational information maximization and noiseless autoencoders. Maximization of the conditional likelihood in noiseless autoencoders \mathcal{M}_C reduces to maximization of the variational lower bound $\tilde{I}(x, y)$ in the corresponding noiseless channels \mathcal{M}_I , where the variational decoder is given by $q(x = s|y) = p(\tilde{x} = s|y)$. Thick arrows show deterministic mappings; the dashed arrow corresponds to the variational distribution.

Specifically, for $\mathcal{F}(q) \equiv \log(q)$, we get

$$\langle \log p(\tilde{x}|y) \rangle_{p(y|x)\tilde{p}(x,\tilde{x})} = \langle \log q(x|y) \rangle_{p(y|x)\tilde{p}(x)}, \quad (3.34)$$

where we have defined $q(x = s|y) \stackrel{\text{def}}{=} p(\tilde{x} = s|y)$ for all patterns $s \in \mathcal{R}_x \equiv \mathcal{R}_{\tilde{x}}$, $y \in \mathcal{R}_y$. We will now use the results (3.33) and (3.34) in order to relate the information-theoretic learning in channel models and the conditional learning in autoencoders.

3.3.2.1 Noiseless Autoencoders

In conventional approaches to autoencoder training, it is typically presumed that the encoding part of the autoencoder is deterministic (Hinton (1989), Baldi and Hornik (1989), Bishop (1995)), i.e.

$$p(y|x^{(i)}) \sim \delta(y - f(x^{(i)})) = \delta(y - y^{(i)}), \quad (3.35)$$

where $y^{(i)} \stackrel{\text{def}}{=} f(x^{(i)})$. With a slight abuse of terminology, we will refer to autoencoders satisfying (3.35) as being *noiseless*, rather than *stochastic*, in order to stress the determinism of the encoding mapping. It is easy to see that for the specific case (3.35), optimization of the exact conditional likelihood $\mathcal{L}_{\tilde{x}|x}$ in the autoencoder \mathcal{M}_C is equivalent to maximizing the generic lower bound on the mutual information in the corresponding *noiseless* channel \mathcal{M}_I . The model of the noiseless channel is defined similarly to expression (3.2), i.e. $\mathcal{M}_I = \tilde{p}(x)p(y|x)$.

Proposition 3.4. For i.i.d. patterns $\{\mathbf{x}\}$, exact conditional likelihood learning in **noiseless autoencoders** \mathcal{M}_C is equivalent to maximizing the generic lower bound on mutual information $\tilde{I}(\mathbf{x}, \mathbf{y}) = \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})} + H_{\tilde{p}}(\mathbf{x})$ in **noiseless channels** \mathcal{M}_I , where the variational decoder $q(\mathbf{x}|\mathbf{y})$ is constrained to be equivalent to the decoding distribution of the autoencoder.

Proof. Let $\tilde{\mathbf{x}}$ denote the reconstructed variables of the autoencoder model, so that $\mathcal{R}_{\mathbf{x}} = \mathcal{R}_{\tilde{\mathbf{x}}}$ and $\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}}) = \tilde{p}(\mathbf{x})\delta(\tilde{\mathbf{x}} - \mathbf{x})$. From expressions (3.28) and (3.35) it is clear that for noiseless autoencoders the exact posterior of the conditional model \mathcal{M}_C is given by

$$p(\mathbf{y}|\mathbf{x}^{(i)}, \tilde{\mathbf{x}}) = \frac{p(\mathbf{y}|\mathbf{x}^{(i)})p(\tilde{\mathbf{x}}|\mathbf{y})}{\int_{\mathbf{y}} p(\mathbf{y}|\mathbf{x}^{(i)})p(\tilde{\mathbf{x}}|\mathbf{y}) d\mathbf{y}} = \frac{\delta(\mathbf{y} - \mathbf{y}^{(i)})p(\tilde{\mathbf{x}}|\mathbf{y}^{(i)})}{p(\tilde{\mathbf{x}}|\mathbf{y}^{(i)})}, \quad (3.36)$$

where $\mathbf{y}^{(i)}$ is defined as in expression (3.35). Then we immediately obtain

$$p(\mathbf{y}|\mathbf{x}^{(i)}, \tilde{\mathbf{x}}) = \delta(\mathbf{y} - \mathbf{y}^{(i)}) = p(\mathbf{y}|\mathbf{x}^{(i)}), \quad (3.37)$$

i.e. for deterministic encoders, the exact posterior in \mathcal{M}_C is conditionally independent of the source reconstructions. By substituting (3.37) into the general expression (B.3) of the conditional likelihood of feed-forward models (see proposition 3.3), we obtain

$$\begin{aligned} \mathcal{L}_{\tilde{\mathbf{x}}|\mathbf{x}} &= \langle \log p(\tilde{\mathbf{x}}|\mathbf{x}, \mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}})\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}})} - \langle KL(p(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}})||p(\mathbf{y}|\mathbf{x})) \rangle_{\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}})} \\ &= \langle \log p(\tilde{\mathbf{x}}|\mathbf{x}, \mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}})\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}})} = \langle \log p(\tilde{\mathbf{x}}|\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}})}, \end{aligned} \quad (3.38)$$

as the Kullback Leibler term cancels for all training patterns $\mathbf{x}^{(i)} \sim \tilde{p}(\mathbf{x})$. Note that in the last identity in (3.38) we have used the fact that $p(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}}) = p(\mathbf{y}|\mathbf{x})$ as the consequence of the deterministic assumption (3.37), and $p(\tilde{\mathbf{x}}|\mathbf{x}, \mathbf{y}) = p(\tilde{\mathbf{x}}|\mathbf{y})$ from the chain structure of the model \mathcal{M}_C (see Appendix C.3, proposition C.1).

Let us now define the conditional distribution $q(\mathbf{x}|\mathbf{y})$, which is constrained to be equivalent to the decoding mapping of the autoencoder, i.e. $q_{\mathbf{x}|\mathbf{y}}(\mathbf{x} = \mathbf{s}|\mathbf{y}) = p_{\tilde{\mathbf{x}}|\mathbf{y}}(\tilde{\mathbf{x}} = \mathbf{s}|\mathbf{y})$ for all $\mathbf{s} \in \mathcal{R}_{\mathbf{x}} \equiv \mathcal{R}_{\tilde{\mathbf{x}}}$, $\mathbf{y} \in \mathcal{R}_{\mathbf{y}}$. Then from lemma 3.1 and (3.34), the exact conditional likelihood $\mathcal{L}_{\tilde{\mathbf{x}}|\mathbf{x}}$ in the considered noiseless autoencoders reduces to

$$\mathcal{L}_{\tilde{\mathbf{x}}|\mathbf{x}} = \langle \log p(\tilde{\mathbf{x}}|\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}})} = \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})}, \quad (3.39)$$

which is effectively the varying part of the generic variational lower bound (2.2) on the mutual information $I(\mathbf{x}, \mathbf{y})$ for the encoder model \mathcal{M}_I . Again, here $\tilde{p}(\mathbf{x})$ is the empirical distribution of the sources, and $p(\mathbf{y}|\mathbf{x})$ is the deterministic encoder given by (3.37). Therefore, from (2.2) and (3.39) we get

$$I(\mathbf{x}, \mathbf{y}) \geq \tilde{I}(\mathbf{x}, \mathbf{y}) = \mathcal{L}_{\tilde{\mathbf{x}}|\mathbf{x}} + H_{\tilde{p}}(\mathbf{x}). \quad (3.40)$$

□

Thus, optimization of the exact conditional likelihood $\mathcal{L}_{\tilde{\mathbf{x}}|\mathbf{x}}$ in conventional noiseless autoencoders may be seen as a special case of the variational IM algorithm for the corresponding noiseless channel, where the variational decoder

is identical to the autoencoder's decoding distribution (see Figure 3.5). From proposition 3.4 we can see that there is only a constant gap between the surfaces defined by $\tilde{I}(\mathbf{x}, \mathbf{y})$ and the conditional likelihood $\mathcal{L}_{\tilde{\mathbf{x}}|\mathbf{x}}$, which for the considered case does not affect the extrema. (Note that this contrasts with the surfaces defined by log-likelihoods \mathcal{L} of generative models \mathcal{M}_L , see (3.8)). Moreover, it is easy to see that the bound on $I(\mathbf{x}, \mathbf{y})$ is saturated if the decoding distribution of the autoencoder is constrained to be identical to the exact posterior of \mathcal{M}_I , i.e.

$$p_{\tilde{\mathbf{x}}|\mathbf{y}}(\tilde{\mathbf{x}} = \mathbf{s}|\mathbf{y}) \propto p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x} = \mathbf{s})\tilde{p}(\mathbf{x} = \mathbf{s}). \quad \forall \mathbf{s} \in \mathcal{R}_{\mathbf{x}} \equiv \mathcal{R}_{\tilde{\mathbf{x}}}, \mathbf{y} \in \mathcal{R}_{\mathbf{y}}. \quad (3.41)$$

Clearly, for this specific case, minimization of the *reconstruction error in noiseless autoencoders* and maximization of the *exact mutual information in noiseless channels* are formally equivalent, if the decoding distribution is the exact posterior.

3.3.2.2 Stochastic Autoencoders with Bayesian Decoders

We have shown that the conventional conditional likelihood training of noiseless autoencoders \mathcal{M}_C may indeed be viewed as a special case of variational information maximization in noiseless channels \mathcal{M}_I , where the variational posterior is given by the \mathcal{M}_C 's decoding distribution. Additionally, for the special case when the decoder is given by the exact Bayesian posterior, the conditional training of noiseless autoencoders is equivalent to maximizing the *exact* mutual information $I(\mathbf{x}, \mathbf{y})$ under the deterministic encoding constraint. Here we relax the noiseless assumption and briefly discuss a general link between $I(\mathbf{x}, \mathbf{y})$ and $\mathcal{L}_{\tilde{\mathbf{x}}|\mathbf{x}}$ maximization. Specifically, we will motivate maximization of the mutual information as an approximate method for reducing the reconstruction error.

Proposition 3.5. *If the decoding distribution of an autoencoder model $\mathcal{M}_C = p_{\tilde{\mathbf{x}}|\mathbf{y}}(\tilde{\mathbf{x}}|\mathbf{y})p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})$ is given by the exact posterior $p_{\tilde{\mathbf{x}}|\mathbf{y}}(\tilde{\mathbf{x}}|\mathbf{y}) \propto \tilde{p}(\mathbf{x})p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})$, then for i.i.d. patterns the conditional likelihood is **lower-bounded** by the exact mutual information of the corresponding model of a stochastic memoryless channel $\mathcal{M}_I = \tilde{p}(\mathbf{x})p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})$ (up to irrelevant constants).*

Proof. The proof follows from Bayes rule and Jensen's inequality (Jensen (1906) and e.g. Korn and Korn (1968)). By definition, we obtain

$$\begin{aligned} \mathcal{L}_{\tilde{\mathbf{x}}|\mathbf{x}} &= \langle \log \langle p_{\tilde{\mathbf{x}}|\mathbf{y}}(\tilde{\mathbf{x}}|\mathbf{y}) \rangle_{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})} \rangle_{\tilde{p}(\mathbf{x})} = \frac{1}{M} \sum_{i=1}^M \log \int_{\mathbf{y}} p_{\tilde{\mathbf{x}}|\mathbf{y}}(\tilde{\mathbf{x}} = \mathbf{s}^{(i)}|\mathbf{y}) p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x} = \mathbf{s}^{(i)}) \, d\mathbf{y} \\ &= \frac{1}{M} \sum_{i=1}^M \log \tilde{p}(\mathbf{x} = \mathbf{s}^{(i)}) \left[\int_{\mathbf{y}} \frac{p_{\mathbf{y}|\mathbf{x}}^2(\mathbf{y}|\mathbf{x} = \mathbf{s}^{(i)})}{p(\mathbf{y})} \right] d\mathbf{y}, \end{aligned} \quad (3.42)$$

where $p(\mathbf{y}) = \langle p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \rangle_{\tilde{p}(\mathbf{x})}$. Then from Jensen's inequality we may transform (3.42) to

$$\begin{aligned} \mathcal{L}_{\tilde{\mathbf{x}}|\mathbf{x}} &\geq \frac{1}{M} \sum_{i=1}^M \langle \log p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x} = \mathbf{s}^{(i)}) + \log \tilde{p}(\mathbf{x} = \mathbf{s}^{(i)}) - \log p(\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x}=\mathbf{s}^{(i)})} \\ &= \langle \log p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})/p(\mathbf{y}) \rangle_{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})} + \langle \log \tilde{p}(\mathbf{x}) \rangle_{\tilde{p}(\mathbf{x})}. \end{aligned} \quad (3.43)$$

From the definition of the mutual information in the encoder model \mathcal{M}_I , we easily get the lower bound on the conditional likelihood

$$\mathcal{L}_{\tilde{x}|x} \geq I(x, y) - H_{\tilde{p}}(x). \quad (3.44)$$

□

Unlike propositions 3.3 and 3.4, which effectively define information-theoretic upper bounds on the conditional likelihood in a way which is vaguely reminiscent of Fano’s inequality (see expression (3.1) and discussion in Appendix B.1), proposition 3.5 provides a general motivation for maximizing the mutual information as a method of decreasing the reconstruction error. Indeed, inequality (3.44) may be viewed as a proper **lower bound** on the average probability of the correct reconstructions in stochastic autoencoders with the specific Bayesian setting of the decoding distribution. As the exact conditional likelihood (3.42) will typically be intractable due to a generally non-trivial form of the log posterior, one could naively hope to simplify the process of maximizing $\mathcal{L}_{\tilde{x}|x}$ by considering optimization of its proper information-theoretic lower bound $I(x, y)$. However, as we have discussed in previous chapters, in many cases of interest the bound (3.44) would involve computing intractable entropic terms. To handle the intractability we may consider optimizing tractable lower bounds on the mutual information, such as the generic bound (2.2), which would correspond to further relaxations of (3.44) and still define proper bounds on $\mathcal{L}_{\tilde{x}|x}$.

Of course, in principle other kinds of lower bounds on the conditional likelihood may be considered. It is interesting and potentially enlightening to see how these objectives and the solutions obtained by the conditional training in completely general, intractable autoencoders could relate to mutual information maximization in stochastic channels. We will now show the fundamental similarity between both approaches for the general variational reformulations.

3.3.3 Variational Information Maximization and Variational Training of Autoencoders

In proposition 3.3 we have shown that for the general case of stochastic autoencoders, the exact conditional likelihood $\mathcal{L}_{\tilde{x}|x}$ defines a relaxation of the variational lower bound on the conditional mutual information $I(\tilde{x}, y|x)$. We have also shown that conventional approaches to training deterministic autoencoders may indeed be viewed as special cases of variational mutual information maximization in noiseless communication channels. While these results may be useful for understanding a general relation between the exact conditional likelihood and mutual information maximization, they are limited to the cases when the computations are tractable. Of course, in general it may be intractable to optimize the exact conditional likelihood $\mathcal{L}_{\tilde{x}|x} = \langle \log \langle p(\tilde{x}|y) \rangle_{p(y|x)} \rangle_{\tilde{p}(x, \tilde{x})}$, and one needs to consider approximations. Here we handle the intractability of integrating over the hidden states by maximizing the variational Jensen’s bound on the objective $\mathcal{L}_{\tilde{x}|x}$, and show that variational approaches to the conditional likelihood training in stochastic autoencoders indeed reduce to optimizing the generic lower bound on $I(x, y)$ in stochastic channels.

Clearly, by analogy with the standard variational approaches of training generative models (see expression (3.17)), we can define the variational lower bound on the conditional likelihood $\mathcal{L}_{\tilde{x}|x}$ in \mathcal{M}_C as

$$\mathcal{L}_{\tilde{x}|x} \geq \tilde{\mathcal{L}}_{\tilde{x}|x} = \underbrace{\langle \log p(\tilde{x}, y|x) \rangle_{q(y|x, \tilde{x})\tilde{p}(x, \tilde{x})}}_{\text{Average energy}} - \underbrace{\langle \log q(y|x, \tilde{x}) \rangle_{q(y|x, \tilde{x})\tilde{p}(x, \tilde{x})}}_{\text{Entropy}}. \quad (3.45)$$

We can now expand the energy term to get

$$\tilde{\mathcal{L}}_{\tilde{x}|x} = \langle \log p(\tilde{x}|x, y) \rangle_{q(y|x, \tilde{x})\tilde{p}(x, \tilde{x})} - \langle KL(q(y|x, \tilde{x}) || p(y|x)) \rangle_{\tilde{p}(x, \tilde{x})}, \quad (3.46)$$

where $q(y|x, \tilde{x})$ is an arbitrary variational posterior, which we constrain to ensure the tractability of computing (3.46). Clearly, (3.46) is saturated for $q(y|x, \tilde{x}) \equiv p(y|x, \tilde{x})$. The standard variational EM algorithm optimizes the objective (3.46) with respect to the encoder $p(y|x)$, decoder $p(\tilde{x}|x, y) = p(\tilde{x}|y)$, and the variational distribution $q(y|x, \tilde{x})$, subject to the imposed constraints. Our goal here is to compare the fixed points of the variational EM algorithm for stochastic autoencoders with the fixed points of the IM algorithm for the corresponding encoder model

$$\tilde{\mathcal{M}}_I = q_{y|x}(y|x)\tilde{p}(x), \quad (3.47)$$

where

$$q_{y|x}(y|x = \mathbf{s}) \stackrel{\text{def}}{=} q(y|x = \mathbf{s}, \tilde{x} = \mathbf{s}) \quad (3.48)$$

for all $\mathbf{s} \in \mathcal{R}_x = \mathcal{R}_{\tilde{x}}$. Again, during inference in \mathcal{M}_C we approximate the exact posterior $p(y|x, \tilde{x})$ by the variational distribution $q(y|x, \tilde{x})$. Note that since the channel encoder $q_{y|x}(y|x)$ corresponds to the specific instance of the tractable variational posterior $q(y|x, \tilde{x})$, it also lies in a tractable family.

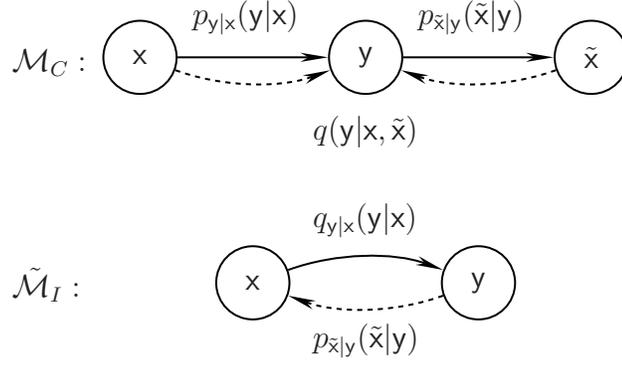
Proposition 3.6. *For i.i.d. patterns $\{x, \tilde{x}\}$, optimization of the standard variational Jensen's bound on the conditional likelihood $\tilde{\mathcal{L}}_{\tilde{x}|x}$ in **stochastic autoencoders** $\mathcal{M}_C = p(y|x)p(\tilde{x}|y)$ reduces to maximization of a specific variational lower bound on the mutual information $I(\tilde{x}, y)$ in the **stochastic memoryless channel** $\tilde{\mathcal{M}}_I = q_{y|x}(y|x)\tilde{p}(x)$.*

Proof. We will prove the proposition by expressing the fixed point updates of the variational EM on the bound $\tilde{\mathcal{L}}_{\tilde{x}|x}$ and comparing them with the variational information maximizing algorithm for the encoder model (3.47).

The autoencoder formulation of the conditional training implies the constraint on the empirical distribution, namely $\tilde{p}(\tilde{x}|x) \sim \delta(\tilde{x}-x)$. This transforms the bound on the conditional likelihood (3.46) to

$$\begin{aligned} \tilde{\mathcal{L}}_{\tilde{x}|x} &= \frac{1}{M} \sum_{m=1}^M \langle \log p_{\tilde{x}|y}(\tilde{x} = \mathbf{s}^{(m)} | y) \rangle_{q(y|x=\mathbf{s}^{(m)}, \tilde{x}=\mathbf{s}^{(m)})} \\ &\quad - \frac{1}{M} \sum_{m=1}^M KL(q(y|x = \mathbf{s}^{(m)}, \tilde{x} = \mathbf{s}^{(m)}) || p_{y|x}(y|x = \mathbf{s}^{(m)})), \end{aligned} \quad (3.49)$$

where $\mathbf{s}^{(m)} \in \mathcal{X}$ is the m^{th} training pattern. Note that due to the specific symmetric constraint on $\tilde{p}(x, \tilde{x})$, the variational posteriors $q(y|x = \mathbf{s}^{(m)}, \tilde{x} \neq \mathbf{s}^{(m)})$ do



$$\boxed{\tilde{I}(x, y) = \tilde{\mathcal{L}}_{\tilde{x}|x} + H_{\tilde{p}}(x) \text{ for } q_{y|x}(y|x) = q(y|x = \tilde{x})}$$

Figure 3.6: Variational information maximization and stochastic autoencoders. Maximization of the standard lower bound on the conditional likelihood $\tilde{\mathcal{L}}_{\tilde{x}|x}$ in \mathcal{M}_C reduces to maximization of the generic lower bound $\tilde{I}(x, y)$ in the corresponding encoder models $\tilde{\mathcal{M}}_I \stackrel{\text{def}}{=} q_{y|x}(y|x)\tilde{p}(x)$ if $q_{y|x=s}(y|x) = q(y|x = s, \tilde{x} = s)$. Here $q(y|x, \tilde{x})$ is the variational posterior of $\tilde{\mathcal{L}}_{\tilde{x}|x}$. Dashed arrows show the graphical structure of the variational distributions.

not affect the bound $\tilde{\mathcal{L}}_{\tilde{x}|x}$; in other words, for the optimization surface (3.49), the cardinality of $p(y|x)$ is greater or equal to the effective cardinality of $q(y|x, \tilde{x})$.

At the first step of the variational EM algorithm we will optimize the bound on the conditional likelihood $\tilde{\mathcal{L}}_{\tilde{x}|x}$ with respect to the encoder distribution $p_{y|x}(y|x)$. From (3.49), it is easy to see that the optimal encoder must satisfy

$$\forall \mathbf{s}^{(m)} \in \mathcal{X}, p_{y|x}(y|x = \mathbf{s}^{(m)}) = q(y|x = \mathbf{s}^{(m)}, \tilde{x} = \mathbf{s}^{(m)}) \equiv q_{y|x}(y|x = \mathbf{s}^{(m)}). \quad (3.50)$$

Clearly, (3.50) is achievable as long as the family $\mathcal{F}_{q_{y|x}}$ of the variational distributions $q(y|x = \tilde{x})$ is a subset of the family $\mathcal{F}_{p_{y|x}}$ of the exact conditionals $p_{y|x}(y|x)$, which will typically be the case for the simplifying variational approximations.

Substituting the optimal encoder (3.50) into the objective (3.49), we obtain

$$\tilde{\mathcal{L}}_{\tilde{x}|x}^{(1)} = \frac{1}{M} \sum_{m=1}^M \langle \log p_{\tilde{x}|y}(\tilde{x} = \mathbf{s}^{(m)}|y) \rangle_{q(y|x=\mathbf{s}^{(m)}, \tilde{x}=\mathbf{s}^{(m)})} = \langle \log p_{\tilde{x}|y}(\tilde{x}|y) \rangle_{q(y|x, \tilde{x})\tilde{p}(x, \tilde{x})} \quad (3.51)$$

which is formally just the non-constant part of the bound on the conditional mutual information $I(y, \tilde{x}|x)$, where the conditional encoder is given by $q(y|x, \tilde{x})$ (see proposition 3.3). The objective (3.51) will now need to be optimized for $q(y|x, \tilde{x})$ and $p_{\tilde{x}|y}(\tilde{x}|y)$, subject to the specific constraints on the decoder and the variational distribution. Without loss of generality, we may assume that the optimum for $q(y|x = \tilde{x})$ is achieved at $q(y|x = \mathbf{s}, \tilde{x} = \mathbf{s}) = r(y|x = \mathbf{s}) \in \mathcal{F}_{q_{y|x}}$, where $r(y|x = \mathbf{s})$ is some distribution in the family of the variational posteriors, and

$\mathbf{s} \in \mathcal{R}_x$. This transforms (3.51) into

$$\tilde{\mathcal{L}}_{\tilde{x}|x}^{(2)} = \langle \log p_{\tilde{x}|y}(\tilde{x}|y) \rangle_{r(y|x)\tilde{p}(x,\tilde{x})} = \langle \log p_{\tilde{x}|y}(x|y) \rangle_{r(y|x)\tilde{p}(x)}, \quad (3.52)$$

where we have applied the results of lemma 3.1. Finally, at the last step of the variational EM we optimize the objective (3.52) for $p_{\tilde{x}|y}$, assuming that $r(y|x)$ is fixed.

It is already easy to see that the iterative optimization of $\tilde{\mathcal{L}}_{\tilde{x}|x}$ is strongly related to maximization of the generic lower bound on the mutual information (2.2). Indeed, we may explicitly write the fixed point solutions for the t^{th} iteration of the algorithm⁹ as

$$\begin{aligned} p_{y|x}^{(t-1)}(y|x = \mathbf{s}^{(m)}) &= q^{(t-1)}(y|x = \mathbf{s}^{(m)}, \tilde{x} = \mathbf{s}^{(m)}), \\ q^{(t)}(y|x = \mathbf{s}^{(m)}, \tilde{x} = \mathbf{s}^{(m)}) &= r^{(t)}(y|x = \mathbf{s}^{(m)}) \\ &\stackrel{\text{def}}{=} \arg \max_q \langle \log p_{\tilde{x}|y}^{(t-1)}(x = \mathbf{s}^{(m)}|y) \rangle_{q(y|x=\mathbf{s}^{(m)}, \tilde{x}=\mathbf{s}^{(m)})}, \end{aligned} \quad (3.53)$$

$$p_{\tilde{x}|y}^{(t)}(\tilde{x} = \mathbf{s}^{(m)}|y) = \arg \max_{p_{\tilde{x}|y}} \langle \log p_{\tilde{x}|y}(x = \mathbf{s}^{(m)}|y) \rangle_{r^{(t)}(y|x=\mathbf{s}^{(m)})}, \quad (3.54)$$

where we implied the relevant constraints on the optimized functional parameters. We may further combine (3.53) and (3.54) to get the optimal encoder for the next iteration

$$\begin{aligned} p_{y|x}^{(t)}(y|x = \mathbf{s}^{(m)}) &= q^{(t)}(y|x = \mathbf{s}^{(m)}, \tilde{x} = \mathbf{s}^{(m)}) \in \mathcal{F}_{q_{y|x}} \subseteq \mathcal{F}_{p_{y|x}} \\ &= \arg \max_{q_{y|x}} \langle \log p_{\tilde{x}|y}^{(t-1)}(x = \mathbf{s}^{(m)}|y) \rangle_{q_{y|x}(y|x=\mathbf{s}^{(m)})}. \end{aligned} \quad (3.55)$$

Clearly, the fixed points (3.55) and (3.56) for the autoencoder's decoding and encoding mappings are equivalent to the ones obtained by the iterative maximization of

$$\tilde{I}(x, y) = \langle \log p_{\tilde{x}|y}(x|y) \rangle_{q_{y|x}(y|x)\tilde{p}(x)} + H_{\tilde{p}}(x), \quad (3.56)$$

where $q_{y|x}(y|x) \equiv q(y|x = \tilde{x}) \in \mathcal{F}_{q_{y|x}}$ for all $x, \tilde{x} \in \mathcal{R}_x$. Note that $\tilde{I}(x, y)$ is the generic lower bound on the mutual information in the stochastic channel $\tilde{\mathcal{M}}_I$, with the variational decoder given by $p_{\tilde{x}|y}(\tilde{x}|y)$. The iterative optimization of (3.57) for $p_{\tilde{x}|y}$ and $q_{y|x}$ defines the simplest form of the IM algorithm (see Section 2.1.2). Therefore, the variational conditional likelihood training in autoencoders may be viewed as a special case of the variational IM algorithm for the stochastic channel $\tilde{\mathcal{M}}_I \stackrel{\text{def}}{=} \tilde{p}(x)q_{y|x}(y|x)$ (see Figure 3.6). \square

It is intuitive that the tightness of $\tilde{\mathcal{L}}_{\tilde{x}|x}$ (and the equivalent objective (3.57)) will strongly depend on the constraints on the family of the variational distributions $\mathcal{F}_{q_{y|x}} \subseteq \mathcal{F}_{p_{y|x}}$. Importantly, we may note that in the case when the computations of $\langle \log p_{\tilde{x}|y}(x|y) \rangle_{p_{y|x}(y|x)}$ are tractable, we may be able to choose the family of the

⁹For the reasons of notational clarity, we will also express the optimal encoder $p_{y|x}^{(t-1)}(y|x)$ for the previous iteration of the algorithm.

variational distributions to satisfy $\mathcal{F}_{q_{y|x}} = \mathcal{F}_{p_{y|x}}$. This more general choice of the variational distributions transforms (3.57) into the equivalent objective

$$\tilde{I}(\mathbf{x}, \mathbf{y}) = \langle \log p_{\tilde{x}|y}(\mathbf{x}|\mathbf{y}) \rangle_{p_{y|x}(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})} + H_{\tilde{p}}(\mathbf{x}) \quad (3.58)$$

with $p_{y|x}(\mathbf{y}|\mathbf{x}) \in \mathcal{F}_{q_{y|x}} = \mathcal{F}_{p_{y|x}}$. From (2.2) it is clear that (3.58) is in fact the generic lower bound¹⁰ on $I(\mathbf{x}, \mathbf{y})$ in the channel model $\mathcal{M}_I = \tilde{p}(\mathbf{x})p_{y|x}(\mathbf{y}|\mathbf{x})$, which is effectively the encoding part of the autoencoder \mathcal{M}_C . Thus, we may state the *equivalence of the variational conditional likelihood training of stochastic autoencoders to the simplest form of the variational information maximization* in the corresponding stochastic channels (\mathcal{M}_I or $\tilde{\mathcal{M}}_I$, depending on the tractability of computing (3.58)). This agrees with the intuitive argument of Section 3.2.2, where we hypothesized that our generic approximate approach to information maximization should be relatable to learning in stochastic autoencoders.

3.4 Summary

Maximum likelihood learning in generative models $\mathcal{M}_{\mathcal{L}}$ and mutual information maximization in encoder models of communication channels \mathcal{M}_I may be viewed as alternative frameworks for finding (unknown) informative representations $\{\mathbf{y}\}$ of the source patterns $\{\mathbf{x}\}$. These methods are fundamentally different in several respects. One of the most important differences is parameterization of the corresponding models. Generative models are parameterized by specifying the data generating process and the prior on the latent variable representations. While this parameterization may often be useful and relevant (for example, when there are reasons to believe that the higher-dimensional observations $\{\mathbf{x}\}$ are indeed generated from lower-dimensional latent variable representations by applying specific noisy transformations), the generative framework is arguably more difficult to apply when there is a need to impose specific constraints on the posteriors $p(\mathbf{y}|\mathbf{x})$. Indeed, as the posterior $p(\mathbf{y}|\mathbf{x})$ of a generative model will typically be a highly non-linear function of $\mathcal{M}_{\mathcal{L}}$'s parameters, explicit constraints on its moments will typically be difficult to impose. Clearly, this contrasts with the encoder models \mathcal{M}_I , where the channel encoder distribution $p(\mathbf{y}|\mathbf{x})$ is a part of the model's specification. In many cases, the choice of the encoder in \mathcal{M}_I may be intuitive (for example, for clustering applications); alternatively, in some cases it may be determined by the environment (for example, for a known type of a neuro-physiological or communication channel).

The second conceptual difference of the likelihood and information maximization methods follows from the definitions of optimality. As the objective of the generative learning in $\mathcal{M}_{\mathcal{L}}$ is quantified as the cross entropy between the model and the empirical distribution (for i.i.d. patterns), a completely unconstrained

¹⁰Moreover, it is easy to see that up to irrelevant constants, expression (3.58) is just the Jensen's bound on the conditional likelihood $\mathcal{L}_{\tilde{x}|\mathbf{x}}$. Since for $\mathcal{F}_{q_{y|x}} \subseteq \mathcal{F}_{p_{y|x}}$ the bound (3.57) will typically be weaker than (3.58), there are no apparent conceptual gains of using the variational extensions of the Jensen's bound on the conditional likelihood for stochastic autoencoders. However, if the optimization is performed numerically, the specifics of the learning may be different.

\mathcal{L} -optimal model would mimic the set of the observations. Since this global optimum is clearly of little interest for any practical inference problem, there is a need for an inductive bias, i.e. constraints on the model distribution $p(\mathbf{x})$ which would hopefully be useful for producing sensible generalizations. Latent variable representations of the observed data, as well as parametric or structured specifications of the models, are useful ways of introducing meaningful constraints, though the resulting models are doomed to produce weaker likelihoods than what one would get by using the empirical distribution $\tilde{p}(\mathbf{x})$ as a model. Divergences and known numerical instabilities and degeneracies of likelihood solutions in certain under-constrained models may themselves be intrinsic artifacts of this definition of optimality. Usually, one may introduce additional constraints by applying Bayesian approaches, which optimize type-2 (marginal) likelihoods.

On the other hand, in the information-maximizing framework, optimal mappings to the hidden space are obtained by maximizing the certainty of reconstructing the source vectors from their latent variable representations. In this case, non-observability of some of the variables is not an inductive modeling assumption; it is a consequence of the channel definition where the information-theoretic learning can make sense. (Indeed, in the complete data case, i.e. when both the sources and the codes are visible, the reduction of uncertainty is specified by the empirical distribution). In general, optimization of mutual information tends to lead to other forms of degeneracies. Indeed, it is intuitive that unconstrained encoding distributions of encoder models will tend to be noiseless, and will tend to produce maximally spread-out representations in the code space. However, in many practical situations, the noise of the encoder $p(\mathbf{y}|\mathbf{x})$ is unavoidable and intrinsic to the environment, which in practice may often simplify model specifications.

Our purpose here was to try to understand possible relations between these approaches, with the specific focus on our variational method for information maximization. We explored the general relation of the generic IM algorithm to maximum likelihood learning in generative models and conditional likelihood learning in stochastic chains. In contrast to much of the previous work which relates the likelihood and the mutual information approaches for relatively simple special cases (e.g. Oja (1989), Pearlmutter and Parra (1996), Cardoso (1997), MacKay (1999b)), we tried to relate the methods for the general variational settings independently of the specific model parameterizations. Specifically, we showed that the likelihood of a generative model in $\mathcal{M}_{\mathcal{L}}$ may be viewed as a lower bound on the mutual information in the corresponding model of the noisy channel \mathcal{M}_I , where the encoding distribution is the exact posterior of the generative model. Moreover, specific tractable lower bounds $\hat{I}(\mathbf{x}, \mathbf{y})$ optimized by the information-maximizing algorithm are formally tighter than the corresponding likelihoods in $\mathcal{M}_{\mathcal{L}}$. The generally non-constant gap between \mathcal{L} and $\hat{I}(\mathbf{x}, \mathbf{y})$ suggests the fundamental differences between the induced optimization surfaces and the solutions obtained by both approaches. A practical side-effect of this study is an information-theoretic objective for training generative models (which we discuss further in Section 5.2.1).

We also demonstrated a close relation between optimization of the simple variational bound on mutual information (Section 2.1.1) and conditional likelihood training in stochastic autoencoders. Specifically, we showed that the conventional approaches to maximizing the exact conditional likelihood in noiseless autoencoders $x \rightarrow y \rightarrow \tilde{x}$ may be viewed as special instances of the generic IM algorithm for the corresponding noiseless channels $x \rightarrow y$. Optimization of the conditional log-likelihood for *stochastic* autoencoders \mathcal{M}_C is a more difficult computational task, since marginalization of the hidden codes may potentially be intractable. Arguably, one of the most straight-forward and rigorous approaches for training such models is by maximizing the variational Jensen’s bound on the conditional likelihood, e.g. by applying the variational EM algorithm. Interestingly, we can show that this procedure for training stochastic autoencoders reduces to the simplest form of the variational IM for a specific noisy channel. Specifically, this happens when the IM is applied to maximizing the generic bound on $I(x, y)$, with the variational decoder defined by the decoding mapping of the conditionally trained stochastic autoencoder \mathcal{M}_C , and the encoding distribution defined by \mathcal{M}_C ’s variational posterior. Thus, the common methods for training the conditional models may in fact be viewed as special cases of the simple variational IM framework (note that optimization of richer bounds on mutual information may potentially be considered (Section 2.3)). Finally, a curious side-product of our exploration of the general properties of the IM algorithm is a tractable model-specific upper bound on the conditional likelihood, which does not ignore the information about the reconstructing distributions (*cf* Fano’s inequality).

Chapter 4

Variational Information Maximization for Linear Dimensionality Reduction

In Chapter 2 and Chapter 3 we described the variational approach to information maximization and discussed how it relates to other methods of training probabilistic graphical models. Similarly to other variational methods, the principal idea was to transform the computationally intractable problem of maximizing the exact mutual information to optimizing tractable bounds on the objective. Our specific focus in this chapter is on applying the *generic* and the *auxiliary variational* bounds on mutual information to extracting informative lower-dimensional projections $\{y\}$ of higher-dimensional data $\{x\}$. In order to get an insight into analytical properties of the IM, we will discuss the case when the projections are stochastic, with the encoding distribution $p(y|x)$ being an isotropic linear Gaussian¹. This case may be formulated from the communication-theoretic viewpoint, where the goal would be to find the compressed representations $\{t\}$ of the data $\{x\}$, for the subsequent transmission of these representations over a non-zero noise Gaussian channel. In this formulation, the goal would be to learn the projections $x \mapsto t$ in such a way that maximizes the amount of information which the received patterns $\{y\}$ contain about the *original* sources $\{x\}$. Obviously, the problem reduces to maximizing mutual information for an isotropic linear Gaussian channel and a generally non-Gaussian distribution of the source patterns.

We will start the discussion by considering optimization of Linsker's *as-if Gaussian* objective (Linsker (1992)), which corresponds to a specific form of our variational formulation for the case of a linear Gaussian variational decoder (see Chapter 2). We show that for the considered channel, optimization of Linsker's *as-if Gaussian* objective criterion cannot improve on PCA projections. Then we show that by considering a richer family of the auxiliary variational bounds (see Section 2.3), we may significantly improve on achievable lower bounds on mutual

¹Note that despite the apparent similarity to linear autoencoders (Baldi and Hornik (1989), Oja (1989), Roweis and Ghahramani (1999)), our formulation is significantly more general in several respects. Specifically, in our case the encoder $p(y|x)$ is constrained to be an isotropic linear Gaussian (rather than a noiseless linear projection), and the variational decoder $q(x|y)$ is an arbitrary distribution in a tractable family (rather than an isotropic linear Gaussian).

information (under the identical encoding constraints). This result is encouraging, as it suggests a simple way to produce tighter lower bounds on the mutual information (compared with PCA) *without* altering the channel specification or increasing the length of the communicated codewords.

Finally, we discuss a simple and practical reformulation of the dimensionality reduction problem, and show that by storing an additional auxiliary variable z (given by the generalized linear projection to the auxiliary space), we may facilitate reconstructions of the sources from noisy lower-dimensional representations without making a recourse to the stored dataset. Effectively, for this case the objective to optimize would be given by $I(\mathbf{x}, \{\mathbf{y}, z\})$, which results in a simplification of the auxiliary variational bound on $I(\mathbf{x}, \mathbf{y})$. Strictly speaking, the resulting reformulation of the optimization problem does not have a direct mapping to a data transmission problem in a simple Gaussian channel; however, it does give rise to an efficient compression and decompression mechanism. Indeed, we show that by a moderate increase in the size of the compressed representations, we may significantly improve on reconstructions from simple *constrained* encodings.

4.1 Introduction

One of the principal goals of dimensionality reduction is to produce a lower-dimensional representation \mathbf{y} of a high-dimensional source vector \mathbf{x} , so that useful information which the codes \mathbf{y} contain about the sources \mathbf{x} is maximally preserved. If it is not known a priori which coordinates of \mathbf{x} may be relevant for a specific task, it is sensible to maximize the amount of information which \mathbf{y} contains about all the coordinates, for all possible source vectors. As discussed in Chapter 2, the fundamental measure of informativeness in this context is the mutual information

$$I(\mathbf{x}, \mathbf{y}) \equiv H(\mathbf{x}) - H(\mathbf{x}|\mathbf{y}),$$

which quantifies the decrease of uncertainty in the pattern \mathbf{x} due to the knowledge of \mathbf{y} . Again, $H(\mathbf{x}) \equiv -\langle \log p(\mathbf{x}) \rangle_{p(\mathbf{x})}$ and $H(\mathbf{x}|\mathbf{y}) \equiv -\langle \log p(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{x}, \mathbf{y})}$ are marginal and conditional entropies respectively, and the angled brackets represent the averages.

A principal motivation for applying information theoretic techniques for learning informative lower-dimensional representations is the general intuition that the lower dimensional codes should preserve useful information about the higher-dimensional data. Moreover, the information maximizing framework of encoder models is particularly convenient for addressing the problem of *constrained* dimensionality reduction. To demonstrate this, suppose that we are interested in learning an optimal undercomplete *orthonormal projection* from the data in the presence of irreducible Gaussian noise (the stochastic subspace selection view). As discussed in Section 3.4, it would generally be difficult to impose specific constraints on the posterior $p(\mathbf{y}|\mathbf{x})$ in the usual generative formulation. For example, if the considered generative model is a factor analyzer $p(\mathbf{y})p(\mathbf{x}|\mathbf{y})$ with $p(\mathbf{y}) \sim \mathcal{N}(0, s\mathbf{I})$, $p(\mathbf{x}|\mathbf{y}) \sim \mathcal{N}(\mathbf{W}\mathbf{y}, \mathbf{\Psi})$, and $\mathbf{\Psi} = \{\psi_{ij}\delta_{ij}\} \neq c\mathbf{I}$, the orthonormal constraints on $\mathbf{A} : \mathbf{x} \rightarrow \langle \mathbf{y}|\mathbf{x} \rangle$ would imply orthonormality of $\mathbf{A} \stackrel{\text{def}}{=} (s\mathbf{I} + \mathbf{W}^T \mathbf{\Psi}^{-1} \mathbf{W})^{-1} \mathbf{W}^T \mathbf{\Psi}^{-1}$

(e.g. von Mises (1964), Bartholomew (1987)). Clearly, by explicitly enforcing (e.g. Arfken (1985), Riley et al. (2002)) the orthonormality constraint, we may significantly complicate the resulting optimization surface. On the other hand, a natural way to address the problem of constrained dimensionality reduction is by considering the information-maximizing paradigm for an encoder model $\tilde{p}(\mathbf{x})p(\mathbf{y}|\mathbf{x})$. In this case we could easily impose the orthonormality constraints (or in fact any other requirements the noisy projections $\mathbf{x} \rightarrow \mathbf{y}$ need to satisfy) by explicitly parameterizing the stochastic mapping $p(\mathbf{y}|\mathbf{x})$. For example, for this specific case we could set $p(\mathbf{y}|\mathbf{x}) \sim \mathcal{N}(\mathbf{W}\mathbf{x}, \Sigma_{\mathbf{y}|\mathbf{x}})$, where $\mathbf{W}\mathbf{W}^T = \mathbf{I}_{|\mathbf{y}|}$. Effectively, this parameterization would be analogous to specifying the conditionals of the discriminative models; however, in contrast to discriminative models which presume observability of the outputs, the lower-dimensional vectors $\{\mathbf{y}\}$ will in our case be hidden. Our focus in this chapter is on applying the variational information maximizing framework to dimensionality reduction, with specific focus on the linear case.

4.1.1 Optimization of Linsker’s Criterion

The principled information theoretic approach to dimensionality reduction would maximize the exact mutual information $I(\mathbf{x}, \mathbf{y})$ with respect to parameters of the encoder $p(\mathbf{y}|\mathbf{x})$. Despite the fact that the dimensionality of the reduced space $|\mathbf{y}|$ will be lower than the dimensionality of the original space, the exact evaluation of $I(\mathbf{x}, \mathbf{y})$ will generally be computationally intractable if $|\mathbf{y}| < |\mathbf{x}|$ is still large. As we mentioned in Section 1.4, the key difficulty lies in the computation of the entropic term $H(\mathbf{x}|\mathbf{y})$, which is tractable only in a few special cases. A computationally tractable alternative to maximizing $I(\mathbf{x}, \mathbf{y})$ is the IM applied to maximizing proper lower bounds on the mutual information.

We will start the discussion by considering optimization of the simple generic lower bound on the mutual information

$$I(\mathbf{x}, \mathbf{y}) \geq \tilde{I}(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} H(\mathbf{x}) + \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})}, \quad (4.1)$$

where $\tilde{p}(\mathbf{x})$ is the empirical distribution, and the variational approximation $q(\mathbf{x}|\mathbf{y})$ of the exact posterior $p(\mathbf{x}|\mathbf{y})$ is given by a linear Gaussian $q(\mathbf{x}|\mathbf{y}) \sim \mathcal{N}(\mathbf{U}\mathbf{y}, \Sigma)$. As we showed in Section 1.4, optimization of the bound (2.2) for this specific choice of the variational decoder reduces to maximization of Linsker’s *as-if Gaussian* criterion

$$2I_G(\mathbf{x}, \mathbf{y}) = \log |\Sigma_{xx}| - \log |\Sigma_{xx} - \Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_{xy}^T| + c, \quad (4.2)$$

where Σ_{xx} , Σ_{xy} , and Σ_{yy} are the partitions of decoder covariance Σ , and c is an irrelevant constant. In other words, maximization of Linsker’s criterion I_G may be seen as a special case of the variational information-maximization formulation for linear Gaussian decoders, *independently* of the specific encoder parameterization. Clearly, the objective (4.2) may be expressed as a function of the encoder parameters. After training, the learned parameters may be used for producing lower-dimensional representations \mathbf{y} for the given sources \mathbf{x} by forward-sampling from the encoder $p(\mathbf{y}|\mathbf{x})$ (or computing the conditional mean $\langle \mathbf{y}|\mathbf{x} \rangle_{p(\mathbf{y}|\mathbf{x})}$).

4.1.1.1 Nature of Optimal Solutions

In the considered special case of a linear Gaussian channel, the encoder is given by $p(y|x) \sim \mathcal{N}(Wx, \Sigma_{y|x})$, and $p(x)$ is the empirical distribution of the training patterns. It is now easy to show that the left singular vectors of the optimal projection weights W^T give rise to the $|y|$ -PCA solution on the sample covariance $\mathbf{S} \stackrel{\text{def}}{=} \langle xx^T \rangle$ (for clarity, we assume that the data is centered). The remaining moments may be trivially expressed as

$$\langle xy^T \rangle = \mathbf{S}W^T, \quad (4.3)$$

$$\langle yy^T \rangle = \Sigma_{y|x} + \mathbf{W}\mathbf{S}\mathbf{W}^T. \quad (4.4)$$

Then by substituting into (4.2), we get

$$\begin{aligned} \tilde{I}(x, y) &= \log |\mathbf{S} - \mathbf{S}\mathbf{W}^T(\mathbf{W}\mathbf{S}\mathbf{W}^T + \Sigma_{y|x})^{-1}\mathbf{W}\mathbf{S}|^{-1} + c \\ &= \log |\mathbf{S}^{-1} + \mathbf{W}^T\Sigma_{y|x}\mathbf{W}| + c, \end{aligned} \quad (4.5)$$

where we used the matrix inversion lemma (e.g. Press et al. (1992)) and let c be a constant.

To explore spectral properties of the optimal solutions, we will now consider the singular value decomposition of the weights, i.e. $\mathbf{W}^T = \mathbf{V}\mathbf{L}\mathbf{R}^T \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{y}|}$. Here $\mathbf{V} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{y}|}$ is a matrix of orthonormal columns, i.e. $\mathbf{V}^T\mathbf{V} = \mathbf{I}_{|\mathbf{y}|}$, $\mathbf{L} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{y}|}$ is a diagonal, and $\mathbf{R} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{y}|}$ is a rigid rotation matrix. By assuming that the irreducible channel noise is white (i.e. $\Sigma_{y|x} = \sigma^2\mathbf{I}$), we may transform the objective (4.5) into

$$\tilde{I}(x, y) = \log |\mathbf{S}^{-1} + \sigma^{-2}\mathbf{V}\mathbf{L}^2\mathbf{V}^T| + c. \quad (4.6)$$

Clearly, the orthonormality constraint on the projection weights $\mathbf{W}\mathbf{W}^T = \mathbf{I}_{|\mathbf{y}|}$ implies $\mathbf{L} = \mathbf{I}$ (which also ensures convergence of the objective (4.6)). This leads to

$$\tilde{I}(x, y) = \log |\mathbf{S}^{-1} + \sigma^{-2}\mathbf{V}\mathbf{V}^T| - \text{tr} \{ \mathbf{M}(\mathbf{V}^T\mathbf{V} - \mathbf{I}) \}. \quad (4.7)$$

A straight-forward matrix optimization for $\mathbf{V} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{y}|}$ leads to

$$(\mathbf{S}^{-1} + \sigma^{-2}\mathbf{V}\mathbf{V}^T)\mathbf{V} = \mathbf{V}\mathbf{M}^{-1}\sigma^{-2}, \quad (4.8)$$

i.e.

$$\mathbf{S}^{-1}\mathbf{V} = \mathbf{V}\tilde{\mathbf{M}}, \quad \text{where } \tilde{\mathbf{M}} \stackrel{\text{def}}{=} (\mathbf{M}^{-1} - \mathbf{I})\sigma^{-2}, \quad (4.9)$$

where $\mathbf{M} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{y}|}$ is a matrix of Lagrange multipliers. Without loss of generality, we may assume that $\tilde{\mathbf{M}}$ is symmetric (since $\text{tr} \{ \mathbf{M} \} = \text{tr} \{ \mathbf{M}^T \}$), i.e. the weights \mathbf{V} satisfying the extremum criterion (4.9) correspond to eigenvectors of \mathbf{S}^{-1} and their rigid rotations. Since the objective (4.6) is not influenced by the rotation factor, we may as well ignore them in our analysis.

Finally, by substituting (4.9) into (4.6) and expressing the objective in terms of the eigenspectra of the sample covariance \mathbf{S} , it is easy to see that the optimal weights \mathbf{V} correspond to *principal* components of \mathbf{S} . Indeed, (4.9) implies that the objective (4.6) may be expressed as

$$\tilde{I}(x, y) = \log \left| \mathbf{V}(\Lambda_{\mathbf{S}}^{-1} + \sigma^{-2}\mathbf{I})\mathbf{V}^T + \tilde{\mathbf{V}}\tilde{\Lambda}_{\mathbf{S}}^{-1}\tilde{\mathbf{V}}^T \right| + c. \quad (4.10)$$

Here $\tilde{\mathbf{V}} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}| - |\mathbf{y}|}$ are the eigenvectors of \mathbf{S} which are orthonormal to the space spanned by $\mathbf{V} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{y}|}$, i.e.

$$\mathbf{V}\mathbf{V}^T + \tilde{\mathbf{V}}\tilde{\mathbf{V}}^T = \mathbf{I}_{|\mathbf{x}|}, \quad (4.11)$$

and $\Lambda_S \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{y}|}$, $\tilde{\Lambda}_S \in \mathbb{R}^{|\mathbf{x}| - |\mathbf{y}| \times |\mathbf{x}| - |\mathbf{y}|}$ are the eigenvalues of \mathbf{S} corresponding to \mathbf{V} and $\tilde{\mathbf{V}}$ respectively. From (4.9) it is clear that $\tilde{\Lambda}_S$ and $\tilde{\mathbf{V}}$ define the spectrum discarded by the projection weights \mathbf{W} . The new objective (4.10) may be equivalently expressed as

$$\tilde{I}(\mathbf{x}, \mathbf{y}) = \sum_{i \in \mathcal{R}} \log(\lambda_i^{-1} + \sigma^{-2}) + \sum_{j \notin \mathcal{R}} \log \lambda_j^{-1} + c, \quad (4.12)$$

where λ_i is the i^{th} eigenvalue of the sample covariance, and \mathcal{R} specifies the spectrum retained in \mathbf{W} . Equivalently, we may get

$$\begin{aligned} \tilde{I}(\mathbf{x}, \mathbf{y}) &= \sum_{i \in \mathcal{R}} \log(\lambda_i^{-1} + \sigma^{-2}) - \sum_{j \notin \mathcal{R}} \log \lambda_j + c \\ &= \sum_{i \in \mathcal{R}} \log(\lambda_i + \sigma^2) - |\mathbf{y}| \log \sigma^2 - \log |\mathbf{S}| + c. \end{aligned} \quad (4.13)$$

From (4.13) it is clear that the bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ is maximized when the spectrum \mathcal{R} retained in the projection weights \mathbf{W} indeed corresponds to the principal components of the sample covariance. This may be further related to the special case of noiseless linear autoencoders discussed in the early work of Baldi and Hornik (1989) and Oja (1989), which may be obtained from (4.13) by computing the limit at $\sigma^2 \rightarrow 0$ and using the result of proposition 3.4.

4.2 Optimization of the Auxiliary Variational Bound

Importantly, the result of Section 4.1.1 suggests that optimization of Linsker’s criterion for isotropic linear Gaussian channels cannot improve on the simple $|\mathbf{y}|$ -PCA projections. Since optimization of Linkser’s criterion may be viewed as a special instance of the IM algorithm with linear Gaussian variational decoders, a natural question to explore is how the linear projections could be affected by using more complex decoder types. A principal conceptual difficulty of applying the bound (2.2) is in specifying a powerful yet tractable variational decoder $q(\mathbf{x}|\mathbf{y})$. Here we consider a richer family of tractable auxiliary variational bounds on $I(\mathbf{x}, \mathbf{y})$ (see Section 2.3), which may overcome the fundamental limitations of Linsker’s criterion.

The key idea of the auxiliary variational method (Agakov and Barber (2005a)) is to introduce mappings to the auxiliary space $\{\mathbf{z}\}$ in a way which does not affect the original channel $p(\mathbf{y}|\mathbf{x})$, and learn the resulting joint distribution $p(\mathbf{x}, \mathbf{y}, \mathbf{z}) = p(\mathbf{x}, \mathbf{y})p(\mathbf{z}|\mathbf{x}, \mathbf{y})$ of the sources \mathbf{x} , encodings \mathbf{y} , and auxiliary (“feature”) variables \mathbf{z} (see Section 2.3 for a detailed discussion). By applying straight-forward algebraic manipulations, we may obtain a tractable lower bound on $I(\mathbf{x}, \mathbf{y})$

$$I(\mathbf{y}, \mathbf{x}) \geq \tilde{I}(\mathbf{y}, \mathbf{x}) \stackrel{\text{def}}{=} H(\mathbf{x}) + H(\mathbf{z}|\mathbf{x}, \mathbf{y}) + \langle \log q(\mathbf{x}|\mathbf{y}, \mathbf{z}) \rangle_{p(\mathbf{x}, \mathbf{y}, \mathbf{z})} + \langle \log q(\mathbf{z}|\mathbf{y}) \rangle_{p(\mathbf{y}, \mathbf{z})} \quad (4.14)$$

The variational distributions $q(\mathbf{x}|\mathbf{y}, z)$ and $q(z|\mathbf{y})$, as well as the auxiliary conditional $p(z|\mathbf{x}, \mathbf{y})$ are chosen to ensure tractability of the objective (4.14). Then it is tractable to optimize $\tilde{I}(\mathbf{x}, \mathbf{y})$ for the channel encoder, variational decoder, and the auxiliary conditional distributions. Effectively, we will still be learning an optimal constrained encoder ($p(\mathbf{y}|\mathbf{x}) \sim \mathcal{N}(\mathbf{W}\mathbf{x}, \mathbf{\Sigma})$ in the special case we consider here), but for a richer family of variational distributions. Note that if the auxiliary space is discrete (for example, if a single auxiliary variable z is multinomial), the variational decoder $q(\mathbf{x}|\mathbf{y}) = \langle q(\mathbf{x}|\mathbf{y}, z) \rangle_{q(z|\mathbf{y})}$ is indeed defined as a multi-modal distribution.

4.2.1 Representations

Here we discuss a tractable choice of the variational parameters of $\tilde{I}(\mathbf{x}, \mathbf{y})$ for a linear Gaussian channel $p(\mathbf{y}|\mathbf{x})$ with constrained projection weights. In the model which we consider, the auxiliary space is given by the multinomial variable z which takes one of $|z|$ possible states $\{z_1, \dots, z_{|z|}\}$. The encoder and the auxiliary conditional distributions are given by $p(\mathbf{y}|\mathbf{x}) \sim \mathcal{N}_{\mathbf{y}}(\mathbf{W}\mathbf{x}; \mathbf{\Sigma})$ and

$$p(z_j|\mathbf{x}, \mathbf{y}) = p(z_j|\mathbf{x}) \propto \exp\{-\mathbf{v}_j^T \mathbf{x} + b_j\} \quad (4.15)$$

respectively. Here $\mathbf{W} \in \mathbb{R}_C^{|\mathbf{y}| \times |\mathbf{x}|}$, $\mathbf{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_{|z|}\} \in \mathbb{R}^{|\mathbf{x}| \times |z|}$, and $\mathbf{b} \in \mathbb{R}^{|z|}$ are weights and biases to be learned, where we have assumed that $\mathbb{R}_C^{|\mathbf{y}| \times |\mathbf{x}|} \subseteq \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$ is a subspace of encoder weights satisfying specific constraints (with $\mathbb{R}_C^{|\mathbf{y}| \times |\mathbf{x}|} \equiv \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$ if the weights are unconstrained). For computational convenience, we assumed that the auxiliary variables z are conditionally independent of the codes \mathbf{y} , i.e. $I(z, \{\mathbf{x}, \mathbf{y}\}) = I(z, \mathbf{x})$. Also, during learning we constrained the variational distribution used for reconstruction of the data \mathbf{x} and auxiliary variables z to satisfy

$$q(\mathbf{x}, z|\mathbf{y}) \propto q(\mathbf{x}|\mathbf{y}, z)q(z), \quad (4.16)$$

which facilitates computations of the integrals in $\tilde{I}(\mathbf{x}, \mathbf{y})$.

Clearly, for the considered parameterization the objective (4.14) is transformed into

$$I(\mathbf{x}, \mathbf{y}) \geq H(z, \mathbf{x}) + \langle \log q(z|\mathbf{y}) \rangle_{p(z, \mathbf{y})} + \langle \log q(\mathbf{x}|\mathbf{y}, z) \rangle_{p(\mathbf{x}, \mathbf{y}, z)}. \quad (4.17)$$

Note that for an arbitrary mixture decoder $q(\mathbf{x}|\mathbf{y})$, the computational complexity of evaluating the first two terms in the bound is linear in the number of states $|z|$. In general, computation of $\langle \log q(\mathbf{x}|\mathbf{y}, z) \rangle_{p(\mathbf{x}, \mathbf{y}, z)}$ is more problematic, since it requires averaging of a non-factorized function of the codes over the channel distribution $p(\mathbf{y}|\mathbf{x})$. For the special case when each component is a Gaussian with a constrained mean $q(\mathbf{x}|\mathbf{y}, z_j) \sim \mathcal{N}_{\mathbf{x}}(\mathbf{U}_j \mathbf{y}, \mathbf{S}_j)$, the rightmost term in (4.17) is expressed as

$$\begin{aligned} \langle \log q(\mathbf{x}|\mathbf{y}, z) \rangle_{p(\mathbf{x}, \mathbf{y}, z)} &= -\frac{1}{2M} \sum_{j=1}^{|z|} \sum_{i=1}^M p(z_j|\mathbf{x}^{(i)}) \text{tr} \left\{ \mathbf{S}_j^{-1} \left(\mathbf{d}_j^{(i)} \mathbf{d}_j^{(i)T} + \mathbf{U}_j \mathbf{\Sigma} \mathbf{U}_j^T \right) \right\} \\ &\quad - \frac{1}{2M} \sum_{j=1}^{|z|} \log |\mathbf{S}_j| \sum_{i=1}^M p(z_j|\mathbf{x}^{(i)}). \end{aligned} \quad (4.18)$$

Here we ignored the irrelevant constants and defined

$$\mathbf{d}_j^{(i)} \stackrel{\text{def}}{=} \mathbf{x}^{(i)} - \mathbf{U}_j \mathbf{W} \mathbf{x}^{(i)} \in \mathbb{R}^{|\mathbf{x}|} \quad (4.19)$$

to be the distortion between the i^{th} pattern and its reconstruction from a noiseless code at the mean of $q(\mathbf{x}|\mathbf{y}, z_j)$. From (4.17) and (4.18) it is easy to see that small values of the distortion terms $\mathbf{d}_j^{(i)}$ lead to improvements in the bound on the mutual information. This agrees with the intuition that the trained model should favour accurate reconstructions of the source patterns from their compressed representations passed through a noisy channel.

4.2.2 Learning Optimal Parameters

From (4.17) it is clear that the optimal settings of the auxiliary conditionals are given by

$$q(z|\mathbf{y}) = p(z|\mathbf{y}) \propto \sum_{i=1}^M p(z|\mathbf{x}^{(i)}) p(\mathbf{y}|\mathbf{x}^{(i)}). \quad (4.20)$$

It is now possible to derive an iterative learning rule for the parameters of $q(\mathbf{x}|\mathbf{y}, z)$, $p(\mathbf{y}|\mathbf{x})$ and $p(z|\mathbf{x})$. These results are obtained by computing matrix derivatives of $\tilde{I}(\mathbf{x}, \mathbf{y})$ and (where possible) deriving the closed-form fixed-point updates. The updates are performed iteratively, assuming parameter independence at each iteration. We will state and briefly discuss the results, omitting the straight-forward derivations.

Optimal decoder: It is easy to see that the considered variational decoder defines a constrained mixture of Gaussians $q(\mathbf{x}|\mathbf{y}) = \langle q(\mathbf{x}|\mathbf{y}, z) \rangle_{q(z|\mathbf{y})}$, where $q(\mathbf{x}|\mathbf{y}, z_j) \sim \mathcal{N}_{\mathbf{x}}(\mathbf{U}_j \mathbf{y}, \mathbf{S}_j)$. For each component j , the optimal weights \mathbf{U}_j parameterizing the component's mean are given by

$$\mathbf{U}_j^{(new)} = \left(\sum_{i=1}^M p(z_j|\mathbf{x}^{(i)}) \mathbf{x}^{(i)} \mathbf{x}^{(i)T} \mathbf{W}^T \right) \left(\sum_{i=1}^M p(z_j|\mathbf{x}^{(i)}) (\mathbf{W} \mathbf{x}^{(i)} \mathbf{x}^{(i)T} \mathbf{W}^T + \mathbf{\Sigma}) \right)^{-1} \quad (4.21)$$

where we have assumed that at the current iteration of the algorithm the encoder weights \mathbf{W} are fixed. Similarly, the update for the components' covariances is given by

$$\mathbf{S}_j^{(new)} = \sum_{i=1}^M p(z_j|\mathbf{x}^{(i)}) \left(\mathbf{d}_j^{(i)} \mathbf{d}_j^{(i)T} + \mathbf{U}_j^{(new)} \mathbf{\Sigma} (\mathbf{U}_j^{(new)})^T \right) \frac{1}{\sum_{i=1}^M p(z_j|\mathbf{x}^{(i)})}. \quad (4.22)$$

Here $\mathbf{\Sigma}$ is the covariance of the channel noise (which is presumed to be fixed and independent of \mathbf{x} and \mathbf{y}), and the distortion $\mathbf{d}_j^{(i)} \in \mathbb{R}^{|\mathbf{x}|}$ is given by (4.19) computed for the new decoder weights $\mathbf{U}_j^{(new)}$. Perhaps not very surprisingly, the fixed point equations (4.21) and (4.22) resemble maximum-likelihood updates for mixtures of constrained Gaussian distributions. The constraints are influenced by the specific choice of encoder and decoder distributions, as well as the mapping to the auxiliary space. Note that for noisy channels and non-singular decoder weights \mathbf{U}_j ,

computation of the inverses in (4.21) should *not* be numerically unstable. However, in order to ensure numerical stability of optimization for the case of small datasets or near-singular decoder weights \mathbf{U}_j , it may be practical to increment \mathbf{S}_j by a multiple of the identity matrix with a small positive scaling factor (which would correspond to learning the constrained covariances $\mathbf{S}_j + \epsilon \mathbf{I}_{|y|}$ for $\epsilon > 0$ and $j = 1, \dots, |z|$).

Optimal auxiliary mappings:

One way to learn the optimal auxiliary conditional $p(z|\mathbf{x})$ is by performing numerical ascent on $\tilde{I}(\mathbf{x}, \mathbf{y})$ with respect to the parameters of the conditional. The gradients are computed from (4.17) as

$$\frac{\partial \tilde{I}}{\partial \mathbf{v}_j} = \frac{1}{M} \sum_{i=1}^M \mathbf{x}^{(i)} p(z_j | \mathbf{x}^{(i)}) (\langle e(z, \mathbf{x}^{(i)}) \rangle_{p(z|\mathbf{x}^{(i)})} - e(z_j, \mathbf{x}^{(i)})), \quad (4.23)$$

where

$$e(z_j, \mathbf{x}^{(i)}) \stackrel{\text{def}}{=} \tilde{q}_{ji} + \log q(z_j) - (1 + \log p(z_j | \mathbf{x}^{(i)})), \quad (4.24)$$

and $\tilde{q}_{ji} \stackrel{\text{def}}{=} \langle \log q(\mathbf{x}^{(i)} | y, z_j) \rangle_{p(y|\mathbf{x}^{(i)})}$ is given by expression (4.18) computed for the new settings of the decoder parameters \mathbf{U}_j and \mathbf{S}_j . Again, z_j is the state of the auxiliary variable z , and i is an index of a training pattern. The gradients for the biases $\partial \tilde{I} / \partial b_j$ have a form similar to (4.23), with an omitted pre-multiplication by $\mathbf{x}^{(i)}$ in the summation. Note that throughout the iterations for $\mathbf{V} \in \mathbb{R}^{|\mathbf{x}| \times |z|}$ and $\mathbf{b} \in \mathbb{R}^{|z|}$, the posterior $p(z|\mathbf{x})$ and the average \tilde{q}_{ji} may be kept fixed. For this case, evaluation of the gradient (4.23) is computationally efficient, as the complexity of computing the averages is just $O(|z|M)$. Having computed the updated parameters of the auxiliary conditional, we can use (4.15) to obtain $p^{(\text{new})}(z_j | \mathbf{x}^{(i)})$ for all $j = \{1, \dots, |y|\}$ and $i = \{1, \dots, M\}$.

Optimal encoder:

Throughout the learning, we assumed that the covariance Σ of the Gaussian encoder was fixed (which corresponds to a fixed channel noise distribution). In this case, optimization for the encoder $p(y|\mathbf{x})$ reduces to learning the projection weights \mathbf{W} . From (4.17), it is easy to get

$$\frac{\partial \tilde{I}}{\partial \mathbf{W}} = \frac{1}{M} \sum_{j=1}^{|z|} \sum_{i=1}^M \left[\left(\mathbf{U}_j^{(\text{new})} \right)^T \left(\mathbf{S}_j^{(\text{new})} \right)^{-1} \left(\mathbf{I}_{|x|} - \mathbf{U}_j^{(\text{new})} \mathbf{W} \right) \right] \mathbf{x}^{(i)} \mathbf{x}^{(i)T} p^{(\text{new})}(z_j | \mathbf{x}^{(i)}), \quad (4.25)$$

where we *implied* that the weights $\mathbf{W} \in \mathbb{R}_C^{y \times |x|} \subseteq \mathbb{R}^{|y| \times |x|}$ lie in the matrix space $\mathbb{R}_C^{|y| \times |x|}$ satisfying any of the required constraints on the encoding weights. In general, if the number of mixture components $|z| > 1$, it is not easy to find a closed form expression for \mathbf{W} . Instead, we may perform numerical optimization of the bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ for \mathbf{W} (subject to the constraints).

Generally, as we mentioned in Section 4.1, the encoder formulation facilitates the choice of constraints on the encoder distribution $p(y|\mathbf{x})$. For example, in our case it is easy to incorporate norm constraints on \mathbf{W} into the objective function

(4.17). In the simplest case when the constraints are soft, (i.e. when the Lagrange multipliers are fixed), this would result in the regularization penalty $-DW$ on the gradient (4.25), where $D = \{D_{ij}\delta_{ij} | D_{ii} \geq 0\} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{Y}|}$ is a fixed diagonal matrix. Alternatively, we can impose hard constraints on the encoder weights by considering a specific construction in a way which always results in the desired singular spectrum. For example, we may impose the orthonormality constraints on the rows of W by parameterizing the weights as

$$W = (\tilde{W}\tilde{W}^T)^{-1/2}\tilde{W}, \quad (4.26)$$

where $\tilde{W} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{X}|}$ is an arbitrary rank- $|\mathcal{Y}|$ real-valued matrix. Clearly, (4.26) implies the desired orthonormality constraint on the projection weights, as $WW^T = I_{|\mathcal{Y}|}$. Other kinds of constraints on the encoder parameters may potentially be considered. For example, if $W = \{w_{ij}\} \in \mathbb{R}_C^{|\mathcal{Y}| \times |\mathcal{X}|}$ is constrained to a hyper-cube (cf expression (4.26)), we may parameterize each weight component as $w_{ij} = f(\tilde{w}_{ij})$, where $f : \mathbb{R} \mapsto [-\omega, \omega]$ defines a mapping from the real space to a closed line segment. Learning the encoder parameters $W \in \mathbb{R}_C^{|\mathcal{Y}| \times |\mathcal{X}|}$ would then involve unconstrained optimization of $\tilde{I}(\mathbf{x}, \mathbf{y})$ for $\tilde{W} = \{\tilde{w}_{ij}\} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{X}|}$, with the corresponding gradients obtained from (4.14) and (4.25) by the chain rule. Note that while the mean of the considered encoder $p(\mathbf{y}|\mathbf{x})$ is linear in the constrained encoder parameters $W \in \mathbb{R}_C^{|\mathcal{Y}| \times |\mathcal{X}|}$, it is generally nonlinear in $\tilde{W} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{X}|}$.

Finally, we note that at the end of each iteration of the optimization procedure, it is necessary to re-compute the moments \tilde{q}_{ji} and the prior on the auxiliary variables $q(z)$ according to (4.18) and (4.20) respectively. The iterations of (4.20) – (4.25) continue until convergence or meeting a termination criterion.

4.2.2.1 Role of the Auxiliary Variables

By analogy with Section 2.3, we may interpret the information maximization framework from the communication-theoretic viewpoint. As the objective function which we optimize is the lower bound on $I(\mathbf{x}, \mathbf{y})$ (bounding the channel capacity in the stochastic communication channel $\mathbf{x} \rightarrow \mathbf{y}$), we presume that the auxiliary variables z are not transmitted across the channel. Their purpose in this context is to define a richer family of variational bounds on the true mutual information in $\mathbf{x} \rightarrow \mathbf{y}$. This auxiliary variational family of lower bounds includes the generic bound (4.1) as a special case. (Indeed, we can easily see that by allowing a flexibility in the mappings to the auxiliary space, we may obtain the bounds which are at least as tight as (4.1). Generally, the simple bound (4.1) arises as a special case of (4.14) when \mathbf{z} takes a single state). The role of the auxiliary variables \mathbf{z} in this context is to capture regularities in the training data and introduce additional dependencies and multi-modality to the structure of the variational decoder.

4.2.3 Learning Optimal Representations in the Augmented $\{\mathbf{y}, \mathbf{z}\}$ -space

Now suppose that the auxiliary variables \mathbf{z} are actually observable at the receiver's end of the channel. Under this assumption, we may consider optimizing an alternative bound $\tilde{I}_H(\mathbf{x}, \{\mathbf{y}, \mathbf{z}\}) \geq I(\mathbf{x}, \mathbf{y})$, defined by analogy with (2.2). (We will use the index I_H to indicate that the channel $\mathbf{x} \rightarrow \{\mathbf{y}, \mathbf{z}\}$ is generally heterogeneous; for example, \mathbf{z} may be a vector of generally unknown class labels, while $\mathbf{y} \in \mathbb{R}^{|\mathbf{y}|}$ may define a lower-dimensional projection). It will lead to a slight simplification of (4.14), which effectively reduces to

$$\tilde{I}_H(\mathbf{x}, \{\mathbf{y}, \mathbf{z}\}) = H(\mathbf{x}) + \langle \log q(\mathbf{x}|\mathbf{y}, \mathbf{z}) \rangle_{\tilde{p}(\mathbf{x})p(\mathbf{y}|\mathbf{x})p(\mathbf{z}|\mathbf{x})}, \quad (4.27)$$

where the cross-entropic term is given by (4.18). For our specific case, this leads to a change in the updates for parameters of the auxiliary conditional $p(\mathbf{z}|\mathbf{x})$, which in this case leads to

$$\frac{\partial \tilde{I}}{\partial \mathbf{v}_j} = \frac{1}{M} \sum_{i=1}^M \mathbf{x}^{(i)} p(z_j|\mathbf{x}^{(i)}) (\langle \log q(\mathbf{x}^{(i)}|\mathbf{y}, \mathbf{z}) \rangle_{p(\mathbf{y}, \mathbf{z}|\mathbf{x}^{(i)})} - \langle \log q(\mathbf{x}^{(i)}|\mathbf{y}, z_j) \rangle_{p(\mathbf{y}|\mathbf{x}^{(i)})}) \quad (4.28)$$

(again, the gradients $\partial \tilde{I} / \partial b_j$ are given by expression (4.28) without the $\mathbf{x}^{(i)}$ factor inside the summation). Clearly, the averages in (4.28) are easy to compute, as both $q(\mathbf{x}^{(i)}|\mathbf{y}, z_j)$ and $p(\mathbf{y}|\mathbf{x}^{(i)})$ in this case are Gaussians. The updates (4.21), (4.22), and (4.25) for the remaining parameters are not affected by the change in the channel definition.

Note that in the communication-theoretic interpretation of the considered case, the auxiliary variables \mathbf{z} will need to be communicated over the channel (generally, at a small increase in the communication cost, which in this case is of the order of $|\mathbf{y}| + |\mathbf{z}|$). For the model parameterization described in Section 4.2.2, this would correspond to sending an additional natural number z , which would effectively index the decoder used at the reconstruction. Generally, the compressed representations of $\{\mathbf{x}\}$ will include not only the codes $\{\mathbf{y}\}$, but also the auxiliary labels z . Finally, we may note that unless $p(\mathbf{z}|\mathbf{x})$ is strongly constrained, the mapping $\mathbf{x} \rightarrow \mathbf{z}$ will typically tend to be nearly noiseless, as this would decrease $H(\mathbf{z}|\mathbf{x})$ and maximize $I(\mathbf{x}, \{\mathbf{y}, \mathbf{z}\})$.

4.2.3.1 Comparison with Mixtures of Probabilistic Principal Component Analyzers

It is interesting to compare the IM framework with the likelihood-based training for mixtures of latent variable models (where the auxiliary variables \mathbf{z} are the mixture components, and \mathbf{y} are the latent variable representations of the data patterns). For such models the likelihood may be expressed as

$$\mathcal{L} \stackrel{\text{def}}{=} \langle \log \langle q(\mathbf{x}|\mathbf{y}, \mathbf{z}) \rangle_{q(\mathbf{y})q(\mathbf{z})} \rangle_{\tilde{p}(\mathbf{x})}, \quad (4.29)$$

where $\tilde{p}(\mathbf{x})$ is the empirical distribution. Effectively, $q(\mathbf{x})$ defines a mixture of latent variable models. While the objective (4.29) is generally different from

$\tilde{I}_H(\mathbf{x}, \{\mathbf{y}, \mathbf{z}\})$ (see expression (4.27)), for several special cases we may observe a close relation between maximizing (4.29) in a mixture model and learning the optimal variational *decoders* in the context of variational information maximization.

In particular, for the special case of learning linear projections of the data patterns to a lower-dimensional space, we can compare the I-step of the IM algorithm in our formulation with the M-step of the EM algorithm applied to training mixtures of probabilistic PCA models $\mathcal{M}_L \stackrel{\text{def}}{=} q(\mathbf{y})q(\mathbf{z})q(\mathbf{x}|\mathbf{y}, \mathbf{z})$, where $q(\mathbf{z})$ is a multinomial distribution, $q(\mathbf{y}) \sim \mathcal{N}_{\mathbf{y}}(0, \mathbf{1})$, and $q(\mathbf{x}|\mathbf{y}, z_j) \sim \mathcal{N}_{\mathbf{x}}(\mathbf{U}_j \mathbf{y}, s_j^2)$ (Tipping and Bishop (1999a), Tipping and Bishop (1999b)). For mixtures of PPCAs, the updates for parameters of the mixture components may be expressed as

$$\mathbf{U}_j^{(new)} = \left(\sum_{i=1}^M p(z_j|\mathbf{x}^{(i)}) \mathbf{x}^{(i)} \langle \mathbf{y}^T \rangle_{q_j^{(i)}} \right) \left(\sum_{i=1}^M p(z_j|\mathbf{x}^{(i)}) \left(\langle \mathbf{y} \rangle_{q_j^{(i)}} \langle \mathbf{y}^T \rangle_{q_j^{(i)}} + \Sigma \right) \right)^{-1} \quad (4.30)$$

and

$$\left(s_j^{(new)} \right)^2 = \sum_{i=1}^M \frac{p(z_j|\mathbf{x}^{(i)})}{|\mathbf{x}| \sum_{i=1}^M p(z_j|\mathbf{x}^{(i)})} \text{tr} \left\{ (\mathbf{x} - \mathbf{U}_j \langle \mathbf{y} \rangle_{q_j^{(i)}}) (\mathbf{x} - \mathbf{U}_j \langle \mathbf{y} \rangle_{q_j^{(i)}})^T + \mathbf{U}_j^{(new)} \Sigma \left(\mathbf{U}_j^{(new)} \right)^T \right\} \quad (4.31)$$

(see Tipping and Bishop (1999a)), where the expectations of the codes are computed over the exact component-based posteriors $\langle \mathbf{y} \rangle_{q_j^{(i)}} \equiv \langle \mathbf{y} \rangle_{q(\mathbf{y}|\mathbf{x}^{(i)}, z_j)}$. The posterior is easily expressed from the mixture model \mathcal{M}_L by Bayes rule $q(\mathbf{y}|\mathbf{x}^{(i)}, z_j) \propto q(\mathbf{x}^{(i)}|\mathbf{y}, z_j)q(\mathbf{y})$, which for the considered case gives a Gaussian with the mean $\langle \mathbf{y} \rangle_{q_j^{(i)}} = (s_j^2|\mathbf{y}| + \mathbf{U}_j^T \mathbf{U}_j)^{-1} \mathbf{U}_j^T \mathbf{x}^{(i)}$ (see e.g. von Mises (1964)).

It is easy to see that under the assumption of isotropic Gaussian mixture components, learning the decoder at the I-step of the variational information-maximizing algorithm given by (4.21) and (4.22) has the same form as the PPCA updates (4.30) and (4.31), with the difference that expectations of the codes $\langle \mathbf{y} \rangle$ are computed over the explicitly *constrained encoder* $p(\mathbf{y}|\mathbf{x}^{(i)}, z_j)$ rather than the posterior $q(\mathbf{y}|\mathbf{x}^{(i)}, z_j)$ expressed from the generative model \mathcal{M}_L by applying Bayes rule. These distributions are generally different; particularly, as we mentioned in Section 4.1, the explicit parameterization of $p(\mathbf{y}|\mathbf{x}^{(i)}, z_j)$ makes it easier to impose the required constraints on the communication channel, while the exact analytical form of $q(\mathbf{y}|\mathbf{x}^{(i)}, z_j)$ might not necessarily satisfy such constraints. Generally, the information-theoretic learning of the optimal encoder involves optimization of the bound $\tilde{I}_H(\mathbf{x}, \{\mathbf{y}, \mathbf{z}\})$ in the space of constrained encoder parameters, rather than probabilistic inference in the generative model \mathcal{M}_L .

Interestingly, we can demonstrate that while the considered variational framework for maximizing the bound on mutual information is generally different from maximizing the likelihood in a mixture of probabilistic PCA models, under the assumption of a fixed Gaussian channel noise it gives rise to the same fixed points as a variational EM applied to fitting a mixture of constrained Gaussians with the components $q(\mathbf{x}|\mathbf{y}, z_j) \sim \mathcal{N}_{\mathbf{x}}(\mathbf{U}_j \mathbf{y}, s_j^2)$ and the uniform (rather than the Gaussian) distribution of the hidden variables (i.e. $q(\mathbf{y}) \sim \mathcal{U}_{\mathbf{y}}$). In other words, this specific

application of our framework leads to the same solutions as a variational approach to fitting a mixture of factor analysis-like models with the uniform (rather than the Gaussian) distribution of the hidden factors. We discuss the link between these methods in Appendix B.2.

4.2.3.2 Comparison with Mixtures of Latent Variable Models

By analogy with the results of Chapter 3, we may point out general differences between maximizing the likelihood in mixture models and maximizing the bound on the mutual information in hybrid channels. In particular, the variational IM algorithm applied to training the encoder model of a hybrid channel $\mathcal{M}_I \stackrel{\text{def}}{=} \tilde{p}(\mathbf{x})p(\mathbf{y}, z|\mathbf{x})$ may be compared with the conventional fitting of a mixture of latent variable models $\mathcal{M}_L = p(\mathbf{y})p(z)p(\mathbf{x}|\mathbf{y}, z)$ to a data set given by the empirical distribution $\tilde{p}(\mathbf{x})$. Not surprisingly, the solutions obtained by training the generative and the encoding models would generally be different. By analogy with proposition 3.1, we may show that for the special case when the posteriors of \mathcal{M}_L and \mathcal{M}_I are constrained to be equivalent, we may define tractable variational lower bounds on $I(\mathbf{x}, \{\mathbf{y}, z\})$ which are at least as tight as (and sometimes significantly tighter than) the bound given by the exact likelihood. Effectively, this means that if the goal is to maximize information content which the hidden codes \mathbf{y} and component labels z contain about the data patterns \mathbf{x} , the variational information-maximizing framework should generally be more preferable than maximum-likelihood approaches (at least, in terms of the resulting bounds on $I(\mathbf{x}, \{\mathbf{y}, z\})$). Some of the intersections of the *approximate* approaches to likelihood and mutual information maximization are described in Appendix B.2.

Furthermore, we note once again that the encoder and the generative frameworks are different conceptually. Fundamentally, the goal of maximizing the mutual information in the hybrid channel $\mathbf{x} \rightarrow \{\mathbf{y}, z\}$ is to learn the *encoder* model (a specific tractable choice of the variational decoder is something which facilitates the generally intractable computations). In contrast, the conceptual goal of likelihood training would be to fit the generative model to data, which in our framework would correspond to learning the variational *decoder* (and the marginal distribution of the hidden variables). As discussed in Chapter 3 and Section 4.1, the fundamental feature of the encoder framework is the possibility of specifying explicit constraints on the encoder model $p(\mathbf{y}|\mathbf{x})$. Thus, the suggested variational information-maximizing framework is particularly convenient in situations when our goal is to learn unknown *constrained* encodings of the visible patterns (in practice, such constraints might be artifacts of biophysical or engineering requirements). In general, this is different from generative models, where the desired constraints on the posterior (encoding) mapping $p(\mathbf{y}|\mathbf{x})$ would need to follow from the explicit parameterization of the priors $p(\mathbf{y})$ and conditionals $p(\mathbf{x}|\mathbf{y})$. It is therefore intuitive that it may be rather difficult to apply generative models for addressing constrained encoding problems exactly. For example, apart from relatively simple cases (e.g. when \mathcal{M}_L is the probabilistic PCA model (Tipping and Bishop, 1999b)), it is difficult to apply an exact generative framework for learning an \mathcal{L} -optimal orthonormal projection of the data. As

we showed, the problem may be conveniently addressed within the information-maximization framework, where the encoder constraints may be satisfied as a part of the model’s specification.

4.3 Demonstrations

Here we demonstrate a few applications of the method to extracting optimal orthonormal subspaces for the digits dataset (which is a sub-sampled lower-dimensional version of MNIST, LeCun and Cortes (1998)). As mentioned in Section 4.2.3, apart from simple cases the problem cannot be easily addressed by exact applications of generative models, as we would require specific orthonormal constraints on the posterior distributions. In all cases, it was assumed that $|y| < |x|$. We also assumed that $p(x) = \tilde{p}(x)$ is the empirical distribution.

4.3.1 Hand-Written Digits: Comparing the Bounds

In the first set of experiments, we compared optimal lower bounds on the mutual information $I(x, y)$ obtained by maximizing the as-if Gaussian $I_G(x, y)$ and the auxiliary variational $\tilde{I}(x, y)$ objectives for hand-written digits. The dataset contained $M = 30$ gray-scaled instances of 14-by-14 digits 1, 2, and 8 (10 of each class), which were centered and normalized. The goal was to find an orthogonal projection of the $|x| = 196$ -dimensional training data into a $|y| = 6$ -dimensional space, so that the bounds $I_G(x, y)$ and $\tilde{I}(x, y)$ were maximized. By analogy with Section 4.1.1, we considered a linear Gaussian channel with an irreducible white noise, which in this case leads to the encoder distribution $p(y|x) \sim \mathcal{N}_y(\mathbf{W}y, s^2\mathbf{I})$ with $\mathbf{W} \in \mathbb{R}^{6 \times 196}$. Our interest was in finding optimal orthogonal projections, so the weights were normalized to satisfy $\mathbf{W}\mathbf{W}^T = \mathbf{I}_{|y|}$ (by considering the parameterization $\mathbf{W} = (\tilde{\mathbf{W}}\tilde{\mathbf{W}}^T)^{-1/2}\tilde{\mathbf{W}}$ with $\tilde{\mathbf{W}} \in \mathbb{R}^{|y| \times |x|}$). Effectively, this case corresponds to finding the most informative compressed representations of the source vectors for improving communication of the *non-Gaussian* data over a noisy Gaussian channel (by maximizing lower bounds on the channel capacity). Our specific interest here was to find whether we may indeed improve on Linsker’s as-if Gaussian bound on the mutual information (with the optima given in this case by the PCA projection) by considering a richer family of auxiliary variational bounds with multi-modal mixture-type decoders.

Figure 4.1 shows typical changes in the auxiliary variational bound $\tilde{I}(x, y)$ as a function of the IM’s iterations T for $|z| \in \{2, \dots, 5\}$ states of the discrete auxiliary variable. (On the plot, we ignored the irrelevant constants $H(x)$ identical for both $\tilde{I}(x, y)$ and $I_G(x, y)$, and interpolated $\tilde{I}(x, y)$ for the consecutive iterations). The mappings were parameterized as described in Section 4.2, with the random initializations of the parameters \mathbf{v}_j and \mathbf{b}_j around zero, and the initial settings of the variational prior $q(z) = 1/|z|$. The encoder weights \mathbf{W} were initialized at 6 normalized principal components $\mathbf{W}_{pca} \in \mathbb{R}^{6 \times 196}$ of the sample covariance $\langle \mathbf{x}\mathbf{x}^T \rangle$, and the variance of the channel noise was fixed at $s^2 = 1$. For each choice of the auxiliary space dimension $|z|$, Figure 4.1 (a) shows the results averaged over 30 random initializations of the IM algorithm. As we see from the plot, the IM

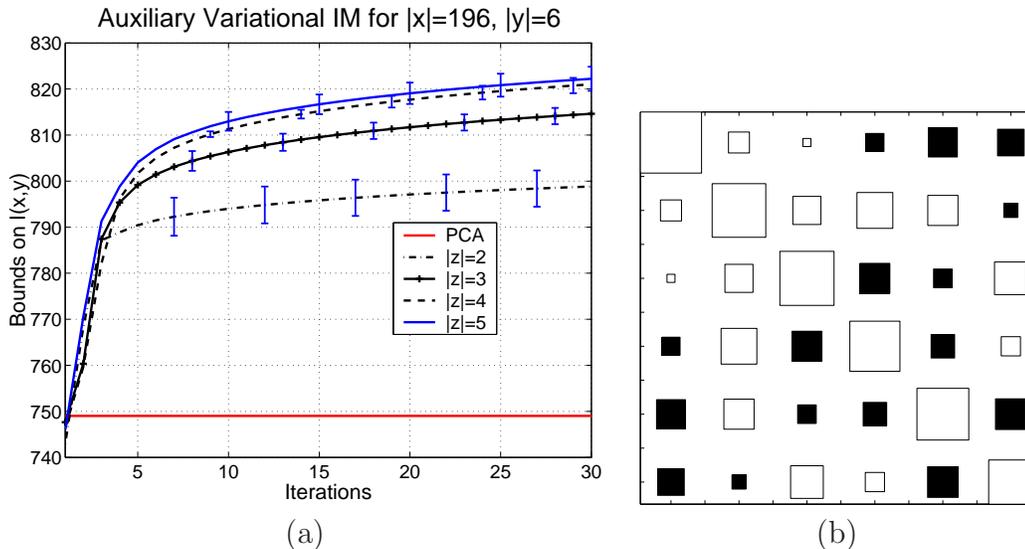


Figure 4.1: Variational information maximization for noisy constrained dimensional-reduction. (a): *Top curves*: Average values of the variational auxiliary bounds $\tilde{I}(x,y)$, obtained by the IM algorithm started at 10 random model initializations (shown for $|z| = 2, \dots, 5$); *bottom line*: the *as-if* Gaussian $I_G(x,y)$ bound (computed numerically). The results are shown for the digits data with $|x| = 196, |y| = 6$ for $M = 30$ patterns and $T = 30$ iterations of the IM. (b): Hinton diagram for $WW_{pca}^T (WW_{pca}^T)^T \in \mathbb{R}^{6 \times 6}$ for $|z| = 3, T = 30$. For orthonormal weights spanning identical subspaces, we would expect to see the identity matrix.

learning leads to a consistent improvement in the auxiliary variational bound, which (on average) varies from $\tilde{I}_0(x,y) \approx 745.7$ to $\tilde{I}_T(x,y) \approx 822.2$ at $T = 30$ for $|z| = 5$. Small variances in the obtained bounds ($\sigma_T \approx 2.6$ for $T = 30, |z| = 5$) indicate a stable increase in the information content independently of the model initializations. From Figure 4.1 (a) we can also observe a consistent improvement in the average $\tilde{I}(x,y)$ with $|z|$, changing as $\tilde{I}_{10}(x,y) \approx 793.9, \approx 806.3, \approx 811.2$, and ≈ 812.9 for $|z| = 2, \dots, 5$ after $T = 10$ IM's iterations. In comparison, the PCA projection weights W_{pca} result in $I_G(x,y) \approx 749.0$, which is visibly worse than the auxiliary bound with the optimized parameters, and is just a little better than $\tilde{I}(x,y)$ computed at a random initialization of the variational decoder for $|z| \geq 2$.

Importantly, we stress once again that the auxiliary variables z are not passed through the channel. In the specific case which we considered here, the auxiliary variables were used to define a more powerful family of variational bounds which we used to extract the \tilde{I} -optimal orthonormal subspace. The results are encouraging, as they show that for a specific constrained channel distribution we may indeed obtain tighter lower bounds on the mutual information $I(x,y)$ *without* communicating more data than in the conventional case. Specifically, for Gaussian channels with orthonormal projections to the code space, we do improve on simple *as-if Gaussian* solutions (leading to the PCA projections) by considering optimization of the auxiliary variational bounds (4.14).

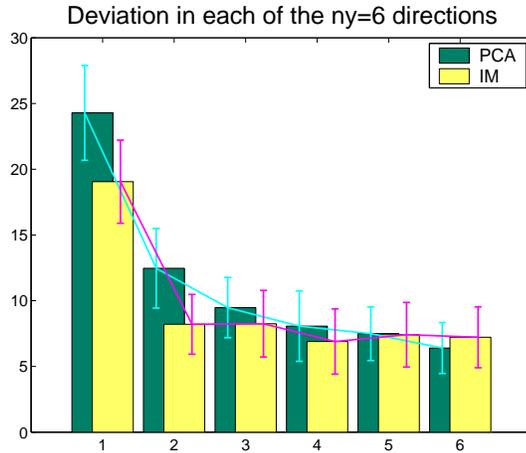


Figure 4.2: *Histogram bars*: marginal variances of the orthonormal noisy projections of the data patterns $\langle (y_i - \langle y_i \rangle)^2 \rangle_{p(y_i|x)\tilde{p}(x)}$ along each of $i = 1, \dots, 6$ dimensions of the code space. The code space $\mathcal{R}^{|y|}$ is spanned by W_{pca} and W , for Linsker’s and the auxiliary variational objectives respectively. *Vertical lines*: variances of the average distances for $M = 30$ data patterns. The results are shown for $|x| = 196$, $|y| = 6$; the auxiliary space size $|z| = 3$; the number of iterations $T = 30$.

As expected, we may also note that the \tilde{I} -optimal *encoder* weights W are in general different from rotations of W_{pca} . This is easy to see by computing $WW_{pca}^T(WW_{pca}^T)^T$, which in our case is visibly different from the identity matrix (see Fig. 4.1 (b) for $|y| = 6$ and $|z| = 3$), which we would have expected to obtain otherwise. Effectively, this means that by allowing a greater flexibility in the choice of the *variational decoder* distributions, the $\tilde{I}(x, y)$ -optimal constrained *encoders* become different from the optimal encoders of simpler models. This result is intuitive: it is natural to expect that a richer structure of the decoder model may change our notion of the optimal codes (at least, at the stage of optimizing $\tilde{I}(x, y)$). And vice versa, a choice of simple variational decoders (such as linear Gaussians) may impose severe constraints on the types of codes which they can decode efficiently, which may lead to a loss in the coding efficiency and result in a general reduction in the retained information content (see Chapter 5.3 for a more detailed discussion).

Finally, figure 4.2 shows average distances of the noisy linear projections of the testing data from the y -space mean (i.e. the marginal variance of the codes) for $M = 30$ patterns. The histogram indicates the variances of the projections for each of the $|y| = 6$ dimensions of the code space, spanned by W_{pca} and W (for Linsker’s and the auxiliary variational bounds respectively). The results are shown for $M = 30$ digit patterns after $T = 30$ iterations of learning (again, we assumed that the size of the auxiliary space $|z| = 3$). We can see that the encoded representations produced by the auxiliary variational method result in a more uniform spectrum of the projection variances, with roughly equal error bars. This result is not unexpected, as under the fixed channel noise assumption, the information content between the data and the codes increases with an increase

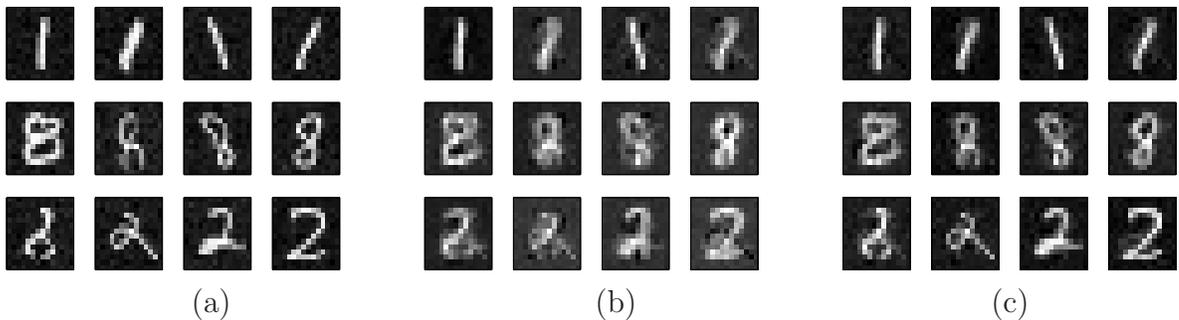


Figure 4.3: Reconstructions of the source patterns from encoded representations. (a): A subset of the generic patterns used to generate the source vectors; (b): the corresponding reconstructions from 6 principal components; (c): the corresponding \tilde{I}_H -optimal reconstructions at $\langle \mathbf{x} \rangle_{q(\mathbf{x}|\mathbf{y},z)} = \mathbf{U}_z \mathbf{y}$ for the hybrid $\{\mathbf{y}, z\}$ representations ($|\mathbf{y}| = 6$, $|z| = 3$).

in the entropy of the encoded representations. Whilst it is not easy to compute the exact entropy of the mixture distribution $p(\mathbf{y})$, qualitatively it is natural to expect that more uniform encodings would typically result in higher values of the mutual information between the encodings and the source patterns.

4.3.2 Hand-Written Digits: Reconstructions

Additionally, for the problem settings described in Sec. 4.3.1, we have computed reconstructions of the source patterns $\{\mathbf{x}\}$ from their noisy encoded representations. First, we generated source vectors by adding an isotropic Gaussian noise to the generic patterns (see Fig. 4.3 (a)), where the variance of the source noise was set as $s_s^2 = 0.5$. Then we computed noisy linear projections $\{\mathbf{y}\}$ of the source vectors by using the I_G - and the \tilde{I}_H - optimal encoder weights (in the latter case, we also computed the auxiliary label z by sampling from the learned $p(z|\mathbf{x})$). This stage corresponds to passing encoded representations over the noisy channels, where the noise variance for the Gaussian part of the channel was fixed at $s^2 = 1$. Finally, we have used the optimal *approximate* decoders to perform the reconstructions from $\{\mathbf{y}\}$ (for I_G -optimal PCA projections) and $\{\mathbf{y}, z\}$ (for \tilde{I}_H -optimal hybrid channels).

As we see from Figure 4.3 (b), (c), by a slight modification of the channel (due to encoding and communicating a multinomial variable z), we may achieve a visible improvement in the reconstruction of the sources by using the \tilde{I}_H - optimal projections². The results are shown for $|\mathbf{y}| = 6$, $|z| = 3$ after $T = 3$ iterations, and the reconstructions are computed at the analytical mean of the decoder’s component $q(\mathbf{x}|\mathbf{y}, z)$ indexed by the auxiliary variable z . Even though the resulting hybrid channel may be difficult to justify from the communication viewpoint, the

²Similar reconstructions could be obtained by maximizing the auxiliary bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ *without* communicating z . However, the approximate decoder for this case would be given as $q(\mathbf{x}|\mathbf{y}) = \sum_z q(\mathbf{x}|\mathbf{y}, z) \frac{\langle p(z|\mathbf{x})p(\mathbf{y}|\mathbf{x}) \rangle_{p(\mathbf{x})}}{\langle p(z|\mathbf{x}) \rangle_{p(\mathbf{x})}}$, which requires knowing $p(\mathbf{x})$.

results suggest that maximization of the bound on $I(\mathbf{x}, \{\mathbf{y}, z\})$ provides a sensible way to reduce dimensionality of the sources for the purpose of reconstructing inherently noisy non-Gaussian patterns. Importantly, the variational decoder $q(z|\mathbf{x}, \mathbf{y})$ which maximizes $\tilde{I}_H(\mathbf{x}, \{\mathbf{y}, z\})$ makes no recourse to $p(\mathbf{x})$. Therefore, just like in the PCA case, we do not need to store the training instances in order to perform an accurate reconstruction from noisy lower-dimensional projections. We note once again that the weights of the optimal encoder were chosen to satisfy the specific orthonormality constraint (though other kinds of constrained encoders may easily be considered). This contrasts with the exact approaches to training generative models, where encoding constraints may be more difficult to enforce.

4.4 Summary

Here we considered an application of the variational information maximizing approach to linear orthonormal dimensionality reduction in the presence of irreducible Gaussian noise. We showed that the well-known *as-if* Gaussian approximation of the mutual information (Linsker (1992)), which may be seen as a special case of the variational bound for correlated linear Gaussian variational distributions, leads to the PCA solution for isotropic linear Gaussian channels. Importantly, this means that by using linear Gaussian variational decoders under the considered Gaussian channel, maximization of the generic lower bound (2.2) on the mutual information cannot improve on the PCA projections.

The situation becomes strikingly different if we consider a richer family of variational auxiliary lower bounds on $I(\mathbf{x}, \mathbf{y})$ under the same encoding constraints. In particular, we showed that in the cases when the source distribution was non-Gaussian, we could significantly improve on the PCA projections by considering multi-modal variational decoders. This confirms the conceptually simple idea that by allowing a greater flexibility in the choice of variational decoders, we may get significant improvements over simple bounds on the mutual information at a limited increase in the computational cost. This result is interesting from the communication-theoretic perspective, as it demonstrates a simple and computationally efficient way to produce better bounds on the capacity of communication channels without altering channel properties (e.g. without communicating more data across the channels).

Finally, we discussed a simple information-theoretic approach to constrained dimensionality reduction for hybrid representations $\mathbf{x} \rightarrow \{\mathbf{y}, z\}$, which may significantly improve reconstructions of the sources $\{\mathbf{x}\}$ from their lower-dimensional representations $\{\mathbf{y}\}$ at a small increase in the transmission cost (given by $|z|$). We applied the hybrid framework for extracting an informative orthonormal projection subspace of the data. While being vaguely related to the exact maximum-likelihood fitting of mixtures of latent variable models (at least, in terms of the specified variable domains), our variational information-maximizing framework is fundamentally different in terms of model specifications. One of the important features of the IM framework is the explicit parameterization of the encoder model, which is particularly convenient in situations when the goal is to learn unknown encodings of the visible patterns for a known family of constrained

encoder distributions. Usually, such constraints would be difficult to impose in generative models. On the other hand, they may be easily introduced in the suggested information-maximizing framework, while a tractable choice of variational decoders simplifies the computations for large-scale stochastic channels.

Generally, we have confirmed the intuition that we may improve on simple generic bounds on $I(x, y)$ by considering a richer family of the *auxiliary variational* bounds, which effectively increase the power of the variational decoders. On the other hand, it is natural to expect that we may maximize the information content by considering richer families of *encoder* distributions. In the next chapter we consider specific applications of the information-maximizing framework to the case of nonlinear channels. We will describe theoretical properties of \tilde{I} -optimal nonlinear Gaussian encoder distributions, and focus specifically on tractable practical applications to information-theoretic clustering.

Chapter 5

Variational Information Maximization for Nonlinear Dimensionality Reduction

In previous chapters we introduced a simple variational lower bound on the mutual information, which resolves some of the computational difficulties of computing the exact mutual information. Then we compared our framework with the conventional learning in generative and autoencoder-type models, and considered specific applications of the variational information maximization to the constrained dimensionality reduction. As we demonstrated in Chapter 4, the variational lower bounds on the mutual information can be made tighter by considering more powerful families of variational decoders. It is also intuitive that one may hope to obtain tighter bounds on the mutual information by increasing the power of the stochastic *encoder*. Intuitively, by appropriately choosing nonlinear encoder distributions so that the encoder satisfies specific local constraints, we may hope to obtain better reconstructions of the transmitted sources, more anthropomorphically sensible visualizations in the compressed variable spaces, etc. Here we explore these matters by considering the problem of maximizing information content for stochastic non-linear channels.

We will focus primarily on the problem of information-theoretic clustering, where the encoder distribution is defined by a generally stochastic, nonlinear mapping from the source patterns to discrete cluster labels. For this case, we will consider two principally different learning techniques. First, we will consider optimizing the specific lower bound on the mutual information, where the encoder distribution of the channel is given by the exact posterior of the corresponding generative model (see expression (3.11) and the discussion in Section 3.2.1). Effectively, this approach may be viewed as an information-theoretic method of training generative models, which we study for the case of Gaussian mixtures. Then we consider a different information-theoretic clustering technique, where we maximize the exact mutual information in explicitly parameterized encoder models. For most of the clustering applications, this may be done reasonably easily, since the cardinality of the code space will typically be low, so that maximization of the exact objective will typically be computationally tractable. For

this case, we describe a simple and practical algorithm for unsupervised discriminative learning of cluster allocations. By allowing the flexibility in the choice of the encoder structure, we may consider a variety of encoder model parameterizations, including those involving nonlinear projections of the source vectors into high-dimensional feature spaces. Empirically, we demonstrate that the resulting information-theoretic clustering approach favorably compares with the common generative clustering methods.

In the second part of the chapter we will briefly review some of the theoretical properties of the IM for higher-dimensional code spaces, and show that some of the popular dimensionality reduction techniques may be seen as special instances of the variational information maximization procedure.

5.1 Introduction

The limiting case of a dimensionality reduction problem, where the code y of a high-dimensional pattern \mathbf{x} is defined by a single multinomial hidden variable, may be viewed as a simple instance of *clustering*. The codes $\{y\}$ in this case may be interpreted as the unknown cluster labels corresponding to the data patterns. If the family of the generally stochastic mappings $\mathbf{x} \rightarrow y$ satisfies specific local constraints, the resulting cluster allocations may provide a simple anthropomorphically sensible visualization of generally high-dimensional patterns.

Arguably, there are at least two simple requirements which a cluster allocation procedure should satisfy. First of all, we may require the resulting clusters of data patterns to be locally smooth. There may be several ways to interpret local smoothness in the clustering context; for example, it may be sensible to require that each two patterns have a high probability of being assigned to the same cluster if the corresponding vectors satisfy specific geometric constraints. Secondly, we may require the clustering method to avoid assigning unique cluster labels to outliers (or other constrained regions in the data space), so that local regions in the data space are not over-represented in the code space.

Generally, we may expect the clusters to be well-separated if patterns \mathbf{x} are predictive of cluster labels y , for all the data patterns $\{\mathbf{x}\}$. Here we argue that a reasonable way of learning optimal cluster allocations is by maximizing the exact mutual information $I(\mathbf{x}, y) = H(y) - H(y|\mathbf{x})$. Indeed, in the clustering context $I(\mathbf{x}, y)$ may be interpreted as a measure of predictability of cluster allocations from training patterns. Intuitively, it is clear that in contrast to the generative mixture model clustering, optimization of $I(\mathbf{x}, y)$ does not intrinsically favour outliers. (Indeed, for the fixed cardinality $|y| < M$ of the code space, assignments of unique cluster labels to under-represented training patterns would have resulted in a decrease in the marginal entropy $H(y)$. In other words, regions under-represented in the data space would have been over-represented in the code space, which for the fixed noise levels would have led to a reduction in the mutual information). Generally, we may expect that the mutual information maximization favours hard assignments of cluster labels to regions of roughly identical sizes, resulting in the growth in $H(y)$ and reduction in $H(y|\mathbf{x})$.

In the first part of the chapter we will discuss simple and efficient information-theoretic clustering algorithms, which may be used for *discriminative unsupervised*¹ learning of nonlinear cluster allocations. First we will consider optimizing the lower bound $\hat{I}(\mathbf{x}, y) = \langle \log p(\mathbf{x}|y) \rangle_{p(y|\mathbf{x})\hat{p}(\mathbf{x})}$, where the encoder distribution of the channel is given by the exact posterior of the corresponding Gaussian mixture model $\mathcal{M}_{\mathcal{L}} = p(y)p(\mathbf{x}|y)$, where $p(y)$ are mixture coefficients, and $p(\mathbf{x}|y)$ are the Gaussian components. Effectively, this approach may be viewed as an information-theoretic method of training Gaussian mixture models. We compare this approach with the standard likelihood maximization for Gaussian mixtures, and show that by analogy with the EM training, the IM algorithm reduces to the k-means clustering in the limiting case of zero-variance components.

Then we consider a different information-theoretic clustering technique, where we impose specific constraints on the encoding distribution of the channel model \mathcal{M}_I , and maximize $I(\mathbf{x}, y)$ directly. Unlike most of the existing information-theoretic approaches applicable in the clustering context (see e.g. Fisher and Principe (1998), Torkkola and Campbell (2000), Gokcay and Principe (2002), Dhillon et al. (2002), Corduneanu and Jaakkola (2003), Jenssen et al. (2003)), we consider optimization of the exact mutual information in an encoder model, which for the purpose of learning cluster allocations may be computed exactly. Furthermore, as a straight-forward extension of our method, we consider clustering in the feature space, with a constrained kernelized representation of the encoder for nonlinearly transformed source patterns. We show that for this case we may apply the information-maximization framework to learn the optimal kernel parameters, which often leads to sensible solutions of the clustering problem.

In the last part of the chapter we extend the discussion of the information-maximizing framework by considering other types of nonlinear encoders. We will focus particularly on the discussion of *non-linear Gaussian* channels with the fixed isotropic noise. Since in this case maximization of the mutual information cannot be performed exactly, we consider optimizing the generic lower bound (2.2) on the mutual information with several specific choices of the variational decoder, and describe properties of the obtained optimal solutions. While our analysis of the analytical properties of the IM solutions is mainly theoretical, it shows interesting links of our variational procedure to other popular dimensionality reduction techniques, inducing PCA, kernel PCA (Schoelkopf et al. (1998), Mika et al. (1999)), and Gaussian Process Latent Variable Models (Lawrence (2003)) as special cases. Some of these theoretical findings extend previous results for noiseless autoencoders (Bourlard and Kamp (1988), Baldi and Hornik (1989), Diamantaras and Kung (1996)); however, they are derived for *stochastic* communication channels and general nonlinear mappings. Most of the derivations for this chapter can be found in Appendix C.

¹As mentioned above, the encoder models are specified by explicitly parameterizing the stochastic mapping $p(y|\mathbf{x})$ to the code space, which corresponds to a similar parameterization of a discriminative model. However, in contrast to discriminative models, the encodings y of the sources \mathbf{x} are presumed to be hidden in our case.

5.2 Nonlinear Channels for Information-Theoretic Clustering

One of the practical side-effects of Section 3.2 is a practical information-theoretic method of training generative models (see expression (3.11)). Here we apply it to training Gaussian mixture models. With a slight abuse of terminology, we will refer to the method as the IM for Gaussian mixtures; by this we mean that the encoder distribution of the channel model trained by the IM has the exact form of the posterior $p(y|x)$ of a Gaussian mixture model. We show that the IM applied for training the considered models indeed provides a sensible clustering technique. Moreover, the results help to empirically verify the theoretical findings of Chapter 3.

5.2.1 Variational Information Maximization for Gaussian Mixtures

In proposition 3.1 we showed that maximization of the exact likelihood \mathcal{L} , expressed from the generative model $\mathcal{M}_L \stackrel{\text{def}}{=} p(y)p(x|y)$ for i.i.d. patterns, may be viewed as maximization of a lower bound on the mutual information in the encoder model $\mathcal{M}_I \stackrel{\text{def}}{=} \tilde{p}(x)p(y|x)$, where y is a latent code, x is a source vector, and $\tilde{p}(x)$ is the empirical distribution. However, as we noted in Section 3.2.1, the likelihood bound on $I(x, y)$ may potentially be weak, and much tighter bounds on $I(x, y)$ may be obtained by considering a special instance of the variational information-maximizing framework. As we showed in Section 3.2.1, the bound provided by the likelihood is weaker or as tight as

$$\hat{I}(x, y) = \langle \log p(x|y) \rangle_{\tilde{p}(y|x)\tilde{p}(x)} + H_{\tilde{p}(x)} = \langle \log p(x|y) \rangle_{p(y|x)\tilde{p}(x)} + H_{\tilde{p}(x)}, \quad (5.1)$$

where we defined the encoding distribution of the encoder model

$$\tilde{p}(y|x) \propto p(y)p(x|y) \quad (5.2)$$

to be the exact posterior of the generative model \mathcal{M}_L .

As in the conventional approaches to latent variable modeling, our goal would be to produce informative latent variable representations of the sources, so optimization of the tighter bound (5.1) on $I(x, y)$ may be a reasonable strategy to consider. Moreover, as discussed in Section 3.2.1, under the considered constraint on the encoder, expression (5.1) may be viewed as an information-theoretic objective for training both the generative \mathcal{M}_L and the encoder \mathcal{M}_I models². Effectively, this is our main motivation for considering optimization of the bound $\hat{I}(x, y)$. If it is indeed the case that optimization of objective (5.1) helps to avoid common degeneracies of the likelihood solutions (such as overfitting to local data segments), we may hope that the data generated from \mathcal{M}_L (trained by maximizing the bound $\hat{I}(x, y)$) could in some sense be more representative of the underlying process than the data sampled from the likelihood-trained mixture.

²This contrasts with the exact likelihood and the exact mutual information training which are commonly applied to training generative and encoder models respectively.

Here we will consider a simple case when \mathcal{M}_L is a Gaussian mixture, where $y \in \{y_j | j = 1, \dots, |y|\}$ is the mixture label, and $p(\mathbf{x}|y_j) \sim \mathcal{N}_{\mathbf{x}}(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$, where $\boldsymbol{\mu}_j \in \mathbb{R}^{|\mathbf{x}|}$, $\boldsymbol{\Sigma}_j \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|}$ are the mean and the covariance corresponding to the j^{th} component y_j . The objective $\hat{I}(\mathbf{x}, y)$ may be easily expressed as

$$\hat{I}(\mathbf{x}, y) = \frac{1}{M} \sum_{m=1}^M \sum_{j=1}^{|y|} \frac{p(y_j)p(\mathbf{x}^{(m)}|y_j)}{p(\mathbf{x}^{(m)})} \log p(\mathbf{x}^{(m)}|y_j), \quad (5.3)$$

which needs to be optimized for $p(y_j)$ and $\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j$ for all $j = 1, \dots, |y|$. For simplicity, we will start the discussion of the resulting optimization procedure by *implying* the constraints on the optimized parameters. Effectively, this implication allows us to treat all the parameters (including the mixture coefficients) as unconstrained vectors in high-dimensional real spaces; we will impose meaningful construction constraints shortly.

5.2.1.1 Learning Optimal Parameters

By computing the matrix derivatives of (5.3) for $\boldsymbol{\Sigma}_j$ and $\boldsymbol{\mu}_j$ (see e.g. Magnus and Neudecker (1999), Minka (2000)), it is easy to find that

$$\frac{\partial \hat{I}(\mathbf{x}, y)}{\partial \boldsymbol{\mu}_j} = -\frac{1}{M} \sum_{m=1}^M \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}^{(m)} - \boldsymbol{\mu}_j) p(y_j | \mathbf{x}^{(m)}) (\hat{\alpha}_j^{(m)} + 1) \quad (5.4)$$

$$\frac{\partial \hat{I}(\mathbf{x}, y)}{\partial \boldsymbol{\Sigma}_j} = -\frac{1}{2M} \sum_{m=1}^M \boldsymbol{\Sigma}_j^{-1} ((\mathbf{x}^{(m)} - \boldsymbol{\mu}_j)(\mathbf{x}^{(m)} - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}_j^{-1} - \mathbf{I}_{|\mathbf{x}|}) p(y_j | \mathbf{x}^{(m)}) (\hat{\alpha}_j^{(m)} + 1), \quad (5.5)$$

where

$$\begin{aligned} \hat{\alpha}_j^{(m)} &\stackrel{\text{def}}{=} \log p(\mathbf{x}^{(m)}|y_j) - \sum_{l=1}^{|y|} p(y_l | \mathbf{x}^{(m)}) \log p(\mathbf{x}^{(m)}|y_l) \\ &= \log \frac{p(\mathbf{x}^{(m)}|y_j)}{p(\mathbf{x}^{(m)})} - KL(p(y|\mathbf{x}^{(m)}) || p(y)), \end{aligned} \quad (5.6)$$

and $p(y_j | \mathbf{x}^{(m)})$ is the exact posterior expressed from the generative model \mathcal{M}_L . Note that in the trivial case when $\hat{\alpha}_j^{(m)} = \text{const}$ for all $m = 1, \dots, M$ and $j = 1, \dots, |y|$, the gradients (5.4) and (5.5) are identical to those obtained by maximizing the log-likelihood of a Gaussian mixture model (up to irrelevant constant pre-factors). We may further express the functional derivatives (e.g. Gelfand and Fomin (1963), Weinstock (1974), Smith (1998)) for the mixture coefficients $p(y_j)$ to get

$$\begin{aligned} \frac{\partial \hat{I}(\mathbf{x}, y)}{\partial p(y_j)} &= \frac{1}{M} \sum_{m=1}^M \frac{p(y_j | \mathbf{x}^{(m)})}{p(y_j)} \left(\log p(\mathbf{x}^{(m)}|y_j) - \sum_{l=1}^{|y|} p(y_l | \mathbf{x}^{(m)}) \log p(\mathbf{x}^{(m)}|y_l) \right) \\ &= \frac{1}{M} \sum_{m=1}^M \frac{p(y_j | \mathbf{x}^{(m)})}{p(y_j)} \hat{\alpha}_j^{(m)}. \end{aligned} \quad (5.7)$$

(again, we have implied that the mixing coefficients $p(y_j)$ lie in a convex space (e.g. Boyd and Vandenberghe (2004)), i.e. there is no need to explicitly incorporate Lagrange multipliers into the objective (5.3)). If the coefficients $\hat{\alpha}_j^{(m)}$ were constant, we could use (5.7) to derive the fixed-point updates for $p(y_j)$ (subject to the normalizing constraints), which in this case would lead to the updates of the iterative EM algorithm for Gaussian mixtures (e.g. Bishop (1995)).

From (5.4) – (5.7) we see that for Gaussian mixture models, optimization of the likelihood \mathcal{L} and the bound $\hat{I}(\mathbf{x}, y)$ result in generally different learning rules, as specified by the model-dependent weighting coefficients $\hat{\alpha}_j^{(m)}$. Note that $\hat{\alpha}_j^{(m)}$ includes the Kullback-Leibler divergence between the posteriors $p(y|\mathbf{x})$ and the prior $p(y)$, which quantifies the difference between the likelihood and the information-theoretic objectives (see expression (3.4)). It is easy to see that in the degenerate case when the latent variables are independent of the data (i.e. for $p(y|\mathbf{x}) = p(y)$), the factor $\hat{\alpha}_j^{(m)}$ will become irrelevant, and both EM and IM will result in simple data averaging for $\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j$ with arbitrary setting of $p(y_j)$ – clearly, this is not an interesting case to consider. Generally, however, $\hat{\alpha}_j^{(m)}$ is a function of the optimized mixture coefficients $p(y)$ and the parameters $\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j$, which makes it awkward to derive the closed-form iterative updates for the model parameters. This contrasts with the standard EM algorithms (Dempster et al. (1977)), which may be viewed as maximizing the expected complete data log-likelihood (see e.g. Bishop (1995)).

A simple alternative to the closed-form iterative optimization could be given by any of the known non-linear optimization methods (e.g. Bishop (1995), Dennis and Schnabel (1996), Bertsekas (1999), Galeev and Tihomirov (2000)), where the optimized parameters lie in appropriately constrained spaces. In this specific case we would need to constrain the covariance matrices to ensure their positive-definiteness; we would also need to ensure that the mixing coefficients $p(y_j)$ follow the probability requirements. Such constraints are easy to impose by considering meaningful non-restrictive constructions. For instance, to ensure positive-definiteness, we may parameterize the covariances $\boldsymbol{\Sigma}_j$ as

$$\boldsymbol{\Sigma}_j = \mathbf{A}_j \mathbf{A}_j^T + I_{|x|} \epsilon, \quad (5.8)$$

where $\mathbf{A}_j \in \mathbb{R}^{|x| \times |x|}$ is an arbitrary real-valued matrix and $\epsilon \gtrsim 0$ is a small constant, which ensures that $\boldsymbol{\Sigma}_j$ is non-singular. Analogously, we may parameterize the mixing coefficients as

$$p(y_j) = \exp\{a_j\} / \sum_{l=1}^{|y|} \exp\{a_l\}, \quad (5.9)$$

which ensures satisfaction of the probability constraints. From the existence of the singular value decomposition of $\boldsymbol{\Sigma}_j$ (see e.g. Golub and Loan (1996)) and continuity of $a_j \in \mathbb{R}$, it is clear that parameterizations (5.8), (5.9) are non-restrictive³.

³In what follows we assume that $\epsilon \rightarrow 0^+$, which implies that eigenvalues of $\boldsymbol{\Sigma}_j$ are the squares of the non-negative singular vectors of \mathbf{A}_j . Clearly, in the considered limit the construction (5.8)

From the definition of the objective (5.1) and the construction (5.8), (5.9), it is easy to see that the constrained optimization of $\hat{I}(\mathbf{x}, \mathbf{y})$ with respect to the covariances Σ_j and mixture coefficients $p(y_j)$ may be transformed to the unconstrained optimization problem for \mathbf{A}_j and a_j as

$$\frac{\partial \hat{I}(\mathbf{x}, \mathbf{y})}{\partial a_j} = p(y_j) \left(\frac{\partial \hat{I}(\mathbf{x}, \mathbf{y})}{\partial p(y_j)} - \sum_l \frac{\partial \hat{I}(\mathbf{x}, \mathbf{y})}{\partial p(y_l)} p(y_l) \right), \quad \frac{\partial \hat{I}(\mathbf{x}, \mathbf{y})}{\partial \mathbf{A}_j} = 2 \frac{\partial \hat{I}(\mathbf{x}, \mathbf{y})}{\partial \Sigma_j} \mathbf{A}_j, \quad (5.10)$$

where $\partial \hat{I}(\mathbf{x}, \mathbf{y}) / \partial \Sigma_j$ and $\partial \hat{I}(\mathbf{x}, \mathbf{y}) / \partial p(y_j)$ are given by (5.5), (5.7), and $j = 1, \dots, |\mathbf{y}|$. The gradients (5.4) and (5.10) for $\boldsymbol{\mu}_j$, \mathbf{A}_j , and a_j may then be used by a numerical optimization procedure performing an ascent on the bound $\hat{I}(\mathbf{x}, \mathbf{y})$.

Finally, we note once again that we can view the procedure of optimizing the bound (5.3) as a way of training the generative model \mathcal{M}_L . Once the model is trained, we can use it for generating new training data. If optimization of the bound $\hat{I}(\mathbf{x}, \mathbf{y})$ helps to avoid degeneracies of the \mathcal{L} -optimal solutions, we may expect that the generated samples may be more representative of the underlying clusters than samples from the mixture model trained by maximizing the likelihood. Specifically, we may hope that the information-theoretic training may limit over-representations of the local segments of the training data, which may occur due to possible singularities of the likelihood solutions.

5.2.1.2 Relation to the K-means Algorithm

Interestingly, we may note that in the limiting case of deterministic *decoders* $p(\mathbf{x}|y_j) \sim \mathcal{N}_{\mathbf{x}}(\boldsymbol{\mu}_j, \Sigma_j)$ with $\Sigma_j = \sigma^2 \mathbf{I} \rightarrow 0$, optimization of $\hat{I}(\mathbf{x}, \mathbf{y})$ reduces to the well-known k-means algorithm. Indeed, for this case the bound on the mutual information may be expressed as

$$\begin{aligned} \hat{I}(\mathbf{x}, \mathbf{y}) &= \langle \log p(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})\hat{p}(\mathbf{x})} \propto - \sum_{m=1}^M \sum_{j=1}^{|\mathbf{y}|} \|\mathbf{x}^{(m)} - \boldsymbol{\mu}_j\|^2 p(y_j|\mathbf{x}^{(m)}) \\ &= - \sum_{m=1}^M \sum_{j=1}^{|\mathbf{y}|} \|\mathbf{x}^{(m)} - \boldsymbol{\mu}_j\|^2 \frac{p(y_j)}{p(y_j) + \sum_{k \neq j}^{|\mathbf{y}|} p(\mathbf{x}^{(m)}|y_k)p(y_k)/p(\mathbf{x}^{(m)}|y_j)}, \end{aligned} \quad (5.11)$$

where we have ignored the entropy of the empirical distribution since it has no effect on the optimization surface for $p(\mathbf{y}|\mathbf{x})$. By analogy with the well-known limiting case of the Gaussian mixtures (see e.g. Bishop (1995)), we can see that in the considered limit the second factor in the summation (5.11) reduces to the Kronecker delta $\delta(j - \arg \min_k \|\mathbf{x}^{(m)} - \boldsymbol{\mu}_k\|)$, which leads to the k-means updates for $\boldsymbol{\mu}_j$ (Hartigan and Wong (1979)). In other words, both the likelihood and the variational information maximization reduce to the k-means algorithm in the limit of zero-variance spherical components of a Gaussian mixture model. This is also a limiting case of a specific Information Bottleneck formulation (e.g. Tishby

does not restrict the covariance space. Non-restrictiveness of the parameterization of $p(y_j)$ trivially follows from the range of (5.9).

et al. (1999), Slonim (2002)) for the model $\mathbf{x} \leftarrow i \rightarrow y$, where $i = 1, \dots, M$ is the training pattern indicator, $\mathbf{x} : i \rightarrow \mathbf{x}$ is a one-to-one deterministic “pattern-extracting” mapping, and $p(y|i) = p(y|\mathbf{x}^{(i)})$ is the exact posterior of a flat spherical Gaussian mixture model (Still et al. (2004), Still and Bialek (2004)). We may therefore conclude that while likelihood maximization, information bottleneck, and the variational IM define very different methodologies for training different graphical models, their point of intersection for the limiting case of Gaussian mixture models in the noiseless limits is the k-means algorithm.

5.2.2 Information Maximization for Clustering with Encoder Models

We will now discuss a more intuitive approach to information-theoretic clustering by considering encoder models $\mathcal{M}_I \stackrel{\text{def}}{=} \tilde{p}(\mathbf{x})p(y|\mathbf{x})$, where $\{\mathbf{x}\}$ is the set of data patterns, and $\{y\}$ is the set of the corresponding (and generally unknown) cluster labels. As in the standard approaches to information maximization in intrinsically tractable channels, we will maximize the exact mutual information with respect to parameters of the explicitly parameterized encoder $p(y|\mathbf{x})$. In many practical clustering applications, we will have the flexibility of choosing the encoder’s parameterization. For clustering, we may be particularly interested in encoders which favour consistent cluster allocations to locally smooth neighborhoods of the training patterns (by this we mean that data patterns which are close to each other according to some distance measure would need to be clustered similarly). As discussed in Chapter 3, one of the fundamental advantages of the information-maximization principle for the encoder models \mathcal{M}_I is the simplicity of imposing such *smoothness* constraints on the encoder distribution; moreover, by aiming to maximize the capacity of the resulting channels, we obtain a proper information-theoretic framework for performing fully unsupervised training. Another potential advantage of clustering with encoder (rather than generative) models is the observation that we do not need to specify the data-generating process to be able to cluster the data properly. Indeed, generative approaches suffer from the fact that $p(\mathbf{x}|y)$ needs to be normalized in \mathbf{x} , which in high dimensions restricts the class of the generative models to mixtures of simple distributions (such as Gaussians). Usually data will lie on low dimensional curved manifolds embedded in the high dimensional \mathbf{x} -space. If we are restricted to using mixtures of Gaussians to model this curved manifold, typically a very large number of mixture components will be required. No such restrictions apply in the information maximization case, so that the mappings $p(y|\mathbf{x})$ may be very complex, subject only to sensible local constraints (Agakov and Barber (2005b)).

Before discussing applications of the IM to specific nonlinear channels, it is important to outline fundamental differences of the suggested formulation from the existing approaches to information-theoretic clustering and feature extraction. We stress that our approach to clustering is a special case of the general formulation of the information-maximization problem for a specific choice of the code variables. In particular, our goal here is to obtain informative codes (cluster labels) by optimizing the information content $I(\mathbf{x}, y)$ in the undercom-

plete stochastic channel. While the information-maximizing framework seems to be somewhat related to the existing information-theoretic approaches to feature extraction (Principe et al. (2000), Torkkola and Campbell (2000), Gokcay and Principe (2002)), it is in many ways fundamentally different. Specifically, these methods presume a complete observability of the data patterns and the corresponding class labels $\{\mathbf{x}^{(i)}, y^{(i)} | i = 1, \dots, M\}$, and suggest to maximize the information content $I(\mathbf{z}, y)$ between the *known* cluster labels $\{y\}$ and the hidden transformations $\{\mathbf{z}\}$ of the training patterns $\{\mathbf{x}\}$. Since the resulting mutual information is generally computationally intractable, it is typically approximated by a somewhat heuristic “*information potential*”, which may be related to the quadratic Renyi entropy $H_R(\mathbf{z}, y) \stackrel{\text{def}}{=} -\log\langle p(\mathbf{z}, y) \rangle_{p(\mathbf{z}, y)}$. The joint distribution $p(\mathbf{z}, y)$ is usually evaluated by applying Parzen density estimation techniques (see e.g. Principe et al. (1998), Fisher and Principe (1998), Torkkola (2000)). Effectively, these methods may be interpreted as a way of training discriminative models $\mathbf{x} \rightarrow \mathbf{z} \rightarrow y$ by maximizing an approximation of $I(\mathbf{z}, y)$, which is computationally and conceptually different from the information-maximizing formulation for encoder models $\mathbf{x} \rightarrow y$.

Other clustering techniques, which make a recourse to information theory, presume partial observability of the cluster labels (Szummer and Jaakkola (2002), Corduneanu and Jaakkola (2003)). The key idea there is to learn the encoder distribution $p(y|\mathbf{x})$ by maximizing the conditional likelihood for the labeled part of the training set $\{\mathbf{x}^{(i)}, y^{(i)} | i = 1, \dots, L\}$, *penalized* by the mutual information $I(\mathbf{x}, y)$ for all (labeled and unlabeled) source patterns $\{\mathbf{x}^{(i)} | i = 1, \dots, M\}$, where $M \geq L$. Clearly, this framework is also different from maximizing mutual information in encoder models; specifically, when applied to clustering of the unlabeled data (i.e. $L = 0$), it would result in *minimization* of the information content between the clusters and the patterns. Of course, this would give rise to independent cluster assignments $p(y|\mathbf{x}) = p(y)$ for unconstrained non-parametric settings. Other methods applied in the vaguely related regularized contexts (Tishby et al. (1999), Slonim et al. (2001), Chechik and Tishby (2002), Dhillon et al. (2002)) are also different from the information-maximizing framework in terms of the conceptual definitions of optimality, channel definitions, observability of the modeled domain, and availability of the additional observable *variables of relevance*.

Our motivation here is to explore whether optimization of the exact mutual information for nonlinear channels may give rise to an anthropomorphically sensible clustering technique. As a simple choice of the encoder for the information-theoretic clustering method, we could trivially consider

$$p(y_j|\mathbf{x}^{(i)}) \propto \exp\{-\|\mathbf{x}^{(i)} - \mathbf{w}_j\|^2/s_j + b_j\}, \quad (5.12)$$

where the cluster centers \mathbf{w}_j , the dispersions s_j , and the biases b_j are the encoder parameters to be learned, and we use the notation y_j to indicate that the code variable y is in state j . Clearly, the encoder constraint (5.12) favors local smoothness of the underlying clusters, so that patterns \mathbf{x} lying close to specific centers \mathbf{w}_j in the *data* space will tend to be clustered together. We may note that parameterization of the channel (5.12) has a strong relation to the posteriors expressed from the isotropic Gaussian mixture models. Indeed, for the specific setting of

$b_j = \log p(y_j) - (|x|/2) \log s_j$, expression (5.12) would reduce to the channel defined by the posterior of the corresponding Gaussian mixture (see Section 5.2.1). In general, for encoder distributions parameterized by (5.12) we will assume that $\mathbf{b} \stackrel{\text{def}}{=} \{b_j\} \in \mathbb{R}^{|y|}$ is an unconstrained free parameter, and the optimized objective is the exact mutual information rather than the specific bound $\hat{I}(\mathbf{x}, \mathbf{y})$ discussed in Section 5.2.1. In principle, we could consider other choices of the encoder distributions; however (5.12) will prove to be particularly convenient for the kernelized representations which we will introduce at a later stage.

As the dimensionality of the code space $|y|$ in clustering problems is usually relatively small, the summations over the $|y|$ states may be performed exactly, and we may proceed by optimizing the exact mutual information $I(\mathbf{x}, \mathbf{y})$. In our case it may be expressed analytically as

$$I(\mathbf{x}, \mathbf{y}) \propto \sum_{m=1}^M \sum_{j=1}^{|y|} p(y_j | \mathbf{x}^{(m)}) \log p(y_j | \mathbf{x}^{(m)}) - \sum_{k=1}^M \sum_{j=1}^{|y|} p(y_j | \mathbf{x}^{(k)}) \log \frac{1}{M} \sum_{l=1}^M p(y_j | \mathbf{x}^{(l)}). \quad (5.13)$$

Clearly, $I(\mathbf{x}, \mathbf{y})$ may be computed in $O(M^2|y|)$.

5.2.2.1 Learning Optimal Parameters

Objective (5.13) needs to be optimized with respect to the encoder parameters (\mathbf{w}_j , s_j , and b_j for the channels defined by (5.12)). By computing the functional derivatives of (5.13) for the encoders $p(y|\mathbf{x})$, we get

$$\frac{\partial I(\mathbf{x}, \mathbf{y})}{\partial p(y_j | \mathbf{x}^{(m)})} = \frac{1}{M} \log \frac{p(y_j | \mathbf{x}^{(m)})}{p(y_j)}. \quad (5.14)$$

This implies

$$\frac{\partial I(\mathbf{x}, \mathbf{y})}{\partial \theta_k} = \frac{1}{M} \sum_{m=1}^M \sum_{j=1}^{|y|} \log \frac{p(y_j | \mathbf{x}^{(m)})}{p(y_j)} \frac{\partial p(y_j | \mathbf{x}^{(m)})}{\partial \theta_k}, \quad (5.15)$$

where θ_k parameterizes $p(y_k | \mathbf{x})$ (again, by construction we have presumed that $p(y|\mathbf{x})$ lies in the convex space of the conditional distributions). Note that due to the normalization, the parameters θ_k would occur in the normalizing constants of $p(y_j | \mathbf{x})$ for all $j = 1, \dots, |y|$, which leads to

$$\frac{\partial p(y_j | \mathbf{x}^{(m)})}{\partial \theta_k} = -p(y_j | \mathbf{x}^{(m)}) \left[\frac{\partial f_j(\mathbf{x}^{(m)})}{\partial \theta_k} - \sum_{l=1}^{|y|} p(y_l | \mathbf{x}^{(m)}) \frac{\partial f_l(\mathbf{x}^{(m)})}{\partial \theta_k} \right]. \quad (5.16)$$

Here we assumed that $p(y_j | \mathbf{x}^{(m)}) \propto \exp\{-f_j(\mathbf{x}^{(m)})\}$ with the potentials defined as $f_j \stackrel{\text{def}}{=} \|\mathbf{x}^{(i)} - \mathbf{w}_j\|^2 / s_j - b_j$. Now we can easily express the gradients of the exact mutual information $I(\mathbf{x}, \mathbf{y})$ for the means \mathbf{w}_j and the isotropic variances s_j as

$$\frac{\partial I(\mathbf{x}, \mathbf{y})}{\partial \mathbf{w}_j} = \frac{1}{M} \sum_{m=1}^M p(y_j | \mathbf{x}^{(m)}) \frac{(\mathbf{x}^{(m)} - \mathbf{w}_j)}{s_j} \alpha_j^{(m)}, \quad (5.17)$$

$$\frac{\partial I(\mathbf{x}, \mathbf{y})}{\partial s_j} = \frac{1}{M} \sum_{m=1}^M p(y_j | \mathbf{x}^{(m)}) \frac{\|\mathbf{x}^{(m)} - \mathbf{w}_j\|^2}{2s_j^2} \alpha_j^{(m)}, \quad (5.18)$$

where $\alpha_j^{(m)}$ is again the increment of the specific Kullback-Leibler term

$$\alpha_j^{(m)} \stackrel{\text{def}}{=} \alpha_j(\mathbf{x}^{(m)}) \stackrel{\text{def}}{=} \log \frac{p(y_j|\mathbf{x}^{(m)})}{p(y_j)} - KL(p(y|\mathbf{x}^{(m)}) \| \langle p(y|\mathbf{x}) \rangle_{\tilde{p}(\mathbf{x})}). \quad (5.19)$$

Note that in contrast to the weighting coefficients (5.6) of the IM training of Gaussian mixture models, the KL divergence in (5.19) is computed between the posterior $p(y|\mathbf{x})$ and their empirical (rather than model-based) average. Analogously, we get

$$\frac{\partial I(\mathbf{x}, y)}{\partial b_j} = \frac{1}{M} \sum_{m=1}^M p(y_j|\mathbf{x}^{(m)}) \alpha_j^{(m)}. \quad (5.20)$$

It is easy to see that apart from the variable factors, the fixed point updates (5.17) and (5.18) are identical to (5.4) and (5.5) with the assumed isotropic covariances. A simple construction constraint ensuring that $s_j > 0$ may be given by $s_j \stackrel{\text{def}}{=} \exp\{\tilde{s}_j\}$ where $\tilde{s}_j \in \mathbb{R}$. For this case, we may re-express the gradients for the variances as

$$\partial I(\mathbf{x}, y) / \partial \tilde{s}_j = s_j \partial I(\mathbf{x}, y) / \partial s_j. \quad (5.21)$$

Expressions (5.17) – (5.21) could then be used by a numerical optimization procedure performing an ascent on $I(\mathbf{x}, y)$ for \mathbf{w}_j , \tilde{s}_j , and b_j , where $j = 1, \dots, |y|$, as discussed in Agakov and Barber (2005c).

5.2.3 Information Maximization for Clustering with Kernelized Encoder Models

We will now extend (5.12) by considering kernelized (Aizerman et al. (1964), Boser et al. (1992)) parameterizations of encoders $p(y_j|\mathbf{x})$. Let us assume that the source patterns $\mathbf{x}^{(i)}$, $\mathbf{x}^{(j)}$ should have a high probability of being assigned to the same cluster, if they lie close to a specific cluster center in some *feature* space. A reasonable choice of the encoder distribution for this case is given by

$$p(y_j|\mathbf{x}^{(i)}) \propto \{-\|\phi(\mathbf{x}^{(i)}) - \mathbf{w}_j\|^2 / s_j + b_j\}, \quad (5.22)$$

where $\phi(\mathbf{x}^{(i)}) \in \mathbb{R}^{|\phi|}$ is the feature vector corresponding to the source pattern $\mathbf{x}^{(i)}$, and $\mathbf{w}_j \in \mathbb{R}^{|\phi|}$ is the (unknown) cluster center in the feature space. Typically, we will presume intractability of explicit computations in feature spaces (which may occur, for example, when $|\phi| \rightarrow \infty$). We will also presume that it is impossible to store vectors in the feature space, which enforces constraints on the encoder parameters $\mathbf{w}_j \in \mathbb{R}^{|\phi|}$.

Note that since each cluster center \mathbf{w}_i in the feature space $\mathbb{R}^{|\phi|}$ has the same dimensionality as the projected source patterns $\phi(\mathbf{x}^{(i)})$, it is representable in the basis of the projections as

$$\mathbf{w}_j = \sum_{m=1}^M \alpha_{mj} \phi(\mathbf{x}^{(m)}) + \mathbf{w}_j^\perp, \quad (5.23)$$

where $\tilde{\mathbf{w}}_j^\perp \in \mathbb{R}^{|\phi|}$ is orthogonal to the span of $\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_M)$, and $\{\alpha_{mj}\}$ is a set of coefficients (here j and m index $|y|$ codes and M patterns respectively). Then we may transform the encoder distribution (5.22) to

$$\begin{aligned} p(y_j|\mathbf{x}^{(m)}) &\propto \exp\left\{-\left(K_{mm} - 2\mathbf{k}^T(\mathbf{x}^{(m)})\mathbf{a}_j + \mathbf{a}_j^T\mathbf{K}\mathbf{a}_j + c_j\right)/s_j\right\} \\ &\stackrel{\text{def}}{=} \exp\{-f_j(\mathbf{x}^{(m)})\}, \end{aligned} \quad (5.24)$$

where $\mathbf{k}(\mathbf{x}^{(m)})$ corresponds to the m^{th} column (or row) of the Gram matrix $\mathbf{K} \stackrel{\text{def}}{=} \{K_{ij}\} \stackrel{\text{def}}{=} \{\phi(\mathbf{x}^{(i)})^T\phi(\mathbf{x}^{(j)})\} \in \mathbb{R}^{M \times M}$, $\mathbf{a}_j \in \mathbb{R}^M$ is the j^{th} column of the matrix of the coefficients $\mathbf{A} \stackrel{\text{def}}{=} \{a_{mj}\} \in \mathbb{R}^{M \times |y|}$, and $c_j = (\mathbf{w}_j^\perp)^T \mathbf{w}_j^\perp - s_j b_j$. Without loss of generality, we may assume that $\mathbf{c} = \{c_j\} \in \mathbb{R}^{|y|}$ is a free unconstrained parameter, as the assumption does not limit the parameter's domain. Also, by analogy with the previous case (see expressions (5.12), (5.21)), we will ensure positivity of the dispersions s_j by assuming $s_j = \exp\{\tilde{s}_j\}$. Learning in the encoder model may now be seen as unconstrained optimization of $I(\mathbf{x}, y)$ for the encoder parameters.

5.2.3.1 Learning Optimal Parameters

It is easy to see that for a *fixed* and *known* Gram matrix $\mathbf{K} \in \mathbb{R}^{M \times M}$, the considered construction (5.23) helps us to avoid computations in high-dimensional feature spaces $\mathbb{R}^{|\phi|}$. Indeed, if we know \mathbf{K} , the encoder (5.24) is a function of the coefficients $\mathbf{A} \stackrel{\text{def}}{=} \{a_{jm}\} \in \mathbb{R}^{M \times |y|}$, $\mathbf{c} \in \mathbb{R}^{|y|}$, and s_j for $j = 1, \dots, |y|$, i.e. evaluation of the encoder potential $f_j(\mathbf{x}^{(m)})$ does not require explicit computations in high-dimensional feature spaces.

The exact mutual information (5.13) should be optimized with respect to the log-dispersions $\tilde{s}_j \equiv \log(s_j)$, biases c_j , and coordinates \mathbf{A} in the space spanned by the feature vectors $\{\phi(\mathbf{x}^{(i)}) | i = 1, \dots, M\}$. From expressions (5.15) and (5.24) we get

$$\frac{\partial I(\mathbf{x}, y)}{\partial \mathbf{a}_j} = \frac{1}{s_j} \langle p(y_j|\mathbf{x}) (\mathbf{k}(\mathbf{x}) - \mathbf{K}\mathbf{a}_j) \alpha_j(\mathbf{x}) \rangle_{\tilde{p}(\mathbf{x})} \in \mathbb{R}^M, \quad (5.25)$$

$$\frac{\partial I(\mathbf{x}, y)}{\partial \tilde{s}_j} = \frac{1}{2s_j} \langle p(y_j|\mathbf{x}) f_j(\mathbf{x}) \alpha_j(\mathbf{x}) \rangle_{\tilde{p}(\mathbf{x})}. \quad (5.26)$$

Also, by analogy with (5.20) we obtain

$$\partial I(\mathbf{x}, y) / \partial c_j = \langle \alpha_j(\mathbf{x}) \rangle_{\tilde{p}(\mathbf{x})}, \quad (5.27)$$

where the coefficients $\alpha_j(\mathbf{x})$ are given by (5.19). For a known Gram matrix \mathbf{K} , the gradients $\partial I / \partial \mathbf{a}_j$, $\partial I / \partial \tilde{s}_j$, and $\partial I / \partial c_j$ given by expressions (5.25) – (5.27) may be used for numerical optimization of the model parameters (e.g. Luenberger (1973), Bertsekas (1999)).

Note that so far we have assumed that the Gram matrix $\mathbf{K} \in \mathbb{R}^{M \times M}$ is fixed and known. Clearly, in general it needs to be recomputed for each new dataset. Moreover, in addition to learning the feature space coordinates \mathbf{A} , we could potentially consider learning optimal features, which would generally lead to changes in the matrices of scalar products. Since we presume that explicit computations

in $\mathbb{R}^{|\phi|}$ are expensive, we cannot compute the Gram matrix by trivially applying its definition $\mathbf{K} = \{\phi(x_i)^T \phi(x_j)\}$. Instead, we may interpret scalar products in feature spaces as *kernel functions*

$$\phi(x^{(i)})^T \phi(x^{(j)}) = \mathcal{K}_{\Theta}(x^{(i)}, x^{(j)}; \Theta), \quad \forall x^{(i)}, x^{(j)} \in \mathcal{R}_x, \quad (5.28)$$

where $\mathcal{K}_{\Theta} : \mathcal{R}_x \times \mathcal{R}_x \rightarrow \mathbb{R}$ satisfies Mercer's kernel properties (symmetry and non-negative definiteness, see e.g. Mercer (1909), Courant and Hilbert (1953), Vapnik (1998), Smola (1998), Cristianini and Shawe-Taylor (2000)). We may therefore evaluate scalar products in feature spaces $K_{ij} = \mathcal{K}_{\Theta}(x^{(i)}, x^{(j)})$ by performing computations in the data space $\mathbb{R}^{|\mathbf{x}|}$; the projections into $\mathbb{R}^{|\phi|}$ will in this case be implied.

A number of nonlinear extensions of common models (applicable in a variety of contexts: from unsupervised dimensionality reduction to supervised classification) typically consider specific kernel functions with fixed parameters (e.g. see a review in Smola (1998), Cristianini and Shawe-Taylor (2000), Scholkopf and Smola (2002)). This raises a natural question of finding optimal settings of kernel parameters, which would implicitly correspond to finding optimal nonlinearities for constrained function spaces. The problem of learning optimal kernels may be conveniently addressed in the (supervised) context of probabilistic classification or regression, with Gaussian process priors over regression functions (Williams and Rasmussen (1996), Williams (1997), Williams (1998)) or the functions' arguments (Barber and Williams (1997), Williams and Barber (1998), Gibbs and MacKay (2000)). Optimal kernel parameters for these cases are typically obtained by maximizing the likelihoods with respect to parameters of covariance functions. Generally, the likelihoods will be computed for a set of labeled patterns. An important difference of our information-theoretic formulation is that it offers a convenient and simple way to learn parameters of kernel functions \mathcal{K}_{Θ} in the unsupervised context. (Additionally, in contrast to Bishop, Svensen and Williams (1998a), Bishop, Svensen and Williams (1998b), we do not need to impose explicit constraints on the code space representations). For the considered channel definitions (5.22), learning the kernel function parameters is particularly computationally convenient; moreover, choosing appropriate constraints on the kernel matrices in this case is generally less complicated than for some other communication channels (*cf* Agakov and Barber (2004c)).

5.2.3.2 Learning Optimal Kernels

We can apply our information-maximizing framework to learning of optimal kernel parameters. Indeed, from (5.13) and (5.15) we can derive the gradients of the exact mutual information with respect to the parameters Θ of the kernel function \mathcal{K}_{Θ} . After some algebraic manipulations, we get

$$\begin{aligned} \frac{\partial I(x, y)}{\partial \Theta} &= \frac{1}{M} \sum_{m=1}^M KL(p(y|x^{(m)})||p(y)) \sum_{k=1}^{|y|} \frac{\partial f_k(x^{(m)})}{\partial \Theta} p(y_k|x^{(m)}) - \\ &\quad \frac{1}{M} \sum_{m=1}^M \sum_{j=1}^{|y|} \frac{\partial f_j(x^{(m)})}{\partial \Theta} p(y_j|x^{(m)}) \log \frac{p(y_j|x^{(m)})}{p(y_j)}, \end{aligned} \quad (5.29)$$

where $f_k(\mathbf{x}^{(m)})$ is given by (5.24). Note that the computational complexity of computing the updates for Θ is $O(M|y|^2)$, where M is the number of training patterns and $|y|$ is the number of clusters. It is easy to see that generally (5.29) does not require further approximations. Also, we may note that neither computation of the objective (5.13), nor computation of its gradients (5.25) – (5.27), (5.29) requires inversion of the Gram matrix $\mathbf{K} \in \mathbb{R}^{M \times M}$.

In the special case of the radial basis function (RBF) kernels

$$\mathcal{K}_\beta(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp\{-\beta\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2\}, \quad (5.30)$$

the gradients of the encoder potentials are given by

$$\frac{\partial f_j(\mathbf{x}^{(m)})}{\partial \beta} = \frac{1}{s_j} \left(\mathbf{a}_j^T \tilde{\mathbf{K}} \mathbf{a}_j - 2\tilde{\mathbf{k}}^T(\mathbf{x}^{(m)}) \mathbf{a}_j \right), \quad (5.31)$$

where $\tilde{\mathbf{K}} \stackrel{\text{def}}{=} \{\tilde{K}_{ij}\} \stackrel{\text{def}}{=} K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})(1 - \delta(\mathbf{x}^{(i)} - \mathbf{x}^{(j)}))$, and δ is the Kronecker delta. By substituting (5.31) into the general expression (5.29), we obtain the gradient of the mutual information with respect to the RBF kernel parameters (Agakov and Barber (2005b)). As usual, the kernel parameter may be learned by performing a numerical ascent on the objective (5.13). Similarly, we may learn continuous parameters of other kernel functions, or use the mutual information (5.13) as a criterion for kernel comparison.

The discussed framework suggests a proper information-theoretic approach to clustering by projecting the data to potentially high-dimensional feature spaces. Importantly, we stress that in contrast to other techniques performing clustering in feature spaces (such as the kernelized k-means, e.g. Dhillon et al. (2004), Wang et al. (2004)), the information-maximizing framework suggests a principled way of learning optimal kernels. Moreover, the proper information-theoretic interpretation may facilitate extensions of the method to richer channel distributions.

5.3 Nonlinear Gaussian Channels

In Section 5.2.3 we assumed that the code space was low-dimensional and discrete (i.e. $\mathbf{y} \in \{y_1, \dots, y_{|y|}\}$), which made it possible to optimize the exact mutual information $I(\mathbf{x}, \mathbf{y})$ analytically. For that case, we considered several kinds of nonlinear Gaussian encoders (see (5.2), (5.12), (5.22)) and outlined the corresponding IM learning rules. We also showed that if the nonlinear channel was defined by the exact posterior of a Gaussian mixture model, the limiting noiseless case of the bound (5.3) gave rise to the well-known k-means clustering algorithm. Here we consider a different family of communication channels, where the code space is high-dimensional and continuous (i.e. $\mathbf{y} \in \mathbb{R}^{|y|}$). Specifically, we focus on applying the variational information-maximizing framework to high-dimensional nonlinear channels $\mathbf{x} \rightarrow \mathbf{y}$ with the independent isotropic Gaussian noise. The encoder distribution we consider is given by $p(\mathbf{y}|\mathbf{x}) \sim \mathcal{N}_{\mathbf{y}}(\mathbf{W}\phi(\mathbf{x}), s^2\mathbf{I}_{|y|})$, where $\mathbf{W} \in \mathbb{R}^{|y| \times |\phi|}$. By analogy with Section 5.2.3, we presume that explicit computations in $\mathbb{R}^{|\phi|}$ are expensive, and work with the kernelized encoders instead. For this case, we

consider optimizing the generic bound (2.2) for several choices of the variational decoder distributions $q(\mathbf{x}|\mathbf{y})$.

Our objective here is to outline several important properties of the IM approach for Gaussian channels and tractable decoding distributions. Specifically, as the first obvious choice of $q(\mathbf{x}|\mathbf{y})$, we consider using linear Gaussian decoders, which simplifies computations of the bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ and facilitates the analysis of optimal solutions for encoder parameters. We show that for the case of isotropic channel noise, nothing is gained by using nonlinear encoders and linear decoders in the context of variational information maximization. This result extends the work of Bourlard and Kamp (1988) and Bourlard (2000) to stochastic communication channels with arguably more general choices of encoding nonlinearities.

Then we consider variational information maximization for the case of nonlinear Gaussian variational distributions. Generally, such choices of $q(\mathbf{x}|\mathbf{y})$ may significantly increase the complexity of computing the generic lower bound $\tilde{I}(\mathbf{x}, \mathbf{y})$, which motivates further reformulations of the variational procedure. We find that the generic lower bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ may indeed be formally modified to ensure tractable computations, which leads to kernel PCA (Schoelkopf et al. (1998)) as the optimal solution for *encoder* weights $\mathbf{W} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$. By analogy with a simpler case of discrete nonlinear channels (see Section 5.2.3), the IM suggests a formal procedure for learning kernel parameters. However, we point out that in order to avoid degenerate solutions for parameters of the encoder model, it may generally be important to impose constraints on the Gram matrix $\mathbf{K} \in \mathbb{R}^{M \times M}$, or carefully choose the feature-to-data decoding mappings.

Finally, we outline a simple relation of our framework to the recent work on Gaussian Process Latent Variable Models (Lawrence (2003)), which may be interpreted as the variational information-maximization procedure in the *noiseless* limit of a nonlinear channel. The presentation in this section summarizes the obtained theoretical results. For their derivations and extended discussions we refer the reader to Appendix C.

5.3.1 Analytic Properties of IM Solutions

As discussed in Section 1.4, optimization of the exact mutual information for high-dimensional code spaces is generally computationally intractable. Here we describe optimization of the bound

$$I(\mathbf{x}, \mathbf{y}) \geq \tilde{I}(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} H(\mathbf{x}) + \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{x})p(\mathbf{y}|\mathbf{x})}$$

for the case of a nonlinear Gaussian encoder $p(\mathbf{y}|\mathbf{x}) \sim \mathcal{N}_{\mathbf{y}}(\mathbf{W}\boldsymbol{\phi}(\mathbf{x}), s^2\mathbf{I}_{|\mathbf{y}|})$. Note that for the data patterns $\{\mathbf{x}^{(m)} | m = 1, \dots, M\}$, the set of encodings $\{\mathbf{y}^{(m)}\}$ is given by a noisy linear projection from the (potentially high-dimensional) feature space $\{\boldsymbol{\phi}(\mathbf{x}^{(m)})\}$. We also assume that $|\mathbf{y}| \leq |\mathbf{x}| \ll |\boldsymbol{\phi}|$ and $|\mathbf{x}| \leq M$, so that \mathbf{y} is a compressed representation of \mathbf{x} , and the number of training points is sufficient to ensure invertibility of the sample covariance (for linearly independent patterns). We will focus primarily on the discussion of the IM framework for the cases when variational decoders $q(\mathbf{x}|\mathbf{y})$ are linear and nonlinear Gaussians.

5.3.1.1 Linear Gaussian Decoders

Let us assume that $p(y|x) \sim \mathcal{N}_y(\mathbf{W}\phi(\mathbf{x}), s^2\mathbf{I}_{|y|})$, and the variational decoder is a linear Gaussian $q(\mathbf{x}|y) \sim \mathcal{N}_x(\mathbf{U}y, \sigma^2\mathbf{I}_{|x|})$. By analogy with (5.23), we note that the encoder weights are representable in the basis of feature vectors $\{\phi(\mathbf{x}^{(i)})^T | i = 1, \dots, M\}$ as

$$\mathbf{W} = \mathbf{A}\mathbf{F}^T + \mathbf{W}^\perp \in \mathbb{R}^{|y| \times |\phi|}, \quad \mathbf{F} \stackrel{\text{def}}{=} [\phi(\mathbf{x}^{(1)}), \dots, \phi(\mathbf{x}^{(M)})] \in \mathbb{R}^{|\phi| \times M}, \quad (5.32)$$

where $\mathbf{A} = \{\alpha_{ij}\} \in \mathbb{R}^{|y| \times M}$ is the matrix of coefficients, and the rows of \mathbf{W}^\perp are the corresponding orthogonal compliments (see expression (5.23)). Also, by analogy with Section 5.2.3 we may define the Gram matrix of scalar products in $\mathbb{R}^{|\phi|}$ as

$$\mathbf{K} \stackrel{\text{def}}{=} \{K_{ij}\} \stackrel{\text{def}}{=} \{\phi(\mathbf{x}^{(i)})^T \phi(\mathbf{x}^{(j)})\} = \mathbf{F}^T \mathbf{F} \in \mathbb{R}^{M \times M}. \quad (5.33)$$

After some straight-forward algebraic manipulations (see Appendix C.2), we may express the bound (2.2) on $I(\mathbf{x}, \mathbf{y})$ as

$$\tilde{I}(\mathbf{x}, \mathbf{y}) \propto \frac{1}{M\sigma^2} \text{tr} \{\mathbf{U}\mathbf{A}\mathbf{B}^T\} - \frac{s^2}{\sigma^2} \text{tr} \{\mathbf{U}^T \mathbf{U}\} - \frac{1}{2M\sigma^2} \text{tr} \{\mathbf{U}^T \mathbf{U} \mathbf{A} \mathbf{K}^2 \mathbf{A}^T\} - \frac{1}{2\sigma^2} \text{tr} \{\mathbf{S}\} + c, \quad (5.34)$$

where $\mathbf{S} \stackrel{\text{def}}{=} \langle \mathbf{x}\mathbf{x}^T \rangle = \sum_m \mathbf{x}^{(m)}(\mathbf{x}^{(m)})^T / M \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|}$ is the sample covariance of the centered data, $\mathbf{B} \stackrel{\text{def}}{=} \sum_{m=1}^M \mathbf{x}^{(m)} \mathbf{k}(\mathbf{x}^{(m)})^T \in \mathbb{R}^{|\mathbf{x}| \times M}$, $\mathbf{k}(\mathbf{x}^{(m)}) \in \mathbb{R}^{M \times 1}$ is the m^{th} column of the Gram matrix, and c is an irrelevant constant. If $\mathbf{K} \in \mathbb{R}^{M \times M}$ is fixed, the objective (5.34) needs to be optimized for the encoder coefficients and decoder weights.

By optimizing the bound (5.34) with respect to the encoder coefficients $\mathbf{A} \in \mathbb{R}^{|y| \times M}$ and plugging the optimal values back into (5.34), we may express the objective $\tilde{I}(\mathbf{x}, \mathbf{y})$ as a function of the *decoder* weights $\mathbf{U} \in \mathbb{R}^{|\mathbf{x}| \times |y|}$. Interestingly, in our case this leads to

$$\tilde{I}(\mathbf{x}, \mathbf{y}) = \frac{1}{2\sigma^2} \text{tr} \{\mathbf{U}(\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T \mathbf{S}\} - \frac{s^2}{2\sigma^2} \text{tr} \{\mathbf{U}\mathbf{U}^T\} - \frac{1}{2\sigma^2} \text{tr} \{\mathbf{S}\} + c, \quad (5.35)$$

which is exactly the expression for the optimal bound for *linear* Gaussian communication channels (see Appendix C.1 and expression (C.7)). The result is derived in the context of variational information maximization for a stochastic channel $\mathbf{x} \rightarrow \mathbf{y}$, and it holds independently of the specific form of the nonlinearity $\phi(\mathbf{x})$. Indeed, we have made no assumptions about the mappings to the feature space. Hence, we reach an important conclusion: in the considered variational framework, *nothing is gained by using a nonlinear encoder and a linear Gaussian variational decoder for isotropic channel noise*. This extends the related work of Boulard and Kamp (1988) and Boulard (2000), who showed that for noiseless 1-layer autoencoders, there are no gains of using specific nonlinearities⁴.

By analogy with Appendix C.1 we may note that the optimal left singular vectors of $\mathbf{U} \in \mathbb{R}^{|\mathbf{x}| \times |y|}$ span the subspace defined by the $|y|$ principal components

⁴To show this, Boulard (2000) considers noiseless autoencoders with the encoding functions approximately linear around $x = 0$, i.e. $\phi(x) \approx a_0 + a_1 x$.

of \mathbf{S} . Moreover, for both linear and nonlinear Gaussian encoders, the bound (5.35) is maximized when the weights are unconstrained (see the discussion in Appendix C.1). By imposing norm constraints on $\mathbf{U} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{y}|}$, both channels result in the same optimization surface for the decoder weights independently of the choice of nonlinearity.

The result suggests that in order to improve the power of the method, we need to consider both nonlinear encoders and decoders. However, from (2.2) it is clear that in the stochastic context, the naive approach of using a nonlinear decoder will typically result in intractable averages over \mathbf{y} in the expression for the variational bound $\tilde{I}(\mathbf{x}, \mathbf{y})$. In order to avoid the computational problems, we derive a modified bound on the mutual information by considering further relaxations of the generic bound and performing decoding in the feature space.

5.3.1.2 Nonlinear Gaussian Decoders

The considered nonlinear Gaussian channel may be represented by the Markov chain $\mathbf{x} \rightarrow \mathbf{f} \rightarrow \mathbf{y}$, where $\mathbf{f} \in \mathbb{R}^{|\phi|}$ and $p(\mathbf{f}|\mathbf{x}) \sim \delta(\mathbf{f} - \phi(\mathbf{x}))$, $p(\mathbf{y}|\mathbf{f}) \sim \mathcal{N}_y(\mathbf{W}\mathbf{f}, \mathbf{\Sigma}_y)$. Intuitively, since projections to the feature space are deterministic, the codes \mathbf{y} are as predictable from the feature vectors $\mathbf{f} \stackrel{\text{def}}{=} \phi(\mathbf{x}) \in \mathbb{R}^{|\phi|}$ as they are from the source variables \mathbf{x} (see Appendix C.3 and proposition C.1). Then the bound on $I(\mathbf{x}, \mathbf{y})$ may be expressed as

$$I(\mathbf{x}, \mathbf{y}) = I(\mathbf{f}, \mathbf{y}) \geq \tilde{I}(\mathbf{f}, \mathbf{y}), \text{ where } \tilde{I}(\mathbf{f}, \mathbf{y}) \stackrel{\text{def}}{=} \langle \log q(\mathbf{f}|\mathbf{y}) \rangle_{\tilde{p}(\mathbf{x})p(\mathbf{f}|\mathbf{x})p(\mathbf{y}|\mathbf{f})} + H(\mathbf{f}). \quad (5.36)$$

Evaluation of the bound (5.36) is complicated by the need of computing the entropy of the features $H(\mathbf{f})$. Despite the fact that the mapping to the feature space is deterministic, we have assumed that we do not generally know explicit feature space representations of the training patterns, i.e. numeric approximations due to Brunel and Nadal (1998), Shriki et al. (2002), Corduneanu and Jaakkola (2003) are not directly applicable. Moreover, even if explicit computations in the feature space were possible, the existing numerical methods would not generally retain a proper bound on $I(\mathbf{x}, \mathbf{y})$, i.e. other ways of handling the intractability may need to be considered.

In Appendix C.3.1 and C.3.2 we discuss several ways of addressing the problem of evaluating the feature space entropy in the context of the optimization problem. As one of such methods, we suggest optimizing a proper variational relaxation of (5.36) given by

$$I(\mathbf{x}, \mathbf{y}) \geq \tilde{I}(\mathbf{x}, \mathbf{y}) \geq \langle \log q(\mathbf{f}|\mathbf{y}) \rangle_{\tilde{p}(\mathbf{x})p(\mathbf{f}|\mathbf{x})p(\mathbf{y}|\mathbf{f})} + \langle \log q(\mathbf{x}|\mathbf{f}) \rangle_{\tilde{p}(\mathbf{x})p(\mathbf{f}|\mathbf{x})} + c, \quad (5.37)$$

where c is an irrelevant constant (note that it corresponds to a sequential application of the generic bound for the chain $\mathbf{x} \rightarrow \mathbf{f} \rightarrow \mathbf{y}$). As in the conventional applications of the variational IM algorithm (see Section 2.1.2), we will be optimizing (5.37) with respect to parameters of the encoder $p(\mathbf{y}|\mathbf{f})$, the *feature decoder* $q(\mathbf{f}|\mathbf{y})$, and the *data decoder* $q(\mathbf{x}|\mathbf{f})$. Note that objective (5.37) is a proper bound on $I(\mathbf{x}, \mathbf{y})$ independently of the specific parameterization of the encoder and the variational decoder distributions.

By assuming that the *feature space* decoder is an isotropic linear Gaussian $q(\mathbf{f}|\mathbf{y}) \sim \mathcal{N}_f(\mathbf{U}\mathbf{y}, \sigma_f^2 \mathbf{I})$ and considering kernelized representations of $p(\mathbf{y}|\mathbf{f})$ and $q(\mathbf{f}|\mathbf{y})$, it is easy to see that the optimal encoder and decoder weights $\mathbf{W}^T, \mathbf{U} \in \mathbb{R}^{|\phi| \times |\mathbf{y}|}$ span the principal eigen-subspace of

$$\mathbf{S}_F = \frac{1}{M} \sum_{i=1}^M \phi(\mathbf{x}^{(i)}) \phi(\mathbf{x}^{(i)})^T. \quad (5.38)$$

Hence, for a fixed kernel function \mathcal{K}_Θ and data decoder $q(\mathbf{x}|\mathbf{f})$ of the considered nonlinear Gaussian channel, the variational lower bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ is maximized by the kernel PCA solutions for encoder and feature decoder weights (see Appendix C.3.3 for derivations and discussions of non-centered representations). As expected, we obtain the linear PCA for the special case of linear mappings $\phi(\mathbf{x}) \equiv \mathbf{x}$.

A few further comments about the objective (5.37) are in order. First of all, we note that analytical properties of the *data decoder* depend on the specific parameterization of $q(\mathbf{x}|\mathbf{f})$. Since integration over \mathbf{x} and \mathbf{f} reduces to evaluations of the empirical averages, the average $\langle \log q(\mathbf{x}|\mathbf{f}) \rangle_{p(\mathbf{x}, \mathbf{f})}$ may typically be easily computed (provided that $q(\mathbf{x}|\mathbf{f})$ is appropriately kernelized, so that no explicit computations in the feature space are performed). Having specified the data decoder, we may optimize (5.37) to learn parameters of the kernel function \mathcal{K}_Θ , which parameterizes scalar products in $\mathbb{R}^{|\phi|}$. While this strategy may indeed lead to nontrivial optimization surfaces for kernel parameters, our current experience shows that the results may be strongly influenced by the specific definitions of the feature-to-data mappings⁵. Nevertheless, while any practical application of the framework may require a careful consideration of appropriate choices for $q(\mathbf{x}|\mathbf{f})$ and \mathcal{K}_Θ , the considered objective (5.2.3) indeed defines a theoretically rigorous tractable way of optimizing the bound on the mutual information in large scale nonlinear Gaussian channels.

5.3.1.3 Noiseless Channels and Gaussian Process Latent Variable Models

The discussed variational framework of maximizing lower bounds on $I(\mathbf{x}, \mathbf{y})$ in nonlinear Gaussian channels is particularly useful in situations when the channel $\mathbf{x} \rightarrow \mathbf{y}$ is intrinsically noisy. However, in many practical situations when the stochasticity of $p(\mathbf{y}|\mathbf{x})$ is not a necessary modeling requirement, simplifications of the IM approach may be possible. One practical example of such situation is an application of the variational information-maximizing framework to nonlinear dimensionality reduction, where the code space $\{\mathbf{y}\}$ is continuous, and the encoder is deterministic, i.e. $p(\mathbf{y}|\mathbf{x}) \sim \delta$. Interestingly, for a specific parameterization of the variational decoder distribution, this noiseless limit of the variational information-maximizing framework is closely related to the recently introduced (Lawrence (2003)) Gaussian Process Latent Variable Models (GPLVMs).

⁵Intuitively, it is clear that unless the feature vectors are constrained, the bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ may diverge due to the contributions of the feature decoder (see expression (5.37) and Appendix C.3.5). The situation is analogous to the case of linear Gaussian channels, where the divergent weights lead to the diminishing noise effects (see e.g. expression (C.8)), which may generally result in the divergence of $I(\mathbf{x}, \mathbf{y})$.

In order to show this, let us consider a conventional noiseless autoencoder $\mathbf{x} \mapsto \mathbf{y} \rightarrow \tilde{\mathbf{x}}$, with $p(\mathbf{y}|\mathbf{x}) \sim \delta(\mathbf{y} - \mathbf{y}(\mathbf{x}))$ and the empirical distribution $\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}}) \propto \sum_m \delta(\mathbf{x} - \mathbf{x}^{(m)})\delta(\mathbf{x} - \tilde{\mathbf{x}})$. For i.i.d. patterns the exact conditional likelihood training in such models reduces to maximizing

$$\mathcal{L}_{\tilde{\mathbf{x}}|\mathbf{x}} = \langle \log \langle q(\tilde{\mathbf{x}}|\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})} \rangle_{\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}})} = \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})}, \quad (5.39)$$

where $p(\mathbf{y}|\mathbf{x})$ and $q(\mathbf{x}|\mathbf{y})$ define the encoding and decoding mappings respectively (also see lemma 3.1 and proposition 3.4). If $\{\mathbf{x}\} \stackrel{\text{def}}{=} \{\mathbf{x}^{(m)} | m = 1, \dots, M\}$ defines the dataset of training patterns with the corresponding codes $\{\mathbf{y}\} \stackrel{\text{def}}{=} \{\mathbf{y}(\mathbf{x}^{(m)}) | m = 1, \dots, M\}$, the objective (5.39) may be transformed to

$$\mathcal{L}_{\tilde{\mathbf{x}}|\mathbf{x}} = \langle \log \hat{q}(\{\mathbf{x}\}|\{\mathbf{y}\}) \rangle_{\hat{p}(\{\mathbf{y}\}|\{\mathbf{x}\})}. \quad (5.40)$$

Here the decoder $\hat{q}(\{\mathbf{x}\}|\{\mathbf{y}\})$ and the encoder

$$\hat{p}(\{\mathbf{y}\}|\{\mathbf{x}\}) = \prod_{m=1}^M \delta(\mathbf{y}^{(m)} - \mathbf{y}(\mathbf{x}^{(m)})) \quad (5.41)$$

are defined for the whole set of the observed patterns $\{\mathbf{x}\}$ and their deterministic projections $\{\mathbf{y}\}$.

We can now see that if the decoding mapping is defined as

$$\hat{q}(\{\mathbf{x}\}|\{\mathbf{y}\}) \propto \prod_{i=1}^{|\mathbf{x}|} \exp \left\{ -\frac{1}{2} (\tilde{\mathbf{x}}^{(i)})^T \mathbf{K}^{-1} \tilde{\mathbf{x}}^{(i)} \right\} = \exp \left\{ -\frac{1}{2} \text{tr} \{ \mathbf{X} \mathbf{K}^{-1} \mathbf{X}^T \} \right\}, \quad (5.42)$$

the conditional likelihood (5.40) reduces to Neil Lawrence's Gaussian Process Latent Variable Models, optimizing the objective

$$\mathcal{L}_{\tilde{\mathbf{x}}|\mathbf{x}} = -\frac{1}{2} \text{tr} \{ \mathbf{X} \mathbf{K}^{-1} \mathbf{X}^T \} - \frac{|\mathbf{x}|}{2} \log |\mathbf{K}| + \text{const}. \quad (5.43)$$

In (5.42) we assumed that $\mathbf{K} \stackrel{\text{def}}{=} \{ \mathcal{K}_{\Theta}(\mathbf{y}^{(i)}, \mathbf{y}^{(j)}) \} \in \mathbb{R}^{M \times M}$ is the symmetric positive-semidefinite covariance defined by the covariance function \mathcal{K}_{Θ} , and $(\tilde{\mathbf{x}}^{(i)})^T \in \mathbb{R}^{1 \times M}$ is the i^{th} row of the data matrix $\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}\} \in \mathbb{R}^{|\mathbf{x}| \times M}$. In this formulation learning corresponds to optimizing the conditional likelihood (5.40) for the encoded representations $\{\mathbf{y}^{(i)} | i = 1, \dots, M\}$ and parameters of the covariance function \mathcal{K}_{Θ} .

On the other hand, it is easy to see that the objective (5.40) may be interpreted as the variational bound on $I(\{\mathbf{x}\}, \{\mathbf{y}\})$ for the **noiseless encoder** (5.41) and the variational decoder $\hat{q}(\{\mathbf{x}\}|\{\mathbf{y}\})$ given by expression (5.42). Effectively, this interpretation offers an alternative, information-theoretic justification of GPLVMs for noiseless encoder models. Note that in contrast to Lawrence (2003), our derivation of (5.40) is independent of the specific assumptions about the model – the GPLVM may be viewed as a special case of the variational information-maximizing framework for noiseless channels with the specific choice of the variational decoder.

5.4 Demonstrations

Here we demonstrate experimental results of maximizing the information content for the discussed discrete nonlinear channels. We will focus specifically on the discussion of information-theoretic clustering (see Section 5.2.1 and Section 5.2.2). First we consider clustering by maximizing the alternative objective function $\hat{I}(\mathbf{x}, \mathbf{y})$ for training Gaussian mixture models (see expression (5.1)). As we showed in proposition 3.1, this objective corresponds to a generic lower bound on $I(\mathbf{x}, \mathbf{y})$ for the encoder model $\mathcal{M}_I = \tilde{p}(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})$, where the encoder distribution is given by the exact posterior of a Gaussian mixture model $\mathcal{M}_L = p(\mathbf{y})p(\mathbf{x}|\mathbf{y})$, i.e.

$$\tilde{p}(\mathbf{y}|\mathbf{x}) \propto p(\mathbf{y})p(\mathbf{x}|\mathbf{y}), \quad \mathbf{y} \in \{y_1, \dots, y_{|y|}\}, \quad p(\mathbf{x}|y_j) \sim \mathcal{N}_{\mathbf{x}}(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j).$$

For this case we will compare maximization of the bound $\hat{I}(\mathbf{x}, \mathbf{y})$ with maximization of the exact likelihood for Gaussian mixtures, and show that the variational approach may indeed favor more uniform cluster assignments. As the method may be used for training both encoder and generative models (due to the specifics of the channel definition), we will refer to it as IM for Gaussian mixtures.

Then we will apply the information-theoretic clustering method to several intrinsically non-Gaussian datasets. For these tasks, we consider the specific definition of the encoder distribution (5.12) and its kernelized extension (5.24) described in Sections 5.2.2, 5.2.3. We compare both approaches to Gaussian mixture, k-means, and the kernelized k-means clustering, and show that for the considered datasets our method may indeed lead to visible anthropomorphic improvements over the common clustering techniques. Interestingly, we show that by *learning* parameters of kernel functions (see expression (5.31)), we may indeed obtain better visualizations of cluster assignments.

5.4.1 IM for Gaussian Mixture Models

5.4.1.1 \mathcal{L} - and $\hat{I}(\mathbf{x}, \mathbf{y})$ -maximization for Gaussian mixtures

Figure 5.1 shows typical soft cluster allocations produced by a mixture model trained by maximizing the exact likelihood (*left*) and the bound $\hat{I}(\mathbf{x}, \mathbf{y})$ on the mutual information (*middle*). The dataset consisted of $M = 150$ training patterns, generated at uniform random from the convex area bounded by a flat triangle (for $|\mathbf{x}| = 2$). The number of mixture components was $|y| = 3$. The initial mixture coefficients were set as $p_0(y_j) = 1/|y|$, the initial covariances were set to be spherical with $\boldsymbol{\Sigma}_j = \mathbf{I}_{|\mathbf{x}|}$, and the components' means were initialized by applying the k-means algorithm. Then the model was trained until convergence by applying the two methods (with the learning rate $\eta = 0.01$ for the IM). Not surprisingly, maximization of the likelihood for the mixture model often resulted in near-singular components' covariances, which led to locally constrained cluster segments (Figure 5.1 (*left*)). On the other hand, under the identical initializations, the IM on $\hat{I}(\mathbf{x}, \mathbf{y})$ typically led to more uniform cluster allocations with non-degenerate covariance matrices (Figure 5.1 (*middle*)).

Figure 5.1 (*right*) shows the means and the variances for the proportions of source patterns allocated to each of the $|y|$ clusters for 10 runs of the EM and

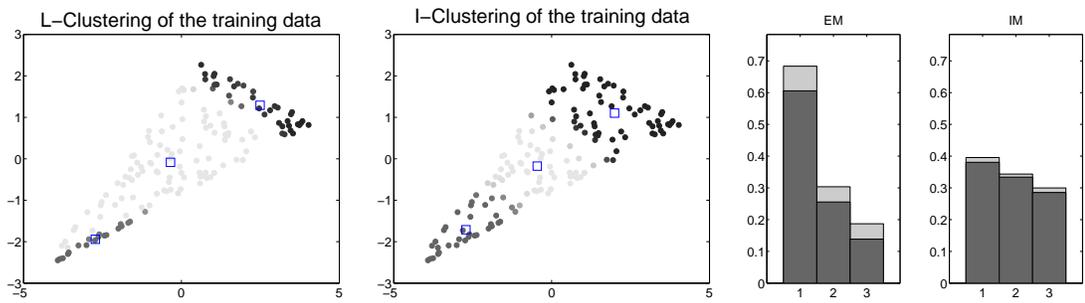


Figure 5.1: \mathcal{L} - and \hat{I} - maximization for a three-component Gaussian mixture ($|y| = 3$). *Left*: clustering with the EM algorithm on \mathcal{L} (the squares show the cluster centers); *Middle*: clustering with the IM on $\hat{I}(x, y)$; *Right*: Proportion of the testing points assigned to each of $|y| = 3$ clusters under $T = 10$ different runs of the optimization procedure (light gray bars indicate the variance). The exact mutual information computed by the EM- and IM-trained models was $I_L(x, y) \approx 0.75$ and $I_I(x, y) \approx 1.00$ respectively; the corresponding log-likelihoods are $\mathcal{L}_L \approx -4.74$ and $\mathcal{L}_I \approx -4.99$.

IM algorithms, where the training set was re-sampled from the underlying distribution at each run. For both EM and IM-trained⁶ models, the cluster allocation $y(x)$ for a pattern x was given by $y(x) = \arg \max_y p(y|x)$. Again, the results suggest that on average the models trained by the IM typically result in more uniform class assignments and smaller variances on cluster sizes. This empirically confirms the informal argument of section 3.2.1 that maximization of the bound $\hat{I}(x, y)$ may favor more spread-out representations in the latent space. We also note that the likelihoods \mathcal{L}_L computed for EM-trained models typically exceed the likelihoods \mathcal{L}_I of Gaussian mixture models trained by maximizing the bound $\hat{I}(x, y)$ (for the illustrated case we had $\mathcal{L}_L \approx -4.74$ and $\mathcal{L}_I \approx -4.99$). At the same time, the exact mutual information computed for a conventionally trained mixture was typically inferior to that of an IM-trained mixture ($I_L \approx 0.7467$ vs. $I_I \approx 0.9950$ for the considered case). Additionally, we see that cluster allocations produced by both training procedures suggest that the models with higher likelihoods do not necessarily lead to better representations of the underlying distribution, which once again illustrates the conceptual problems of maximizing the likelihoods of under-constrained models for learning informative representations of the data.

5.4.1.2 Samples from the Trained Models

As mentioned in Section 5.2.1, our motivation for maximizing the bound $\hat{I}(x, y)$ was to produce informative latent variable representations y of the sources $\{x\}$. Importantly, due to the specific parameterization of the encoder distribution, the suggested framework of maximizing $\hat{I}(x, y)$ could be used for training both encoder and generative models. Our hope was that by maximizing the bound,

⁶Throughout the discussion in this section, we consider a special instance of the IM algorithm for the specific definition of the encoder distribution (5.2).

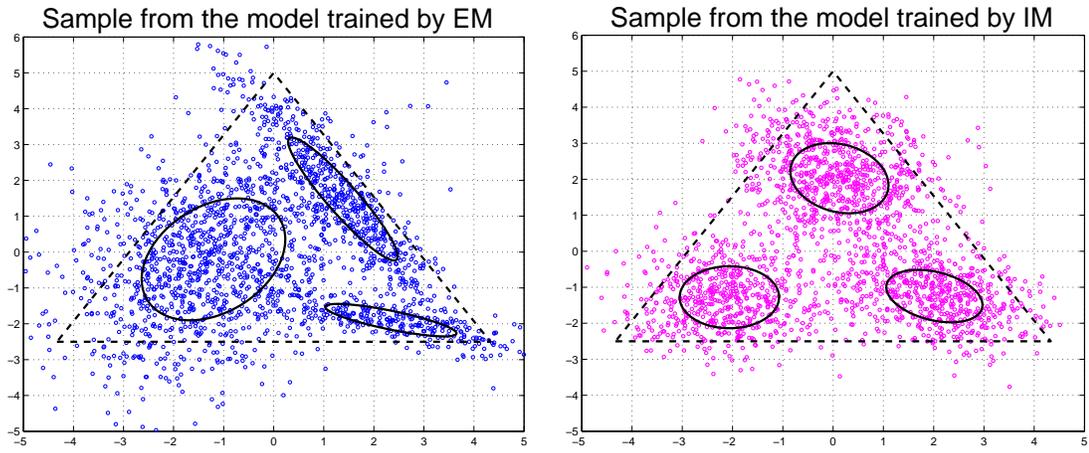


Figure 5.2: Samples from the Gaussian mixture model \mathcal{M}_L trained by the EM- and IM-algorithms with $|x| = 2$, $|y| = 3$. The number of training patterns was $M = 250$. *Left:* \mathcal{M}_L is trained by maximizing the likelihood \mathcal{L} (resulting in $I_L(x, y) \approx 0.654$); *Right:* \mathcal{M}_L is trained by maximizing the bound $\hat{I}(x, y)$ (resulting in $I_I(x, y) \approx 0.911$). The ellipses show probability contours of the mixture components at $1/2$ Mahalanobis distance from the means μ_j for $j = 1, \dots, |y|$. The dashed lines indicate the boundary of the convex region used to generate the underlying data.

the resulting generative models could potentially give rise to samples which were somewhat more representative of the underlying data distribution. Figure 5.2 illustrates typical samples produced by the three-component Gaussian mixture model trained by maximizing the likelihood (*left*) and the bound $\hat{I}(x, y)$ (*right*). As in the previous case, we assumed that $|x| = 2$, $|y| = 3$, and both models were trained by starting from identical initializations at the flat component weights, the initial means given by the k-means algorithm, and the initial covariances set as $\Sigma_j = I_{|x|}$. For illustration purposes, the training data was sampled from the equiangular triangle, which is more difficult to model by a Gaussian mixture distribution than the data shown on Figure 5.1 (*top left*) (for example, it is clear that in the limiting case when one of the angles approaches zero, the data could be trivially modeled by a single near-singular Gaussian).

We can see that in the considered case the Gaussian mixture model trained by maximizing the likelihood leads to the components fitting local segments of the dataset, with a wide flat component responsible for the remaining training patterns (see Figure 5.2 (*left*)). While the resulting samples have a vague resemblance to the underlying equiangular area, they are arguably less representative of the convex region than the samples produced by the IM-trained mixture (Figure 5.2 (*right*)). Note that the latter model is characterized by flatter spectra of the components' covariances (allowing more uniform space coverage within the constrained area). Moreover, for the IM-trained model the components' means lie close to the corners of the underlying shape, which in this case preserves the higher-level relations between the angles.

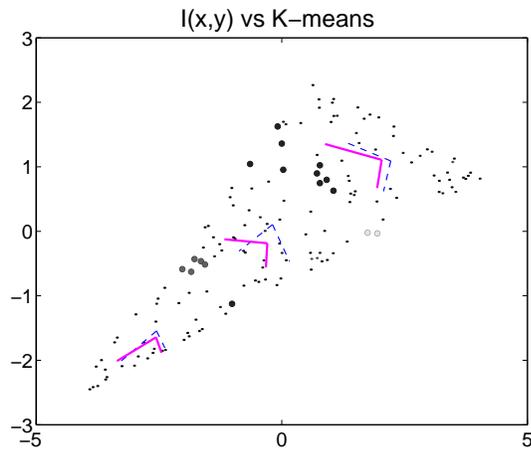


Figure 5.3: $\hat{I}(x, y)$ -maximization vs. K-means. *Solid lines*: eigenvectors of the covariances Σ_j obtained by the IM. *Dashed lines*: eigenvectors of the sample covariances Σ_j^{KM} associated with each of the clusters obtained by the deterministic k-means algorithm. All the eigenvectors are weighted by the corresponding eigenvalues and centered at the components' means (μ_j and μ_j^{KM} for the IM and k-means respectively). The IM was initialized at μ_j^{KM} and Σ_j^{KM} . The filled circles show the training patterns which are classified differently by the k-means and the IM started at the considered initialization.

5.4.1.3 $\hat{I}(x, y)$ -maximization for Gaussian Mixtures and the K-means Algorithm

In Section 5.2 we showed that in the limiting noiseless case of the decoder distributions $p(x|y)$, optimization of the specific bound $\hat{I}(x, y)$ for the considered models reduces to the k-means algorithm. However, in the non-limiting cases, maximization of the specific bound $\hat{I}(x, y)$ for Gaussian mixtures is different from the k-means, which may be seen analytically from (5.11) or illustrated by the deviations of the IM from the k-means initializations. Figure 5.3 compares the means and covariances μ_j, Σ_j of the IM-trained mixture model with the cluster centers and the associated sample covariances $\mu_j^{KM}, \Sigma_j^{KM}$ obtained by the k-means algorithm. In the illustrated case, the IM-trained mixture was initialized at μ_j^{KM} and Σ_j^{KM} , with the initial mixing coefficients set according to the cluster sizes. Figure 5.3 indicates the changes in the parameters and the resulting cluster allocations. We see that for nonzero covariances, the IM diverges from the initial settings, though for the considered dataset the cluster assignments produced by the k-means and the IM are not very different.

Generally, the results of our experiments in this section have confirmed the intuition that optimization of the bound $\hat{I}(x, y)$ may result in more uniform cluster assignments than the likelihood-based training of generative models. We also showed that samples from the mixture models trained by maximizing the bound may indeed be qualitatively more representable of the underlying distributions. Extensions of this form of the IM to training other kinds of generative models may potentially be considered, provided that the encoder may be computed exactly.

As demonstrated in this section, the resulting models may be used for clustering and data generation.

5.4.2 Kernelized Information Theoretic Clustering

5.4.2.1 Gaussian Mixture, K-means, and Kernelized Information-Theoretic Clustering for Spiral Data

Here we consider a direct application of the encoder models (5.12) and their kernelized extensions (5.24) to clustering of intrinsically non-Gaussian data. In contrast to the bound $\hat{I}(\mathbf{x}, y)$ optimized in Section 5.4.1, the encoder models considered here were trained by maximizing the exact mutual information $I(\mathbf{x}, y)$ as discussed in Section 5.2.2. For the kernelized encoders (5.22), we also learned optimal kernel parameters (focusing specifically on the RBF kernel, see expression (5.30)). We compare both information-theoretic approaches to the k-means and the Gaussian mixture clustering (assuming the standard EM training).

As discussed earlier in Section 5.4.1, clustering by maximizing the mutual information corresponds to learning of an optimal encoder, where the code space $\{y\}$ defines the generally unknown cluster labels. In the first set of experiments, we considered the discrete code space with $|y| = 3$ states of the output variables. (Effectively, the size of the code space $|y|$ gives an upper bound on the number of clusters; in what follows we assume that this number is fixed). The unlabeled training data was generated as $x_1(t) = t \cos(t)/4$, $x_2(t) = t \sin(t)/4$, with t changing uniformly in $[0, 3.4\pi]$ (here x_1 and x_2 correspond to different coordinates of the sources $\mathbf{x} \in \mathbb{R}^{|\mathbf{x}|}$, $|\mathbf{x}| = 2$). The total number of patterns was $M = 70$.

Figure 5.4 illustrates typical allocations of cluster labels to the considered dataset. The data was clustered by a Gaussian mixture model trained by maximizing the likelihood (Figure 5.4 *left*), the k-means algorithm (Figure 5.4 *medium*), and the kernelized encoder model (5.24) trained by maximizing $I(\mathbf{x}, y)$ (Figure 5.4 *right*). For the considered RBF kernel (5.31), we learned the inverse variance parameter β . Color intensity of each pattern shown on Figure 5.4 is given as an average of the cluster colors weighted by the probability of cluster allocations $p(y_j|\mathbf{x})$. The “pure” cluster colors corresponding to deterministic cluster assignments are shown by light-, medium-, and dark-gray squares.

As we see from Figure 5.4 (*left*), Gaussian mixture models trained by maximizing the likelihood resulted in largely stochastic cluster allocations; from the corresponding plot, it is qualitatively difficult to distinguish between patterns which belong to “medium-gray” and “dark-gray” clusters. Moreover, for the illustrated case we see a large disproportion in cluster sizes. Specifically, we see that 3 or 4 patterns (out of 70) form a unique cluster (shown in light-gray on Figure 5.4 (*left*)). This observation is further confirmed by Figure 5.5 (*left*), which shows responsibilities of the three mixture components for all 70 data patterns. Figures 5.4 (*middle*) and 5.5 (*middle*) show the corresponding cluster assignments for the k-means algorithm. We see that while the similarly clustered points indeed lie close to each other in \mathbb{R}^2 according to the L_2 -norm, the allocated clusters are *not* locally smooth in t . (We also note that maximization of the bound (5.1) on $I(\mathbf{x}, y)$ for the encoding distribution corresponding to the exact posterior of a

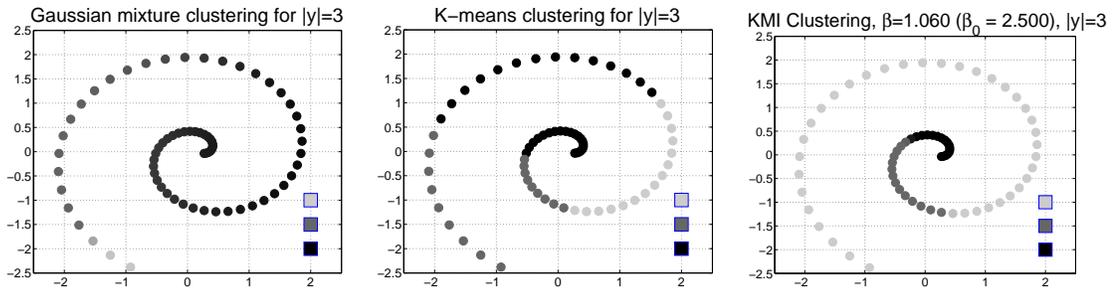


Figure 5.4: Clustering for $|y| = 3$. *Left*: Gaussian mixtures; *Middle*: K-means; *Right*: information-maximization for the (RBF-)kernelized encoder. Light, medium, and dark-gray squares show the cluster colors corresponding to deterministic cluster allocations. The color intensity of each training point $x^{(m)}$ is the average of the pure cluster intensities, weighted by the responsibilities $p(y_j|x^{(m)})$. Nearly indistinguishable dark colors for the majority of patterns under the Gaussian mixture clustering indicate soft cluster assignments (see also Figure 5.5). Note that by applying the kernelized IM algorithm, we obtain nearly deterministic cluster assignments to locally smooth data regions.

Gaussian mixture model led to the results similar to a soft form of the k-means clustering. This agrees with previous observations that the cluster allocations produced by a model maximizing (5.1) are qualitatively not very different from the k-means (see Section 5.4.1)).

Figures 5.4 and 5.5 (*right*) illustrate cluster allocations produced by maximizing the mutual information for the kernelized encoder (5.22). As mentioned above, in the considered set of experiments we have applied the RBF kernels (5.31) with adaptable parameters β . The kernel parameters β were initialized at $\beta_0 = 2.5$. The initial settings of the coefficients $A \in \mathbb{R}^{M \times |y|}$ in the feature space were samples from $\mathcal{N}_{A_{ij}}(0, 0.1)$. The log-variances $\tilde{s}_1, \dots, \tilde{s}_{|y|}$ were initialized at zeros. The encoder parameters A and $\{\tilde{s}_j | j = 1, \dots, |y|\}$ (along with the RBF kernel parameter β) were optimized by applying the scaled conjugate gradients procedure (see e.g. Bishop (1995)). The training stopped at the 46th iterations of the algorithm, after the changes in the consecutive evaluations of the objective $I(x, y)$ remained lower than 10^{-15} for three consecutive iterations. We see that in the considered case we indeed obtain non-degenerate well-separated clusters which are locally smooth in t (see Figure 5.4 (*right*)). The local smoothness of cluster allocations is further confirmed by Figure 5.5 (*right*) (note a slight decrease in coding certainty for the patterns lying close to the cluster boundaries). The results are shown for $\beta \approx 0.825$ obtained by the non-linear ascent on $I(x, y)$.

Note that in the experiments described here we have made no strong assumptions about the choice of the kernel function which could be particularly suitable for clustering this specific dataset. Undoubtedly, a careful choice of the kernel could potentially lead to a better visualization of the locally smooth, non-degenerate structure. Instead, we considered the commonly used RBF kernel with adaptable kernel parameters. The results suggest that maximization of mu-

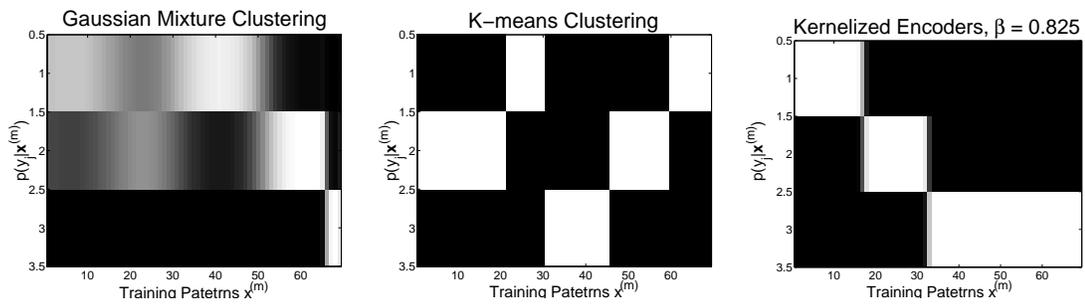


Figure 5.5: Probabilities $p(y_j|x^{(m)})$ for the spiral data, $x \in \mathbb{R}^2$, $y \in \{y_1, y_2, y_3\}$. The total number of patterns $M = 70$. *Left*: Gaussian mixtures; *Middle*: K-means; *Right*: information-maximization for the kernelized encoder model (RBF kernel with the learned parameter $\beta = 0.825$).

tual information in discrete nonlinear channels may indeed be a useful technique to consider in unsupervised clustering applications. We will now consider another application of the method, and show that by learning kernel parameters we may indeed obtain better visualizations than by using fixed kernel functions or applying more common clustering techniques.

5.4.2.2 Kernelized Information-Theoretic Clustering for Spatially Translated Letters

Figures 5.6 and 5.7 show an application of the the mutual information maximization to clustering of spatially translated letters, for the code sizes $|y| = 2$ and $|y| = 3$ respectively. As in the previous case shown on Figure 5.4, patterns allocated to different clusters are shown by different color intensities. The data consisted of $M = 210$ patterns, sampled from models of Latin letters⁷ (70 per letter). The training set was constructed in such a way that the distance between the neighboring letters was roughly constant (by this we mean that the distances between the sample means for the neighboring letter models were fixed to be the same). In the considered case, we assumed that the neighboring letters were relatively close to each other (compared to the letter sizes), so that points sampled from different letters could in fact be geometrically closer to each other than points sampled from a single letter model. As in the previous set of experiments, we have compared information-theoretic clustering in kernelized encoder models with the k-means and Gaussian mixture clustering. We have also explored the effects of introducing projections into the feature space and applied the IM procedure for both encoder models described in Section 5.2.2 and 5.2.3 (which we will refer to as *simple*, or *non-kernelized* encoders and *kernelized* encoders respectively).

By analogy with the previously discussed experiments (see Figure 5.4), we can see that for the considered case the IM framework indeed results in anthro-

⁷The generative model for each letter is given by a mixture of constrained uniform line segments with additive spherical Gaussian noise.

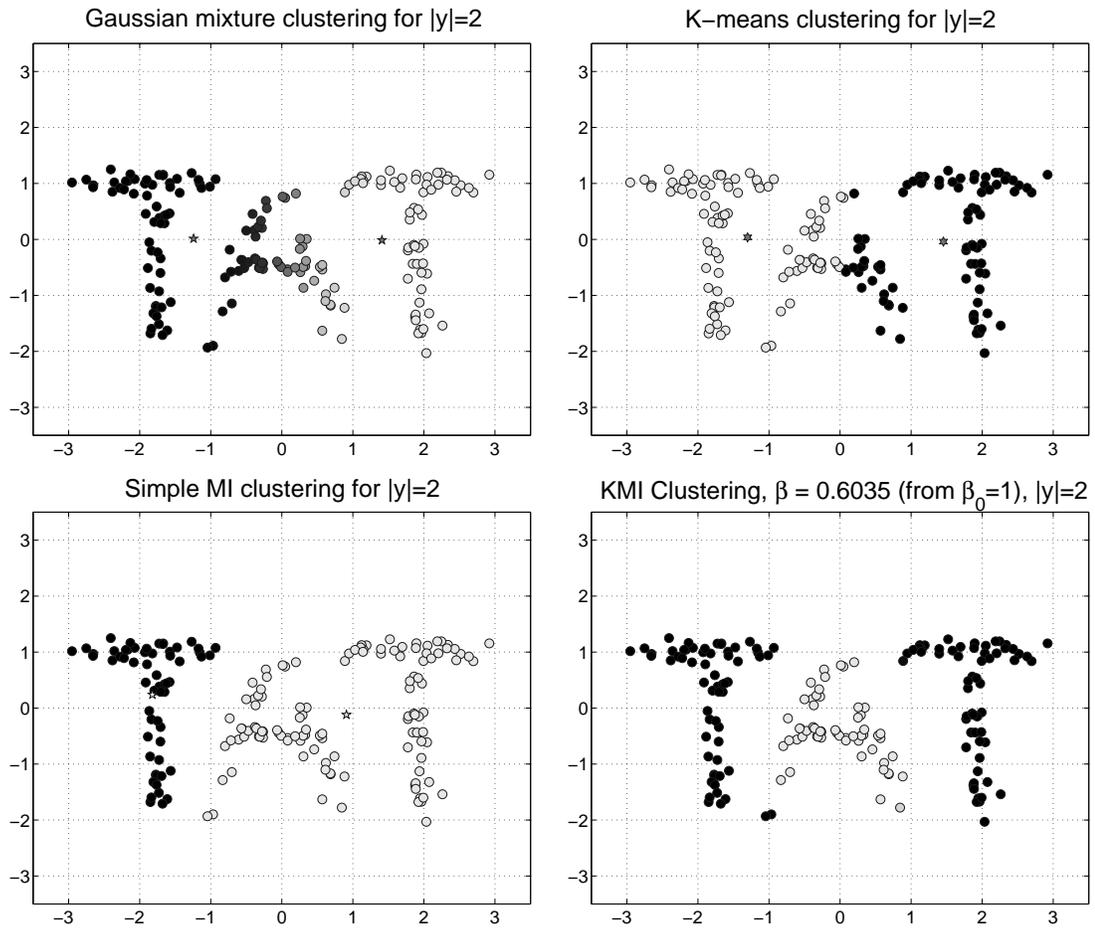


Figure 5.6: Learning cluster allocations for $|y| = 2$. Where appropriate, the stars show the cluster centers. *Top left:* clustering with the two-component Gaussian mixture trained by the EM algorithm on \mathcal{L} ; *Top right:* clustering with the k-means; *Bottom left:* clustering with the encoder model $p(y_j|x) \propto \exp\{-\|x - w_j\|^2/s_j\}$ trained by maximizing mutual information $I(x, y)$; *Bottom right:* clustering with the kernelized encoder model $p(y_j|x) \propto \exp\{-\|\phi(x) - w_j\|^2/s_j\}$ trained by maximizing $I(x, y)$ (the results are shown for the RBF kernel). For the kernelized model, the inverse variance β of the RBF kernel varied from $\beta_0 = 1$ (at the initialization) to $\beta \approx 0.604$ after convergence.

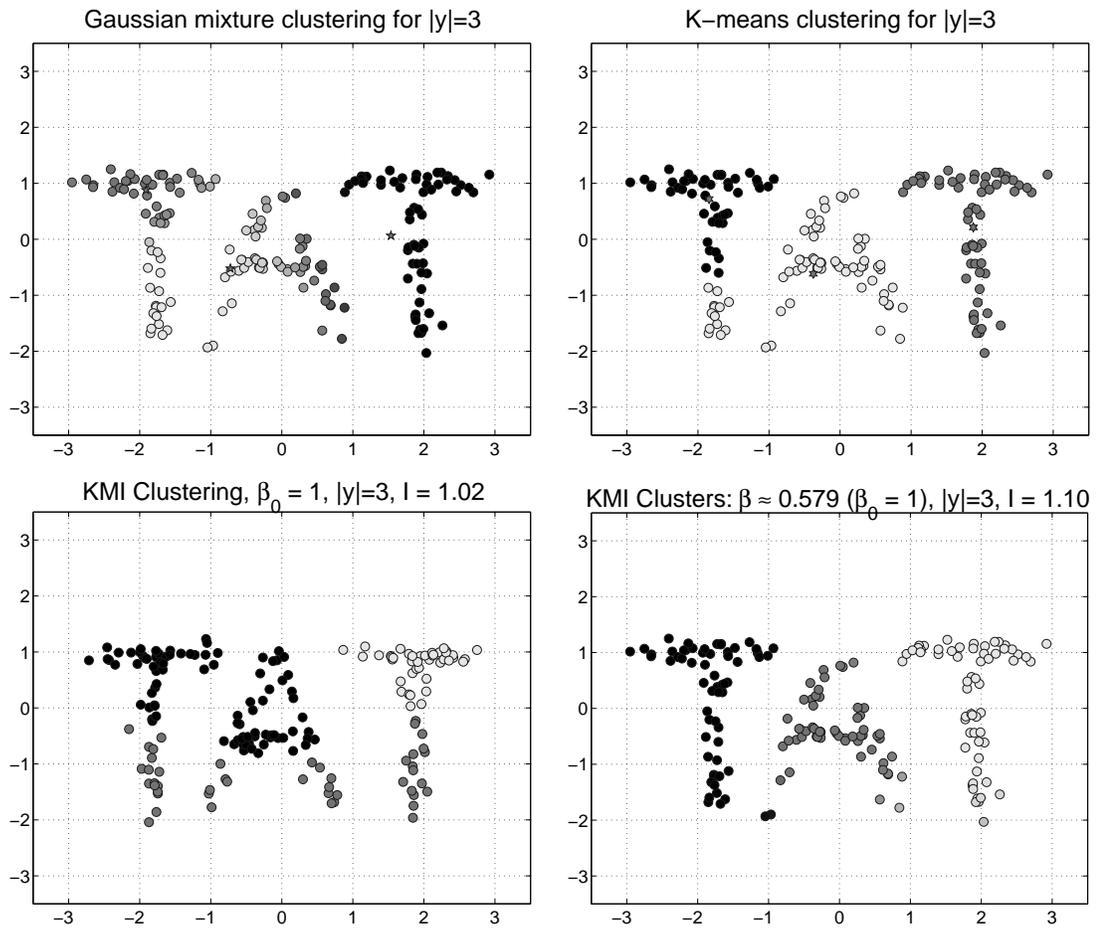


Figure 5.7: Learning cluster allocations for $|y| = 3$. Where appropriate, the stars show the cluster centers. *Top left*: clustering with the three-component Gaussian mixture trained by the EM algorithm on \mathcal{L} ; *Top right*: clustering with the k-means algorithm; *Bottom left*: clustering in the kernelized encoder model (the results are shown for the RBF kernel with the inverse variance fixed at $\beta_0 = 1$); *Bottom right*: clustering with the kernelized encoder model with the adaptable kernel parameter (the inverse variance of the RBF kernel varied from $\beta_0 = 1$ (at the initialization) to $\beta \approx 0.579$ after convergence). Note that by learning the kernel parameter β we obtain higher values of the mutual information ($I \approx 1.10$ vs $I \approx 1.02$) and more intuitive cluster assignments.

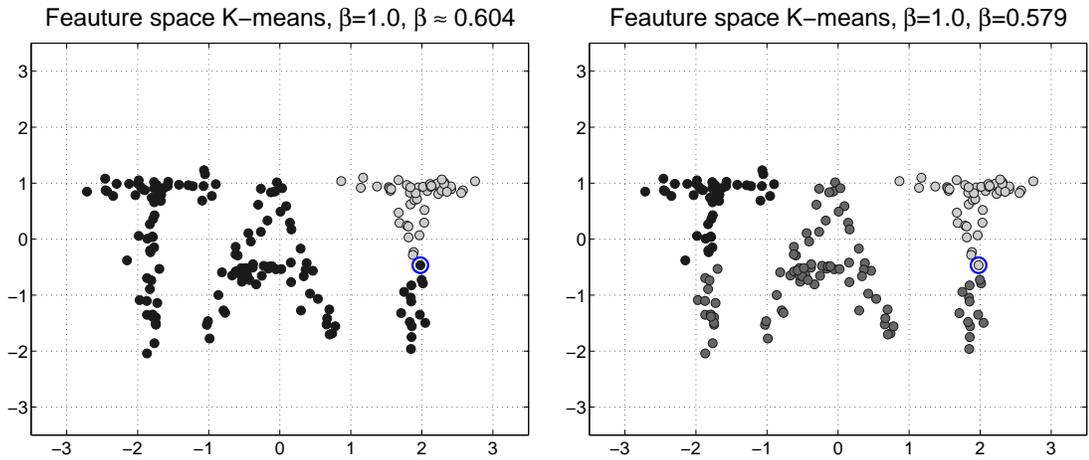


Figure 5.8: Learning cluster allocations by the feature-space k-means. The results are shown for the fixed RBF kernels with the inverse variance parameter β . *Left*: $|y| = 2$, $\beta_1 = 1$ and $\beta_2 = 0.604$; *Right*: $|y| = 3$, $\beta_1 = 1$ and $\beta_2 = 0.579$. The number of training patterns was $M = 210$. The patterns shown by the double circles \odot are the only ones clustered differently for the two settings of the kernel parameters (β_1 and β_2), assuming identical initializations. The parameters β_2 were set to the converged values of β for the corresponding encoder models (cf Figure 5.6 (*bottom left*), Figure 5.7 (*bottom left*)).

pomorphically sensible cluster allocations. Indeed, in contrast to the k-means or Gaussian mixture clustering, the information-theoretic approach appears to favour smoother local representations, which results in arguably more intuitive cluster assignments. Figure 5.6 shows typical cluster assignments produced by the k-means algorithm (*top right*) and Gaussian mixtures trained by the EM (*top left*) for $|y| = 2$. Again, two different clusters are illustrated by different color intensities (light- and dark-gray), and soft cluster assignments are shown by the intermediate intensities. As usual, the EM for the mixture model was started at the k-means initialization, with the initial covariances set to the sample covariances of the associated k-means clusters, and the mixture coefficients proportional to the cluster sizes. As we see from the plots, both methods result in roughly identical cluster allocations, and for this case divide the dataset into two roughly symmetric parts, consisting of the letter *T* and half of the letter *A*. (Not surprisingly, the cluster allocations produced by the Gaussian mixture model are soft – other than that, the resulting clusters are not too different from those given by the k-means). Generally, other initializations of the mixture model resulted in more degenerate cluster assignments.

On the other hand, we can see that for the considered case the clusters produced by the encoder models are arguably more intuitive. Figure 5.6 (*bottom left*) shows two typical clusters generated by the non-kernelized encoder model trained by the IM for $|y| = 2$. The initial weights of the simple encoder model (5.12) were set according to $W_{ij}^{(0)} \sim \mathcal{N}_{W_{ij}}(0, 0.1)$, $W^{(0)} = \{W_{ij}^{(0)}\} \in \mathbb{R}^{|x| \times |y|}$, with the log-variances $\tilde{s}_1^{(0)}, \tilde{s}_2^{(0)} \sim \mathcal{N}(0, 0.1)$. We used a similar initialization for the

kernelized encoder (5.22), where instead of the matrix of the weights \mathbf{W} we have learned the matrix of the feature space coefficients $\mathbf{A}^{(0)} = \{A_{ij}^{(0)}\} \in \mathbb{R}^{M \times |y|}$. For the kernelized model, we considered the RBF kernel with the initial settings of $\beta_0 = 1$. Both the non-kernelized and kernelized encoder models were trained by performing a numerical ascent on $I(\mathbf{x}, y)$ (by the scaled conjugate gradients). As we see from the plot, the simpler model extracts a single letter T as a separate cluster, and allocates the remaining training patterns to the other cluster (Figure 5.6 (*bottom left*)). The clusters produced by the trained kernelized encoder model are typically different (see Figure 5.6 (*bottom right*)). As we see from the plot, the kernelized encoder model clusters both T 's similarly, separating them from the model of A , which may arguably be an intuitive coding scheme. Note that both T s are clustered together despite being located geometrically far from each other in the data space. The results of (Figure 5.6 (*bottom right*)) are shown for $\beta \approx 0.6035$ obtained after convergence of the SCG learning procedure. Importantly, the obtained results appear to be stable for different samples from the underlying distribution and different initializations of the IM learning procedure (provided that the RBF kernel parameter β is not too large or too low, i.e. $\mathbf{K} \in \mathbb{R}^{M \times M}$ has a non-degenerate spectrum).

Figure 5.7 shows typical cluster allocations produced by the three methods for the codesize $|y| = 3$. Again, Gaussian mixtures and the k-means (Figure 5.7 (*top left*), Figure 5.7 (*top right*)) result in an arguably inferior performance to the kernelized IM clustering (Figure 5.7 (*bottom right*)). Moreover, we can see that by learning the RBF kernel parameter β , we may obtain better visualizations than for the case when β is fixed (Figure 5.7 (*bottom left*)). Also, by learning kernel parameters we typically get higher values of the mutual information $I(\mathbf{x}, y)$ ($I \approx 1.10$ vs $I \approx 1.02$ when β is fixed for the illustrated case). The results were confirmed repeatedly for different samples from the underlying distribution and different initializations of the learning algorithms. Generally, the kernelized encoder with the adaptable RBF kernel parameters was the only method (out of the considered ones) which resulted in an almost perfect separation of the letter models.

Finally, we have compared the information-theoretic clustering method with the kernel k-means algorithm (see e.g. Zhang and Rudnicky (2002), Dhillon et al. (2004), Wang et al. (2004)), where for a *fixed kernel function* each pattern is deterministically assigned to a cluster depending on the distance from the cluster mean in the feature space. Figure 5.8 illustrates a typical application of the kernelized k-means to clustering of the considered dataset for $|y| = 2$ (*left*), and $|y| = 3$ (*right*), started at the same initializations as the previous experiments shown on Figure 5.6 and Figure 5.7. Again, for comparison with the kernelized IM clustering, we show the results for the RBF kernels with the inverse variance parameter β . For both choices of the code space sizes $|y|$, the kernel parameter was set to the initial and the converged value of β obtained during the *kernelized encoder* clustering (see Figure 5.6 (*bottom right*), 5.7 (*bottom right*)); for example, for $|y| = 2$ we have considered $\beta = 1$ and $\beta = 0.6035$. We see that while the settings of β could influence the resulting cluster allocations, the clusters were not very different (assuming identical initializations for both choices of β). As

mentioned in Section 5.2.3, the principle advantage of our method is that it may indeed be viewed as an optimization procedure (namely, mutual information maximization) for learning kernel parameters. This contrasts with the majority of other approaches to unsupervised nonlinear clustering in feature spaces, which typically consider fixed similarity measures. As we have shown, learning kernel parameters may indeed lead to more intuitive visualizations of the underlying structure (see also Agakov and Barber (2005b) for an empirical comparison with the spectral clustering method of Ng et al. (2001)).

5.5 Summary

In this chapter we described several possible extensions of the IM framework to nonlinear encoder models. First, we considered optimizing the specific lower bound on the mutual information, where the encoder distribution of the channel was given by the exact posterior of the corresponding generative model. We explored this case for Gaussian mixtures, and compared the EM and the IM learning for a dataset which could not be easily modeled by a Gaussian mixture distribution. For this case we empirically demonstrated that both methods give rise to generally different optimization surfaces, and optimization of the bound $\hat{I}(\mathbf{x}, y)$ may favour more uniform representations in the code space, which may in some cases lead to an arguably more informative representation of the underlying distribution. Additionally, we pointed out that in the limiting special case of asymptotically noiseless spherical Gaussian decoders, the IM algorithm optimizing the considered bound $\hat{I}(\mathbf{x}, y)$ reduces to the k-means algorithm.

Then we focused on the problem of information-theoretic clustering in encoder models, where the generally stochastic nonlinear mapping from the sources to the discrete codes was learned by maximizing the exact mutual information $I(\mathbf{x}, y)$. For this case, we described a simple and practical algorithm applicable to unsupervised clustering, and explored its extensions by considering kernelized encoder models. Empirically, we demonstrated that the resulting information-theoretic clustering approach favorably compares with the common clustering techniques, and the option of learning kernel functions may indeed be of a practical benefit for visualizing the underlying structure of the data.

Finally, we reviewed some of the theoretical properties of the IM for higher-dimensional code spaces, and showed that some of the popular dimensionality reduction techniques may be interpreted as special instances of the variational information maximizing procedure. Specifically, we extended the work of Bourlard and Kamp (1988) and Bourlard (2000) to *arbitrary kernelizable* feature mappings applied in the context of stochastic encoding, and showed the nothing could be gained by using nonlinear encoders and linear variational decoders in the context of variational information maximization in Gaussian channels. To handle the intrinsic constraints of linear Gaussian variational decoders applied in the context of nonlinear encodings, we suggested a proper variational relaxation of the bound on $I(\mathbf{x}, y)$ which could be optimized for the encoder, *data decoder*, and *feature decoder* distributions under the assumption of a nonlinear variational decoding distribution. For *nonlinear Gaussian* encoding distributions, this led to

kernel PCA (Schoelkopf et al. (1998)) as the optimal solution for encoder weights. Additionally, we outlined a simple relation of the variational framework to the recent work on Gaussian Process Latent Variable Models (Lawrence (2003)), which may be interpreted as the variational information-maximization procedure in the *noiseless* limit of a nonlinear channel.

As an extension of the work described in this chapter, we note that a further study of the bound (5.37) on mutual information for nonlinear Gaussian encoders may need to be considered. As discussed in Section 5.3 and Appendix C.3.5, the bound indeed provides a proper information-theoretic objective for learning optimal kernel parameters. However, our current experience suggests that the results may be strongly influenced by the choice of constraints on the Gram matrix $\mathbf{K} \in \mathbb{R}^{M \times M}$ and specific definitions of the feature-to-data mappings. By choosing appropriate constraints on the kernel functions and the variational decoder distributions, one may hope to adapt the nonlinear Gaussian IM framework to visualizing high-dimensional data⁸. Another application field to explore is communication of discrete-valued data over channels with Gaussian noise, where a specific practical application may include code division multiple access in cellular telephony (see e.g. Viterbi (1995)).

Additionally, a comparison of the IM clustering with other families of clustering techniques will need to be considered. Specifically, it is interesting to see how information-maximization in nonlinear channels with the considered definition of the encoding distribution (5.22) could be related to the common spectral clustering approaches (see e.g. Perona and Freeman (1998), Kannan et al. (2000), Ng et al. (2001), Yu and Shi (2003)). Most of such approaches use eigenvalues of the fixed similarity matrices to transform the original data set, and apply any of the known clustering techniques on the transformations (see Weiss (1999) for a unifying discussion). It is intuitive that it may indeed be possible to interpret some of the spectral clustering methods from the information-theoretic viewpoint. Indeed, it is well known that the standard k-means algorithm (Hartigan and Wong (1979)) may be viewed as the trace maximization problem of the Gram matrix of the original data patterns (Zha et al. (2001)). Recently this result has been generalized to show that popular multi-class spectral clustering techniques using *normalized cuts* (e.g. Shi and Malik (2000), Ng et al. (2001)) may in fact be interpreted as a form of the weighted k-means algorithm (Bach and Jordan (2003), Dhillon et al. (2004)). The results of Section 5.2.1 suggest that it may potentially be possible to relate the spectral clustering methods to a specific form of the variational information-maximizing procedure, though a direct relation between the methods currently remains unclear.

Fundamentally, one of the principal advantages of our approach to nonlinear dimensionality reduction is a simple way to learn parameters of the kernel function in the unsupervised context. As we showed, the resulting procedure is numerically and conceptually simple; specifically, in the considered cases it did not require

⁸Note that if the projection noise may be reduced to zero, one may consider optimizing simpler objective functions, such as (5.40) – cf Gaussian Process Latent Variable Models. It is therefore believed that a further study of nonlinear Gaussian channels for dimensionality reduction may be of a practical interest mainly in situations when the encoding noise is unavoidable.

computations of the matrix inversions or eigenvalue decompositions of the Gram matrices. Also, once the channel distributions (and the corresponding kernels) are parameterized, the method does not require complex problem-specific algorithmic heuristics. Furthermore, one may expect that the information-theoretic view of clustering could potentially offer common advantages over the algorithmic approaches; specifically, it may be possible to extend the method to richer encoding distributions.

Chapter 6

Variational Information Maximization for Learning High-Dimensional Discrete Representations

In Chapter 5 we described applications of the information-maximizing framework to learning optimal parameters of nonlinear encoder models. One specific application of the framework which we considered was information-theoretic clustering. Due to the low cardinality of the codes (cluster labels), optimization of mutual information between the discrete codes y and the continuous sources x could in that case be performed exactly, and a number of nonlinear extensions of simple channels could easily be considered (see Section 5.2).

While being practically useful for visualizing unlabeled data, information-theoretic clustering may be viewed as a limiting case of discrete encoding. In many domains much richer encoding schemes must be considered. For example, in the neurophysiological domain it may often be a matter of interest to estimate how much information about the source stimuli may be contained in a (generally unknown) *population* of neural spikes. Additionally, the problem domain may impose constraints on the channel distributions, which may generally lead to the need of learning stochastic high-dimensional encoding mappings. As discussed in Section 1.4, maximization of the exact mutual information in this case may potentially be problematic, as it will generally require an explicit integration over the high dimensional encodings.

The primary goal of this chapter is to explore applicability of the variational information-maximizing framework in the context of learning high-dimensional binary representations of continuous source patterns. While the obtained results are general, we believe that the application area where they may be particularly useful is stochastic neural coding. We show that the simplest instance of our variational information-maximizing formulation provides a convenient framework for population coding in stochastic point neuron models. By analogy with Linsker (1997), we demonstrate that it is possible to derive an information-maximizing procedure which only requires local computations; however, our results are signifi-

cantly more general, as they are applicable in the case of stochastic, non-invertible encoding mappings. Moreover, we demonstrate that for the considered choice of a conditionally factorized encoder model, our variational method favorably compares with two approximate techniques maximizing Fisher Information-based *approximations* of lower bounds on the generally intractable $I(\mathbf{x}, \mathbf{y})$.

6.1 Introduction

The problem of encoding real-valued stimuli $\mathbf{x} \in \mathbb{R}^{|\mathbf{x}|}$ by a population of binary spikes $\mathbf{y} \in \{-1, 1\}^{|\mathbf{y}|}$ may be addressed in many different ways. Essentially, the goal is to adapt the parameters of any mapping $p(\mathbf{y}|\mathbf{x})$ to make a desirable population code for a given set of input stimuli $\{\mathbf{x}\}$. There are many possible ways to address this problem. One could be that *any* reconstruction based on the population should be accurate. This is typically handled by appealing to the Fisher Information (e.g. Cramer (1946), Cover and Thomas (1991)) which, with care, can be used in order to bound mean square reconstruction error (see e.g. Johnson (2003) for an introductory discussion).

Here we consider maximizing the amount of information which the spiking patterns contain about the stimuli. In this framework, a population coding formulation may be viewed as a special case of the information-maximizing problem, where continuous source signals are mapped into a discrete high-dimensional space. While much of the earlier work focuses on simple channels (e.g. Linsker (1988), Pouget et al. (1998), Zhang and Sejnowski (1999), Bethge et al. (2002)) or noiseless invertible mappings (Nadal and Parga (1994), Bell and Sejnowski (1995), Linsker (1997)), it is particularly interesting to address the problem of high-dimensional stochastic coding. Many previous attempts to apply mutual information¹ to population coding have been made (e.g. Brunel and Nadal (1998), Stocks and Mannella (2001), Kang and Sompolinsky (2001), Samengo and Treves (2001)). However, for large population sizes and non-invertible mappings to the code space, maximization of the mutual information is generally a computationally intractable task (see Section 1.4). Most current studies address the problem of computational intractability of maximizing $I(\mathbf{x}, \mathbf{y})$ by considering alternative objective criteria, e.g. those based on numerical approximations of $I(\mathbf{x}, \mathbf{y})$ derived under specific asymptotic assumptions (see e.g. Brunel and Nadal (1998), Kang and Sompolinsky (2001), Hoch et al. (2003)). While these approximations may result in asymptotic efficiency for many encoding channels, their applicability for specific models may be strongly affected by the form of encoder distributions, which may lead to instability of the resulting optimization procedure for specific encoding schemes (Agakov and Barber (2004b)).

In Chapter 2 we described a simple and general variational procedure optimizing a proper bound on mutual information $I(\mathbf{x}, \mathbf{y})$ between the sources \mathbf{x} and the codes \mathbf{y} . In contrast to most of the existing techniques, the procedure maximizes

¹Many other methods suggesting to optimize alternative objectives, e.g. *redundancy* (Barlow (1989), Atick (1992), Redlich (1993), Field (1994)), may be shown to correspond to specific relaxations (Nadal et al. (1998)) of the exact mutual information, which also involves generally intractable computations.

a proper bound on $I(\mathbf{x}, \mathbf{y})$, rather than an asymptotic approximation of a bound (Brunel and Nadal (1998)) or approximation of the exact mutual information for a local region of the source space (Szummer and Jaakkola (2002), Corduneanu and Jaakkola (2003)). The primary goal of this chapter is to apply the variational IM method to learning *high-dimensional stochastic* binary representations of continuous source patterns. As we believe that this procedure offers a convenient framework for addressing sub-goals of predictive population coding, we focus specifically on a biologically inspired channel parameterization.

6.2 Variational Learning of Population Codes

The principled information theoretic approach to learning neural codes involves maximization of the mutual information with respect to parameters of the encoder $p(\mathbf{y}|\mathbf{x})$. In what follows we assume a conditionally factorized decoder, i.e. $p(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^{|\mathbf{y}|} p(y_i|\mathbf{x})$, which is arguably the simplest case of a feed-forward mapping (each unit y_i is defined by a simple *point-neuron* model). This assumption also facilitates comparisons with other approximate information-maximizing techniques (such as the common approximations of Brunel and Nadal (1998), or reformulations of the local criteria derived in a different context by Corduneanu and Jaakkola (2003)).

Since for large-scale stochastic systems exact evaluation of $I(\mathbf{x}, \mathbf{y})$ is in general computationally intractable, we consider optimizing the generic lower bound $\tilde{I}(\mathbf{x}, \mathbf{y}) = \text{const} + \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})}$, where $\tilde{p}(\mathbf{x})$ is the empirical distribution of the high-dimensional continuous stimuli $\mathbf{x} \in \mathbb{R}^{|\mathbf{x}|}$. To learn optimal stochastic representations of the continuous training patterns $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}$ according to the lower bound $\tilde{I}(\mathbf{x}, \mathbf{y})$, we need to choose a continuous density function for the decoder $q(\mathbf{x}|\mathbf{y})$. Computationally, it is convenient to assume that the decoder is given by the isotropic Gaussian $q(\mathbf{x}|\mathbf{y}) \sim \mathcal{N}(\mathbf{U}\mathbf{y}, s^2\mathbf{I})$, where $\mathbf{U} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{y}|}$. Other (e.g. non-linear) variational decoders may potentially be considered and relaxations analogous to the ones described in Section 5.3.1 may potentially be used. However, we show that in situations when the variational decoder is a constrained linear Gaussian, the resulting iterative optimization procedure is particularly convenient, which results in a local (and arguably more biologically plausible) learning scheme.

We note that the proposed method describes a theoretically rigorous framework for maximizing information content which the high-dimensional spikes contain about (generally high-dimensional) continuous stimuli. The encoder distribution $p(\mathbf{y}|\mathbf{x})$ may in this case be interpreted as a stochastic mapping from the physical to the neural domain. One particularly interesting biological example where such mappings may occur is a mammalian retina, where the continuous-valued sources $\mathbf{x} \in \mathbb{R}^{|\mathbf{x}|}$ define activations at the photoreceptor layer, and the high-dimensional outputs $\mathbf{y} \in \{-1, 1\}^{|\mathbf{y}|}$ correspond to the on-off encodings in the ganglion cells (e.g. Watanabe and Rodieck (1989), Lee et al. (1998)). While the neurophysiological interpretation of the variational decoder $q(\mathbf{x}|\mathbf{y})$ is not properly understood, effectively the distribution defines the stochastic perceptual projection of the retinal encoder models (Eckmiller et al. (2005)). Under the considered

parameterization, the local field of each perceptual unit (x_i under the variational distribution $q(x_i|\mathbf{y})$) is a linear combination of the ganglion firings, with an added white noise. Effectively, this choice of the decoder indicates that small changes in the post-synaptic firings do not significantly vary our guesses about the generating stimuli. It is clear that for the considered choice of the decoder distribution, the percepts are conditionally independent (i.e. $q(\mathbf{x}|\mathbf{y}) = \prod_{i=1}^{|\mathbf{x}|} q(x_i|\mathbf{y})$), i.e. the reconstruction at each perceptual unit is independent from the neighboring percepts for a given pattern of ganglion firings. Conveniently, this facilitates derivations of a local learning rule. Note that for the considered case it is straightforward to evaluate the generic bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ exactly, since it only involves computations of the second-order moments of \mathbf{y} over the factorized distribution $p(\mathbf{y}|\mathbf{x})$.

6.2.1 Local Iterative Learning

Here we consider the case of high-dimensional continuous patterns $\mathbf{x} \in \mathbb{R}^{|\mathbf{x}|}$ represented by stochastic firings of the post-synaptic neurons $\mathbf{y} \in \{-1, +1\}^{|\mathbf{y}|}$. For each neuron y_i , we assume the logistic parameterization of the encoder $p(y_i|\mathbf{x})$, so that the probability of firing monotonically increases with an increase in the membrane potential (using any other parameterization of $p(y_i|\mathbf{x})$ will lead to a straight-forward re-formulation of the model). For conditionally independent activations, we obtain

$$p(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^{|\mathbf{y}|} p(y_i|\mathbf{x}) \stackrel{\text{def}}{=} \prod_{i=1}^{|\mathbf{y}|} \sigma(y_i(\mathbf{w}_i^T \mathbf{x} + b_i)) \quad (6.1)$$

where $\mathbf{w}_i \in \mathbb{R}^{|\mathbf{x}|}$ is a vector of the synaptic weights for neuron y_i , b_i is the corresponding threshold, and $\sigma(a) \stackrel{\text{def}}{=}} 1/(1 + e^{-a})$. It is obvious that (6.1) indeed defines a properly normalized conditional distribution.

By utilizing the factorial assumptions $p(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^{|\mathbf{y}|} \sigma(y_i(\mathbf{w}_i^T \mathbf{x} + b_i))$ and $q(\mathbf{x}|\mathbf{y}) \sim \mathcal{N}(\mathbf{U}\mathbf{y}, s^2\mathbf{I})$, the straight-forward substitution into (2.2) leads to

$$\tilde{I}(\mathbf{x}, \mathbf{y}) \propto \sum_{m=1}^M \text{tr} \left\{ \mathbf{U} \langle \mathbf{y} \rangle_{p(\mathbf{y}|\mathbf{x}^{(m)})} \mathbf{x}_m^T - \frac{1}{2} \mathbf{U}^T \mathbf{U} \langle \mathbf{y} \mathbf{y}^T \rangle_{p(\mathbf{y}|\mathbf{x}^{(m)})} \right\} + \text{const}, \quad (6.2)$$

which needs to be optimized for the encoder and decoder weights $\mathbf{W} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$, $\mathbf{U} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{y}|}$ and the biases $\mathbf{b} \in \mathbb{R}^{|\mathbf{y}|}$ (we have assumed that s^2 is a constant incorporated into the proportionality factor). We may now compute the matrix derivatives of (6.2) to obtain

$$\frac{\partial \tilde{I}(\mathbf{x}, \mathbf{y})}{\partial \mathbf{U}} \propto \sum_{m=1}^M [\mathbf{x}^{(m)} \boldsymbol{\lambda}^T(\mathbf{x}^{(m)}) - \mathbf{U} (\boldsymbol{\lambda}(\mathbf{x}^{(m)}) \boldsymbol{\lambda}^T(\mathbf{x}^{(m)}) - \mathbf{D}_{\lambda_x}(\mathbf{x}^{(m)}) + \mathbf{I})] \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{y}|} \quad (6.3)$$

where $\mathbf{x}^{(m)} = \{x_m^1, \dots, x_m^{|\mathbf{x}|}\} \in \mathbb{R}^{|\mathbf{x}|}$ is the m^{th} vector of input stimuli, and

$$\lambda_i(\mathbf{x}^{(m)}) \stackrel{\text{def}}{=} \langle y_i \rangle_{p(y_i|\mathbf{x}^{(m)})} = 2\sigma(\mathbf{w}_i^T \mathbf{x}^{(m)} + b_i) - 1 \quad (6.4)$$

is the conditional mean. The diagonal matrix $D_{\lambda_x}(\mathbf{x}_m) \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|}$ is defined by

$$D_{\lambda_x} \stackrel{\text{def}}{=} \text{diag}(\lambda_1^2(\mathbf{x}), \dots, \lambda_{|\mathbf{y}|}^2(\mathbf{x})) = \mathbf{I} - \text{cov}(\mathbf{y}|\mathbf{x}) \in [0, 1]^{|\mathbf{y}| \times |\mathbf{y}|}, \quad (6.5)$$

which is a measure of consistency of neural firings. Note that evaluation of the gradient (6.3) requires computing the expected firings $\lambda_i(\mathbf{x})$ of the output units $\{y_i | i = 1, \dots, |\mathbf{y}|\}$, which only involves local feed-forward computations. Specifically, we note that computations of the gradients for the decoder weights $\partial \tilde{I} / \partial \mathbf{U}$ do not require global calculations (such as evaluations of matrix inverses, etc.), which would have been difficult to justify biologically.

Analogously, by differentiating the bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ with respect to the weights and biases of the i^{th} encoding unit, we get

$$\frac{\partial \tilde{I}(\mathbf{x}, \mathbf{y})}{\partial w_{ij}} \propto \sum_{m=1}^M x_j^{(m)} [1 - \lambda_i^2(\mathbf{x}^{(m)})] \mathbf{u}_i^T [\mathbf{x}^{(m)} + \mathbf{u}_i \lambda_i(\mathbf{x}^{(m)}) - \mathbf{U} \boldsymbol{\lambda}(\mathbf{x}^{(m)})], \quad (6.6)$$

$$\frac{\partial \tilde{I}(\mathbf{x}, \mathbf{y})}{\partial b_i} \propto \sum_{m=1}^M [1 - \lambda_i^2(\mathbf{x}^{(m)})] \mathbf{u}_i^T [\mathbf{x}^{(m)} + \mathbf{u}_i \lambda_i(\mathbf{x}^{(m)}) - \mathbf{U} \boldsymbol{\lambda}(\mathbf{x}^{(m)})] \quad (6.7)$$

where \mathbf{u}_i corresponds to the i^{th} column of the decoder weights $\mathbf{U} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{y}|}$, and w_{ij} links the j^{th} pre-synaptic and the i^{th} post-synaptic units. Note that the updates for the biases b_i and encoder weights w_{ij} involve computations of the third and fourth power terms $x_k^{(m)} \langle y_j | \mathbf{x}^{(m)} \rangle \langle y_l | \mathbf{x}^{(m)} \rangle$ and $x_i^{(m)} x_k^{(m)} \langle y_j | \mathbf{x}^{(m)} \rangle \langle y_l | \mathbf{x}^{(m)} \rangle$, which may also be performed locally.

An arguably more biologically plausible interpretation of the learning rule may be obtained by performing the stochastic updates. If $\tilde{y}_i^{(m)} = 1$ with probability $\sigma(\mathbf{w}_i^T \mathbf{x}^{(m)} + b_i)$ (and $\tilde{y}_i^{(m)} = -1$ otherwise), the gradients (6.6) and (6.7) may be transformed to

$$\Delta w_{ij}^{(m)} = \eta \text{var}(y_i | \mathbf{x}^{(m)}) x_j^{(m)} \mathbf{u}_i^T [\mathbf{x}^{(m)} + \mathbf{u}_i \tilde{y}_i^{(m)} - \mathbf{U} \tilde{\mathbf{y}}^{(m)}] \quad (6.8)$$

$$\Delta b_i^{(m)} = \eta \text{var}(y_i | \mathbf{x}^{(m)}) \mathbf{u}_i^T [\mathbf{x}^{(m)} + \mathbf{u}_i \tilde{y}_i^{(m)} - \mathbf{U} \tilde{\mathbf{y}}^{(m)}], \quad (6.9)$$

where η is the learning rate, and $\Delta w_{ij}^{(m)}$, $\Delta b_i^{(m)}$ are the encoder updates for the m^{th} observation. The pre-factor $\text{var}(y_i | \mathbf{x}^{(m)}) = 1 - \lambda_i(\mathbf{x}^{(m)}) \in [0, 1]$ indicates that optimally, the training should slow down once the weights saturate and the firings become more deterministic. (Suboptimally and conventionally, the term is ignored, which corresponds to its first-order Taylor expansion around zero field of y_i). Similarly, we may derive the stochastic updates for the parameter of the perceptual mapping, leading to

$$\Delta \mathbf{U}^{(m)} = \eta_U [\mathbf{x}^{(m)} \tilde{\mathbf{y}}^{(m)} - \mathbf{U} \tilde{\mathbf{Y}}^{(m)}] \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{y}|}, \quad (6.10)$$

where η_U is the learning rate, $\tilde{\mathbf{Y}}^{(m)} \stackrel{\text{def}}{=} \{\tilde{y}_i^{(m)} \tilde{y}_j^{(m)} (1 - \delta_{ij}) + \delta_{ij}\} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{y}|}$, and δ_{ij} is Kronecker delta. Note that the updates in (6.8) – (6.10) are easily representable as weighted Hebbian and anti-Hebbian terms, where the receptive field of the i^{th}

post-synaptic unit y_i is affected not only by the activations at the pre-synaptic layer \mathbf{x} , but also (implicitly) by the stochastic firings of the neighboring post-synaptic units \tilde{y}_j ($j \neq i$).

Generally, expressions (6.3), (6.6) – (6.7) define the updates of the IM algorithm on $\tilde{I}(\mathbf{x}, \mathbf{y})$ and may be used in any standard numerical optimization procedure (see e.g. Bishop (1995)). An arguably more biologically plausible alternative is to perform the stochastic ascent (expressions (6.8) – (6.10)). Finally, we note once again that the explicit parameterization of the encoder $p(\mathbf{y}|\mathbf{x})$ and the variational decoder $q(\mathbf{x}|\mathbf{y})$ makes it easy to incorporate additional constraints on the encoder and decoder parameters $\mathbf{W} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$, $\mathbf{b} \in \mathbb{R}^{|\mathbf{y}|}$, and $\mathbf{U} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{y}|}$, which may be used to transform the variational IM learning so that the solutions satisfy additional requirements.

6.2.2 Optimal Encoder Models

In order to simplify the analysis of the variational bound (6.2), we will follow the strategy of Section 2.2.2 and re-define the objective $\tilde{I}(\mathbf{x}, \mathbf{y})$ as a function of the encoder parameters. As mentioned in Chapter 2, for non-convex functions this may lead to a different optimization surface and a generally more complex (non-local) learning rule. However, the analysis is still interesting, as the resulting bound may be more easily compared with the common numerical approximations of the generally intractable mutual information $I(\mathbf{x}, \mathbf{y})$.

Expressing the bound (6.2) as a function of the encoder $p(\mathbf{y}|\mathbf{x})$ alone, we get

$$\mathbf{U} = \langle \mathbf{x}\mathbf{y}^T \rangle \langle \mathbf{y}\mathbf{y}^T \rangle^{-1}, \quad \tilde{I}(\mathbf{x}, \mathbf{y}) \propto \text{tr} \{ \langle \mathbf{x}\mathbf{y}^T \rangle \langle \mathbf{y}\mathbf{y}^T \rangle^{-1} \langle \mathbf{y}\mathbf{x}^T \rangle \} + \text{const} \quad (6.11)$$

(from now on we will ignore the constant which has no effect on the optimization surface). The objective (6.11) is a proper bound on $I(\mathbf{x}, \mathbf{y})$ for any choice of the stochastic mapping $p(\mathbf{y}|\mathbf{x})$. We may therefore² use it for optimizing a variety of channels with continuous source vectors.

For the considered parameterization of the encoding distribution (6.1), we may transform (6.11) to obtain

$$\tilde{I}(\mathbf{x}, \mathbf{y}) \propto \text{tr} \{ \langle \mathbf{x}\boldsymbol{\lambda}_x^T \rangle \langle \boldsymbol{\lambda}_x \boldsymbol{\lambda}_x^T - \mathbf{D}_{\lambda_x} + \mathbf{I} \rangle^{-1} \langle \boldsymbol{\lambda}_x \mathbf{x}^T \rangle \}, \quad (6.12)$$

where the averages are computed over $\tilde{p}(\mathbf{x})p(\mathbf{y}|\mathbf{x})$, and $\tilde{p}(\mathbf{x}) \propto \sum_{m=1}^M \delta(\mathbf{x} - \mathbf{x}^{(m)})$ is the empirical distribution. Again, $\boldsymbol{\lambda}_x \in [-1, 1]^{|\mathbf{y}|}$ is a vector whose elements $\lambda_i(\mathbf{x}) \stackrel{\text{def}}{=} \langle y_i \rangle_{p(\mathbf{y}|\mathbf{x})} = 2\sigma(\mathbf{w}_i^T \mathbf{x} + b_i) - 1$ correspond to expected firings of the i^{th} unit for a fixed stimulus \mathbf{x} , and $\mathbf{D}_{\lambda_x}(\mathbf{x}^{(m)}) \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{y}|}$ is defined by expression (6.5). The expectations in (6.12) are computed over the empirical distribution.

Since the lower bound (6.12) depends only on the thresholds and synaptic weights, the learning rule is easily obtained by differentiating (6.12) with respect

²From (6.11) it is clear that if $\langle \mathbf{y}\mathbf{y}^T \rangle$ is near-singular, the varying part of the objective $\tilde{I}(\mathbf{x}, \mathbf{y})$ may be infinitely large. However, if the mapping $\mathbf{x} \mapsto \mathbf{y}$ is probabilistic and the number of training stimuli M exceeds the dimensionality of the neural codes $|\mathbf{y}|$, the optimized criterion is typically positive and finite.

to the encoder parameters $\mathbf{b} \in \mathbb{R}^{|\mathcal{Y}|}$ and $\mathbf{W} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{X}|}$ (where rows of \mathbf{W} are given by $\mathbf{w}_i^T \in \mathbb{R}^{1 \times |\mathcal{X}|}$). This leads to

$$\Delta \mathbf{W} \propto \sum_{m=1}^M (\mathbf{I} - \mathbf{D}_{\lambda(\mathbf{x}^{(m)})}) \left(\tilde{\mathbf{D}} \boldsymbol{\lambda}(\mathbf{x}^{(m)}) + \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{yx} [\mathbf{x}^{(m)} - \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\lambda}(\mathbf{x}^{(m)})] \right) (\mathbf{x}^{(m)})^T \quad (6.13)$$

where $\Delta \mathbf{W} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{X}|}$ is the weight update, $\boldsymbol{\Sigma}_{yy} \stackrel{\text{def}}{=} \langle \mathbf{y} \mathbf{y}^T \rangle$, $\boldsymbol{\Sigma}_{yx} \equiv \boldsymbol{\Sigma}_{xy}^T \stackrel{\text{def}}{=} \langle \mathbf{y} \mathbf{x}^T \rangle$ are the second-order moments, and $\tilde{\mathbf{D}}$ is the diagonal matrix

$$\tilde{\mathbf{D}} \stackrel{\text{def}}{=} \text{diag} \left\{ \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{yx} (\boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{yx})^T \right\} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{Y}|}, \quad (6.14)$$

where $\text{diag}(\mathbf{A}) = \{a_{ij} \delta_{ij}\}$, and δ_{ij} is Kronecker delta. The update for the threshold $\Delta \mathbf{b} \in \mathbb{R}^{|\mathcal{Y}|}$ has the same form as (6.13) without the post-multiplication of each term by the training stimulus $(\mathbf{x}^{(m)})^T$.

From (6.13) it is clear that the magnitude of each weight update $\Delta \mathbf{w}_i \in \mathbb{R}^{|\mathcal{X}|}$ should decrease with an increase in the corresponding conditional variance $\text{var}(y_i | \mathbf{x}_m)$. Similarly to the case of the local iterative learning (described in Section 6.2.1), the pre-factor

$$(\mathbf{I} - \mathbf{D}_{\lambda(\mathbf{x}^{(m)})}) = \text{cov}(\mathbf{y} | \mathbf{x}^{(m)}) = \text{diag} \left\{ \text{var}(y_1 | \mathbf{x}^{(m)}), \dots, \text{var}(y_{|\mathcal{Y}|} | \mathbf{x}^{(m)}) \right\} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{Y}|} \quad (6.15)$$

may be interpreted as a variable learning rate. It is intuitive that as training continues, the magnitudes of the synaptic weights would typically increase (as this would lead to a decrease in the conditional entropy $H(\mathbf{y} | \mathbf{x})$). For each unit y_i , this would typically result in a decrease in the conditional variance $\text{var}(y_i | \mathbf{x}^{(m)})$, leading to the parameter saturation effect. The sufficient condition for the training to stop is the noiseless limit of the encoding projection $\text{cov}(\mathbf{y} | \mathbf{x}^{(m)}) \approx \mathbf{0}_{|\mathcal{Y}|}$; for the considered parameterization of the encoding distribution, this may only happen for the divergent membrane potentials of the post-synaptic units. By analogy with a discussion of the Gaussian channel in Section 4.1.1, this stipulates the need of additional constraints³ on the encoder weights.

It is clear that by analogy with (6.8) – (6.10) we may consider stochastic approximations of (6.13), which leads to an arguably more biologically plausible formulation. Nevertheless, the fundamental criticism of the learning rule (6.13) still applies: the rule (6.13) has an intrinsically non-local nature, and a careful approximation of the second-order moments $\boldsymbol{\Sigma}_{yy}^{-1}$, $\boldsymbol{\Sigma}_{yx}$ may need to be considered. The learning rule (6.13) may then be decomposed as a weighted combination of Hebbian and anti-Hebbian terms, though the nature of dependencies of the weighting coefficients on the firing patterns and encoder parameters is non-local.

Finally, it is interesting to note that by considering the first-order expansion of each pre-factor $\text{var}(y_i | \mathbf{x}^{(m)})$ around $\mathbf{w}_i^T \mathbf{x}^{(m)} + b_i \approx 0$, we may transform (6.13) to

$$\Delta \mathbf{W} \propto \tilde{\mathbf{D}} \langle \boldsymbol{\lambda}_x \mathbf{x}^T \rangle + \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{yx} (\boldsymbol{\Sigma}_{xx} - \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \langle \boldsymbol{\lambda}_x \mathbf{x}^T \rangle) \quad (6.16)$$

³Possibly the simplest of such constraints could be introduced by considering a soft constraint on the variances of the conditional firings, which amends (6.13) with an anti-Hebbian term $-\mathbf{A} \langle \boldsymbol{\lambda}_x \mathbf{x}^T \rangle \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{X}|}$ such that $\mathbf{A} \succeq \mathbf{0} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{Y}|}$.

where $\Sigma_{xx} \stackrel{\text{def}}{=} \langle \mathbf{x}\mathbf{x}^T \rangle$. Clearly, (6.16) is decomposable as a combination of the stochastic Hebbian and anti-Hebbian terms, with the weighting coefficients determined by the second-order moments of the firings and input stimuli. Additionally, it is easy to see that the parenthesized factor at the correction of the Hebbian term in (6.16) corresponds to the Schur compliment of the joint correlation $\tilde{\Sigma} \stackrel{\text{def}}{=} \langle [\mathbf{x} \ \mathbf{y}] [\mathbf{x} \ \mathbf{y}]^T \rangle_{p(\mathbf{x}, \mathbf{y})}$ (see e.g. von Mises (1964)). This may be viewed as a conditional covariance $\tilde{\Sigma}_{x|y}$ of the decoder expressed from the *joint* Gaussian model $\tilde{p}_{\mathbf{x}\mathbf{y}} \sim \mathcal{N}_{\mathbf{x}\mathbf{y}}(\mathbf{0}, \tilde{\Sigma})$, resulting in

$$\Delta W \propto \tilde{D} \langle \lambda_x \mathbf{x}^T \rangle + \Sigma_{yy}^{-1} \Sigma_{yx} \tilde{\Sigma}_{x|y}, \quad \tilde{\Sigma}_{x|y} \stackrel{\text{def}}{=} \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx}. \quad (6.17)$$

(cf Linsker’s *as-if Gaussian* bound (1.16)). From (6.13) and (6.17) we see that a simple fixed-rate Hebbian update would generally be a suboptimal approximation of both the exact gradient of the bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ and its more conventional fixed-rate approximation.

6.3 Fisher Information and Mutual Information

Now we will briefly review two classes of methods using the Fisher Information criterion in order to approximate the exact mutual information. The first class of methods which we consider proposes to optimize an approximation of a specific bound on $I(\mathbf{x}, \mathbf{y})$, which is conveniently derived by applying the Cramer-Rao and the data processing inequalities (Brunel and Nadal (1998)), and computing the limit for $|\mathbf{y}| \rightarrow \infty$. The same result may be obtained by considering a numerical approximation of the “local mutual information” $I_x = \langle \log p(\mathbf{y}|\mathbf{x})/p(\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})}$ for $|\mathbf{y}| \rightarrow \infty$, and integrating the approximation over the sources $\tilde{p}(\mathbf{x})$ (Kang and Sompolinsky (2001)).

The second method is inspired by Szummer and Jaakkola (2002) and Corduneanu and Jaakkola (2003), who also propose to approximate the local information $I_{\mathcal{R}}(\mathbf{x}, \mathbf{y})$ for each small region of the source space $\mathcal{R} \subset \mathcal{R}_x, \mathcal{R} \rightarrow 0$ (cf Kang and Sompolinsky (2001)), but use different numerical relaxations, which generally leads to a different optimization surface. Specifically, their approximations do not depend on the population sizes $|\mathbf{y}|$, which suggests better convergence properties for low and intermediate values of $|\mathbf{y}|$. Though the original results of Szummer and Jaakkola (2002) and Corduneanu and Jaakkola (2003) are applied in the different context of semi-supervised learning, we show that their approach may be considered as a general alternative for computing an approximate lower bound on the exact mutual information $I(\mathbf{x}, \mathbf{y})$.

It turns out that both approximations we discuss here involve computations of the Fisher Information criterion. To differentiate between the methods, we will refer to the method of Szummer and Jaakkola (2002) and Corduneanu and Jaakkola (2003) as the *local approximation* of mutual information, since the fundamental assumption used in the derivations is the locality of the regions (the Fisher Information matrix arises as a result of a specific expansion of the encoder distribution). In contrast, we will refer to the approach of Brunel and Nadal (1998) as *Fisher approximation* of mutual information (or *Fisher criterion*), since

the term arises as the result of applying the fundamental Cramer-Rao inequality. For both of the considered methods, we outline the corresponding learning rules for the considered case of a sigmoidal network.

6.3.1 Fisher Approximation of Mutual Information

Let us assume that the stochastic firings of the output units are conditionally independent given the source pattern, i.e. $p(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^{|\mathbf{y}|} p_i(y_i|\mathbf{x})$, where all $p_i(y_i|\mathbf{x}) \equiv p(y_i|\mathbf{x})$ are in the same parametric family. If we view $\mathbf{x} \in \mathbb{R}^{|\mathbf{x}|}$ as a fixed parameter, and the vector of binary spikes $\mathbf{y} \in \{\pm 1\}^{|\mathbf{y}|}$ – as $|\mathbf{y}|$ independent identically distributed samples from $p(y_i|\mathbf{x})$, we may apply the results of the theory of statistical parameter estimation (see e.g. Cramer (1946)) to bound the mutual information $I(\mathbf{x}, \mathbf{y})$ (Brunel and Nadal (1998)).

Indeed, let $\hat{\mathbf{x}} \in \mathbb{R}^{|\mathbf{x}|}$ be a statistical estimator of the input stimulus \mathbf{x} obtained from the stochastic neural firings \mathbf{y} . Generally, we may assume that the sources, the encodings, and the estimators form a Markov chain $\mathbf{x} \rightarrow \mathbf{y} \mapsto \hat{\mathbf{x}}$, where $p(\hat{\mathbf{x}}|\mathbf{y}) \sim \delta(\hat{\mathbf{x}} - \hat{\mathbf{x}}(\mathbf{y}))$ and $p(\mathbf{y}|\mathbf{x})$ is factorized in \mathbf{y} . It is well known that if the estimator $\hat{\mathbf{x}}(\mathbf{y})$ is unbiased, i.e. $\langle \hat{\mathbf{x}}(\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})} = \mathbf{x}$, then its covariance may be bounded according to the Cramer-Rao inequality, i.e.

$$\left\langle (\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x}) (\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x})^T \right\rangle_{p(\mathbf{y}|\mathbf{x})} \succeq \mathbf{F}_{\mathbf{x}} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|}, \quad (6.18)$$

where $\mathbf{A}_1 \succeq \mathbf{A}_2$ indicates that $\mathbf{A}_1 - \mathbf{A}_2$ is positive semi-definite, and

$$\mathbf{F}_{\mathbf{x}} = \{F_{ij}(\mathbf{x})\} \stackrel{\text{def}}{=} - \left\{ \langle \partial^2 \log p(\mathbf{y}|\mathbf{x}) / \partial x_i \partial x_j \rangle_{p(\mathbf{y}|\mathbf{x})} \right\} \quad (6.19)$$

$$\equiv \left\{ \langle \partial \log p(\mathbf{y}|\mathbf{x}) / \partial x_i \cdot \partial \log p(\mathbf{y}|\mathbf{x}) / \partial x_j \rangle_{p(\mathbf{y}|\mathbf{x})} \right\} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|} \quad (6.20)$$

is the Fisher Information matrix (see e.g. Cramer (1946), Cover and Thomas (1991), Johnson (2003)). It is easy to see that we may equivalently express the bound (6.18) as

$$\text{cov}(\hat{\mathbf{x}}|\mathbf{x}) \succeq \mathbf{F}_{\mathbf{x}} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|}. \quad (6.21)$$

Here $\text{cov}(\hat{\mathbf{x}}|\mathbf{x})$ is the conditional covariance of $p(\hat{\mathbf{x}}|\mathbf{x})$, and we have used the fact that for unbiased estimators $\langle \hat{\mathbf{x}} \rangle_{p(\hat{\mathbf{x}}|\mathbf{x})} = \langle \hat{\mathbf{x}}(\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})} = \mathbf{x}$.

If the estimator $\hat{\mathbf{x}}(\mathbf{y})$ is also *efficient*, its covariance saturates the Cramer-Rao bound, which results in an upper bound on the entropy of the conditional distribution $H(p(\hat{\mathbf{x}}|\mathbf{x})) \leq H(\mathcal{N}_{\mathbf{x}}(\mathbf{0}, \mathbf{F}_{\mathbf{x}}^{-1}))$. The bound follows from the well-known fact that for a fixed covariance, the maximum entropy distribution is a Gaussian (see e.g. Feller (1971), McEliece (1977)). We may now obtain a lower bound on the mutual information

$$I(\mathbf{x}, \mathbf{y}) \geq H(\hat{\mathbf{x}}) + \frac{1}{2} \langle \log |\mathbf{F}_{\mathbf{x}}| \rangle_{\hat{p}(\mathbf{x})} + \text{const}, \quad (6.22)$$

where we have used the data processing inequality $I(\mathbf{x}, \mathbf{y}) \geq I(\mathbf{x}, \hat{\mathbf{x}})$ for the considered chain $\mathbf{x} \rightarrow \mathbf{y} \mapsto \hat{\mathbf{x}}$.

Unfortunately, despite the fact that the mapping $\mathbf{y} \mapsto \hat{\mathbf{x}}$ is deterministic, exact computation of the entropy of statistical estimates $H(\hat{\mathbf{x}})$ in the objective (6.22) is in general computationally intractable. Brunel and Nadal (1998) show that under some assumptions we may assume $H(\hat{\mathbf{x}}) \approx H(\mathbf{x})$, which leads to the approximation

$$I(\mathbf{x}, \mathbf{y}) \gtrsim \tilde{I}_F(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} H(\mathbf{x}) + \frac{1}{2} \langle \log |\mathbf{F}_{\mathbf{x}}| \rangle_{\tilde{p}(\mathbf{x})} + \text{const} \quad (6.23)$$

(specifically, this approximation applies for $|\mathbf{y}| \rightarrow \infty$ when the estimators are sharply peaked around the mean values). Brunel and Nadal also show that under similar assumptions, (6.23) may be used as an approximation of $I(\mathbf{x}, \mathbf{y})$ independently of the bias of the estimator. Note that since $H(\mathbf{x})$ is independent of the parameters of $p(\mathbf{y}|\mathbf{x})$, maximization of (6.23) is equivalent to maximization of (6.22) where the generally intractable term $H(\hat{\mathbf{x}})$ (entropy of the mixture) is ignored.

6.3.1.1 Sigmoidal Activations

It is straight-forward to see that for the considered sigmoidal activations (6.1), each element $F_{ij}(\mathbf{x})$ of the Fisher Information matrix $\mathbf{F}_{\mathbf{x}} = \{F_{ij}(\mathbf{x})\}$ is given by

$$F_{ij}(\mathbf{x}) = \sum_{l=1}^{|\mathbf{y}|} w_{li} w_{lj} \sigma(\mathbf{w}_l^T \mathbf{x} + b_l) (1 - \sigma(\mathbf{w}_l^T \mathbf{x} + b_l)), \quad i, j = 1, \dots, |\mathbf{y}|, \quad (6.24)$$

where \mathbf{w}_l^T is the l^{th} row of the synaptic weights $\mathbf{W} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$. Then we may express the Fisher approximation of the mutual information (6.23) as

$$\tilde{I}_F(\mathbf{x}, \mathbf{y}) \propto \sum_{m=1}^M \log |\mathbf{W}^T (\mathbf{I} - \mathbf{D}_{\lambda_{x_m}}) \mathbf{W}| + \text{const}, \quad (6.25)$$

where the constant incorporates the remaining terms which have no effect on the optimization. Again, $\mathbf{I} - \mathbf{D}_{\lambda_{x_m}} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{y}|}$ is the conditional covariance of the stochastic spikes (given by expression (6.5)).

By computing the matrix derivatives of (6.25) for $\mathbf{W} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$ and $\mathbf{b} \in \mathbb{R}^{|\mathbf{y}|}$ and performing some straight-forward algebraic manipulations, we get

$$\frac{\partial \tilde{I}_F}{\partial \mathbf{W}} = \frac{1}{4M} \sum_{m=1}^M \left\{ 2(\mathbf{I} - \mathbf{D}_{\lambda_{x_m}}) \mathbf{W} \mathbf{A}_{x_m} - \mathbf{C}_{x_m} (\mathbf{I} - \mathbf{D}_{\lambda_{x_m}}) \boldsymbol{\lambda}(\mathbf{x}^{(m)}) (\mathbf{x}^{(m)})^T \right\} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}. \quad (6.26)$$

Here $\boldsymbol{\lambda}(\mathbf{x}^{(m)}) \stackrel{\text{def}}{=} \langle \mathbf{y} \rangle_{p(\mathbf{y}|\mathbf{x}^{(m)})} \in \mathbb{R}^{|\mathbf{y}|}$ is a vector of the expected firing at the encoding layer (see expression (6.4)), and

$$\mathbf{A}_{x_m} = (\mathbf{W}^T (\mathbf{I} - \mathbf{D}_{\lambda_{x_m}}) \mathbf{W})^+ \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|} \quad (6.27)$$

$$\mathbf{C}_{x_m} = \text{diag} (\mathbf{W} \mathbf{A}_{x_m} \mathbf{W}^T) \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{y}|}, \quad (6.28)$$

where \mathbf{A}^+ is the pseudo-inverse⁴ (so that $\mathbf{A}\mathbf{A}^+\mathbf{A} = \mathbf{A}$). Analogously,

$$\frac{\partial \tilde{I}_F}{\partial \mathbf{b}} = -\frac{1}{4M} \sum_{m=1}^M \mathbf{C}_{x_m} (\mathbf{I} - \mathbf{D}_{\lambda_{x_m}}) \boldsymbol{\lambda}(\mathbf{x}^{(m)}) \in \mathbb{R}^{|\mathbf{y}|} \quad (6.29)$$

with the similar definitions of $\boldsymbol{\lambda}(\mathbf{x}^{(m)}) \in \mathbb{R}^{|\mathbf{y}|}$ and $\mathbf{C}_{x_m} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{y}|}$.

It is interesting to note that for the square model with $|\mathbf{x}| = |\mathbf{y}|$, optimization of (6.25) leads to

$$\Delta \mathbf{W} = 2\mathbf{W}^{-T} - \langle \boldsymbol{\lambda}_x \mathbf{x}^T \rangle, \quad (6.30)$$

which (apart from the coefficient at the inverse weight – *redundancy* term) has the same form as the learning rule of Bell and Sejnowski (1995) derived for the invertible channel with $p(\mathbf{y}|\mathbf{x}) \sim \delta(\mathbf{y} - \sigma(\mathbf{x}))$. If the encoded representations are overcomplete (i.e. $|\mathbf{x}| < |\mathbf{y}|$) and the variable rate $\text{cov}(\mathbf{y}|\mathbf{x}^{(m)}) = \mathbf{I} - \mathbf{D}_{\lambda_{x_m}}$ is approximated by the 1st-order Taylor expansion around zero receptive fields, the redundancy term in (6.30) is replaced by the transposed pseudo-inverse $\mathbf{W}(\mathbf{W}^T\mathbf{W})^{-1} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$. Notably, the resulting learning rule (6.26), (6.29) is non-local (as it involves computations of the inverses), and the weight updates (6.26), (6.30) have no Hebbian terms.

More importantly, one can see that optimization of the Fisher approximation of the mutual information (6.25) may be problematic when $\mathbf{W}^T(\mathbf{I} - \mathbf{D}_{\lambda_{x_m}})\mathbf{W} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|}$ is rank-deficient, which complicates applicability of the method for a variety of tasks involving relatively low-dimensional encodings of high-dimensional input stimuli. Apart from instability of numerical optimization on $\tilde{I}_F(\mathbf{y}, \mathbf{x})$ (which we have partially addressed by computing the pseudo-inverses, rather than the exact inverses in (6.27), (6.28)), the conceptual problem of optimizing the criterion $\tilde{I}_F(\mathbf{y}, \mathbf{x})$ is the weakness of the approximate bound for $|\mathbf{y}| < |\mathbf{x}|$ (in most cases, such channel choices would violate the fundamental assumptions of the approximation). Additionally, from (6.25) we may note that the approximate bound $\tilde{I}_F(\mathbf{x}, \mathbf{y})$ may become loose with the decrease in the conditional variance of the stochastic firings. Intuitively, this happens because the directions of low variations in the code space swamp the optimized volume of the $|\mathbf{x}|$ -dimensional manifold (see (6.25)), thus leading to a drop in the determinant.

6.3.2 Local Approximation of Mutual Information

Here we describe an alternative approximation of mutual information $I(\mathbf{x}, \mathbf{y})$. It is inspired by recent work of Szummer and Jaakkola (2002) and Corduneanu and Jaakkola (2003), who approximate information content $I(\mathbf{x} \in \mathcal{R}, \mathbf{y})$ in infinitely small local regions $\mathcal{R} \subset \mathcal{R}_x$, and apply the approximations as regularizers for semi-supervised classification. Here we show that their method may in fact be used to obtain an alternative approximate lower bound on mutual information.

Consider the model $r \rightarrow \mathbf{x} \rightarrow \mathbf{y}$, where r defines a local region in the data space, so that $\mathcal{R}_r \subset \mathcal{R}_x$. The local information $I_r(\mathbf{x}, \mathbf{y})$ in the region r may be

⁴We make a recourse to computing the pseudo-inverse in order to handle the singularity of the gradient $\partial \tilde{I}_F(\mathbf{x}, \mathbf{y}) / \partial \mathbf{W}$ for rank-deficient Fisher Information matrices (which happens for example when $|\mathbf{x}| > |\mathbf{y}|$). Note that $\mathbf{A}^+ = \mathbf{A}^{-1} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|}$ when $\text{rank}(\mathbf{A}) = |\mathbf{x}|$.

defined as

$$I_r(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} \int_{\mathcal{R}_y} \int_{\mathcal{R}_r} p(\mathbf{x}|r)p(\mathbf{y}|\mathbf{x}) \log \frac{p(\mathbf{y}|\mathbf{x})}{p(\mathbf{y}|r)} \, d\mathbf{x} \, d\mathbf{y}, \quad (6.31)$$

where \mathcal{R}_y is the code space. The conditional $p(\mathbf{x}|r)$ may generally be computed as $p(\mathbf{x}|r) \propto p(\mathbf{x})\mathbb{I}(\mathbf{x} \in \mathcal{R}_r)$ [\mathbb{I} defines an indicator variable], where $p(\mathbf{x})$ is the distribution of the source variables. By computing the 2^{nd} -order Taylor expansion of $p(\mathbf{y}|\mathbf{x})$ around $p(\mathbf{y}|\langle \mathbf{x} \rangle_{p(\mathbf{x}|r)})$ and performing some algebraic manipulations, Corduneanu and Jaakkola (2003) show that the local information (6.31) may be approximated as

$$I_r(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \text{tr} \left\{ \text{cov}(\mathbf{x}|r) \mathbf{F}_{\langle \mathbf{x} \rangle_{p(\mathbf{x}|r)}} \right\} + O(\mathcal{R}_r^3), \quad (6.32)$$

where $\mathbf{F}_{\langle \mathbf{x} \rangle_{p(\mathbf{x}|r)}}$ is the Fisher Information matrix (6.19) computed at the local mean $\langle \mathbf{x} \rangle_{p(\mathbf{x}|r)}$.

If the local regions \mathcal{R}_r are symmetric and centered at the training patterns $\{\mathbf{x}\}$ then we may define $\text{cov}(\mathbf{x}|r) \stackrel{\text{def}}{=} v_{\mathbf{x}|r} \mathbf{I}_{|\mathbf{x}|}$, and up to a constant pre-factor we obtain⁵

$$\langle I_r(\mathbf{x}, \mathbf{y}) \rangle_{p(r)} \approx v_{\mathbf{x}|r} \langle \text{tr} \{ \mathbf{F}_{\mathbf{x}} \} \rangle_{\tilde{p}(\mathbf{x})} \quad (6.33)$$

(approximation is accurate for small regions \mathcal{R}_r , i.e. small variances $v_{\mathbf{x}|r}$). Again, $\tilde{p}(\mathbf{x})$ is the empirical distribution. We also note that when expression (6.33) is coupled with the minimization of the sum-of-squared error, the resulting objective is equivalent (to first order) to minimization of the sum-of-squared error with noise, i.e. when the source distribution is given by $p(\mathbf{x}) = \sum_{m=1}^M \mathcal{N}_{\mathbf{x}}(\mathbf{x}^{(m)}, v_{\mathbf{x}|r} \mathbf{I}_{|\mathbf{x}|})$ for the training set $\{\mathbf{x}^{(m)} | m = 1, \dots, M\}$, and the encoding distribution $p(\mathbf{y}|\mathbf{x})$ is approximated around $p(\mathbf{y}|\mathbf{x}^{(m)})$ (see Bishop (1995)).

We will now show that (6.33) may be used to derive an approximate lower bound on the global mutual information $I(\mathbf{x}, \mathbf{y})$. To do so, we will outline a general relation between $I(\mathbf{x}, \mathbf{y})$ and the local average $\langle I_r(\mathbf{x}, \mathbf{y}) \rangle_{p(r)}$ for the considered model $p(r, \mathbf{x}, \mathbf{y}) = p(r)p(\mathbf{x}|r)p(\mathbf{y}|\mathbf{x})$. From the chain rule on mutual information (see e.g. Cover and Thomas (1991)), we may express the joint mutual information between the code and the source vectors (with the corresponding local partitionings) as

$$I(\{r, \mathbf{x}\}, \mathbf{y}) = I(r, \mathbf{y}) + I(\mathbf{x}, \mathbf{y}|r). \quad (6.34)$$

On the other hand, for the considered chain model we get

$$I(\{r, \mathbf{x}\}, \mathbf{y}) = H(\mathbf{y}) - H(\mathbf{y}|\mathbf{x}, r) = I(\mathbf{x}, \mathbf{y}), \quad (6.35)$$

where $I(\{r, \mathbf{x}\}, \mathbf{y})$ is the amount of information which the codes \mathbf{y} contain about the sources and the regions \mathbf{x}, r jointly. From (6.31), (6.33), and the definition of the conditional mutual information

$$I(\mathbf{x}, \mathbf{y}|r) \stackrel{\text{def}}{=} H(\mathbf{y}|r) - H(\mathbf{y}|\mathbf{x}, r) = \langle I_r(\mathbf{x}, \mathbf{y}) \rangle_{p(r)}, \quad (6.36)$$

⁵This result may be generalized for the case of infinitely small local regions \mathcal{R}_r , provided that the covering of the data space satisfies specific properties. For example, (6.33) applies if \mathcal{R}_x consists of axis-parallel cubes $\mathcal{R}_r \rightarrow 0$ centered at the axis-parallel lattice points, see Corduneanu and Jaakkola (2003) for details.

it is clear that

$$I(\mathbf{x}, \mathbf{y}) = I(r, \mathbf{y}) + \langle I_r(\mathbf{x}, \mathbf{y}) \rangle_{p(r)} \gtrsim v_{\mathbf{x}|r} \langle \text{tr} \{ \mathbf{F}_{\mathbf{x}} \} \rangle_{\tilde{p}(\mathbf{x})} \stackrel{\text{def}}{=} \tilde{I}_L(\mathbf{x}, \mathbf{y}), \quad (6.37)$$

where the slackness of the bound is stipulated by the non-negativity of $I(r, \mathbf{y})$. Note that the variance $v_{\mathbf{x}|r}$ is a constant parameter which does not affect the optimization surface for the encoding distribution $p(\mathbf{y}|\mathbf{x})$. Note that for the special case of one-dimensional input stimuli, the approximation (6.37) reduces to the well-known scalar Fisher criterion optimized by a number of the currently used approaches to population coding (e.g. Pouget et al. (1998), Zhang and Sejnowski (1999), Bethge et al. (2002)).

Despite the apparent similarity of the Fisher approximation criterion $\tilde{I}_F(\mathbf{x}, \mathbf{y})$ (expression (6.23)) and the local criterion $\tilde{I}_L(\mathbf{x}, \mathbf{y})$ (expression (6.37)), one can see that in general the methods result in different fixed points. Indeed, if $\lambda_i(\mathbf{F}_{\mathbf{x}})$ is the i^{th} eigenvalue of $\mathbf{F}_{\mathbf{x}} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|}$ then the extrema of $\tilde{I}_F(\mathbf{x}, \mathbf{y})$ must satisfy

$$\partial \tilde{I}_F(\mathbf{x}, \mathbf{y}) / \partial \mathbf{W} = \left\langle \sum_{i=1}^{|\mathbf{x}|} \lambda_i^{-1}(\mathbf{F}_{\mathbf{x}}) \frac{\partial \lambda_i(\mathbf{F}_{\mathbf{x}})}{\partial \mathbf{W}} \right\rangle_{\tilde{p}(\mathbf{x})} = \mathbf{0} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}. \quad (6.38)$$

On the other hand,

$$\partial \tilde{I}_L(\mathbf{x}, \mathbf{y}) / \partial \mathbf{W} = \left\langle \sum_{i=1}^{|\mathbf{x}|} \frac{\partial \lambda_i(\mathbf{F}_{\mathbf{x}})}{\partial \mathbf{W}} \right\rangle_{\tilde{p}(\mathbf{x})} = \mathbf{0} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}. \quad (6.39)$$

Clearly, both (6.38) and (6.39) hold simultaneously if $\partial \lambda_i(\mathbf{F}_{\mathbf{x}}) / \partial \mathbf{W} = \mathbf{0}$ for all eigenvalues $\{\lambda_i | i = 1, \dots, |\mathbf{x}|\}$, and all training patterns $\{\mathbf{x}^{(m)} \in \mathbb{R}^{|\mathbf{x}|} | m = 1, \dots, M\}$. However, it is clear that despite positive semi-definiteness of $\mathbf{F}_{\mathbf{x}} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|}$, the criteria $\tilde{I}_F(\mathbf{x}, \mathbf{y})$ and $\tilde{I}_L(\mathbf{x}, \mathbf{y})$ will generally give rise to different solutions.

Finally, it is important to note that in contrast to the Fisher approximation $\tilde{I}_F(\mathbf{x}, \mathbf{y})$ of Brunel and Nadal (1998), the local approximation $\tilde{I}_L(\mathbf{x}, \mathbf{y})$ is also defined for $|\mathbf{x}| > |\mathbf{y}|$. This suggests that in general numerical optimization of (6.37) may be more stable than optimization of (6.23) for small population sizes. Optimization of the local approximation is also less computationally demanding, since by optimizing $\tilde{I}_L(\mathbf{x}, \mathbf{y})$ we avoid computing the inverse of the Fisher Information matrix. Moreover, as we will see, for the considered logistic parameterization of the encoder distribution, numerical ascent on $\tilde{I}_L(\mathbf{x}, \mathbf{y})$ may often lead to better values of the exact mutual information $I(\mathbf{x}, \mathbf{y})$ (which we will compute for small models for the purpose of testing).

6.3.2.1 Sigmoidal Activations

By analogy with Section 6.3.1, we may derive the local approximation of $I(\mathbf{x}, \mathbf{y})$ for the considered sigmoidal encoder model (6.1). Up to the constant pre-factor $v_{\mathbf{x}|r}$, the local approximation is given by

$$\tilde{I}_L(\mathbf{x}, \mathbf{y}) \propto \sum_{m=1}^M \text{tr} \{ \mathbf{W}^T (\mathbf{I} - \mathbf{D}_{\lambda_{x_m}}) \mathbf{W} \}, \quad (6.40)$$

where $\mathbf{I} - \mathbf{D}_{\lambda_{x_m}} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{Y}|}$ is given by expression (6.5). Note that in contrast to the previously discussed Fisher approximation for sigmoidal activations (6.25), the objective (6.40) is finite even when $\mathbf{F}_x \in \mathbb{R}^{|\mathcal{X}| \times \mathcal{X}}$ is rank-deficient. This generally leads to a more stable behavior of the numerical optimization procedures performing an ascent on $\tilde{I}_L(\mathbf{x}, \mathbf{y})$.

The gradients for the encoder parameters in this case are given by

$$\frac{\partial \tilde{I}_L}{\partial \mathbf{W}} \propto \frac{1}{4M} \sum_{m=1}^M \left\{ 2(\mathbf{I} - \mathbf{D}_{\lambda_{x_m}}) \mathbf{W} - \text{diag}(\mathbf{W} \mathbf{W}^T) (\mathbf{I} - \mathbf{D}_{\lambda_{x_m}}) \boldsymbol{\lambda}(\mathbf{x}^{(m)}) (\mathbf{x}^{(m)})^T \right\} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{X}|} \quad (6.41)$$

$$\frac{\partial \tilde{I}_L}{\partial \mathbf{b}} \propto -\frac{1}{4M} \sum_{m=1}^M \text{diag}(\mathbf{W} \mathbf{W}^T) (\mathbf{I} - \mathbf{D}_{\lambda_{x_m}}) \boldsymbol{\lambda}(\mathbf{x}^{(m)}) \in \mathbb{R}^{|\mathcal{Y}|}, \quad (6.42)$$

where $\text{diag}(\mathbf{W} \mathbf{W}^T) \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{Y}|}$ is the diagonal matrix of the squares of L2-norms $\|\mathbf{w}_i\|^2$ of the encoder weight vectors. Note that symbolically the only difference of the gradients of the local approximation $\tilde{I}_L(\mathbf{x}, \mathbf{y})$ (expressions (6.41) and (6.42)) from the gradients of the Fisher approximation $\tilde{I}_F(\mathbf{x}, \mathbf{y})$ ((6.26) and (6.29)) is multiplication by the non-local term $\mathbf{A}_{x_m} \in \mathbb{R}^{|\mathcal{X}| \times |\mathcal{X}|}$ in $\partial \tilde{I}_F / \partial \mathbf{W}$ and $\mathbf{C}_{x_m} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{Y}|}$ (see expressions (6.27), (6.28)).

6.3.3 Constraints on the Encoder Parameters

Interestingly, from the definitions (6.25), (6.40) of the Fisher and local approximations of the mutual information for the considered encoder model, one can see that unconstrained optimization for the encoder weights $\mathbf{W} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{X}|}$ does not necessarily lead to their divergence. To see this, consider the simple case of a scalar input stimulus (e.g. $x \in \mathbb{R}$ may correspond to the angle of saccadic eye movements or direction of head motion) encoded by a single spiking neuron $y \in \{-1, 1\}$. The only synaptic weight for this simple model would be given by $w \in \mathbb{R}$ (for clarity, we ignore the bias b). From (6.40) and the definition of the conditional variance (6.5), it is clear that

$$\tilde{I}_L(\mathbf{x}, \mathbf{y}) \propto \langle w^2 / (2 + e^{-wx} + e^{wx}) \rangle_{\tilde{p}(\mathbf{x})} \geq 0, \quad (6.43)$$

where $\tilde{p}(\mathbf{x})$ is the empirical distribution. Assume that $\exists \epsilon > 0$ such that for all training patterns $|x^{(m)}| > \epsilon$. Then by computing the limit of (6.43) for $|w| \rightarrow \infty$ we can see that the approximation $\tilde{I}_L(\mathbf{x}, \mathbf{y})$ is *minimized* for the divergent weight norm $|w|$. (Note that finding the analytical solution for w which maximizes (6.43) would require solving a transcendental equation even for this simple model with $|\mathcal{X}| = |\mathcal{Y}| = 1$). Similar intuition applies in the case of Brunel and Nadal's approximation⁶ $\tilde{I}_F(\mathbf{x}, \mathbf{y})$.

Analogously, for the high-dimensional case we may note that the conditional variances of the stochastic firings are functions of the singular values $\mathbf{L} \in \mathbb{R}^{|\mathcal{X}| \times |\mathcal{X}|}$

⁶Consider the case when $|\mathcal{X}| = |\mathcal{Y}| = 1$. For the Fisher approximation we note that from Jensen's inequality $\tilde{I}_F(x, y) = \langle \log |\mathbf{F}_x| \rangle_{\tilde{p}(x)} \leq \log \langle |\mathbf{F}_x| \rangle_{\tilde{p}(x)} \stackrel{\text{def}}{=} \hat{I}_F(x, y)$. It is easy to see that $\lim_{|w| \rightarrow \infty} \hat{I}_F(x, y) \rightarrow -\infty$. Since $\tilde{I}_F(x, y) \leq \hat{I}_F(x, y)$, we get $\lim_{|w| \rightarrow \infty} \tilde{I}_F(x, y) \rightarrow -\infty$.

of the encoding weights $\mathbf{W} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$. For the considered choice of the encoding distributions, higher weight magnitudes would generally lead to lower conditional variances of the firings, i.e. the growth in $\|\mathbf{W}\|_F \equiv \text{tr}^{1/2}\{\mathbf{W}\mathbf{W}^T\}$ in the approximation $\tilde{I}_L(\mathbf{x}, \mathbf{y})$ would be compensated by a decrease in $\mathbf{I} - \mathbf{D}_{\lambda_{x_m}} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{y}|}$. This suggests generally non-trivial changes in $\hat{I}_L(\mathbf{x}, \mathbf{y})$ with an increase in the weight magnitudes (*cf* Gaussian channels with the isotropic noise).

Despite the observation that unconstrained optimization of $\tilde{I}_F(x, y)$ and $\tilde{I}_L(x, y)$ would not necessarily lead to indefinite growth of the encoder weights \mathbf{W} (at least, for the considered channel), in many cases we may still be interested in constraining the magnitudes of the synaptic weights or the conditional variances of the firings. This may be motivated, for example, by neuro-physiological constraints on the channel, or technical limitations on precision. Possibly one of the simplest ways of introducing implicit constraints on the encoder parameters for the considered channels is to penalize large Frobenius norms $\|\mathbf{W}\|_F$, or small conditional variances $\text{var}(y_j|\mathbf{x}) = 1 - \lambda_j^2(\mathbf{x}^{(m)})$ of the firing units. It is straight-forward to see that the gradients of the penalty terms would in this case be given as $-\mathbf{M}\mathbf{W}$ and $-\mathbf{M}\langle \boldsymbol{\lambda}_x \mathbf{x}^T \rangle \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$ (which is the anti-Hebbian term) respectively, where $\mathbf{M} \succeq \mathbf{0} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{y}|}$ is a fixed positive semi-definite matrix. However, while being conceptually and computationally simple, the choice of the penalty constants \mathbf{M} may be rather heuristic. A more rigorous approach would involve optimization of the dual Lagrangian

$$\tilde{\mathbb{L}} \stackrel{\text{def}}{=} \sup_{\mathbf{W}} \{\mathbb{L}(\mathbf{W}, \mathbf{M})\} \stackrel{\text{def}}{=} \sup_{\mathbf{W}} \{\tilde{I}_F(\mathbf{x}, \mathbf{y}) - \mathbf{m}^T \mathbf{f}(\mathbf{W})\}, \quad (6.44)$$

where $\mathbf{f}(\mathbf{W}) \geq \mathbf{0}$ defines a set of the inequality constraints, and $\mathbf{m} \geq \mathbf{0}$ is a vector of Lagrange multipliers (see e.g. Bertsekas (1996), Boyd and Vandenberghe (2004)), which requires further numerical approximations when the dual cannot be expressed as an analytical function of \mathbf{W} alone (e.g. Rubinov and Yang (2003)).

Alternatively, we may impose explicit construction constraints on the encoder parameters $\mathbf{W} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$, for example by constraining the weights to a unit hypersphere. For example, for $|\mathbf{y}| \geq |\mathbf{x}|$ we may define $\mathbf{W} = \tilde{\mathbf{W}}(\tilde{\mathbf{W}}^T \tilde{\mathbf{W}})^{-1/2}$, where $\tilde{\mathbf{W}} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$ is an arbitrary matrix such that $\text{rank}(\tilde{\mathbf{W}}) = |\mathbf{y}|$, and re-compute the derivatives for $\tilde{\mathbf{W}}$ (the case when $|\mathbf{x}| \geq |\mathbf{y}|$ is analogous). We may also bound each synaptic weight, so that $w_{ij} \in [-\omega, \omega]$, for example by imposing parametric constraints $w_{ij} \stackrel{\text{def}}{=} \phi(a_{ij})$, where $\phi(a) : \mathbb{R} \rightarrow [-\omega, \omega]$. Presumably, such constraints on the individual synaptic weights are more biologically intuitive; moreover, they do not require non-local computations. Optimization of $\tilde{I}_F(\mathbf{x}, \mathbf{y})$ or $\tilde{I}_L(\mathbf{x}, \mathbf{y})$ with respect to $\mathbf{W} = \{w_{ij}\} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$ would then need to be replaced by optimization with respect to $\mathbf{A} = \{a_{ij}\} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$, with the gradients expressed e.g. by

$$\partial \tilde{I}_F(\mathbf{x}, \mathbf{y}) / \partial \mathbf{A} = \partial \tilde{I}_F(\mathbf{x}, \mathbf{y}) / \partial \mathbf{W} \circ \nabla_{\mathbf{A}} \phi(\mathbf{A}) \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}. \quad (6.45)$$

Here $\mathbf{A}_1 \circ \mathbf{A}_2$ denotes the element-wise product, and $\nabla_{\mathbf{A}} \phi(\mathbf{A}) \stackrel{\text{def}}{=} \{\partial \phi(a_{ij}) / \partial a_{ij}\} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$. Interestingly, we can empirically demonstrate that in some cases (e.g. when the input stimuli $\mathbf{x} \in \mathbb{R}^{|\mathbf{x}|}$ are sufficiently close to $\mathbf{0}_{|\mathbf{x}|}$), this choice of the constraint on $\mathbf{W} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$ pushes the optimal weights towards the corners of the

hypercube $[-\omega, \omega]^{|y| \times |x|}$, thus suggesting a straight-forward approach to unsupervised information-theoretic training of ising perceptrons (see e.g. Penney and Sherrington (1993), Rosen-Zvi and Kanter (2001)).

6.4 Variational Lower Bound vs. Fisher and Local Approximations of Mutual Information

Since $\tilde{I}_F(\mathbf{x}, \mathbf{y})$ and $\tilde{I}_L(\mathbf{x}, \mathbf{y})$ are not proper lower bounds on the mutual information, it is difficult to analyze their tightness or compare them with the generic variational bound (2.2). To illustrate a relation between the approaches, we may consider a Gaussian decoder $q(\mathbf{x}|\mathbf{y}) \sim \mathcal{N}_{\mathbf{x}}(\boldsymbol{\mu}_y; \boldsymbol{\Sigma})$, which transforms the variational bound into

$$\tilde{I}(\mathbf{x}, \mathbf{y}) = -\frac{1}{2} \langle \text{tr} \{ \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_y) (\mathbf{x} - \boldsymbol{\mu}_y)^T \} \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})} + \frac{1}{2} \log |\boldsymbol{\Sigma}^{-1}| + \acute{c}. \quad (6.46)$$

Here \acute{c} incorporates $H(\mathbf{x})$ and other constants which do not affect the optimization surface for the encoder parameters, and $\boldsymbol{\Sigma} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|}$ is a function of parameters of the conditional $p(\mathbf{y}|\mathbf{x})$.

Clearly, if the log eigenspectrum of the inverse covariance of the decoder is constrained to satisfy

$$\sum_{i=1}^{|\mathbf{x}|} \log l_i(\boldsymbol{\Sigma}^{-1}) = \sum_{i=1}^{|\mathbf{x}|} \langle \log l_i(\mathbf{F}_{\mathbf{x}}) \rangle_{\tilde{p}(\mathbf{x})}, \quad (6.47)$$

where $\{l_i(\boldsymbol{\Sigma}^{-1})\}$ and $\{l_i(\mathbf{F}_{\mathbf{x}})\}$ are eigenvalues of $\boldsymbol{\Sigma}^{-1}$ and $\mathbf{F}_{\mathbf{x}}$ respectively, then the lower bound (6.46) reduces to the objective (6.23) amended with the average quadratic reconstruction error

$$\tilde{I}(\mathbf{x}, \mathbf{y}) = -\frac{1}{2} \underbrace{\langle \text{tr} \{ \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_y) (\mathbf{x} - \boldsymbol{\mu}_y)^T \} \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})}}_{\text{reconstruction error}} + \frac{1}{2} \underbrace{\langle \log |\mathbf{F}_{\mathbf{x}}| \rangle_{\tilde{p}(\mathbf{x})}}_{\text{Fisher criterion}} + \acute{c}. \quad (6.48)$$

Arguably, it is due to the subtraction of the non-negative quadratic term that (6.46) remains a general lower bound independently of the parameterization of the model and spectral properties of $\mathbf{F}_{\mathbf{x}} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|}$. We may find an analogous relation for the local approximation of mutual information $\tilde{I}(\mathbf{x}, \mathbf{y})$ by constraining the log eigenspectrum of $\boldsymbol{\Sigma}^{-1}$ to satisfy

$$\sum_{i=1}^{|\mathbf{x}|} \log l_i(\boldsymbol{\Sigma}^{-1}) = v_{\mathbf{x}|r} \sum_{i=1}^{|\mathbf{x}|} \langle l_i(\mathbf{F}_{\mathbf{x}}) \rangle_{\tilde{p}(\mathbf{x})}, \quad (6.49)$$

where $v_{\mathbf{x}|r}$ is a constant variance of the symmetric regions (see the discussion in Section 6.3.2). It is easy to see that in this case we get

$$\tilde{I}(\mathbf{x}, \mathbf{y}) = -\frac{1}{2} \underbrace{\langle \text{tr} \{ \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_y) (\mathbf{x} - \boldsymbol{\mu}_y)^T \} \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})}}_{\text{reconstruction error}} + \frac{1}{2} \underbrace{\langle \text{tr} \{ \mathbf{F}_{\mathbf{x}} \} \rangle_{\tilde{p}(\mathbf{x})}}_{\text{local criterion}} + \acute{c}. \quad (6.50)$$

One potential advantage of the variational approach over the Fisher approximation of mutual information is the fact that the optimized objective $\tilde{I}(\mathbf{x}, \mathbf{y})$ remains a proper bound on $I(\mathbf{x}, \mathbf{y})$ independently of the modeling assumptions. This contrasts with the approximations of Brunel and Nadal (1998) and the related methods of Kang and Sompolinsky (2001) and Hoch et al. (2003), where the accuracy is strongly influenced by the size of the code space. While the method based on local approximations of Szummer and Jaakkola (2002), Corduneanu and Jaakkola (2003) helps to avoid some of the constraints on the encoding mappings (specifically, it is applicable for $|\mathbf{y}| < |\mathbf{x}|$), it may be rather difficult to justify from the information-theoretic viewpoint (as by introducing the noise into the empirical distribution we decrease the information content between the codes and the original sources).

Another principal advantage of the variational approach to information maximization is the flexibility in the choice of the variational decoder. Intuitively, if the Fisher Information matrices are nearly singular, both (6.23) and (6.48) may be quite weak. However, by relaxing (6.47) and imposing full-rank constraints on the covariances of the variational decoders, the variational bound may be significantly strengthened. Moreover, as we showed in Section 6.2.1, by imposing specific constraints on the variational decoder, we can derive an iterative learning rule, which only requires local computations (such as evaluations of weighted Hebbian and anti-Hebbian terms). As discussed in Section 6.2.1, this may lead to an arguably more biologically plausible learning.

6.5 Demonstrations

Variational IM vs Fisher and Local Approximations

As we mentioned in Section 6.4, due to the fact that the Fisher and the local approximation criteria are not proper bounds on $I(\mathbf{x}, \mathbf{y})$, it is difficult to justify comparisons of the generic lower bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ with $\tilde{I}_F(\mathbf{x}, \mathbf{y})$ and $\tilde{I}_L(\mathbf{x}, \mathbf{y})$. Indeed, while we can generally justify comparisons of proper bounds (as a tighter bound would imply a smaller gap from the true unknown functional for a specific parameter subspace), it is not always easy to judge on accuracy of approximations of unknown underlying functionals. Generally, this complicates empirical comparisons of the described learning methods for large scale problems.

In order to gain intuition on how the approximate information-maximizing methods influence the true underlying objective $I(\mathbf{x}, \mathbf{y})$, we considered a low-scale problem (so that the mutual information $I(\mathbf{x}, \mathbf{y})$ could be computed exactly). We were particularly interested to see whether or not maximization of the generic lower bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ could lead to consistent improvements in the true mutual information $I(\mathbf{x}, \mathbf{y})$ for the considered conditionally factorized channel (6.1). To answer this question, we computed the exact value of $I(\mathbf{x}, \mathbf{y})$ at each iteration of the variational IM learning, where the encoder parameters $\mathbf{W} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$, $\mathbf{b} \in \mathbb{R}^{|\mathbf{y}|}$ were obtained by maximizing the bound $\tilde{I}(\mathbf{x}, \mathbf{y})$. We compared this with the changes in the exact mutual information under maximization of the Fisher approximation criterion $\tilde{I}_F(\mathbf{x}, \mathbf{y})$ of Brunel and Nadal (1998), and the local approx-

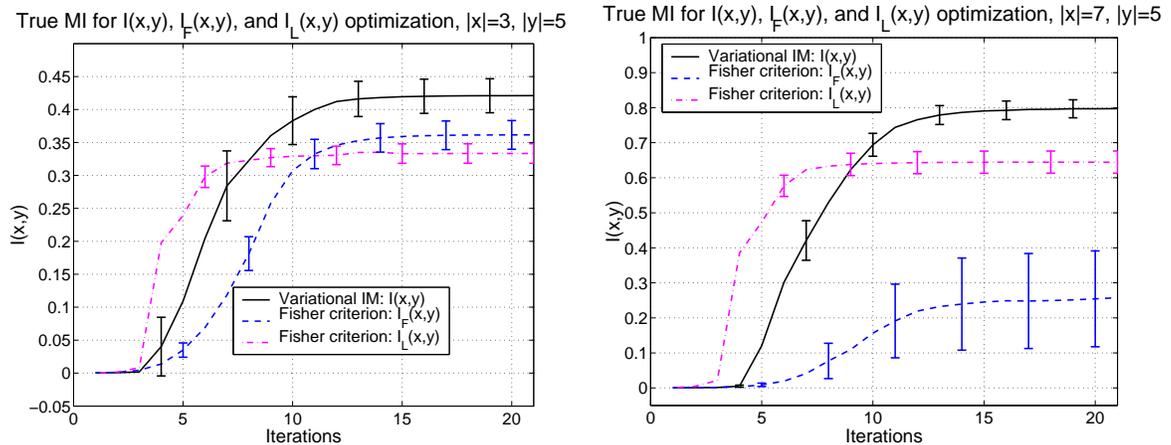


Figure 6.1: Changes in the exact mutual information $I(x, y)$ for parameters of the encoder $p(y|x)$ obtained by maximizing the variational lower bound $\tilde{I}(x, y)$, the Fisher information criterion $\tilde{I}_F(x, y)$, and the local approximation criterion $\tilde{I}_L(x, y)$. The number of training stimuli $M = 20$. The results are averaged over 30 runs from random parameter initializations. *Left*: $|x| = 3$, $|y| = 5$ *Right*: $|x| = 7$, $|y| = 5$.

imation criterion $\tilde{I}_L(x, y)$ motivated in Section 6.2.1 and inspired by Corduneanu and Jaakkola (2003). To ensure that the true mutual information $I(x, y)$ could indeed be computed exactly, we restricted the dimensionality of the response variables $|y|$. In order to illustrate the effects which the codesize $|y|$ could have on the performance, we considered the cases when $|y| < |x|$ and $|y| > |x|$.

Figure 6.1 illustrates changes in the exact mutual information $I(x, y)$ with iterations of the scaled conjugate gradients (SCG) optimization on $\tilde{I}(x, y)$, $\tilde{I}_F(x, y)$, and $\tilde{I}_L(x, y)$. The variational decoder $q(x|y)$ of the generic lower bound was chosen to be an isotropic linear Gaussian with the unconstrained optimal weights (see expression (6.11) and the discussion in Section 6.2.2). The plot shows the mean values and error bars on $I(x, y)$ for 30 runs from different initializations. The initial settings of the encoder parameters were the same for all the considered learning methods. The dimensionality of the code space was chosen to be $|y| = 5$. The training stimuli $\mathbf{x} \in \mathbb{R}^{|\mathbf{x}|}$ were sampled from $\mathbf{x} \sim \mathcal{N}_{\mathbf{x}}(\mathbf{0}, \mathbf{I}_{|\mathbf{x}|})$ and centered, where we considered $|\mathbf{x}| = 3$ (Figure 6.1 (*left*)) and $|\mathbf{x}| = 7$ (Figure 6.1 (*right*)). The results are shown for $M = 20$ training patterns. Moreover, we imposed additional parametric constraints on each synaptic weight, so that $w_{ij} = \sigma(a_{ij}) - 0.5 \in (-0.5, 0.5)$. (The gradients for this case could be easily obtained from (6.45)). The auxiliary parameters $\{a_{ij}\}$ were initialized at random as $a_{ij} \sim \mathcal{N}_a(0, 0.1)$.

As we can see from (Figure 6.1 (*left*)), all three methods tended to lead to consistent improvements in the true mutual information for $|\mathbf{x}| \leq |y|$. The variational approach usually resulted in higher values of $I(x, y)$ after just a few training iterations, and typically converged to higher values of $I(x, y)$. The local approximation criterion described in Section 6.2.1 may be characterized by a rapid convergence and low variance of the mutual information estimates, which may be explained by

the smoothness of the optimized surface defined by $\tilde{I}_L(\mathbf{x}, \mathbf{y})$. For the considered case of overcomplete representations, the Fisher approximation criterion $\tilde{I}_F(\mathbf{x}, \mathbf{y})$ typically resulted in a stable convergence to relatively high values of $I(\mathbf{x}, \mathbf{y})$.

In the case of undercomplete binary representations (i.e. $|\mathbf{x}| > |\mathbf{y}|$), optimization of the Fisher criterion $\tilde{I}_F(\mathbf{x}, \mathbf{y})$ was numerically unstable and often led to no visible improvements of mutual information $I(\mathbf{x}, \mathbf{y})$ over its starting values at initializations. This is not surprising, as the objective $\tilde{I}_F(\mathbf{x}, \mathbf{y})$ (and its gradients) are undefined for singular Fisher Information matrices (note that in the considered case $\text{rank}(\mathbf{F}_x) \leq \min\{|\mathbf{x}|, |\mathbf{y}|\} = |\mathbf{y}|$, i.e. $\mathbf{F}_x \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|}$ is rank-deficient). Further approximations aimed at handling singularities of the gradients (see Section 6.3.1) could lead to slight improvements in $I(\mathbf{x}, \mathbf{y})$, though their dynamics was rather inconsistent. In practice, this could often be characterized by non-monotonic changes in $I(\mathbf{x}, \mathbf{y})$ with the number of SCG iterations on $\tilde{I}_F(\mathbf{x}, \mathbf{y})$. Not surprisingly, after averaging over a number of independent runs, optimization of $\tilde{I}_F(\mathbf{x}, \mathbf{y})$ often led to high variances of the mutual information evaluations (Figure 6.1 (*right*)). In contrast, optimization of the local approximation criterion $\tilde{I}(\mathbf{x}, \mathbf{y})$ and the variational lower bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ typically led to consistent improvements in the exact mutual information. Again, while learning by maximizing the local approximations often resulted in quicker convergence, the variational IM applied to the considered channel usually led to consistently higher values of $I(\mathbf{x}, \mathbf{y})$ after a small number of training iterations.

Finally, we note that the imposed constraint on the individual weights $\{w_{ij}\}$ resulted in a hard upper bound on $\|\mathbf{W}\|_F \leq \omega \sqrt{|\mathbf{y}||\mathbf{x}|}$. The comparison of the methods is fair in the sense that the obtained results are not simply influenced by the weight re-scalings. Other types of constraints (e.g. soft constraints on the conditional variances of the stochastic firings) lead to qualitatively similar relations between the learning methods (see Agakov and Barber (2004b)). While the considered problem is low-scale, the observed relations between the learning methods help to understand the effects which optimization of the approximate objective criteria may have on the true mutual information. Specifically, the results confirm the intuitive limitations of the commonly used Fisher approximations of $I(\mathbf{x}, \mathbf{y})$.

Hypercubic Constraints and Discrete-Valued Weights

It is interesting to note that by imposing hypercubic constraints on the encoder weights, i.e. constraining each synaptic weight w_{ij} to lie in a symmetric region $[-\omega, \omega]$, and carefully re-scaling the input stimuli, we may often observe the situation when the optimal weight parameters $\mathbf{W} \in [-\omega, \omega]^{|\mathbf{y}| \times |\mathbf{x}|}$ obtained by maximizing any of the three objective criteria $\tilde{I}(\mathbf{x}, \mathbf{y})$, $\tilde{I}_F(\mathbf{x}, \mathbf{y})$, and $\tilde{I}_L(\mathbf{x}, \mathbf{y})$ lie very close to the corners of the hypercube. Not surprisingly, this behavior is strongly influenced by the cube's size ω , dimensionality of the source space $|\mathbf{x}|$, and the moments of the source distribution $\tilde{p}(\mathbf{x})$ (as all these entities influence the effective activation fields of the output units). Here we will not aim at characterizing the relation theoretically, but demonstrate the situation when the effect occurs.

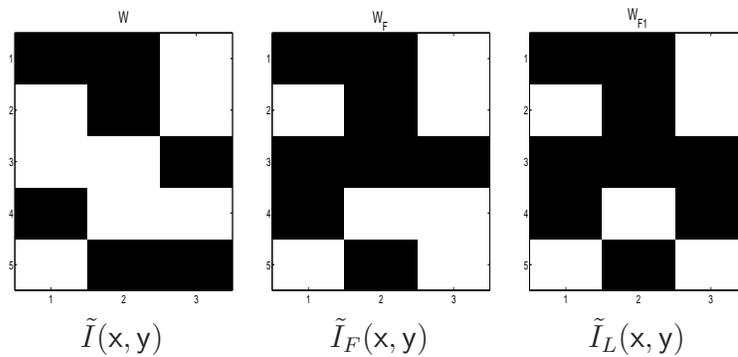


Figure 6.2: Optimal encoder weights $\mathbf{W} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$ obtained by maximizing $\tilde{I}(\mathbf{x}, \mathbf{y})$, $\tilde{I}_F(\mathbf{x}, \mathbf{y})$, and $\tilde{I}_L(\mathbf{x}, \mathbf{y})$. Here $|\mathbf{x}| = 3$, $|\mathbf{y}| = 5$, and $M = 20$. The white squares correspond to the settings $w_{ij} \approx 0.5$. The black squares correspond to the settings $w_{ij} \approx -0.5$. For the illustrated case, $\max_{ij} \{|w_{ij} - 0.5|\} \sim O(10^{-4})$.

Specifically, we will consider the settings of the previous set of experiments, where $|\mathbf{x}| = 3$, $|\mathbf{y}| = 5$, the weights are constrained as $w_{ij} = \sigma(a_{ij}) - 0.5 \in (-0.5, 0.5)$ for some real-valued parameters a_{ij} , and the sources are sampled as $\mathbf{x} \sim \mathcal{N}_{\mathbf{x}}(\mathbf{0}, 0.7\mathbf{I})$. Optimization of the three objectives with respect to the auxiliary parameters $\mathbf{A} = \{a_{ij}\} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$ results in the nearly discrete-valued weights $\mathbf{W} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$ shown on Figure 6.2, where $\max_{i,j} ||w_{i,j} - 0.5| \sim O(10^{-4})$. Figure 6.3 (*left*) shows the changes in the exact mutual information $I(\mathbf{x}, \mathbf{y})$ computed at the parameters obtained by maximizing $\tilde{I}(\mathbf{x}, \mathbf{y})$, $\tilde{I}_F(\mathbf{x}, \mathbf{y})$, and $\tilde{I}_L(\mathbf{x}, \mathbf{y})$. Note that the resulting weights have nearly identical norms, i.e. effects which the norm rescalings have on the objective criteria may safely be ignored. Figure 6.3 (*right*) shows typical changes in the exact mutual information $I(\mathbf{x}, \mathbf{y})$ computed for the discrete-valued weight parameters $\mathbf{W} \in \{-0.5, 0.5\}^{|\mathbf{y}| \times |\mathbf{x}|}$ around \mathbf{W}_I , \mathbf{W}_F , and \mathbf{W}_L (which are the optimal weights obtained by maximizing $\tilde{I}(\mathbf{x}, \mathbf{y})$, $\tilde{I}_F(\mathbf{x}, \mathbf{y})$, and $\tilde{I}_L(\mathbf{x}, \mathbf{y})$ respectively). As we see from the plot, the variational IM maximizing the bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ results in highest values of the exact $I(\mathbf{x}, \mathbf{y})$. (Additionally, it turns out that in the considered case the weight $\mathbf{W}_I \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$ gives rise to one of the equivalent global optima of the exact $I(\mathbf{x}, \mathbf{y})$ computed for the 2^{15} combinations of the discrete weights $\{-0.5, 0.5\}^{|\mathbf{y}| \times |\mathbf{x}|}$).

The result is potentially interesting, as it suggests the existence of encoder models which may be used to transform the combinatorial search over discrete-valued parameters to a continuous optimization problem. It is indeed intuitive that by maximizing the exact $I(\mathbf{x}, \mathbf{y})$, the optimal encoder model would encourage spread-out representations in the code space. However, as we showed in Section 6.3.3, for the considered parameterization of the encoder and the considered constraints on the parameters, the weights cannot grow indefinitely. This suggests a curious relation between dimensionality of the sources $|\mathbf{x}|$, the weight range defined by ω , the scalings of the source patterns, and the observed weight saturation effects, which will be interesting to investigate in the future.

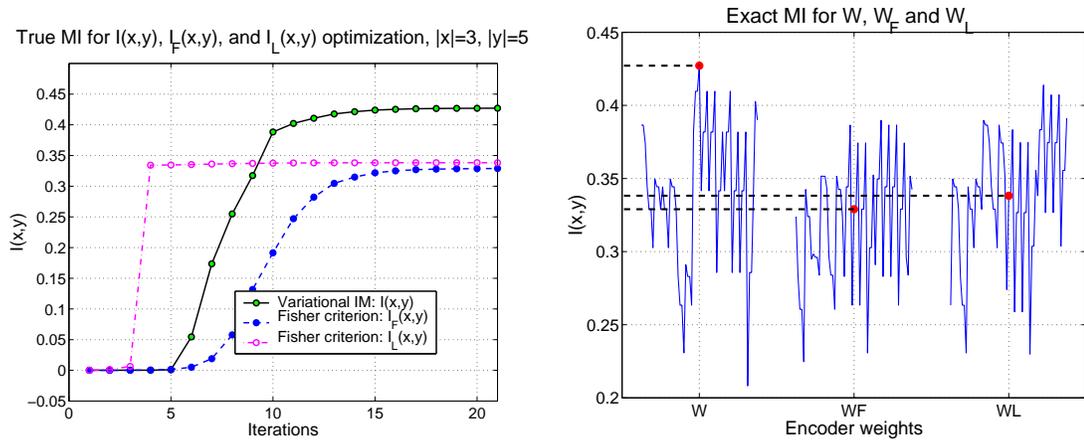


Figure 6.3: Exact mutual information as a function of (approximately) discrete-valued weights. *Left*: typical changes in the exact $I(x,y)$ under maximization of $\tilde{I}(x,y)$, $\tilde{I}_F(x,y)$, and $\tilde{I}_L(x,y)$ for the auxiliary parameters $\{a_{ij}\}$. *Right*: values of the exact mutual information computed at $W = \{-0.5, 0.5\}^{|y| \times |x|}$ around the optimal W_i , W_F , and W_L . The x -axis corresponds to the decimal representations of the binary encodings of the discrete weight vector $w \in \{-0.5, 0.5\}^{|x| \times |y|}$ obtained by flattening $W \in \mathbb{R}^{|y| \times x}$.

Variational IM: Stochastic Representations of the Digit Data

Finally, we have applied the variational IM method to stochastic coding of visual patterns for the case when the number of the input stimuli $|x|$ significantly exceeds the number of the spiking units $|y|$. (The problem of coding in this case may be viewed in the context of stochastic compression). Our motivation here was to check whether the variational approach to information maximization could lead to at least vaguely interpretable firing frequencies of small groups of neurons for the considered biologically-inspired channel parameterization; specifically, we were interested in checking whether frequency of stochastic firings of the encoding units could in some sense be representative of different aspects of the visual stimuli. (The stochasticity was an artifact of the constrained channel parameterization, as specified by the encoding distribution (6.1)). We were also interested to see how different firing frequencies could affect reconstructions of input stimuli obtained by applying the variational decoder.

Note that for the considered case the information-theoretic formulation is particularly attractive, as it allows to impose explicit biologically-inspired constraints on parameters of the encoding mappings. Additionally, in the exact formulation learning of encodings is entirely unsupervised, which is particularly attractive from the biological perspective. Our method is an approximation of the exact case, with a specific choice of the variational distribution. In our experiments we used the simplest form of the Gaussian decoder discussed in Section 6.2.2. After numerical optimization for $p(y|x)$ and $q(x|y)$ with an explicit constraint on the variance of the conditional firings (enforced by the penalty term $m = 0.2$, see Section 6.3.3), we applied the variational decoder to perform reconstruction of

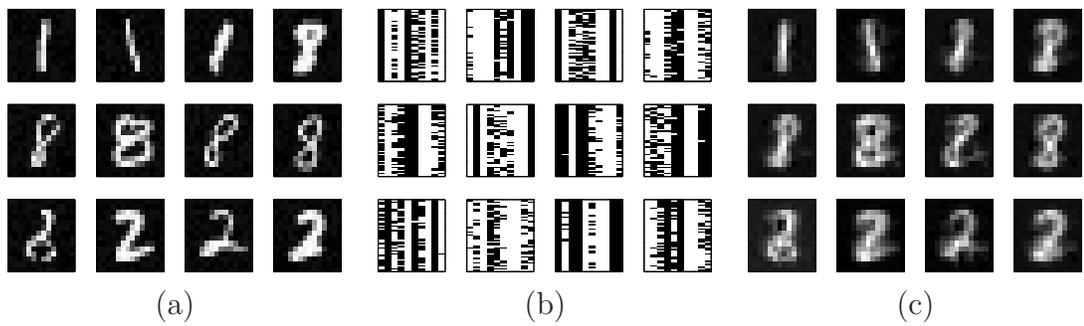


Figure 6.4: Stochastic binary encodings of the continuous input stimuli. (a): a subset of the original visual stimuli. (b): 30 samples from the corresponding encoding distribution $p(y|x)$ for $|y| = 10$ spiking units (black and white dashes correspond to silent and spiking output units). (c): Reconstructions of the original sources from 50 samples of neural spikes. Note that we imposed soft constraints on the variances of firings.

196-dimensional continuous visual stimuli from 10 spiking neurons (i.e. $|x| = 196$ and $|y| = 10$). The training stimuli consisted of 30 instances of digits 1, 2, and 8 (10 of each class). The source variables were reconstructed from 50 stochastic spikes at the mean of the optimal approximate decoder $q(x|y)$. Note that since $|x| > |y|$, the problem could not be efficiently addressed by optimizing the Fisher approximation criterion (6.25). Clearly, the deterministic approaches (see e.g. Bell and Sejnowski (1995) or Shriki et al. (2002)) are not applicable either, since the considered channel is noisy and undercomplete. On the other hand, the variational IM algorithm is applicable and numerically stable.

Figure 6.4 illustrates a subset of the original source signals (a), samples of the corresponding binary responses for units (b), and reconstructions of the source data from the population of binary spikes (c). (For Figure 6.4 (b), the x -axis corresponds to each of the 10 units, and the y -axis corresponds to each of 30 samples from the encoding distribution). We see that while stochastic encodings corresponding to distinct patterns are generally different, there are some visibly distinct input stimuli characterized by very similar firing patterns (see e.g. the rightmost patterns in the first and the second rows), and it is arguably the frequency of firings of a small number of neurons which accounts for differences in the corresponding inputs and reconstructions. For example, the firing patterns of the rightmost stimuli in the first and the second rows are quite similar, apart from the noisy firings of the 2nd, 9th and the 10th units. However, the resulting reconstructions are visibly different (see the sizes of the bottom loops of the digits “8”). The result is consistent with observations that receptive fields of post-synaptic cells in the neocortex are influenced by distinct features of input stimuli, which in turn determines firing frequencies of individual neurons (e.g. Markram et al. (1998)).

Table 6.1: Objective functions for approximate information maximization

1. <i>Invertible channels:</i>	$I(\mathbf{x}, \mathbf{y}) = \langle \log \mathbf{J}_{\mathbf{x}} \rangle_{\tilde{p}(\mathbf{x})}$
2. <i>Overcomplete deterministic channels:</i>	$I(\mathbf{x}, \mathbf{y}) = \langle \log \mathbf{J}_{\mathbf{x}}^T \mathbf{J}_{\mathbf{x}} \rangle_{\tilde{p}(\mathbf{x})}$
3. Fisher approximation:	$\tilde{I}_F(\mathbf{x}, \mathbf{y}) = \langle \log \mathbf{F}_{\mathbf{x}} \rangle_{\tilde{p}(\mathbf{x})}$
4. Local approximation:	$\tilde{I}_F(\mathbf{x}, \mathbf{y}) = \langle \log \exp\{\text{tr}\{\mathbf{F}_{\mathbf{x}}\}\} \rangle_{\tilde{p}(\mathbf{x})}$
5. Variational lower bound:	$\tilde{I}(\mathbf{x}, \mathbf{y}) = \langle \log q(\mathbf{x} \mathbf{y}) \rangle_{p(\mathbf{y} \mathbf{x})\tilde{p}(\mathbf{x})}$

6.6 Summary

The primary goal of this chapter was to explore applicability of the variational information-maximizing framework in the context of learning *high-dimensional* binary representations of continuous source patterns. We described an application of the variational approach to information maximization for the case when continuous source stimuli were represented by conditionally independent stochastic binary responses. We also showed that for the considered encoding distribution it was possible to derive a local iterative learning rule, which gives rise to a form of weighted variable-rate Hebbian learning (interestingly, the receptive fields of the post-synaptic units were influenced not only by the activations at the pre-synaptic layers, but also (implicitly) by the activations of the neighboring post-synaptic units). This result generalizes the work of Linsker (1997), who derived a local approximation of Bell and Sejnowski’s information-maximizing rule for deterministic invertible channels (Bell and Sejnowski (1995)).

Moreover, we have applied the numerical approximations used for deriving the local regularizers of Szummer and Jaakkola (2002), Corduneanu and Jaakkola (2003), and showed that the results may be used to approximate a lower bound on $I(\mathbf{x}, \mathbf{y})$. While the resulting approximation typically has better convergence properties than the common Fisher approximation of mutual information (see e.g. Brunel and Nadal (1998)), our empirical results indicate that the variational approach may be more attractive. Additionally, our results indicate that the considered methods addressing approximate information maximization may be viewed as approximations of our variational approach; however, generally they do not preserve a proper bound on the mutual information, and may be less computationally and numerically appealing.

Table 6.1 summarizes effective criteria optimized by Bell and Sejnowski (1995), Shriki et al. (2002), Brunel and Nadal (1998), the local approximation based on the work of Szummer and Jaakkola (2002) and Corduneanu and Jaakkola (2003), and our generic variational approach. Here the symmetric matrix $\mathbf{J}_{\mathbf{x}} = \{J_{ij}(\mathbf{x})\} \stackrel{\text{def}}{=} \{\partial y_i(\mathbf{x})/\partial x_j\} \in \mathbb{R}^{|\mathbf{x}|\times|\mathbf{x}|}$ is the Jacobian of the deterministic invertible mapping $\mathbf{x} \mapsto \mathbf{y}$ (with $|\mathbf{y}| = |\mathbf{x}|$), $\mathbf{F}_{\mathbf{x}} \in \mathbb{R}^{|\mathbf{x}|\times|\mathbf{x}|}$ is the Fisher Information matrix (6.19), and $q(\mathbf{x}|\mathbf{y})$ is an approximate decoder lying in a tractable family. In this chapter we focused primarily on the discussion of the last three approaches (shown in **bold**), since they are applicable for approximate maximization of mutual information in stochastic, rather than deterministic, channels. For this reason we have ignored the comparisons of our method with those of Bell and Sejnowski (1995) and Shriki et al. (2002), which presume noiseless mappings from the code to the data space.

Note that in contrast to the other techniques, the suggested variational method optimizes a proper lower bound independently of the choice of the decoder, dimensionality of the input stimuli, number of post-synaptic neurons, or noise of the stochastic firings. This extends applicability of the variational approach to dimensionality reduction, compression, syndrome decoding, and facilitates applications of the method to undercomplete and overcomplete stochastic coding. Of course, more biologically realistic channels and applications should potentially be considered; our results here mainly serve to illustrate the potential advantages of our variational information-maximizing formulation over the common (Brunel and Nadal (1998), Kang and Sompolinsky (2001)) and less common (Szummer and Jaakkola (2002), Corduneanu and Jaakkola (2003) and Section 6.2.1) approaches to population coding of high-dimensional input stimuli.

Chapter 7

Auxiliary Variational Inference and Variational Mutual Information Maximization

In Chapter 2 we discussed a simple and general variational approach to maximizing the generic lower bound on $I(\mathbf{x}, \mathbf{y})$ for a stochastic channel. In the subsequent chapters we considered generalizations and applications of the variational information-maximizing approach (see Chapter 4, Chapter 5, and Chapter 6), and discussed the relation of the variational IM to conditional likelihood training (Chapter 3). Here we change the perspective rapidly and demonstrate that the generic bound on the mutual information may be obtained in the context of *auxiliary variational statistical inference* (Agakov and Barber (2004a)), where we aim to lower-bound a generally intractable normalizing constant of a Markov network¹. Generally, the approach we describe here defines a completely different family of variational bounds; however, maximization of the generic lower bound on mutual information (2.2) may be seen as its fundamentally important subgoal.

While little work appears to have been done on developing a simple and easily generalizable variational framework for approximate information maximization, variational methods have proved popular and effective for inference and maximum-likelihood learning in intractable graphical models. In this context they are often applied for bounding the likelihoods and the normalizing constants (see e.g. Jaakkola (1997), Jaakkola and Jordan (1998), Barber and Wiergerinck (1998), Bishop et al. (1998), Lawrence (2000), Wainwright et al. (2001), Wainwright et al. (2002), Beal (2003)). The majority of the standard approaches to lower-bounding the normalizing constants are based on non-negativity of the Kullback-Leibler divergence $KL(q(\mathbf{x})||p(\mathbf{x}))$ between a tractable variational distribution $q(\mathbf{x})$ and the original distribution $p(\mathbf{x})$ (see e.g. Jordan et al. (1998) for a general introduction to variational inference and learning). Here we show that by expressing the bound on the normalizing constant from the Kullback-Leibler divergence $KL(q(\mathbf{x}, \mathbf{y})||p(\mathbf{x}, \mathbf{y}))$ in a specifically defined augmented variable space $\{\mathbf{x}, \mathbf{y}\}$, we may improve the standard bounds. It turns out that the improvement

¹We can use a fundamentally similar method for inference and learning in arbitrary distributions, but focus on Markov networks for clarity.

of the proposed *auxiliary variational* method over a convex combination of simple variational bounds is given by a specific form of the generic lower bound on mutual information $I(\mathbf{x}, \mathbf{y})$ derived in Chapter 2 (see expression (2.2)).

In the context of approximate probabilistic inference, we may use the suggested approach for improving on the lower bounds of standard factorized approximations. Indeed, we show that the method described here forms a more powerful class of approximations than any *structured mean field* technique. The existing lower bounds of the variational mixture models (Lawrence et al. (1998), Jaakkola and Jordan (1998)) can be viewed as computationally expensive special cases of our method. A byproduct of our work is an efficient way to calculate a set of mixture coefficients for any set of tractable distributions, which principally improves on the flat combination (Agakov and Barber (2004a)).

7.1 Introduction

Probabilistic graphical models provide a convenient framework for graphical representation of joint probability distributions, and facilitate computation of many quantities of interest required for both inference and learning. Generally, probabilistic treatment of uncertainty offers a consistent and principled framework for inference in complex domains (Chapter 1). However, many distributions used for modeling practical domains are inherently intractable, which motivates the need for accurate and efficient approximations. In this chapter we focus on approximate inference, specifically on computation of lower bounds on normalization constants of undirected graphical models, which can also be used to approximate marginals of a formally intractable distribution. Fundamentally similar methods can be applied for inference and learning in arbitrary distributions; we will focus on inference in undirected models for clarity.

To be explicit, we will consider distributions $p(\mathbf{x})$ of the (Boltzmann) form

$$p(\mathbf{x}) = \exp\{-E(\mathbf{x})\}/Z, \quad Z = \sum_{\mathbf{x}} \exp\{-E(\mathbf{x})\}. \quad (7.1)$$

The complexity of carrying out the summation required to compute Z depends on the graphical structure of $p(\mathbf{x})$. In (poly)trees the normalisation constant can be computed in time linear in the number of variables in the distribution. In a general graph, the time is exponential in the size of the largest clique in the associated junction tree (Lauritzen and Spiegelhalter (1988), Jordan (2005)). Here we are interested in cases where exact computation by the junction tree algorithm is infeasible.

Variational approximations have been extensively used in physics and engineering and more recently applied to approximate inference and learning in intractable graphical models (see e.g. Saul et al. (1996), Jordan et al. (1998); Barber and Wiergerinck (1998); Wainwright et al. (2002)). In this context they were shown to result in relatively simple representations of the induced optimization problems; at the same time, it was shown that they often led to accurate approximations of the generally intractable quantities of interest. Their other

advantage for graphical models is availability of rigorous bounds on the normalizing constant (e.g. Jaakkola and Jordan (1996)), which contrasts with other (e.g. Monte-Carlo) approximations (see e.g. Neal (1993)). Availability of such bounds on the underlying intractable objectives facilitates comparisons of variational approximation techniques, where the tightness of a bound on an unknown quantity is often used as a relative measure of performance of an approximate method.

A popular class of lower bounds on $\log Z$ is based on the non-negativity of the Kullback-Leibler divergence

$$KL(q(\mathbf{x})\|p(\mathbf{x})) = \langle \log q(\mathbf{x}) \rangle_{q(\mathbf{x})} - \langle \log p(\mathbf{x}) \rangle_{q(\mathbf{x})} \geq 0. \quad (7.2)$$

Here $\langle \dots \rangle_{q(\mathbf{x})}$ denotes an average over $q(\mathbf{x})$, and the bound is saturated if and only if $q(\mathbf{x}) \equiv p(\mathbf{x})$. In the case of the Boltzmann distribution (7.1), non-negativity of the KL divergence (7.2) yields the well-known class of lower bounds

$$\log Z \geq \langle -\log q(\mathbf{x}) \rangle_{q(\mathbf{x})} - \langle E(\mathbf{x}) \rangle_{q(\mathbf{x})}, \quad (7.3)$$

where $q(\mathbf{x})$ is typically restricted to a class of tractable distributions and varied to obtain the tightest bound within the tractable class. Coupled with an upper bound on the normalizing constant (Wainwright et al., 2002), expression (7.3) may be used for bounding expectations in the original model. A further use for this procedure is to provide a lower bound on the marginal likelihood in situations of observed and unobserved variables, which is a natural derivation and extension (Neal and Hinton, 1998) of the EM procedure (Dempster et al. (1977)).

7.1.1 Existing Variational Approximations

The computational tractability of the bound (7.3) depends on the choice of the approximating distribution $q(\mathbf{x})$. Possibly one of the simplest choices is given by the factorized *mean field* model (see Figure 7.1 (a), (b)) with $q_{MF}(\mathbf{x}) = \prod_i q(x_i)$, which discards all the edges from the original graph. The factorized assumption often results in a tractably computable bound (7.3). However, the bound may be inaccurate when the variables in the true distribution are strongly correlated. Moreover, it can be shown that the mean field distribution is log-concave (in the space of functional parameters), which implies a fundamental limitation of the mean field approximation in the case when $p(\mathbf{x})$ is locally multi-modal, since significant mass contributing to the partition function may be missed.

One way to go beyond the factorized assumption for $q(\mathbf{x})$ is to consider a *structured mean field* approximation (Ghahramani and Jordan, 1995; Barber and Wiegerinck, 1998) which retains some of the structure of $p(\mathbf{x})$. In this case it is often assumed that $q(\mathbf{x})$ has a sparse graphical representation (e.g. it is a (poly)tree, see Figure 7.1 (c)), which typically leads to an improvement on the bound at a moderate increase in computational cost. However, we may note that discarded edges in $q(\mathbf{x})$ introduce conditional independencies which may not exist in the original distribution $p(\mathbf{x})$.

Other approaches (Lawrence et al., 1998; Jaakkola and Jordan, 1998) bound the log partition function by considering mixtures of mean field type models

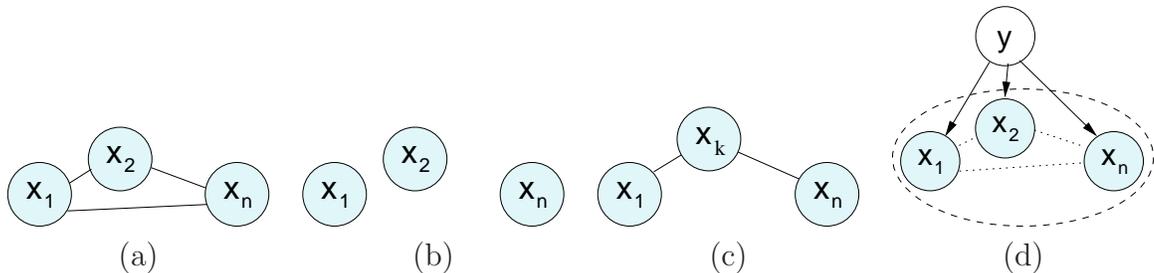


Figure 7.1: Undirected models and their approximations (a schematic plot). (a) A fully connected pairwise Markov network representing the intractable distribution $p(\mathbf{x})$; (b) a standard mean field approximation $q_{MF}(\mathbf{x})$; (c) a structured mean field approximation $q_{SMF}(\mathbf{x})$; (d) a mixture of mean field models. (All the variables \mathbf{x} are coupled through the mixture label y . The dotted lines serve to indicate that the marginal $q_{MMF}(\mathbf{x})$ expressed from $q_{MMF}(\mathbf{x}, y)$ is in general fully connected). The transparent node y shows the auxiliary variable. The shaded nodes indicate the variables forming the space $\{\mathbf{x}\}$ of the original model $p(\mathbf{x})$.

(see Figure 7.1 (d)). This formally extends the standard factorized approximations, since the resulting approximating distribution $q(\mathbf{x})$ is generally multi-modal and not factorized in \mathbf{x} . However, in general optimization of the bound (7.3) in this case requires minimization of the KL divergence between two fully connected distributions $p(\mathbf{x})$ and $q(\mathbf{x})$, which requires a formally computationally intractable evaluation of the entropy of the mixture $H(\mathbf{x})$. Further approximations suggested by Lawrence et al. (1998), Jaakkola and Jordan (1998), El-Hay and Friedman (2002) circumvent the computational intractability of variational mixture approximations by relaxing (7.3) at the cost of introducing additional variational relaxations. Unfortunately, unless *all* the mixture components $q(\mathbf{x}|y)$ have the same tractable structure, optimization of the resulting bound on $\log Z$ may become numerically unstable, which may complicate generalizations of the existing results to variational mixtures of *arbitrarily structured* tractable experts. Additionally, it is difficult to formally analyze the optimal solution induced by the resulting optimization procedure, and the existing procedure is not easily generalizable to richer families of variational bounds.

7.2 Auxiliary Variational Method

We will now explore the idea of using augmented variable spaces in the context of bounding the normalizing constant, and show that the formulation gives rise to the generic lower bound on mutual information.

Note that the computational problems of the variational mixture approximations arise from the fact that marginalization of the mixture labels y from the joint distribution $q(\mathbf{x}, y)$ results in a fully connected marginal $q(\mathbf{x})$ (see Figure 7.1 (d)). Informally, computation of the bound on $\log Z$ in this case requires minimization of the KL divergence between two fully-connected distributions (see expression (7.2)). It is intuitive that from the computational viewpoint it would

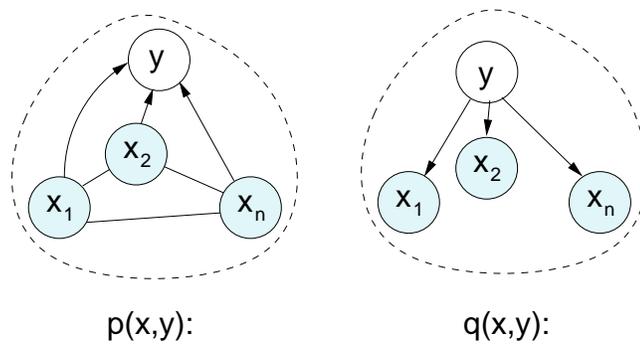


Figure 7.2: An auxiliary mean field approximation. The target distribution $p(x, y)$ is approximated by $q(x, y)$, which is structured in the *augmented space*. Note that the mapping $p(y|x)$ was introduced in a way which does not affect the original fully connected pairwise target distribution $p(x)$. The transparent node y shows the auxiliary variable. The shaded nodes indicate the variables forming the space $\{x\}$ of the original model $p(x)$.

be significantly more beneficial to retain the structural form of the joint distribution and use $q(x, y)$ as an approximation. The variational distribution would in this case be defined over the augmented variable space $\{x, y\}$, while the target $p(x)$ is defined over the space of the original variables $\{x\}$. In order for the bound (7.3) to be well-defined, we introduce *auxiliary variables* y to the target distribution. This can be readily done in such a way that the marginal $\tilde{p}(x)$ of the joint distribution $p(x, y)$ has the same graphical structure as the original target $p(x)$ (see Figure 7.2). Then we minimize the KL divergence between $q(x, y)$ and $p(x, y)$ in the joint variable spaces. Note that in contrast to standard structured approximations, all the variables x of the marginal $q(x)$ remain coupled, which agrees with the graphical structure of the fully connected Markov network $p(x)$. However, similarly to structured mean field techniques, our *auxiliary variational* method does not require a direct evaluation of the computationally intractable entropy of the mixture $H_q(x)$, as by minimizing the KL divergence in the augmented variable space we would effectively be fitting a structured distribution. Note that the idea of introducing “dummy” variables y is conceptually similar to the auxiliary variational bound on the mutual information discussed in Section 2.3, though the projections and the bounds are defined differently.

Another motivation for introducing the auxiliary variables is the reported success of auxiliary sampling techniques, such as the *Swendsen-Wang* (Swendsen and Wang, 1987), *partial decoupling* (Higdon, 1998) algorithms. It has been shown that by augmenting the original variable space with auxiliary variables and drawing samples from the joint distribution $p(x, y)$ in the augmented space, one can achieve a significant improvement over standard MCMC approaches. The purpose of the auxiliary variables in this context is to capture (structural) information about clusters of strongly correlated variables. Our hope was that by performing approximations in the augmented space, where the auxiliary variables capture useful regularities about the data, we may improve over standard

variational approaches.

Indeed, we demonstrate that the auxiliary variational technique forms a more powerful class of approximations than any structured mean field approach. Moreover, the method offers an efficient way of calculating a set of mixture coefficients for any choice of tractable approximators (for example, trees with different structures). These coefficients may be used to form a mixture which is principally better than a single best tractable approximator or their flat combination. Finally, we show that our approach provides a computationally and conceptually simple alternative to the existing bounds optimized by the variational mixture methods.

7.2.1 Optimizing the Auxiliary Variational Bound

Let $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y}|\mathbf{x})$ define the joint distribution of the original variables \mathbf{x} and auxiliary variables \mathbf{y} in the augmented $\{\mathbf{x}, \mathbf{y}\}$ space. From the divergence $KL(q(\mathbf{x}, \mathbf{y})||p(\mathbf{x}, \mathbf{y}))$ in the joint space it is easy to obtain an expression for the lower bound on the log partition function of $p(\mathbf{x}) = \exp\{-E(\mathbf{x})\}/Z$, which is given by

$$\log Z \geq -\langle \log q(\mathbf{x}, \mathbf{y}) \rangle_{q(\mathbf{x}, \mathbf{y})} - \langle E(\mathbf{x}) \rangle_{q(\mathbf{x})} + \langle \log p(\mathbf{y}|\mathbf{x}) \rangle_{q(\mathbf{x}, \mathbf{y})}. \quad (7.4)$$

Here $p(\mathbf{y}|\mathbf{x}) \stackrel{\text{def}}{=} p(\mathbf{y}|\mathbf{x}; \Psi)$ is an *auxiliary conditional* distribution parameterized by Ψ . (In principle, we do not need to impose parametric constraints on $p(\mathbf{y}|\mathbf{x})$; for sparse discrete models the auxiliary conditional may be defined by the conditional probability table). Equivalently, (7.4) may be written as

$$\begin{aligned} \log Z &\geq -\langle KL(q(\mathbf{x}|\mathbf{y})||p(\mathbf{x})) \rangle_{q(\mathbf{y})} + \tilde{I}(\mathbf{x}, \mathbf{y}) \\ &= \sum_{\mathbf{y}} q(\mathbf{y}) [\langle -E(\mathbf{x}) - \log q(\mathbf{x}|\mathbf{y}) \rangle_{q(\mathbf{x}|\mathbf{y})}] + \tilde{I}(\mathbf{x}, \mathbf{y}), \end{aligned} \quad (7.5)$$

where $\tilde{I}(\mathbf{x}, \mathbf{y})$ is defined as

$$\tilde{I}(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} \sum_{\mathbf{x}} \sum_{\mathbf{y}} q(\mathbf{x}, \mathbf{y}) \log \frac{p(\mathbf{y}|\mathbf{x})}{q(\mathbf{y})} = \langle \log p(\mathbf{y}|\mathbf{x}) \rangle_{q(\mathbf{x}, \mathbf{y})} + H(\mathbf{y}). \quad (7.6)$$

Note that up to the notational invariance, (7.6) is exactly the generic lower bound on mutual information $I(\mathbf{x}, \mathbf{y})$ expressed for the encoder model $q(\mathbf{x}, \mathbf{y}) = q(\mathbf{y})q(\mathbf{x}|\mathbf{y})$ (see expression (2.2)). In this context, the mapping to the auxiliary space $p(\mathbf{y}|\mathbf{x})$ may be interpreted as the variational decoder of the IM framework. The leftmost term in the r.h.s. of expression (7.6) defines a convex summation of the standard lower bounds (7.3) on $\log Z$ for the set of $\{\mathbf{x}\}$ -variables, which cannot improve on the single best bound in the set. Clearly, (7.5) may improve on the standard bounds only if $\tilde{I}(\mathbf{x}, \mathbf{y}) > 0$. In the case that the auxiliary variables \mathbf{y} contain no information about \mathbf{x} , i.e. $p(\mathbf{y}|\mathbf{x}) = p(\mathbf{y})$, it is straightforward to show that the method reproduces the standard variational bound which uses the single best approximation $q(\mathbf{x}|\mathbf{y})$. It is intuitive that by considering less trivial auxiliary conditional mappings we may improve on the generic lower bound $\tilde{I}(\mathbf{x}, \mathbf{y})$.

If $p(y|x) \equiv q(y|x)$ then $\tilde{I}(x, y)$ defines the exact mutual information between the original variables \mathbf{x} and the auxiliary variables \mathbf{y} . However, this specification leads to computational difficulties of evaluating or bounding the intractable entropy of the mixture $H(y)$ (see Section 1.4). By analogy with the variational IM method, we may constrain $p(y|x)$ so that it lies in a tractable family (we will denote the constrained distribution as $p(y|x; \Psi)$, where Ψ is a set of parameters; structural constraints may be considered similarly). Then a rigorous variational approach to maximizing $\log Z$ would involve maximizing the bound (7.6) with respect to the parameters (or clique potentials) of the auxiliary conditional distribution $p(y|x, \Psi)$ (which corresponds to the variational decoder in the IM terminology), the marginal $q(y)$, and the conditionals $q(x|y)$. The general iterative optimization algorithm in this case is given as follows:

1. Choose the auxiliary conditional $p(y|x)$. For the remainder, we choose

$$p(y|x) = \exp\{\Psi(y; x)\} / Z_{y|x}, \quad Z_{y|x} = \sum_y \exp\{\Psi(y; x)\}, \quad (7.7)$$

though more general distributions may potentially be considered.

2. Initialize $q(x|y)$, $q(y)$, and parameters Ψ of $p(y|x)$.
3. For the fixed $q(y, x)$, obtain Ψ^{new} by solving for zeros of

$$\partial \log Z / \partial \Psi = \langle \partial \log p(y|x) / \partial \Psi \rangle_{q(x,y)}, \quad (7.8)$$

or performing numerical ascent on $\log Z$ for Ψ (see the bound (7.6)). Clearly, an unconstrained optimization would result in $p(y|x) = q(y|x)$, as this would maximize the generic bound $\tilde{I}(x, y)$.

4. For the fixed $p^{new}(y, x) \stackrel{\text{def}}{=} p(x)p(y|x; \Psi^{new})$ and $q(y)$ set

$$q^{new}(x|y) \propto p^{new}(y, x) \quad (7.9)$$

for all instances \mathbf{y} .

5. For the fixed $p^{new}(y, x)$ and $q^{new}(x|y)$, set

$$q^{new}(y) \propto \exp \left\{ - \sum_x q^{new}(x|y) \log \frac{q^{new}(x|y)}{p(x)p(y|x; \Psi^{new})} \right\}. \quad (7.10)$$

6. Iterate steps 3–5 until a termination criterion is met.

Note that in the case of a constrained auxiliary conditional $p(y|x)$, step 3 is analogous to the M-step of the generalized EM algorithm (Neal and Hinton, 1998), while steps 4 and 5 are analogous to the E-step. An update rule for each term was obtained from (7.6) by computing the corresponding functional derivatives while keeping other terms fixed.

Up to this point, the results are general and applicable (at least, in theory) for arbitrarily parameterized distributions. However, as discussed in Section 1.4, for computational reasons it is fundamentally important to impose constraints on the auxiliary conditional $p(\mathbf{y}|\mathbf{x})$ and the approximate joint $q(\mathbf{x}, \mathbf{y})$ (effectively, both distributions define a set of variational parameters). If both distributions are in a tractable family (for example, if each node y_i has a small number of \mathbf{x} -parents), the bound (7.6) may be computed and optimized exactly (see the discussion in Section 2.2.1). Another case leading to exact computations is when $p(\mathbf{y}|\mathbf{x}) \sim \mathcal{N}(\Psi\mathbf{x}, \Sigma)$, which would correspond to Linsker’s *as-if Gaussian* bound on mutual information (see Section 2.2.2). Note that fundamentally the variational bound on $\log Z$ is tractable only in cases when $\tilde{I}(\mathbf{x}, \mathbf{y})$ may be exactly computed and optimized. Optimization of the remaining terms in the bound (7.6) should not be problematic (at least, this is the case when the variational components $q(\mathbf{x}|\mathbf{y})$ are structured and the underlying distribution $p(\mathbf{x})$ is a pairwise Markov network $p(\mathbf{x}) \propto \exp\{-(\mathbf{x}^T \mathbf{W} \mathbf{x} + \mathbf{x}^T \mathbf{b})\}$ for some $\mathbf{W} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|}$, $\mathbf{b} \in \mathbb{R}^{|\mathbf{x}|}$).

If one wishes to have a large number of parental variables \mathbf{x} influencing \mathbf{y} , further approximations may need to be employed. For the special case of distributions of the form (7.7), we can utilize the standard linear upper bound $\log x \leq mx - \log m - 1$. This transforms the objective (7.6) to

$$\log Z \geq 1 + \sum_{\mathbf{y}} q(\mathbf{y}) \langle -E(\mathbf{x}) - \log q(\mathbf{x}, \mathbf{y}) \rangle_{q(\mathbf{x}|\mathbf{y})} + \langle \Psi(\mathbf{y}; \mathbf{x}) + \mu(\mathbf{x}; \mathbf{y}) - e^{\mu(\mathbf{x}; \mathbf{y})} Z_{\mathbf{y}|\mathbf{x}} \rangle_{q(\mathbf{x}, \mathbf{y})} \quad (7.11)$$

where $e^{\mu(\mathbf{x}; \mathbf{y})}$ is an additional variational functional of the exponential form (see e.g. Jaakkola and Jordan (1998)). In general, optimization of the bound (7.11) is numerically unstable and computationally expensive. However, by analogy with the general formulation of the IM framework, we may avoid computational difficulties of computing the bound (7.6) by constraining $p(\mathbf{y}|\mathbf{x}, \Psi)$ to lie in a tractable family. Generally, such constraints would obviate a recourse to (7.11).

7.2.2 Specific Auxiliary Representations

Effectively, the problem of choosing an appropriate (tractable yet general) mapping to the auxiliary space is strongly related to the problem of choosing an appropriate variational decoder (see e.g. Sections 1.5, 2.2.1). Generally, any kind of tractable constraints on $p(\mathbf{y}|\mathbf{x})$ may be used, such as the Gaussian or a factorized approximation (see e.g. the discussion in Section 2.2.1). Additionally, we note that by analogy with Section 2.3 we may use the auxiliary variational lower bound on mutual information² (2.38). Here we will briefly outline other choices of the auxiliary conditional $p(\mathbf{y}|\mathbf{x})$, which in practice often lead to accurate approximations.

²This should not be confused with the auxiliary variational lower bound on $\log Z$ (expression (7.5)).

Discrete Spaces with Parametric Auxiliary Distributions

If the auxiliary space is given by a single multinomial variable $y \in \{1, \dots, M\}$, a natural choice for $p(y|\mathbf{x})$ is to use a *softmax* type representation

$$p(y_k|\mathbf{x}) \propto \exp \{f(\mathbf{x}; \mathbf{u}^{(k)})\}, \quad \mathbf{U} = \{\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(M)}\} \in \mathbb{R}^{|\mathbf{x}| \times M}. \quad (7.12)$$

Here $f(\mathbf{x}; \mathbf{U})$ is some differentiable function and $p(y_k|\mathbf{x})$ is the probability of the auxiliary variable y being in state $k = 1, \dots, M$. Although the expectation $\langle \log p(y_k|\mathbf{x}) \rangle_{q(\mathbf{x}|y_k)}$ is in general intractable, we can make a use of the transformed bound (7.11). A significantly cheaper alternative is to perform optimization of (7.6) by approximating the term at $\log p(y_k|\langle \mathbf{x} \rangle_{q(\mathbf{x}|y_k)})$, which is exact as long as the mapping from the auxiliary to the data space is deterministic, i.e. $q(\mathbf{x}|y) \sim \delta(\mathbf{x} - \langle \mathbf{x} | y \rangle)$. (Note that in the IM terminology this choice of the variational component distribution would corresponds to a noiseless encoding model, which may greatly simplify the computations). Generally, this approximation may be accurate as long as each component $q(\mathbf{x}|y)$ is sharply peaked around the mode.

Additionally, in the case when $|\mathbf{x}|$ is large, $p(y|\mathbf{x}) = \prod_k p(y_k|\mathbf{x})$, and each factor utilizes a simple (generalized) linear type dependency

$$p(y_k|\mathbf{x}) = p(y_k|a_k), \quad a_k \stackrel{\text{def}}{=} f_k(\mathbf{x}^T \mathbf{u}^{(k)} + b^{(k)}), \quad (7.13)$$

the *Gaussian field* approximation (Barber and Sollich (2000)) can be employed. From the Central Limit Theorem we can assume approximate Gaussianity of the field a_k and approximate the expectation $\langle \log p(y_i|\mathbf{x}^i) \rangle_{q(\mathbf{x}|y_k)}$ by performing 1-D Gaussian integration of the general form $\int_a f(a)p(a)$, where $p(a) \sim \mathcal{N}(\mu_a, \sigma_a^2)$. The means and variances of the fields are readily relatable to the first and second order moments of $q(\mathbf{x}|y_k)$.

In practice, such approximations do not lead to significant deviations and are shown to be both accurate and efficient (Barber and Sollich (2000), Agakov and Barber (2003)). Probably the greatest disadvantage of these relaxations is due to the fact that for any realistic limit of $|\mathbf{x}|$ the bound will no longer be strict. One way to address this problem is to impose additional structural constraints on the conditionals by limiting the number of parental variables for each factor.

Discrete Spaces with Structured Auxiliary Distributions

If $\boldsymbol{\pi}_x(y_i)$ and $\boldsymbol{\pi}_y(y_i)$ correspond to \mathbf{x} - and \mathbf{y} -parents of the auxiliary variable y_i in the graph of $p(\mathbf{y}|\mathbf{x})$, the generic bound on the mutual information $\tilde{I}(\mathbf{x}, \mathbf{y})$ in the bound (7.6) is expressed as

$$\begin{aligned} \tilde{I}(\mathbf{x}, \mathbf{y}) &= \langle \log p(\mathbf{y}|\mathbf{x}) \rangle_{q(\mathbf{x}, \mathbf{y})} + H(\mathbf{y}) \\ &= \sum_{i=1}^{|\mathbf{y}|} \langle \log p(y_i | \boldsymbol{\pi}_x(y_i), \boldsymbol{\pi}_y(y_i)) \rangle_{q(y_i, \boldsymbol{\pi}_y(y_i), \boldsymbol{\pi}_x(y_i))} + H(\mathbf{y}). \end{aligned} \quad (7.14)$$

We can always choose a mapping to the auxiliary space in such a way that computation of the entropic term $H(\mathbf{y})$ is not problematic. Here we will focus

mainly on the discussion of the computational complexity of $\langle \log p(\mathbf{y}|\mathbf{x}) \rangle_{q(\mathbf{x},\mathbf{y})}$. Note that the representational complexity of each conditional $p(y_i|\boldsymbol{\pi}_x(y_i), \boldsymbol{\pi}_y(y_i))$ is $\sim O(s^{|\boldsymbol{\pi}_x(y_i)|+|\boldsymbol{\pi}_y(y_i)|})$, where s is the number of states (for simplicity assumed to be equal for all variables y_i, x_j). Since we are free to choose a form of the distribution $p(\mathbf{y}|\mathbf{x})$, we can limit its parental structure, so that $|\boldsymbol{\pi}_x(y_i)|+|\boldsymbol{\pi}_y(y_i)|$ is low. For discrete variables this allows an exact representation of the conditionals.

In the special case when $q(\mathbf{y}, \mathbf{x}) = \prod_{i=1}^{|\mathbf{y}|} q(y_i) \prod_{j=1}^{|\mathbf{x}|} q(x_j|\boldsymbol{\pi}_y(x_j))$ is an irregular sparse bipartite graph (i.e. the number of \mathbf{y} -parents $\boldsymbol{\pi}_y(x_j)$ of each variable x_j is low), and the auxiliary conditional $p(\mathbf{y}|\mathbf{x})$ is a sparse structured distribution, one may perform the marginalization (7.14) explicitly, so that

$$\begin{aligned} \langle \log p(y_i|\boldsymbol{\pi}_x(y_i), \boldsymbol{\pi}_y(y_i)) \rangle_{q(y_i, \boldsymbol{\pi}_y(y_i), \boldsymbol{\pi}_x(y_i))} &= \sum_{y_i, \boldsymbol{\pi}_y(y_i)} q(y_i, \boldsymbol{\pi}_y(y_i)) \times \\ &\sum_{\boldsymbol{\pi}_x(y_i)} q(\boldsymbol{\pi}_x(y_i)|y_i, \boldsymbol{\pi}_y(y_i)) \log p(y_i|\boldsymbol{\pi}_x(y_i), \boldsymbol{\pi}_y(y_i)), \end{aligned} \quad (7.15)$$

where

$$q(\boldsymbol{\pi}_x(y_i)|y_i, \boldsymbol{\pi}_y(y_i)) \equiv \sum_{\tilde{\boldsymbol{\pi}}_i} q(\tilde{\boldsymbol{\pi}}_i) \prod_{j \in \boldsymbol{\pi}_x(y_i)} q(x_j|\tilde{\boldsymbol{\pi}}_i, y_i, \boldsymbol{\pi}_y(y_i)) \quad (7.16)$$

$$= \sum_{\tilde{\boldsymbol{\pi}}_i} q(\tilde{\boldsymbol{\pi}}_i) \prod_{j \in \boldsymbol{\pi}_x(y_i)} q(x_j|\boldsymbol{\pi}_y(x_j)), \quad (7.17)$$

and $\tilde{\boldsymbol{\pi}}_i \stackrel{\text{def}}{=} \boldsymbol{\pi}_y(\boldsymbol{\pi}_x(y_i)) \setminus \{\boldsymbol{\pi}_y(y_i) \cup y_i\}$. Here we extended the definition of the \mathbf{y} -parents in the conditional $q(\mathbf{x}|\mathbf{y})$, so that $\boldsymbol{\pi}_y(\boldsymbol{\pi}_x(y_i))$ defines all the \mathbf{y} -parents (expressed from $q(\mathbf{x}|\mathbf{y})$) of all the \mathbf{x} -parents (expressed from $p(\mathbf{y}|\mathbf{x})$) of the auxiliary variable y_j . The explicit marginalization (7.17) is $\sim O(s^{|\tilde{\boldsymbol{\pi}}_i|})$, and the complexity of explicit computations of the bound (7.15) is $\sim O(s^{|\boldsymbol{\pi}_x(y_i)|+|\boldsymbol{\pi}_y(\boldsymbol{\pi}_x(y_i))|})$. Clearly, it is acceptable as long as both $p(\mathbf{y}|\mathbf{x})$ and $q(\mathbf{x}|\mathbf{y})$ are sparse. For specific sparse structures of the auxiliary conditional $p(\mathbf{y}|\mathbf{x})$ and the variational distribution $q(\mathbf{x}, \mathbf{y})$, one may also consider evaluating the bound (7.15) by applying algorithms of the exact inference (e.g. Jensen (1996), Jordan (2005)).

Clearly, if there are no other (e.g. parametric) constraints on the auxiliary conditional $p(y_i|\boldsymbol{\pi}_x(y_i), \boldsymbol{\pi}_y(y_i))$, optimization of the bound (7.5) for $p(\mathbf{y}|\mathbf{x})$ results in the optimal gain

$$\max_{p(\mathbf{y}|\mathbf{x})} \left\{ \tilde{I}(\mathbf{x}, \mathbf{y}) \right\} = H(\mathbf{y}) - \sum_{i=1}^{|\mathbf{y}|} H(y_i|\boldsymbol{\pi}_x(y_i), \boldsymbol{\pi}_y(y_i)), \quad (7.18)$$

where $H(y_i|\boldsymbol{\pi}_x(y_i), \boldsymbol{\pi}_y(y_i))$ are the structural relaxations of $H_q(\mathbf{y}|\mathbf{x})$ expressed from the variational distribution $q(\mathbf{x}, \mathbf{y})$. Computations of each entropic term $H(y_i|\boldsymbol{\pi}_y(y_i), \boldsymbol{\pi}_x(y_i))$ are of the same order of complexity $\sim O(s^{|\boldsymbol{\pi}_x(y_i)|+|\boldsymbol{\pi}_y(\boldsymbol{\pi}_x(y_i))|})$, and under the sparsity constraints may be performed exactly (see also the discussion in Section 2.2.1).

7.3 Relation to Variational Mixture Models

In Section 7.2 we described optimization of the auxiliary variational bound on the log partition function subject to parametric constraints on the auxiliary conditional distribution. Fundamentally, it is exactly due to the choice of a constrained conditional $p(\mathbf{y}|\mathbf{x})$ (the variational decoder in the IM formulation) that computationally efficient evaluations of the bound (7.5) were possible. The *variational mixture* methods of Jaakkola and Jordan (1998), Bishop et al. (1998) may be viewed as maximizing the same objective criterion (7.5) for the case when there are no utilizable structural or parametric constraints on $p(\mathbf{y}|\mathbf{x})$. In this case, optimization of (7.5) for large-scale models requires arguably less general relaxations of (7.6), such as the one given by (7.11). Our argument here is that by constraining the auxiliary conditional (variational decoder) $p(\mathbf{y}|\mathbf{x})$ to lie in a tractable family of distributions, we may avoid some of the problems of the existing variational mixture formulations.

The existing variational mixture approaches (Jaakkola and Jordan (1998); Bishop et al. (1998); Lawrence et al. (1998); El-Hay and Friedman (2002)) express the mutual information term in the objective criterion (7.5) as

$$\begin{aligned} \tilde{I}(\mathbf{x}, y) &\geq \left\langle \log \frac{\tilde{q}(\mathbf{x}|y)}{q(y)} \right\rangle_{q(\mathbf{x}, y)} + \langle \log \lambda(y) \rangle_{q(y)} + 1 - \sum_y \lambda(y) \sum_{\mathbf{x}} \tilde{q}(\mathbf{x}|y) q_{mix}(\mathbf{x}) \\ &\stackrel{\text{def}}{=} \acute{I}(\mathbf{x}, y), \end{aligned} \quad (7.19)$$

which is obtained from (7.6) by applying the upper bound on the logarithmic function

$$\log x \leq \lambda x - \log \lambda - 1 \quad (7.20)$$

(see e.g. Jordan et al. (1998)). Commonly, the bound (7.19) is optimized with respect to the variational functionals $\lambda(y)$, “smoothing” conditionals $\tilde{q}(\mathbf{x}|y)$, and parameters of the mixture distribution $q(y)$ and $q(\mathbf{x}|y)$, where $q_{mix}(\mathbf{x}) \stackrel{\text{def}}{=} \langle q(\mathbf{x}|y) \rangle_{q(y)}$. In this case it is usually presumed that $y \in \{y_1, \dots, y_{|y|}\}$ is the space of mixing coefficients, so that the marginal $q_{mix}(\mathbf{x})$ may be computed exactly³. In order for the computations in (7.19) to be tractable, the smoothing distributions $\tilde{q}(\mathbf{x}|y)$ will need to be factorized (e.g. Jaakkola and Jordan (1998)). While theoretically one could use (7.19) to variationally fit a mixture of trees of different structures (as opposed to the mixture of completely factorized (Lawrence et al. (1998)) or identically structured (El-Hay and Friedman (2002)) models), it may be numerically difficult, as optimization would involve computations of non-factorized ratios of summations of exponentially small terms (see e.g. Bishop et al. (1998)), which follows from the need of computing the high-dimensional convolutions $\langle q_{mix}(\mathbf{x}) \rangle_{\tilde{q}(\mathbf{x}|y)}$.

Moreover, it is not easy to interpret the optimal solutions of (7.19). Effectively, by optimizing (7.19) with respect to the variational parameters $\lambda(y)$, we obtain

$$\max_{\lambda(y)} \acute{I}(\mathbf{x}, y) = \left\langle \log \frac{\tilde{q}(\mathbf{x}|y)q(y)}{\langle \tilde{q}(\mathbf{x}|y) \rangle_{q_{mix}(\mathbf{x})}} \right\rangle_{q(\mathbf{x}, y)} + H(y) \leq I(\mathbf{x}, y). \quad (7.21)$$

³Richer representations of the y -space may potentially be considered, provided that $q(\mathbf{x}, y)$ is in a tractable family, and the integrals in (7.19) may be computed exactly.

The bound (7.21) is saturated if and only if

$$\tilde{q}(\mathbf{x}|y)q(y)/\langle\tilde{q}(\mathbf{x}|y)\rangle_{q_{mix}(\mathbf{x})} \equiv q(y|\mathbf{x}). \quad (7.22)$$

However, for a general choice of $\tilde{q}(\mathbf{x}|y)$ the term in the l.h.s. of (7.22) does not define a proper distribution in y , which complicates the analysis of effects which relaxations of the structure of the conditional $\tilde{q}(\mathbf{x}|y)$ may have on the bound (cf variational information maximization). Probably the key conceptual problems of the existing variational mixture approaches are analytical difficulties of interpreting the optimal solutions, and practical difficulties of generalizing the existing bound (7.19) to richer, non-factorized families of the variational distributions. Some other practical limitations include general numerical instability of the direct computations of the average $\langle q_{mix}(\mathbf{x}) \rangle_{\tilde{q}(\mathbf{x}|y)}$ in large-scale models, and the iterative nature of optimization for individual factors $\tilde{q}(x_i|y)$ of the smoothing distribution $\tilde{q}(\mathbf{x}|y) = \prod_{i=1}^{|\mathbf{x}|} \tilde{q}(x_i|y)$ (which may sometimes result in a low convergence speed).

Our method addresses the subgoal of optimizing the generic lower bound $\tilde{I}(\mathbf{x}, y)$ by applying the variational IM. Apart from being conceptually and computationally simpler than the existing variational mixture approaches (Jaakkola and Jordan (1998); Bishop et al. (1998); Lawrence et al. (1998); El-Hay and Friedman (2002)), our method is also arguably more general. Specifically, by considering an unconstrained mapping to the auxiliary space and applying the bound on the logarithm (7.20), we can use our formulation (7.5) to arrive at (7.19), but not vice versa (also, note that (7.11) is a formal generalization of (7.19)). Importantly, we are not confined to using the specific relaxations (7.19), and may handle the intractability of computing $I(\mathbf{x}, y)$ by imposing constraints on the auxiliary conditional $p(y|\mathbf{x})$. Generally, the optimal auxiliary conditional distributions $p(y|\mathbf{x})$ approximate the posteriors $q(y|\mathbf{x})$ expressed from the variational model $q(y, \mathbf{x})$, which leads to a consistent improvement in the bound on $I(\mathbf{x}, y)$. Also, due to the flexibility of choosing the form of the mapping to the auxiliary space, our method suggests extensions of the variational mixture approaches to factorial and structured auxiliary spaces. Specifically, we note that our bound on $\log Z$ may be used in conjunction with the auxiliary variational bound on mutual information (2.38) to improve on simple (e.g. factorized) choices of the auxiliary conditional $p(y|\mathbf{x})$; it may also be used in situations when the auxiliary space is high-dimensional.

7.4 Demonstrations

Here we demonstrate the results comparing performance of the auxiliary variational framework with the standard factorized approaches. Also, for discrete variable spaces we apply our method to computing reweightings of fixed approximations $q(\mathbf{x}|y)$ of the original model $p(\mathbf{x})$, where each approximating distribution $q(\mathbf{x}|y_i)$, $i = 1, \dots, |y|$ had a unique structure. We also look at the qualitative effects which the choice of the structure of the auxiliary conditional $p(y|\mathbf{x})$ may have on the theoretically achievable improvements over a mixture of simple variational approximators.

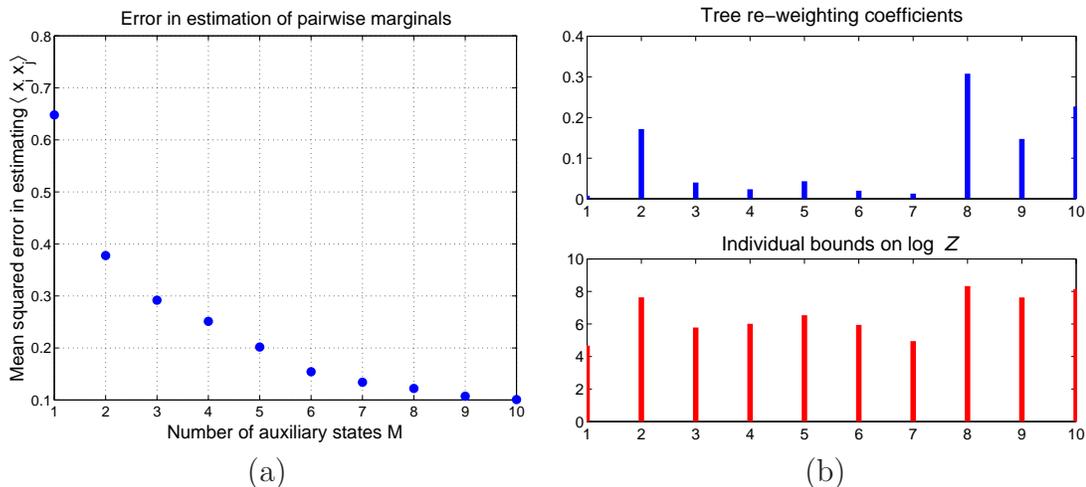


Figure 7.3: An auxiliary variational framework for approximate inference. (a) Systematic changes in the mean squared error for the estimates of the second-order moments with the growth in the number of mixture states M ; (b) *Top*: re-weighting coefficients for a set of fixed structured approximators (each $q(x|y)$ is a uniquely structured spanning tree); *Bottom*: the lower bounds on $\log Z$ obtained by each individual tree. Note that greater mixing coefficients were assigned to the tree-structured distributions which individually had resulted in higher values of lower bounds on $\log Z$.

Inference in Discrete Markov Networks

Here we demonstrate systematic changes in the auxiliary variational estimates of the second-order moments for discrete variable spaces. Throughout the simulations, it was assumed that $p(\mathbf{x})$ was a pairwise Markov network with the energy

$$E(\mathbf{x}) = -(\mathbf{x}^T \mathbf{W} \mathbf{x} + \mathbf{x}^T \mathbf{b}), \quad \mathbf{W} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|}, \quad \mathbf{b} \in \mathbb{R}^{|\mathbf{x}|}, \quad (7.23)$$

and $\mathbf{x} \in \{-1, 1\}^{|\mathbf{x}|}$. In the following experiments it is assumed that $y \in \{1, \dots, M\}$ is multinomial and $p(y|\mathbf{x})$ has a softmax form (7.13). Initially, we did not impose structural (sparsity) constraints on the auxiliary conditional distribution $p(y|\mathbf{x})$, and approximated the averages in the bound (7.5) at the means (i.e. $\langle \log p(y|\mathbf{x}) \rangle_{q(\mathbf{x}|y)} \approx \log p(y|\langle \mathbf{x} \rangle_{q(\mathbf{x}|y)})$). (Strictly speaking, the approximation is only accurate in the limit of noiseless projections from the auxiliary space, and one way to ensure its validity would be to explicitly constrain $q(x_i|y) \in \{0, 1\}$; we did not impose such constraints, hoping that maximization of the bound on mutual information $I(\mathbf{x}, y)$ in (7.5) would tend to favor low conditional entropies $H(\mathbf{x}|y)$, which would imply sharply peaked distributions). Analogous experiments were repeated for the case of conditionally factorized auxiliary state representations, i.e. $p(y|\mathbf{x}) = \prod_{i=1}^{|y|} p(y_i|\mathbf{x})$, $|y| \sim O(\log_2 M)$, which led to almost identical results.

To demonstrate the effect of the cardinality of the auxiliary space on the marginals, we generated 100 biases $\mathbf{b} \in \mathbb{I}^{10}$ and matrices $\mathbf{W} \in \mathbb{I}^{10 \times 10}$ of the 10-D pairwise Markov network $p(\mathbf{x})$, where \mathbb{I} defines the uniform range $[-1, 1]$. Then by analogy with Lawrence et al. (1998) we computed the squared errors $\epsilon(M)$

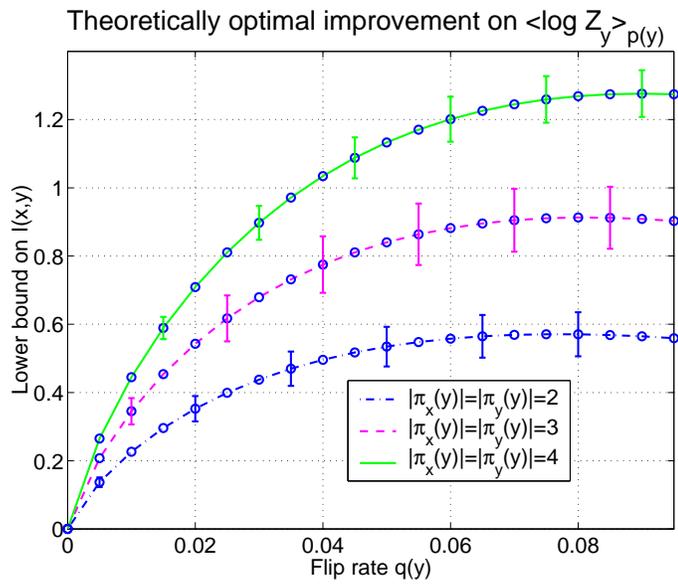


Figure 7.4: Influence of the structure of the auxiliary conditional $p(y|x)$ and the variational $q(x|y)$ distributions on the theoretically achievable improvement of the auxiliary variational lower bound on $\log Z$ over a mixture of simple bounds. The x -axis correspond to the “flip rate” $q(y)$. The y -axis correspond to the bound $\tilde{I}(x, y)$. The results are computed for $|x| = 20$, $|y| = 40$. The curves correspond to the exact values of the bound $\tilde{I}(x, y)$ for random structures of $q(x|y)$ and $p(y|x)$ which satisfy the specified sparsity constraints ($|\pi_y(x)| = 3$, $|\pi_x(y)|$ and $|\pi_y(y)|$ are shown in the legend box). Note that a choice of a richer auxiliary conditional $p(y|x)$ generally leads to higher values of $\tilde{I}(x, y)$.

between exact and estimated second moments $\langle x_i x_j \rangle$, averaged for all networks and all $i \neq j$. Note that for $M = 1$ the auxiliary variational representation (7.6) is equivalent to the mean field model. As can be seen from Figure 7.3 (a), we obtain a systematic improvement in the accuracy with an increase in M . The scale of the changes in the accuracy is different from that reported by Lawrence et al. (1998), whose mean field error $\epsilon(1)$ was approximately 0.15, though we observe qualitatively similar improvements. This discrepancy is undoubtedly due to details of the optimization approaches, and does not detract from our conclusion that the auxiliary method conveys a systematic improvement.

Reweighting Structured Approximators

Here we demonstrate how we can easily apply a simple reformulation of our method for computing reweightings of fixed structured approximations $q(x|y)$. For the weights $W \in \mathbb{I}^{10 \times 10}$ chosen at uniform random, we generated $K = 10$ random spanning trees with the weights $W^{(m)}$ such that $W_{ij}^{(m)} = W_{ij}$ for all i, j and $m = 1, \dots, K$. Then we re-weighted the trees by recomputing $q(y)$ according

to

$$q(y_k) \propto \exp \left\{ \langle -E(\mathbf{x}) - \log q(\mathbf{x}|y_k) + \log p(y_k|\mathbf{x}) \rangle_{q(\mathbf{x}|y_k)} \right\}, \quad (7.24)$$

$$= p(y_k) \exp \left\{ -KL(q(\mathbf{x}|y_k) \| p(\mathbf{x}|y_k)) \right\} \quad (7.25)$$

which follows directly from (7.10). Note that the mixing coefficient $q(y_k)$ of each component $q(\mathbf{x}|y_k)$ in the variational distribution is proportional to the surrogate marginal in the auxiliary space $p(y_k)$, and the discrepancy between the mixture component $q(\mathbf{x}|y_k)$ and the model $p(\mathbf{x}|y_k)$ for a fixed setting of the component's label y_k .

Figure 7.3 (b) illustrates typical re-weightings of fixed structured approximators $q(\mathbf{x}|y)$ and the induced lower bounds on $\log Z$ for the case of the softmax parameterization of $p(y|\mathbf{x})$ and for *fixed* parameters of the auxiliary conditional (selected at uniform random on \mathbb{I}). The corresponding bounds in this case were $L_r \approx 8.38$, $L_u \approx 8.87$, $L_b \approx 8.33$, and $L_{av} \approx 9.52$ for the random, uniform, best single, and auxiliary variational weightings respectively. Note that the single-step computation of (7.24) has an acceptable order of computational complexity ($\sim O(|\mathbf{x}|^2)$ if we consider the approximation $\log p(y|\langle \mathbf{x} \rangle_{q(\mathbf{x}|y)})$ and a Markov network (7.23)), and can be easily extended to the case of higher dimensional models.

The results are straight-forward but potentially interesting, as they suggest a quick and simple way to compute mixing coefficient for a set of trained structured approximations $q(\mathbf{x}|y)$. So instead of maximizing the bound for K standard intrinsically tractable approximators, choosing a single best one which leads to the tightest lower bound on $\log Z$, and discarding the remaining $K - 1$ variational distributions, we can use all of the approximations to obtain a tighter lower bound on the normalizing constant. Interestingly, we can obtain these results even without learning the optimal mapping $p(y|\mathbf{x})$ to the auxiliary space, and by performing generally suboptimal approximations of the conditional entropies in $\tilde{I}(\mathbf{x}, y)$ at the mean of $q(\mathbf{x}|y)$. An arguably more principled approach to the problem of re-weighting the variational distributions would involve optimization of the auxiliary variational lower bound on $\log Z$ with respect to $q(y)$ and $p(y|\mathbf{x})$ for a set of given structured components $q(\mathbf{x}|y)$. Formally, a more careful choice of constraints on $p(y|\mathbf{x})$ may need to be considered in order to ensure tractability of computing $\tilde{I}(\mathbf{x}, y)$. Generally, however, it may be practically attractive to apply a simple one-step procedure to produce a tighter bound on $\log Z$ for a given choice of standard variational distributions.

Structured Auxiliary Mappings

It is interesting to see how the choice of structured constraints on both the auxiliary mapping $p(y|\mathbf{x})$ and the variational conditional $q(\mathbf{x}|y)$ could affect the theoretically optimal improvement over a convex combination of approximations. In order to check this, we assumed that the auxiliary space was binary and high-dimensional, so that $y \in \{0, 1\}^{|y|}$ (cf the variational mixture methods of Jaakkola and Jordan (1998), Bishop et al. (1998)). We also assumed that different dimensions of the auxiliary vector were identically distributed, so that $q(\mathbf{x}, y) = q(\mathbf{x}|y) \prod_{i=1}^{|y|} q_i(y_i)$ and $q_i(y_i = a) \equiv q_j(y_j = a) \stackrel{\text{def}}{=} q(y = a)$ for all

$i, j \in \{1, \dots, |y|\}$. The variational conditional $q(x|y)$ was constrained as

$$q(x|y) = \prod_{i=1}^{|\mathbf{x}|} q(x_i | \boldsymbol{\pi}_y(x_i)) = \prod_{i=1}^{|\mathbf{x}|} \delta \left(x_i - \sum_{y_j \in \boldsymbol{\pi}_y(x_i)}^{\oplus} y_j \right), \quad (7.26)$$

where $\boldsymbol{\pi}_y(x_i)$ are the y -parents of variable x_i in the graph for $q(x|y)$, and \sum^{\oplus} defines the modulo 2 summation. Equivalently, we can define $q(x|y) \sim \delta(\mathbf{x} - \mathbf{G}\mathbf{y})$, where $\mathbf{G} \in \{0, 1\}^{|\mathbf{x}| \times |\mathbf{y}|}$ is a binary matrix such that $G_{ij} = 1$ if and only if y_j is a parent of x_i in the graph⁴ for the variational conditional $q(x|y)$.

To ensure tractability of the computations, we also imposed structural constraints on the stochastic mapping to the auxiliary space, so that the auxiliary conditional was defined as $p(y|x) = \prod_{i=1}^{|\mathbf{y}|} p(y_i | \boldsymbol{\pi}_x(y_i), \boldsymbol{\pi}_y(y_i))$, where $\boldsymbol{\pi}_x(y_i)$ and $\boldsymbol{\pi}_y(y_i)$ are the x - and y -parents of unit y_i under the model $p(y|x)$, and $|\boldsymbol{\pi}_x(y_i)| \ll |\mathbf{y}|$ (additional care had been taken to ensure that $p(y|x)$ defined a proper distribution, i.e. its graph was directed acyclic). The variables were defined over discrete high-dimensional domains, so that $\mathbf{x} \in \{0, 1\}^{|\mathbf{x}|}$ and $\mathbf{y} \in \{0, 1\}^{|\mathbf{y}|}$, and under the imposed sparsity constraints the conditional $p(y_i | \boldsymbol{\pi}_x(y_i), \boldsymbol{\pi}_y(y_i))$ could be immediately obtained by maximizing the exact value of the bound (7.6).

Figure 7.4 shows the improvement $\tilde{I}(\mathbf{x}, \mathbf{y})$ over the convex combination in (7.5) as a function of the flip rate $q(y) \in [0, 0.1]$ for $|\mathbf{x}| = 20$ and $|\mathbf{y}| = 40$. The results are shown for the optimal settings of the auxiliary conditional $p(y|x)$, obtained by maximizing the bound on $\log Z$ (which immediately leads to the constrained approximations of $q(x|y)$). Throughout the experiments, it was assumed that $|\boldsymbol{\pi}_y(x_i)| = 3$, and the structure of $p(y|x)$ varied as shown on Figure 7.4. For each choice of the structural constraints, the results were averaged over 20 runs with random choices of $q(x|y)$. We see that models with richer structures of the auxiliary conditional typically lead to greater bounds. Note that for the trivial independent setting $p(y|x) = p(y)$ (i.e. $|\boldsymbol{\pi}_x(y_i)| = 0$), there would be no theoretical improvements over a standard approximation, so that $\tilde{I}(\mathbf{x}, \mathbf{y}) = 0$. However, the considered choices do lead to consistent improvements in (7.6) with an increase in the complexity of the auxiliary mappings. This also empirically confirms the analytical results of Section 2.2.1, and agrees with the intuition that richer structures of variational distributions tend lead to consistent improvements over simpler approximations.

We have repeated similar experiments for different sizes of the auxiliary space $|\mathbf{y}|$ (assuming that $|\mathbf{x}|$ was fixed). Our initial hope was that by considering high-dimensional factorial representations in the auxiliary space, we could potentially increase the theoretically achievable improvements over mixtures of simple bounds. It is easy to see that if M is the effective number of mixture states then $\tilde{I}(\mathbf{x}, \mathbf{y}) \leq I(\mathbf{x}, \mathbf{y}) \leq \log M$, with the maximum achieved for flat $q(y)$ and deterministic posteriors $q(y|x)$ (see expression (7.5) and a discussion in Bishop et al.

⁴This formulation is particularly interesting, as it could be shown to be relatable to the settings of a syndrome decoding problem (see e.g. McEliece (1977)), where the marginal $q(y)$ defines the bit error probability (flip rate) of a binary symmetric channel, $\mathbf{x} \in \{0, 1\}^{|\mathbf{x}|}$ is the *syndrome* vector, and $\mathbf{y} \in \{0, 1\}^{|\mathbf{y}|}$ is the (unknown) noise vector. The goal of decoding in this context would be to perfectly determine the unknown noise \mathbf{y} from the syndrome vector \mathbf{x} , which for the half-rate codes ($|\mathbf{y}| = 2|\mathbf{x}|$) could theoretically be done for $q(y) \lesssim 0.11$.

(1998)). By considering the factorial auxiliary state representations, we hoped to massively increase the effective number of mixture states (so that $M \sim O(s^{|y|})$) and achieve nearly linear improvement in the bound (7.5) with $|y|$. Unfortunately, our initial intuition has not yet received strong empirical confirmations. While choosing richer structures of the auxiliary conditionals $p(y|x)$ could indeed be helpful for obtaining consistent improvements in $\tilde{I}(x, y)$, choices of higher-dimensional auxiliary spaces did not always lead to systematic improvements in the bounds, especially for higher flip rates $q(y)$ and random irregular structures of $p(y|x)$ and $q(x|y)$.

This effect may be explained by the sparsity constraints which we needed to impose on the auxiliary conditional. Let us presume that the auxiliary conditional is given by $p(y|x) = \prod_{i=1}^{|y|} p(y_i | \pi_x(y_i), \pi_y(y_i))$, where each factor $p(y_i | \pi_x(y_i), \pi_y(y_i))$. Without loss of generality we may define the exact posterior expressed from the variational model $q(y, x)$ as $q(y|x) = \prod_{i=1}^{|y|} q(y_i | \tilde{\pi}_x(y_i), \tilde{\pi}_y(y_i))$, where $\tilde{\pi}_x(y_i) \supseteq \pi_x(y_i)$ and $\tilde{\pi}_y(y_i) \supseteq \pi_y(y_i)$. Then the difference between the exact mutual information $I(x, y)$ and the tightest lower bound (7.6) obtained at the optimal settings of $p(y|x)$ is given by

$$\sum_{i=1}^{|y|} I(y_i, \{\tilde{\pi}_x(y_i), \tilde{\pi}_y(y_i)\} \setminus \{\pi_x(y_i), \pi_y(y_i)\} | \{\pi_x(y_i), \pi_y(y_i)\}),$$

which quantifies the overcounting effects which occur when we bound the conditional entropy $H(y|x)$ by a summation of the marginals. (For example, if both the variational distribution $q(x, y)$ and the auxiliary conditional $p(y|x)$ are chains, this difference would be negligible only in situations when $q(y_i | y_{i-1}, x_i) \approx q(y_i | y_{i-1}, x)$, leading to $I(y_i, x \setminus x_i | y_{i-1}, x_i) \approx 0$. This is effectively analogous to the cases when we may ignore future observations during the inference in state-space models, which clearly imposes significant limitations). One may expect that for general distributions $q(x, y)$, the growth in $|y|$ and $|x|$ leads to increasing overcounting effects, resulting in a greater volume loss. These effects are not accounted for in the current formulation of our method, and will need to be addressed in further extensions of this work. E.g. the overcounting problem may be addressed by cluster variation methods (see e.g. Kikuchi (1951), Yedidia et al. (2000a), Yedidia et al. (2004)), which approximate entropy over a region as a weighted summation of entropies over smaller regions. Generally, however, our motivation here was to retain a rigorous bound on the intractable entropy (leading to a proper lower bound of $\log Z$), and improving our current method in a rigorous and general way is a challenging problem for future research.

7.5 Summary

Here we described an approach which generalizes the standard Kullback-Leibler variational procedures for evaluating the variational lower bounds on the generally intractable normalizing constants of undirected graphical models. Our work here was partially motivated by the success of auxiliary sampling techniques (e.g.

Swendsen and Wang (1987), Higdon (1998), Neal (1993)), which introduce auxiliary variables and draw samples from a specific joint distribution defined in the augmented variable space. The role of the auxiliary variables in this context is to capture (structural) information about clusters of correlated variables, which proves to be effective for decreasing the time gap between subsequent independent samples.

In this chapter we considered an auxiliary variable extension of common variational methods for bounding the log partition function, which could be useful in the context of approximate probabilistic inference. We showed that it is possible to define a tractable variational framework which leads to systematic improvements over the standard theory for any structured approximation. While the method described here is of a potential interest as a general approach for approximate inference, it demonstrates a curious link to our variational information-maximizing framework discussed in Chapter 2. Specifically, it turns out that the improvement of the proposed bound on $\log Z$ over a convex combination of standard bounds is given by a specific form of the generic lower bound on mutual information (see expression (2.2)).

Our auxiliary variational method for inference is related to the family of variational mixture models (see Jaakkola and Jordan (1998), Lawrence et al. (1998)), which can be seen as special and more computationally expensive cases of our approach. Specifically, we can obtain the existing bounds by considering an unconstrained mapping to the auxiliary space and applying a standard relaxation of the logarithmic function (7.20) (but we cannot generally derive our framework from the existing variational mixture approximations). Importantly, we may avoid computational, numerical, and practical problems of applying the existing variational mixture methods by considering tractable constraints on the mappings from the data space $\{\mathbf{x}\}$ to the auxiliary space $\{\mathbf{y}\}$. In this formulation, we may view our variational information-maximizing framework as an integral subgoal of auxiliary variational inference. The method is attractive both computationally and analytically, as the flexibility of the choice of the auxiliary conditional distribution (the variational decoder in the IM terminology) suggests simple generalizations of the generic approach to richer families of variational decoders (such as chains, polytrees, and mixtures-of-experts, see e.g. Section 2.2.1 and Section 2.3). Additionally, the variational IM algorithm is easy to understand analytically (which becomes particularly attractive when compared with the analysis of the existing variational mixture models, where interpretation of the optimal smoothing factors and their effects on the induced bounds on $I(\mathbf{x}, y)$ at this stage remain, to us, largely unclear). Indeed, the optimization surface of our method is concave in the auxiliary conditional $p(\mathbf{y}|\mathbf{x})$, which for sparse models makes it possible to find the optimal projections analytically (see Section 2.2 and Section 7.2.2).

In the context of probabilistic variational inference, we may use the suggested approach for improving on the lower bounds on $\log Z$ produced by the standard approximations. Indeed, we showed that the method described here forms a more powerful class of approximations than any structured mean field technique. We also showed that a richer choice of structures of the auxiliary conditional

distribution may indeed lead to tighter lower bounds on the log partition function compared with simpler approximations. While these results are potentially interesting (especially when compared with standard variational approaches), the theoretically achievable gains over mixtures of standard approximations are bounded to be logarithmic in the size of the auxiliary space. Our initial hope was that by considering factorial representations of the auxiliary space, we could increase the effective number of the states, which could potentially lead to super-logarithmic improvements. Unfortunately, our current experience suggests that by naively increasing dimensionality of the auxiliary space (and imposing sparsity constraints on $p(\mathbf{x}|y)$), we do not always get significant improvements over the lower-dimensional formulations. This may potentially be explained by the fact that the simple formulation of the method leads to the bounds which ignore the overcounting problem in bounding the conditional entropy $H_q(y|\mathbf{x})$. Addressing this matter in a systematic way can make the method potentially more attractive for variational inference.

A practical byproduct of our work is an efficient way to calculate a set of mixture coefficients for any set of tractable distributions, which principally improves on the flat combination (Agakov and Barber (2004a)). An extension of this work would be an application to a real-world problem (possibly in a structured context, e.g. by extending and re-formulating the results of Ghahramani and Hinton (1998), Ghahramani and Hinton (2000)).

Chapter 8

Discussion

Finding regularities in observations of the external world is an important task handled by many biological organisms. It is believed that learning about such regularities involves adapting structural and physiological properties of brain synapses, resulting in internal encodings of the external stimuli. Intuitively, meaningful internal representations should in some sense be informative about the environment, in which case biological learning may be viewed as a process of finding informative representations of the observations. Of course, in biology such internal representations would be intrinsically constrained by the neurophysiological properties of biological networks.

Many applications of machine learning are aiming to address a fundamentally similar task, where the key goal is to automatically find meaningful representations of the observations. For example, a system of automated medical diagnostics may be applied for determining hidden diseases which could have given rise to the observed symptoms; images of human faces taken by a robot's camera may be used by a robot to recognize human emotions, etc. The process of finding unknown informative descriptions (such as diseases and emotions) of the observations (such as symptoms and photographs) is generally referred to as the process of *inference*, which is one of the fundamental tasks of machine learning. Another fundamental task is *learning*, which may be viewed as an automated procedure for finding mathematical formalisms which would result in meaningful inferences. In practice learning often corresponds to finding an optimal *model* which provides a description of how the data relates to its representations. Both learning and inference may in principle be very computationally demanding. In this work we aimed to target computational intractability of a class of learning approaches.

Generative vs encoder models

Generally, one may distinguish two different approaches to unsupervised learning. The first class of methods aims to learn a constrained statistical model of the observations. The key idea there is to find a model of a probability distribution which would be likely to generate the data. Constraints on the distributions are introduced by utilizing the prior knowledge about the modeled domain, for example by choosing a specific parameterization and structure for a class of models. In these methods learning typically corresponds to fitting a constrained model to

the observations. Informally, this may be viewed as minimizing a measure of discrepancy between the empirical distribution and the marginal distribution of the observed variables expressed from the model. Inference, i.e. the process of finding unknown informative representations of the observations, can be addressed by applying fundamental probability rules. An example of such class of methods includes maximization of marginal likelihoods in *generative latent variable* models, where for a fixed set of parameters the model specifies how to generate the observations. Usefulness of the extracted hidden variable representations expressed from the model would in this case be quantified by how well the model fits to the set of observations. A different class of methods aims to extract informative representations directly from the set of observations. Rather than learning the density model of the observations, the methods aim to find a mapping (generally, stochastic) from the observations to the latent variable representations by optimizing other measures of informativeness. A popular class of such methods is based on the idea of maximizing the mutual information between the observations and the unknown representations in the *encoder*, or *recognition* models.

There are fundamental differences in how we parameterize and train generative and recognition models. Specifically, while parameterization of a generative model imposes explicit constraints on the distribution of the observations given latent variable representations (which may require knowledge about the process which generates the data), parameterization of an encoder model imposes explicit constraints on the *encoding* distribution. This may be particularly important in situations when constraints on the encoding distribution are partially known from the physical properties of the channels (which may be the case in engineering or neurobiological domains), or when explicit constraints on the encoding distribution form a part of a specific problem formulation (for example, in problems of constrained subspace selection). In some applications it may be easier to parameterize an encoder, rather than a generative model (one example is clustering with a defined family of similarity measures); in other cases, the situation may be completely opposite (for instance, when the family of distributions generating a specific dataset is known).

Generally, a choice of a model (generative *vs* encoder) may be strongly dependent on a specific task. To illustrate this, consider a problem of constructing a system for automated medical diagnostics, where patients' symptoms are associated with diseases. Many medical experts may arguably find it easier to parameterize a stochastic mapping $D \rightarrow S$ from the diseases (D) to the symptoms (S) (which requires encyclopedic knowledge), rather than $S \rightarrow D$ (which requires medical heuristics). Once the model is specified and trained, it may be applied for inferring combinations of hidden diseases from a combination of the observed symptoms. We see that the former specification of the system defines a standard generative model; the latter parameterization corresponds to an encoder model. If we have a good idea about symptoms which could be "generated" by a given disease, the generative model is more easily parameterizable.

Now let us re-formulate the problem slightly. Assume that we are interested in finding clusters of symptoms which could form a new (and unknown) family of diseases, which we may be aiming to name. By specifying a similarity metric

in $\{S\}$ (i.e. defining what it means for the symptoms to be similar), we may find it easier to describe a mapping from a symptom to a disease family, even though we may not necessarily know much about disease families. This would be an example of an encoder model, where the unknown cluster labels (disease families) carry useful information about the observations (symptoms). Effectively, this formulation ($S \rightarrow D$) defines a *discriminative unsupervised* framework, where the model is parameterized similarly to a conditionally trained classifier (like in the supervised learning formulation), but the representations D are unknown. Such parameterizations may be particularly useful in situations when little is known about the process which generates the data.

Unfortunately, for both generative and encoder models the problem of finding informative latent variable representations $\{y\}$ of the observations $\{x\}$ may become very difficult in the presence of noise, which motivates a need of approximations. In this thesis we considered a class of machine learning methods maximizing the mutual information $I(x, y)$ in encoder models, and addressed the fundamental computational problems of the exact formulation by applying variational approximations. Our focus on the discussion of variational methods for information maximization was partially motivated by their popularity and effectiveness for inference and likelihood training in graphical models, and the apparent lack of understanding of how the methods could in general be applied in the information-maximizing context (indeed, many of the currently existing approximations of mutual information are either heuristic or too specific). The fundamental advantage of the variational methods over other approximations is availability of rigorous bounds on the intractable quantities, which facilitates comparisons of different variational approximation techniques, and in some cases makes them particularly attractive for learning. For this reason, one of our technical subgoals here was to try to retain rigorous bounds on $I(x, y)$, whenever possible.

General results

In this thesis we described a family of variational lower bounds on mutual information $I(x, y)$, which gave rise to a formal and theoretically justified approach to information maximization in noisy channels. While the formulation of the generic bound on the objective criterion is straight-forward, it appears to have attracted little previous attention as a practical tool for approximate maximization of information content. The fundamental idea of the approach is to introduce a *variational decoder* $q(x|y)$ which is constrained to lie in a tractable family. Effectively, an iterative information-maximizing (IM) algorithm optimizing the generic lower bound $\tilde{I}(x, y)$ extends the family of the generally intractable Blahut-Arimoto type algorithms (Arimoto (1972), Blahut (1972)), and reduces to them in the special case when the variational decoder is unconstrained. Qualitatively, this is similar to the variational EM algorithm for likelihood maximization (Neal and Hinton (1998)), which reduces to the standard EM (Dempster et al. (1977)) for unconstrained variational posteriors. The generality of the approach allows a flexibility in the choice of variational decoders or specific optimization procedures, which suggests that the method may naturally generalize other techniques for approx-

imate maximization of mutual information. Indeed, in addition to generalizing the conventional Blahut-Arimoto formulation, the IM generalizes the existing *as-if* Gaussian criterion (Linsker (1992)), which may be seen as a specific way of optimizing the variational lower bound on mutual information for a specific choice of linear Gaussian variational decoders. Additionally, factorial relaxations of the generally intractable conditional entropy $H(\mathbf{x}|\mathbf{y})$, which are sometimes used in approximations of $I(\mathbf{x}, \mathbf{y})$, may be seen as special cases of the generic lower bound on $I(\mathbf{x}, \mathbf{y})$ with structural constraints on the parental structure of the variational decoder distribution.

In our work we also explored general relations of the variational IM algorithm to maximum likelihood learning in generative models and conditional likelihood learning in noiseless and stochastic autoencoders (where by a noiseless autoencoder we mean a self-supervised network with a deterministic encoding distribution $p(\mathbf{y}|\mathbf{x}) \sim \delta$). In contrast to much of the previous work which related the likelihood and the mutual information approaches for relatively simple special cases (e.g. Pearlmutter and Parra (1996), Cardoso (1997), MacKay (1999b)), we aimed at relating the methods for the general variational settings independently of the specific model parameterizations. To make the comparisons possible, we assumed that both the encoder and the generative model led to identical inferences of the latent representations \mathbf{y} for all the source patterns \mathbf{x} . For this case, we showed that the likelihood of a generative model $\mathbf{y} \rightarrow \mathbf{x}$ could be viewed as a relaxation of a specific generic bound on $I(\mathbf{x}, \mathbf{y})$, defined for the corresponding model $\mathbf{x} \rightarrow \mathbf{y}$ of a noisy channel. Interestingly, while generally it is not easy to relate learning by maximizing the exact likelihood and the exact mutual information for generative and encoder models, we may find specific cases when the fixed points of the generic lower bound on $I(\mathbf{x}, \mathbf{y})$ are identical to those of the standard variational Jensen’s bound on the likelihood (with identical constraints on the *encoding mapping* of the encoder model and the *variational posterior* of the generative model). For example, this happens for noiseless channels *independently* of invertibility of the mappings $\mathbf{x} \mapsto \mathbf{y}$ (provided that the generative model is characterized by flat priors on latent variables \mathbf{y}). This result shows interesting intersections of the variational IM and EM approaches, though in general the methods are quite different.

Another curious result of our study is a link of the simple form of the IM algorithm to conditional likelihood training in stochastic autoencoders. While the relation of the probability of correct *deterministic* reconstructions to the mutual information in the corresponding stochastic channels is well-known (Fano (1961)), its stochastic generalizations appear to have attracted little attention within the machine learning community. (This was probably the reason for a large number of heuristic or very specific approximations of the generally intractable mutual information $I(\mathbf{x}, \mathbf{y})$, which to the best of our knowledge are rarely (if at all) compared with the conditional training in stochastic autoencoders). A simple way to maximize the information content which the codes \mathbf{y} contain about the sources \mathbf{x} would be to maximize the conditional likelihood $p(\tilde{\mathbf{x}}|\mathbf{x})$ in a stochastic autoencoder model $\mathbf{x} \rightarrow \mathbf{y} \rightarrow \tilde{\mathbf{x}}$ for some choice of the decoding distribution. It is easy to see that in stochastic models the conditional likelihood is generally intractable,

and it is necessary to consider its approximations. As a result of our study we showed that the standard variational Jensen’s bounds on the conditional likelihood have the same fixed points as the generic lower bound on $I(\mathbf{x}, \mathbf{y})$. This result is rather disappointing, as it shows that by considering a simple *stochastic* model performing *self-supervised* training we could effectively arrive at the fixed points identical to the ones produced by the simplest formulation of the IM approach. However, our method is easier than standard variational approaches of maximizing the conditional likelihood in stochastic autoencoders, which would effectively correspond to a more expensive way of optimizing the simple generic bound on the mutual information. Generally, this results in the computational and representational efficiency and desirable convergence properties of the IM approaches, which optimize the bound for a significantly lower number of variational parameters (lower by the order of magnitude compared with a naive application of the variational Jensen’s bound on $p(\tilde{\mathbf{x}}|\mathbf{x})$).

Another result of the suggested work is the definition of a richer family of tractable *auxiliary variational* lower bounds on $I(\mathbf{x}, \mathbf{y})$, which formally generalizes on the generic lower bounds on $I(\mathbf{x}, \mathbf{y})$. By analogy with the auxiliary sampling techniques (see e.g. Swendsen and Wang (1987), Higdon (1998)), the key idea was to introduce additional variables, which could be used for capturing useful regularities of the source patterns and for introducing global dependencies to the decodings. Importantly, projections to the auxiliary space were defined in a way which did not alter properties of the original channel. We demonstrated that the auxiliary variational bounds may indeed be significantly tighter than the simple generic criteria. This result is potentially interesting from the communication-theoretic perspective, as it demonstrates a simple and computationally efficient way to produce tighter bounds on the capacity of a communication channel without assuming that more data is being transmitted across the channel. Generally, this family of variational methods may be used as a simple and tractable approach for improving on simple bounds on mutual information.

Interestingly, the variational formulation of the information maximizing procedure suggests a relation between maximizing the generic lower bound on $I(\mathbf{x}, \mathbf{y})$ and computing an optimal estimate of an intractable posterior $p(\mathbf{x}|\mathbf{y})$ in a generative model, where the codes \mathbf{y} are visible and the sources \mathbf{x} are hidden. One of the goals of variational inference in graphical models is to approximate moments of the generally intractable posterior $p(\mathbf{x}|\mathbf{y})$ by the moments of a simpler distribution $q(\mathbf{x}|\mathbf{y})$. Standard mean field approaches assume that $q(\mathbf{x}|\mathbf{y}) \propto \prod_{i=1}^{|\mathbf{x}|} q(x_i|\mathbf{y})$, which usually leads to $q(x_i|\mathbf{y})$ approximating any one of a large number of local *modes* of the model-specific posterior $p(x_i|\mathbf{y})$ (where $q(\mathbf{x}|\mathbf{y})$ is assumed to be factorized in \mathbf{x}). On the other hand, by imposing parametric constraints on the variational decoder, so that $q(\mathbf{x}|\mathbf{y}, \Theta) = \prod_{i=1}^{|\mathbf{x}|} q(x_i|\mathbf{y}, \Theta_i)$, and optimizing our bound on $I(\mathbf{x}, \mathbf{y})$ for parameters of the decoder, we obtain posterior mean estimates which are good *on average*.

This outlines a fundamental problem of using the variational decoders $q(\mathbf{x}|\mathbf{y}, \Theta)$ of the IM formulation for approximate inference, even compared with standard mean field methods. For example, this is the case for the problem of *syndrome decoding* in binary symmetric channels. Specifically, if the exact posterior $p(\mathbf{x}|\mathbf{y})$

is sharply peaked (which is the case for specific choices of the encoder distribution $p(\mathbf{y}|\mathbf{x})$, see e.g. Luby et al. (2001)), a mean field theory may indeed provide an accurate approximation of $p(\mathbf{x}|\mathbf{y})$ for each setting of \mathbf{y} . (Generally, however, such accurate approximations may not always be easily found, due to a large number of local modes). In contrast, by learning a factorized variational distribution which is only optimal on average, variational estimates of the exact marginal for each specific codeword \mathbf{y} may in general be very poor. One could hope to obtain improvements by considering significantly richer families of variational decoders (e.g. by having no parametric constraints on each factor $q(x_i|\mathbf{y})$). Unfortunately, unless the parental structure of each $q(x_i|\mathbf{y})$ is very sparse, the computational and representational complexity of using such approximations will generally be prohibitively high, while a choice of sparse structures will still lead to the averaging effects. While it may potentially be possible to modify the IM framework to obtain better estimates of the exact posterior means $\langle x_i \rangle_{p(x_i|\mathbf{y})}$, we do not advocate using simple constrained variational decoders obtained by our method as a competitive technique for error correction in binary symmetric channels.

Generally, we stress that a correct way of viewing our variational approach would be to interpret it as a general framework for learning an optimal *encoding* distribution by maximizing a proper bound on $I(\mathbf{x}, \mathbf{y})$. Once the optimal encoder is learned, one may choose any of a number of approximate inference techniques to perform the reconstructions. Using the variational decoder $q(\mathbf{x}|\mathbf{y})$ for inference is only one of such choices (which may be good or poor, depending on the specific application).

Case studies

As a part of our exploration of the information maximizing framework, we considered applications to constrained dimensionality reduction. Specifically, we discussed several ways of applying the framework to information-theoretic clustering, where the encoding distribution was defined either by the exact posterior of the corresponding latent variable model, or by explicitly choosing a specific nonlinear projection into the code space. The former approach could also be used for training generative models, which typically resulted in a more uniform coverage of the data space compared with the conventional training by maximizing the likelihood in mixture models. The latter approach to information-theoretic clustering could be applied to learning parameters of kernel functions (within a specific family), which was beneficial for visualizing the underlying structure of the data. Empirically, the method favorably compared with the Gaussian mixture, feature-space k-means, and non-kernelized information-theoretic clustering.

We also outlined analytical properties of the optimal IM solutions for the case of real-valued encoded representations \mathbf{y} . Specifically, we extended the work of Boullard and Kamp (1988) and Boullard (2000) to *arbitrary kernelizable* feature mappings in the stochastic information-theoretic context, and showed the nothing could be gained by using nonlinear encoders and linear variational decoders in the context of variational information maximization in noisy Gaussian channels. To handle the intrinsic constraints of linear Gaussian variational decoders applied in the context of nonlinear encodings, we suggested a proper variational

relaxation of the bound on $I(\mathbf{x}, \mathbf{y})$. For *nonlinear Gaussian* encoding distributions, this led to kernel PCA (Schoelkopf et al. (1998)) as the optimal solution for encoder’s weights. Additionally, in the deterministic limit of non-parametric encoding mappings, our framework led to the Gaussian Process Latent Variable Models (Lawrence (2003)) for a specific choice of the variational decoder.

As an immediate extension of our work, we note that a further study of the bound on $I(\mathbf{x}, \mathbf{y})$ for nonlinear Gaussian channels may need to be considered. Our current study indicates that the IM framework may be used to learn non-degenerate kernel parameters, but the obtained visualization and reconstruction results are strongly influenced by the choice of the kernel and the constraints on the feature-to-data decoding mappings. Importantly, we note that if the encoding noise may in principle be reduced to zero, one may consider optimizing simpler objective functions, such as those of Gaussian Process Latent Variable Models (Lawrence (2003)). These models also appear to lead to better visualization results, as they consider noiseless non-parametric encoders and rich families of variational decoders (given in these models by products of Gaussian processes). It is therefore believed that a further study of variational IM in nonlinear Gaussian channels for visualization or dimensionality reduction may be of a practical interest mainly in situations when the noise of the encoding distribution is unavoidable.

An arguably more promising application field to explore would be communication of discrete-valued data over channels with Gaussian noise, where a specific practical application may include code division multiple access in cellular telephony (see e.g. Viterbi (1995)). Additionally, it is interesting to explore the relation of our kernelized information-theoretic clustering approach (Agakov and Barber (2005b), Agakov and Barber (2005c)) to some of the common spectral clustering methods (e.g. Shi and Malik (2000), Yu and Shi (2003)), which were shown to be related to the weighted feature-space k-means (Dhillon et al. (2004)). Some of our preliminary results extending the work of Chapter 5 suggest that by constraining the feature-space coefficients, it may potentially be possible to relate the mentioned spectral clustering methods to a form of the constrained variational information-maximizing procedure (though a direct relation between the methods is not entirely clear at this stage).

As another application of the variational IM approach, we explored applicability of the variational IM framework in the context of learning high-dimensional binary representations of continuous source patterns for a specific biologically inspired parameterization of the encoding distribution (defining a population of point-neuron models). We believe that the application area where the results may be found to be particularly useful is stochastic neural coding (as neurophysiological properties of biological networks suggest a need of explicit constraints on the encoding distribution). For this case, we compared our variational approach with Brunel and Nadal’s Fisher approximation of mutual information (Brunel and Nadal (1998)). Moreover, we compared our approach with the recent results of Szummer and Jaakkola (2002) and Corduneanu and Jaakkola (2003), which may be seen as another approximation of a lower bound on $I(\mathbf{x}, \mathbf{y})$ formally generalizing on the criteria optimized by a variety of common popula-

tion coding methods (Pouget et al. (1998), Zhang and Sejnowski (1999), Bethge et al. (2002)). Our empirical results indicated that for the considered cases, our variational approach was most preferable (though the results here are rather preliminary, since in order to provide a sensible comparison of our bound with the existing non-bounding approximations, we have constrained our models to be low-dimensional). Additionally, we demonstrated that for the considered encoding distribution it was possible to derive a *local* learning rule, which may potentially be attractive from neuro-biological and computational perspectives (this also generalizes work of Linsker (1997), who had derived local approximations of *infomax* learning for invertible channels). In the future, more biologically realistic channels and applications should be considered. Additionally, for the specific biologically interesting channels, we are planning to conduct a tractability study of other lower bounds on mutual information known from statistical mechanics of supervised learning (Oppen and Haussler (1995), Haussler and Oppen (1997)) and compare them with our variational approach. Nevertheless, we can say that at this stage our current results illustrate potential advantages of the variational information-maximizing framework over the common (Brunel and Nadal (1998), Kang and Sompolinsky (2001)) and less common (Szummer and Jaakkola (2002), Corduneanu and Jaakkola (2003) and Section 6.2.1) approaches to population coding of high-dimensional input stimuli.

Finally, we considered a seemingly unrelated problem of lower bounding the normalizing constant (partition function) Z of an undirected graphical model (similar methods may be used to bound the likelihood of any probability distribution). Specifically, we introduced an auxiliary variable extension of any structured mean field theory, which could be useful in the context of approximate probabilistic inference. We showed that by considering a projection to the auxiliary variable space, and expressing the Kullback-Leibler divergence between the distributions defined over the augmented variable spaces, it was possible to define a tractable variational framework which leads to systematic improvements over the standard theories. While the auxiliary variational extension of standard approximation theories is of a potential interest as a general approach to approximate inference, it demonstrates a curious link to our variational information-maximizing framework. In particular, it turns out that the improvement of the proposed bound on $\log Z$ over a convex combination of standard bounds is given by a specific form of the generic lower bound on mutual information. The variational IM framework may therefore be seen as addressing an integral subgoal of auxiliary variational inference.

This part of our study suggests interesting directions for further research, including exploration of the effects which a choice of the mapping to the auxiliary space may have on the bounds for high-dimensional auxiliary spaces. Unfortunately, our current results suggest that the sparsity constraints on the structures of auxiliary conditional (variational decoder) distributions may be too restrictive. (For example, in the case when both the variational distribution $q(\mathbf{x}, \mathbf{y})$ and the auxiliary conditional $p(\mathbf{y}|\mathbf{x})$ have chain structures, the theoretically achievable gain produced by our approximation (over a simple theory) would only be high ($\sim O(|\mathbf{y}|)$) in situations when $q(y_i|y_{i-1}, x_i) \approx q(y_i|y_{i-1}, \mathbf{x})$, where \mathbf{y} and \mathbf{x}

are the auxiliary and latent variables respectively. This is effectively analogous to the cases when we may ignore future observations during inference in state-space models, which may be the case only for a limited number of applications). Nevertheless, we can formally show that by considering richer structures of the auxiliary conditional distributions we formally improve on simple approximations, modeling them as generally suboptimal special cases. Generally, modifying our framework for obtaining even better improvements would be an interesting and challenging direction for future research. We expect, however, that obtaining significantly better improvements over standard theories without sacrificing the rigorous bound may be difficult.

A further study of our auxiliary variational method for inference may also be motivated from a different perspective. Note that even in situations when the exact posterior may be accurately approximated by a factorized uni-modal distribution, finding such good factorized approximations by using the standard techniques may be complicated. While our results suggest that increasing the cardinality of the auxiliary variables, or using richer structures of variational decoders helps to produce better lower bounds on the partition function (by missing less of the probability mass), at this stage it remains largely unclear how helpful the auxiliary variables could be for preventing convergence of each component of the variational distribution to a suboptimal local mode of the exact posterior. A part of our motivation for further exploration of the method is the reported success of the auxiliary sampling techniques, which often result in efficient exploration of the distribution (by allowing potentially large changes in clusters of variables). Motivated by this observation, we are planning to continue studying our framework by empirically applying it to real-world tasks, as well as analyzing the theoretically achievable improvements over simple bounds for other choices of constraints on the auxiliary conditional distribution.

Summary

To summarize, we note that for the purpose of finding informative representations of the data, learning by maximizing the mutual information in encoder models is an alternative to learning probability models of the observations. In some cases, e.g. when specific tasks require explicit parameterization of the encoding distribution, information-maximization in encoder models may be the method of choice. It may be particularly attractive for clustering and constrained dimensionality reduction (for both deterministic and intrinsically noisy encoding mappings). An example of an intrinsically noisy channel is a model of a biological network (for example, a retinal encoder model defining a mapping from photoreceptors to ganglion representations in the retina), where a choice of constraints on encoder parameters is motivated by physiological properties of biological systems.

An important practical complication of maximizing mutual information in stochastic environments is high computational complexity of the exact optimization procedure, which is only tractable in a few simple special cases. In this thesis we suggested a general variational framework for maximizing a family of proper lower bounds on mutual information in stochastic environments. It turned out that in its simplest formulation our framework led to the same fixed points as the

approaches aiming to optimize the variational Jensen’s bounds on the likelihoods of stochastic self-supervised models; however, our method is computationally and representationally simpler than naive applications of the variational EM for such models. Additionally, our approach may be easily generalized to define richer types of bounds on mutual information which formally generalize simpler approximations. We believe that this makes our method potentially useful as a general variational framework for approximate information maximization.

We note, however, that the simple formulation of the suggested framework may in some cases be limited. Specifically, our empirical experience and intuition suggest that using parametric or structurally constrained variational decoders for approximating moments of the exact posterior $p(\mathbf{x}|\mathbf{y})$ may in some cases be too restrictive (for example, this is the case for error-correction in binary symmetric channels). This follows from the fact that the IM-optimal constrained variational decoder computes the marginal estimates which are good on average (for all possible codewords $\{\mathbf{y}\}$). Thus, for any *specific* codeword \mathbf{y} a constrained IM-optimal decoder may often be inferior to standard approximation theories. In future work it will be interesting to see whether or not using the IM framework for learning the optimal *encoder* may improve on standard inferences in the resulting models. Improving our framework to be competitive for error correction applications may be particularly interesting and challenging, as this is one case where a naive application of our approach is not expected to work well.

Despite the apparent limitations of our framework, we believe that the results presented in this work may be potentially interesting from several perspectives. First of all, in contrast to the majority of the existing approximations of $I(\mathbf{x}, \mathbf{y})$ (see e.g. Brunel and Nadal (1998), Kang and Sompolinsky (2001), Torkkola (2000), Gokcay and Principe (2002), Szummer and Jaakkola (2002), Corduneanu and Jaakkola (2003)), our method optimizes a proper lower bound, rather than a surrogate objective criterion or an approximation of $I(\mathbf{x}, \mathbf{y})$ (which may only be accurate under specific asymptotic assumptions, and weak or even undefined when the assumptions are violated). Secondly, the flexibility of the choice of the variational distribution makes it possible to generalize and improve other bounds on mutual information. For example, we may tractably extend our method to the family of auxiliary variational bounds on $I(\mathbf{x}, \mathbf{y})$, which may be used to improve on any simple generic approach without altering properties of the original channel. Moreover, in contrast to other approaches (e.g. Lawrence et al. (1998), Jaakkola and Jordan (1998), Brunel and Nadal (1998)), our method may be easily generalized to the structured context. Finally, we can demonstrate that by applying the IM framework to optimizing the bounds on $I(\mathbf{x}, \mathbf{y})$ in the augmented space of encoder and variational decoder parameters, we can often obtain simpler optimization procedures than those resulting from expressing the bounds as functions of the encoder alone. This leads to an interesting observation that in some cases the IM framework may be used to derive optimization procedures which only require local computations. This may be particularly attractive from the neuro-biological, but also from the computational perspectives.

Possibly the most important contribution of this work is a rigorous and general variational framework for maximizing the mutual information in intrinsically

intractable channels. We show that it leads to simple, stable, and easily generalizable optimization procedures, which (despite some intrinsic constraints) outperform and supersede many of the common approximate information-maximizing techniques.

Appendix A

Linsker's Bound: Centering in the Code Space

A brief comment about centering the data in the code space is in order. Consider the *as-if Gaussian* objective

$$I_G(\mathbf{x}, \mathbf{y}) \propto \log |\langle \mathbf{x}\mathbf{x}^T \rangle - \log \langle \mathbf{x}\mathbf{x}^T \rangle - \langle \mathbf{x}\mathbf{y}^T \rangle \langle \mathbf{y}\mathbf{y}^T \rangle^{-1} \langle \mathbf{y}\mathbf{x}^T \rangle| \quad (\text{A.1})$$

with source vectors \mathbf{x} and codes \mathbf{y} (see Section 2.2.2). Here all the averages are computed over $p(\mathbf{x}, \mathbf{y}) = \tilde{p}(\mathbf{x})p(\mathbf{y}|\mathbf{x})$, where $p(\mathbf{y}|\mathbf{x})$ is the encoder and $\tilde{p}(\mathbf{x})$ is the empirical distribution. In Section 2.2.2 we showed that (A.1) is in fact a special case of the generic lower bound on $I(\mathbf{x}, \mathbf{y})$

$$\tilde{I}(\mathbf{x}, \mathbf{y}) = \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{\tilde{p}(\mathbf{x})p(\mathbf{y}|\mathbf{x})} + \text{const} \quad (\text{A.2})$$

for the variational decoder $q(\mathbf{x}|\mathbf{y}) \sim \mathcal{N}_{\mathbf{x}}(\mathbf{U}\mathbf{y}, \mathbf{\Sigma})$, where we assumed $\langle \mathbf{x} \rangle_{p(\mathbf{x})} = 0$ and $\langle \mathbf{y} \rangle_{p(\mathbf{y})} = 0$. In practice, centering the source vectors $\{\mathbf{x}\}$ should not be problematic, since transformations of the sources are not explicitly affected by parameterization of the encoder $p(\mathbf{y}|\mathbf{x})$. In some cases, $\langle \mathbf{x} \rangle_{p(\mathbf{x})} = 0 \Rightarrow \langle \mathbf{y} \rangle_{p(\mathbf{y})} = 0$ [for example, this is the case for linear Gaussian encoders $p(\mathbf{y}|\mathbf{x}) \sim \mathcal{N}_{\mathbf{y}}(\mathbf{W}\mathbf{x}, \sigma^2 \mathbf{I}_{|\mathbf{y}|})$]. However, in general the assumption of the centered *encodings* $\{\mathbf{y}\}$ implies additional constraints on the encoder distribution $p(\mathbf{y}|\mathbf{x})$, which may in general be difficult to enforce.

We will now consider non-centered codes (i.e. $\langle \mathbf{y} \rangle_{p(\mathbf{y})} \neq 0$) and a more general form of the linear Gaussian variational decoder, namely $q(\mathbf{x}|\mathbf{y}) \sim \mathcal{N}_{\mathbf{x}}(\mathbf{U}\mathbf{y} + \boldsymbol{\mu}, \mathbf{\Sigma})$. A straight-forward substitution into (A.2) then results in

$$\tilde{I}(\mathbf{x}, \mathbf{y}) = -\frac{1}{2} \langle \text{tr} \{ \mathbf{\Sigma}^{-1} (\mathbf{x} - \mathbf{U}\mathbf{y} - \boldsymbol{\mu})(\mathbf{x} - \mathbf{U}\mathbf{y} - \boldsymbol{\mu})^T \} \rangle_{p(\mathbf{x})p(\mathbf{y}|\mathbf{x})} - \frac{1}{2} \log |\mathbf{\Sigma}| + \text{const}. \quad (\text{A.3})$$

Then by maximizing $\tilde{I}(\mathbf{x}, \mathbf{y})$ for $\boldsymbol{\mu} \in \mathbb{R}^{|\mathbf{x}|}$ we obtain $\boldsymbol{\mu} = -\mathbf{U}\langle \mathbf{y} \rangle$. By comparing the mean of the reconstructed sources (computed for the variational distribution $q(\mathbf{x}|\mathbf{y})$) with the empirical mean of the source vectors, we easily see that $\langle \mathbf{x} \rangle_{q(\mathbf{x}|\mathbf{y})p(\mathbf{y})} = \mathbf{U}\langle \mathbf{y} \rangle_{p(\mathbf{y})} + \boldsymbol{\mu} = \langle \mathbf{x} \rangle_{p(\mathbf{x})} = 0$, i.e. the variational decoder does not introduce a bias into the reconstructions. By substituting $\boldsymbol{\mu} = \mathbf{U}\langle \mathbf{y} \rangle \in \mathbb{R}^{|\mathbf{y}|}$ and solving for $\mathbf{U} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{y}|}$ and $\mathbf{\Sigma} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|}$, we obtain

$$\tilde{I}(\mathbf{x}, \mathbf{y}) = -\log |\mathbf{\Sigma}_{xx} - \mathbf{\Sigma}_{xy} \mathbf{\Sigma}_{yy}^{-1} \mathbf{\Sigma}_{yx}|, \quad (\text{A.4})$$

where we ignored the irrelevant constants and let $\Sigma_{xx} = \langle \mathbf{x}\mathbf{x}^T \rangle$, $\Sigma_{xy} = \Sigma_{yx}^T = \langle \mathbf{x}\mathbf{y}^T \rangle$, and $\Sigma_{yy} = \langle \mathbf{y}\mathbf{y}^T \rangle - \langle \mathbf{y} \rangle \langle \mathbf{y} \rangle^T$. This is exactly the general form of Linsker's *as-if Gaussian* criterion (see expression (1.16)).

Appendix B

Analysis of Variational Information Maximization

Here we outline several simple relations of the variational lower bound on mutual information to the standard variational bounds on the likelihood and conditional likelihood. Section B.1 discusses the case of conditional likelihood training in general feed-forward models. A curious result derived there is a simple decoder-specific upper bound on the conditional log-likelihood. Section B.2 discusses a relation of the generic lower bound on $I(x, y)$ for noiseless or Gaussian channels to the standard bound on the likelihood for flat mixtures. Effectively, it illustrates a sufficient condition for the variational EM and IM algorithms to converge to identical fixed points.

B.1 Variational Information Maximization and Feed-Forward Models

Here we will prove proposition 3.3. By analogy with Section 3.3.1, we will define the feed-forward and the conditional encoder models as

$$\mathcal{M}_C \stackrel{\text{def}}{=} p(y|x)p(\tilde{x}|y), \quad \mathcal{M}_{IC} \stackrel{\text{def}}{=} \tilde{p}(x, \tilde{x})p(y|x, \tilde{x}), \quad (\text{B.1})$$

where $p(y|x, \tilde{x})$ is the exact posterior of \mathcal{M}_C .

Proposition B.1. *For i.i.d. patterns $\{x, \tilde{x}\}$, conditional likelihood learning in the feed-forward model \mathcal{M}_C corresponds to maximization of a **lower bound** on the conditional mutual information $I(\tilde{x}, y|x)$ in \mathcal{M}_{IC} . Up to irrelevant constants, this bound is weaker or as tight as $\hat{I}_C(\tilde{x}, y|x) = \langle \log p(\tilde{x}|x, y) \rangle_{p(y|x, \tilde{x})\tilde{p}(x, \tilde{x})}$.*

Proof. It is easy to see that for i.i.d. data, the average conditional log-likelihood (3.29) may be expressed as

$$\begin{aligned} \mathcal{L}_{\tilde{x}|x} &= \langle \log p(\tilde{x}|x) \rangle_{\tilde{p}(\tilde{x}, x)} \\ &= \langle \log p(y|x) + \log p(\tilde{x}|x, y) - \log p(y|x, \tilde{x}) \rangle_{\tilde{p}(x, \tilde{x})}, \end{aligned} \quad (\text{B.2})$$

where $\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}})$ is the empirical distribution, and \mathbf{y} is a latent variable in the conditional $p(\tilde{\mathbf{x}}, \mathbf{y}|\mathbf{x})$. By averaging both parts of (B.2) over the exact posterior $p(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}})$, we obtain

$$\mathcal{L}_{\tilde{\mathbf{x}}|\mathbf{x}} = \langle \log p(\tilde{\mathbf{x}}|\mathbf{x}, \mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}})\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}})} - \langle KL(p(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}})||p(\mathbf{y}|\mathbf{x})) \rangle_{\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}})}. \quad (\text{B.3})$$

We will assume that $p(\tilde{\mathbf{x}}, \mathbf{y}|\mathbf{x})$ is in the tractable family, i.e. no variational relaxations of the conditional likelihood are required to perform the computations in (B.3). Note that the Kullback-Leibler divergence in (B.3) cancels if and only if $p(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}}) \equiv p(\mathbf{y}|\mathbf{x})$, which is generally not the case for stochastic feed-forward models.

By definition, the conditional mutual information $I(\tilde{\mathbf{x}}, \mathbf{y}|\mathbf{x})$ in the corresponding recognition model \mathcal{M}_{IC} may be expressed as

$$\begin{aligned} I(\tilde{\mathbf{x}}, \mathbf{y}|\mathbf{x}) &= H_{\tilde{p}}(\tilde{\mathbf{x}}|\mathbf{x}) + \langle \log \tilde{p}(\tilde{\mathbf{x}}|\mathbf{y}, \mathbf{x}) \rangle_{\tilde{p}(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}})\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}})} \\ &= H_{\tilde{p}}(\tilde{\mathbf{x}}|\mathbf{x}) + \langle \log \tilde{p}(\tilde{\mathbf{x}}|\mathbf{y}, \mathbf{x}) \rangle_{p(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}})\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}})}, \end{aligned} \quad (\text{B.4})$$

where $\tilde{p}(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}}) = p(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}})$ by construction (3.30), $H_{\tilde{p}}(\tilde{\mathbf{x}}|\mathbf{x}) \stackrel{\text{def}}{=} -\langle \log \tilde{p}(\tilde{\mathbf{x}}|\mathbf{x}) \rangle_{\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}})}$ is the entropy of the empirical distribution (independent of the functional parameters), and $\tilde{p}(\tilde{\mathbf{x}}|\mathbf{y}, \mathbf{x})$ is the Bayes-optimal decoder of \mathcal{M}_{IC} given by

$$\tilde{p}(\tilde{\mathbf{x}}|\mathbf{y}, \mathbf{x}) \propto p(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}})\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}}). \quad (\text{B.5})$$

Here $p(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}})$ is the exact posterior of the feed-forward model \mathcal{M}_C . Again, it is easy to see that due to the mixture form of the Bayesian decoder $\tilde{p}(\tilde{\mathbf{x}}|\mathbf{y}, \mathbf{x})$, the exact evaluations of the conditional mutual information $I(\tilde{\mathbf{x}}, \mathbf{y}|\mathbf{x})$ are in general computationally intractable. However, we may re-define the generic lower bound on the mutual information (2.2) for the conditional case, and obtain a tractable bound on (B.4) as

$$I(\tilde{\mathbf{x}}, \mathbf{y}|\mathbf{x}) \geq \hat{I}_C(\tilde{\mathbf{x}}, \mathbf{y}|\mathbf{x}) \stackrel{\text{def}}{=} H_{\tilde{p}}(\tilde{\mathbf{x}}|\mathbf{x}) + \langle \log p(\tilde{\mathbf{x}}|\mathbf{x}, \mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}})\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}})}, \quad (\text{B.6})$$

which is saturated if and only if $p(\tilde{\mathbf{x}}|\mathbf{x}, \mathbf{y}) \equiv \tilde{p}(\tilde{\mathbf{x}}|\mathbf{x}, \mathbf{y})$. For the considered case of tractable distributions $p(\tilde{\mathbf{x}}, \mathbf{y}|\mathbf{x})$, the bound (B.6) will also be tractable. Moreover, from the non-negativity of the KL divergence we get

$$\hat{I}_C(\tilde{\mathbf{x}}, \mathbf{y}|\mathbf{x}) \geq H_{\tilde{p}}(\tilde{\mathbf{x}}|\mathbf{x}) + \langle \log p(\tilde{\mathbf{x}}|\mathbf{x}, \mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}})\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}})} - \langle KL(p(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}})||p(\mathbf{y}|\mathbf{x})) \rangle_{\tilde{p}(\mathbf{x}, \tilde{\mathbf{x}})}, \quad (\text{B.7})$$

which by transitivity leads to

$$I(\tilde{\mathbf{x}}, \mathbf{y}|\mathbf{x}) \geq \hat{I}_C(\tilde{\mathbf{x}}, \mathbf{y}|\mathbf{x}) \geq H_{\tilde{p}}(\tilde{\mathbf{x}}|\mathbf{x}) + \mathcal{L}_{\tilde{\mathbf{x}}|\mathbf{x}}. \quad (\text{B.8})$$

□

The result shows that the conditional likelihood training in chains \mathcal{M}_C may be viewed as maximization of a specific lower bound on $I(\tilde{\mathbf{x}}, \mathbf{y}|\mathbf{x})$ in the corresponding conditional encoder models \mathcal{M}_{IC} . Alternatively (and perhaps more intuitively), we may view the result as *an upper bound* on $\mathcal{L}_{\tilde{\mathbf{x}}|\mathbf{x}}$. From proposition 3.3 we may

immediately obtain an upper bound on the average log-probability of producing the required outputs $\tilde{\mathbf{x}}$ from the encodings \mathbf{y} with the chosen model-specific decoder $p(\tilde{\mathbf{x}}|\mathbf{y})$:

$$\mathcal{L}_{\tilde{\mathbf{x}}|\mathbf{x}} = \frac{1}{M} \sum_{i=1}^M \log p(\tilde{\mathbf{x}} = \tilde{\mathbf{x}}^{(i)} | \mathbf{x} = \mathbf{x}^{(i)}) \leq \underbrace{\frac{1}{M} \sum_{i=1}^M \langle \log p(\tilde{\mathbf{x}}^{(i)} | \mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x}^{(i)}, \tilde{\mathbf{x}}^{(i)})}}_{\hat{I}_C(\tilde{\mathbf{x}}, \mathbf{y}|\mathbf{x}) + \text{const}}, \quad (\text{B.9})$$

where $(\mathbf{x}^{(i)}, \tilde{\mathbf{x}}^{(i)}) \in \mathcal{X}_C$ is the i^{th} source-output pair of the training set $\mathcal{X}_C = \{(\mathbf{x}^{(i)}, \tilde{\mathbf{x}}^{(i)}) | i = 1, \dots, M\}$. For the case when $\tilde{\mathbf{x}}^{(i)} = \mathbf{x}^{(i)}$ for all M training patterns, the feed-forward model \mathcal{M}_C becomes an autoencoder, and the conditional likelihood defines the probability of correct reconstructions of the training set. Note that the bound (B.9) is fundamentally different from Fano’s inequality (3.1) in the sense that (B.9) is a functional of the decoder distribution $p(\tilde{\mathbf{x}}|\mathbf{y})$, while Fano’s result (3.1) ignores any knowledge about the specific distributions used for source reconstructions. Moreover, whilst Fano’s bound on reconstruction error is tractable only for simple channels where it is possible to compute $I(\mathbf{x}, \mathbf{y}) = H(\mathbf{y}) - H(\mathbf{y}|\mathbf{x})$ exactly, our result (B.9) does not require computations of entropies of mixture distributions.

B.2 Variational Information Maximization and Flat Mixture Models

Here we compare variational information maximization (IM) and variational expectation maximization (EM) for flat latent variable models (i.e. generative models where the latent variables are uniformly distributed).

Let $\mathcal{M}_L \stackrel{\text{def}}{=} q(\mathbf{y})q(\mathbf{x}|\mathbf{y})$ define a latent variable model, where \mathbf{x} is a data pattern, and \mathbf{y} is its latent variable representation. (NB: we change the conventional notations for the purpose which will soon become clear). Following the standard variational procedure of maximizing the log-likelihood $\mathcal{L} \stackrel{\text{def}}{=} \langle \log q(\mathbf{x}) \rangle_{\tilde{p}(\mathbf{x})}$ in the generative model \mathcal{M}_L , it is straight-forward to obtain the standard variational Jensen’s bound on \mathcal{L}

$$\begin{aligned} \mathcal{L} &= \left\langle \log \int_{\mathbf{y}} q(\mathbf{x}, \mathbf{y}) d\mathbf{y} \right\rangle_{\tilde{p}(\mathbf{x})} \\ &= \left\langle \log \int_{\mathbf{y}} q(\mathbf{x}, \mathbf{y}) \frac{p(\mathbf{y}|\mathbf{x})}{p(\mathbf{y}|\mathbf{x})} d\mathbf{y} \right\rangle_{\tilde{p}(\mathbf{x})} \\ &\geq \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})} - \langle KL(p(\mathbf{y}|\mathbf{x}) \| q(\mathbf{y})) \rangle_{\tilde{p}(\mathbf{x})} \stackrel{\text{def}}{=} \tilde{\mathcal{L}}, \end{aligned} \quad (\text{B.10})$$

where we have applied Jensen’s inequality (e.g. Jensen (1906), Hardy et al. (1988)). Here $\tilde{p}(\mathbf{x})$ is the empirical distribution, and $p(\mathbf{y}|\mathbf{x})$ is a variational distribution approximating the true posterior $q(\mathbf{y}|\mathbf{x}) \propto q(\mathbf{y})q(\mathbf{x}|\mathbf{y})$ expressed from the generative model \mathcal{M}_L . As usual, in order to simplify computations of (B.10), the variational posterior $p(\mathbf{y}|\mathbf{x})$ needs to be constrained to ensure tractability of

computing the averages. The standard variational extension of the expectation-maximizing algorithm (see e.g. Neal and Hinton (1998)) trains the generative model \mathcal{M}_L by iteratively optimizing (B.10) with respect to the model parameters $q(\mathbf{y})$, $q(\mathbf{x}|\mathbf{y})$, and the variational posterior $p(\mathbf{y}|\mathbf{x})$.

Note that we can equivalently express the bound (B.10) as

$$\tilde{\mathcal{L}} = \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})} + \langle \log q(\mathbf{y}) \rangle_{p(\mathbf{y}|\mathbf{x})\tilde{p}(\mathbf{x})} + \langle H(p(\mathbf{y}|\mathbf{x})) \rangle_{\tilde{p}(\mathbf{x})}, \quad (\text{B.11})$$

where $H(p(\mathbf{y}|\mathbf{x})) \stackrel{\text{def}}{=} -\langle \log p(\mathbf{y}|\mathbf{x}) \rangle_{p(\mathbf{y}|\mathbf{x})}$. It is now straight-forward to see that for the case of uniform code distributions ($q(\mathbf{y}) \sim \mathcal{U}_{\mathbf{y}}$) with deterministic variational posteriors ($p(\mathbf{y}|\mathbf{x}) \sim \delta$), maximization of the bound (B.11) reduces to optimizing

$$\tilde{\mathcal{L}} = \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{\tilde{p}(\mathbf{x})p(\mathbf{y}|\mathbf{x})} + \text{const.} \quad (\text{B.12})$$

Up to irrelevant constants, the bound on the likelihood for this case is in fact the generic lower bound (2.2) on the mutual information for the noiseless channel model $\mathcal{M}_I(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} \tilde{p}(\mathbf{x})p(\mathbf{y}|\mathbf{x})$, which uses the variational posterior of $\tilde{\mathcal{L}}$ as the *encoder* distribution, i.e.

$$\tilde{I}(\mathbf{x}, \mathbf{y}) = H_{\tilde{p}(\mathbf{x})} + \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{\tilde{p}(\mathbf{x})p(\mathbf{y}|\mathbf{x})}, \quad \text{where } p(\mathbf{y}|\mathbf{x}) \sim \delta. \quad (\text{B.13})$$

Here the variational decoder $q(\mathbf{x}|\mathbf{y})$ of the bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ is the conditional of the latent variable model \mathcal{M}_L . From (B.12) and (B.13), we can see that *for i.i.d. patterns $\{\mathbf{x}\}$, maximization of the standard variational lower bound on the likelihood in a **flat** latent variable model with a **deterministic** variational posterior $p(\mathbf{y}|\mathbf{x})$ reduces to maximizing the generic lower bound on the mutual information in the corresponding **noiseless channel**.*

Note that the same result holds if $p(\mathbf{y}|\mathbf{x})$ is a Gaussian with a fixed covariance. Indeed, if $p(\mathbf{y}|\mathbf{x}) \sim \mathcal{N}_{\mathbf{y}}(\mathbf{f}(\mathbf{x}), \mathbf{\Sigma})$, the entropic term in (B.11) is a function of the fixed noise covariance $\mathbf{\Sigma}$, i.e. up to irrelevant constants the bound $\tilde{\mathcal{L}}$ defined by (B.11) is identical to the generic lower bound on the mutual information (B.13) for the corresponding Gaussian channel.

The result may be trivially generalized to flat mixtures of latent variable models $\mathcal{M}_{LH}(\mathbf{x}, \mathbf{y}, z) \stackrel{\text{def}}{=} q(\mathbf{y})q(z)q(\mathbf{x}|\mathbf{y}, z)$ with deterministic variational posteriors $p(\mathbf{x}|\mathbf{y}, z) \sim \delta$ and flat priors on the latent variables $q(\mathbf{y}) \sim \mathcal{U}_{\mathbf{y}}$ and mixture coefficients $q(z) = 1/|z|$ (here we assumed that $\mathbf{x} \in \mathbb{R}^{|\mathbf{x}|}$, $\mathbf{y} \in \mathbb{R}^{|\mathbf{y}|}$, and $z \in \{1, \dots, |z|\}$). By analogy with (B.11)–(B.13), it is easy to see that optimization of the standard variational lower bound on \mathcal{L} for this case reduces to maximizing the generic lower bound on $I(\mathbf{x}, \{\mathbf{y}, z\})$ in the corresponding **noiseless hybrid channel** $\mathbf{x} \rightarrow \{\mathbf{y}, z\}$

$$\tilde{I}_H(\mathbf{x}, \{\mathbf{y}, z\}) = H(\mathbf{x}) + \langle \log q(\mathbf{x}|\mathbf{y}, z) \rangle_{\tilde{p}(\mathbf{x})p(\mathbf{y}|\mathbf{x})p(z|\mathbf{x})}, \quad (\text{B.14})$$

where $p(\mathbf{y}|\mathbf{x}) \sim \delta$, $p(z|\mathbf{x}) \sim \delta$, and $q(\mathbf{x}|\mathbf{y}, z)$ is the conditional of the generative mixture of latent variable models $\mathcal{M}_{LH}(\mathbf{x}, \mathbf{y}, z)$. Similar results hold if $p(z|\mathbf{x}) \sim \delta$ and $p(\mathbf{y}|\mathbf{x}) \sim \mathcal{N}_{\mathbf{y}}(\mathbf{f}(\mathbf{x}), \mathbf{\Sigma})$ where $\mathbf{\Sigma}$ is fixed. In the encoder model formulation, the fixed covariance corresponds to the irreducible independent channel noise.

To summarize, we may note that *while generally the optimization surfaces defined by the standard variational lower bounds on the likelihood and the generic*

variational lower bounds on the mutual information are different, the obtained solutions coincide for flat mixture models \mathcal{M}_L with deterministic variational posteriors, and the corresponding noiseless encoder models \mathcal{M}_I (where the variational decoding distribution corresponds to the conditional of the generative model).

Appendix C

Variational Information Maximization for Isotropic Gaussian Channels

Here we analyze properties of the variational IM algorithm for isotropic Gaussian channels $\mathbf{x} \rightarrow \mathbf{y}$. Specifically, we focus on the theoretical analysis of the bounds on $I(\mathbf{x}, \mathbf{y})$ for *nonlinear* Gaussian encoders $p(\mathbf{y}|\mathbf{x}) \sim \mathcal{N}_{\mathbf{y}}(\boldsymbol{\phi}(\mathbf{x}), \sigma^2\mathbf{I})$. As usual, we will presume intractability of explicit computations in $\mathbb{R}^{|\phi|}$, and derive a kernelized extension of the generic lower bound

$$I(\mathbf{x}, \mathbf{y}) \geq \tilde{I}(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} H(\mathbf{x}) + \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{x}, \mathbf{y})}.$$

As the variational information-maximizing framework offers flexibility in choosing optimal variational decoders $q(\mathbf{x}|\mathbf{y})$, we consider optimizing $\tilde{I}(\mathbf{x}, \mathbf{y})$ for several of such choices. As the first obvious choice, we consider using linear Gaussian decoders, which simplifies computations of the bound and facilitates the analysis of optimal solutions for the encoder parameters. Our motivation here is to check whether there is any gain in using nonlinearities of the codes if the variational decoder distribution remains linear. As expected, we show that linear Gaussian variational decoders may indeed be too restrictive. Specifically, we show that in the case of the isotropic noise of Gaussian channels, nothing is gained by using nonlinear encoders and linear decoders in the context of variational information maximization. This agrees with the more specific result for noiseless one-layer autoencoders with linear output units (Boulevard and Kamp (1988), Boulevard (2000)), but is derived for *stochastic* channels with arguably less specific choices of nonlinearities.

Then we consider variational information maximization for nonlinear variational decoders. Since the choice of nonlinearity may significantly influence the complexity of computing variational lower bounds $\tilde{I}(\mathbf{x}, \mathbf{y})$, relaxations of the generic variational procedure may need to be considered. We show that the generic lower bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ may indeed be formally modified to ensure tractable computations. This naturally relaxes the bound on $I(\mathbf{x}, \mathbf{y})$ to provide an objective function for Kernel PCA. By analogy with a simpler case of discrete nonlinear channels (see Section 5.2.3), the information-theoretic formulation suggests a

proper way for learning kernel parameters of KPCA models; however, the choice of constraints on the kernel matrices will prove to be crucial for avoiding degenerate (asymptotically noiseless) solutions.

Throughout the discussion in this chapter, we will make several references to the optimal bound on $I(\mathbf{x}, \mathbf{y})$ for isotropic linear Gaussian decoders $q(\mathbf{x}|\mathbf{y}) \sim \mathcal{N}_x(\mathbf{U}\mathbf{y}, \sigma^2\mathbf{I})$. This is a more constrained variational decoder than the correlated linear Gaussian (which gives rise to Linsker's *as-if Gaussian* approximation, see Section 1.4, 4.1.1). However, the solutions obtained for this case will prove to be important for considering nonlinear generalizations of the bound $\tilde{I}(\mathbf{x}, \mathbf{y})$, so we will discuss them first. Then in Appendix C.2 we will move to a discussion of nonlinear Gaussian channels.

C.1 Linear Gaussian Channels: Linear Decoders

$$p(\mathbf{y}|\mathbf{x}) \sim \mathcal{N}_y(\mathbf{W}\mathbf{x}, s^2\mathbf{I}), \quad q(\mathbf{x}|\mathbf{y}) \sim \mathcal{N}_x(\mathbf{U}\mathbf{y}, \sigma^2\mathbf{I})$$

Let the encoder and decoder be given by $p(\mathbf{y}|\mathbf{x}) \sim \mathcal{N}_y(\mathbf{W}\mathbf{x}, s^2\mathbf{I})$ and $q(\mathbf{x}|\mathbf{y}) \sim \mathcal{N}_x(\mathbf{U}\mathbf{y}, \sigma^2\mathbf{I})$ respectively. Our goal is to learn optimal settings of $\mathbf{W} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$ and $\mathbf{U} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{y}|}$ (for fixed decoder and encoder variances σ^2 and s^2) by maximizing the simple variational bound (4.1) on the mutual information, which in this case is expressed as

$$\tilde{I}(\mathbf{x}, \mathbf{y}) = \frac{1}{\sigma^2} \text{tr} \{ \mathbf{U}\mathbf{W}\mathbf{S} \} - \frac{1}{2\sigma^2} \text{tr} \{ \mathbf{U}\mathbf{\Sigma}\mathbf{U}^T \} - \frac{1}{2s^2} \text{tr} \{ \mathbf{S} \} + c. \quad (\text{C.1})$$

Here $c = -|\mathbf{x}|/2 \log(2\pi\sigma^2)$ is a constant which does not affect the optimization surface for the encoder and decoder weights,

$$\mathbf{S} = \langle \mathbf{x}\mathbf{x}^T \rangle = \frac{1}{M} \sum_m \mathbf{x}^{(m)} (\mathbf{x}^{(m)})^T \quad (\text{C.2})$$

is the sample covariance of the centered data, and

$$\mathbf{\Sigma} = \mathbf{I}s^2 + \mathbf{W}\mathbf{S}\mathbf{W}^T \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{y}|} \quad (\text{C.3})$$

is the covariance of the marginal distribution of responses $p(\mathbf{y})$. In the following we assume that the weights \mathbf{W} and \mathbf{U} are non-singular, which in our context means that the dimensionality of the codewords $|\mathbf{y}|$ is sufficiently low. Note that we make no parametric assumptions about the distribution of the sources $p(\mathbf{x})$ or non-singularity of the sample covariance \mathbf{S} .

Unsurprisingly, objective (C.1) is closely related to the least squares reconstruction error in a linear autoencoder. What is interesting in this context that it provides a proper bound on the mutual information, independently of the source distribution.

C.1.1 Nature of Optimal Solutions

Unconstrained optimization of (C.1) for the encoder's weights \mathbf{W} leads to the extremum condition

$$\mathbf{U}^T \mathbf{S} = \mathbf{U}^T \mathbf{U} \mathbf{W} \mathbf{S}. \quad (\text{C.4})$$

By assuming that \mathbf{y} is a compressed representation of the source \mathbf{x} (i.e. $|\mathbf{x}| > |\mathbf{y}|$), we obtain $\mathbf{W} \mathbf{S} = (\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T \mathbf{S} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{x}|}$. (It is safe to assume non-singularity; indeed, if $\mathbf{U}^T \mathbf{U}$ is singular, we may ensure non-singularity by reducing dimensionality of the codewords). Then it is clear that

$$\begin{aligned} \text{tr} \{ \mathbf{U} \mathbf{W} \mathbf{S} \mathbf{W}^T \mathbf{U}^T \} &= \text{tr} \{ \mathbf{U} \mathbf{W} \mathbf{S} \mathbf{W}^T \mathbf{U}^T \} = \text{tr} \{ (\mathbf{U}^T \mathbf{U} \mathbf{W} \mathbf{S}) \mathbf{W}^T \} \\ &= \text{tr} \{ \mathbf{U} \mathbf{W} \mathbf{S} \} = \text{tr} \{ \mathbf{U} (\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T \mathbf{S} \}, \end{aligned} \quad (\text{C.5})$$

leading to

$$\text{tr} \{ \mathbf{U} \boldsymbol{\Sigma} \mathbf{U}^T \} = \text{tr} \{ \mathbf{U} (\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T \mathbf{S} \} + s^2 \text{tr} \{ \mathbf{U} \mathbf{U}^T \}. \quad (\text{C.6})$$

By substituting expression (C.6) into (C.1), we may express the bound as a function of the decoder alone as

$$\tilde{I}(\mathbf{x}, \mathbf{y}) = \frac{1}{2\sigma^2} \text{tr} \{ \mathbf{U} (\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T \mathbf{S} \} - \frac{s^2}{2\sigma^2} \text{tr} \{ \mathbf{U} \mathbf{U}^T \} - \frac{1}{2\sigma^2} \text{tr} \{ \mathbf{S} \}, \quad (\text{C.7})$$

where we have ignored the irrelevant constant c .

Let $\mathbf{U} = \mathbf{V} \mathbf{L} \mathbf{R}^T$ be the singular value decomposition of the decoder weights \mathbf{U} (see e.g. Golub and Loan (1996)). By definition, $\mathbf{L} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{y}|}$ is diagonal and invertible, and $\mathbf{V} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{y}|}$, $\mathbf{R} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{y}|}$ are orthogonal, so that $\mathbf{V}^T \mathbf{V} = \mathbf{R}^T \mathbf{R} = \mathbf{R} \mathbf{R}^T = \mathbf{I}_{|\mathbf{y}|}$. From (C.4) it is clear that $\mathbf{W} = \mathbf{R} \mathbf{L}^{-1} \mathbf{V}^T$, i.e. $\mathbf{W} \mathbf{U} = \mathbf{I}_{|\mathbf{y}|}$. Substitution into (C.7) results in

$$\tilde{I}(\mathbf{x}, \mathbf{y}) \propto \text{tr} \{ \mathbf{V}^T \mathbf{S} \mathbf{V} \} - s^2 \text{tr} \{ \mathbf{L}^2 \} - \text{tr} \{ \mathbf{S} \}, \quad (\text{C.8})$$

which is a function of the decoder's singular vectors and scalings. Optimizing (C.8) for \mathbf{V} under the orthonormality constraints on \mathbf{V} , we readily obtain the PCA solution (and its rigid rotations). It is clear that in the case of a noiseless channel ($s^2 \rightarrow 0$), the solutions are invariant with respect to the scalings \mathbf{L} .

For noisy channels (i.e. for $s^2 > 0$), maximization of (C.8) with respect to the scalings of the *decoder* weights results in $\mathbf{L} \rightarrow 0$, which leads to the divergence of the Frobenius norm $\|\mathbf{W}\|_F$. Effectively, this corresponds to an approximately noiseless encoder, since distinct source patterns will be mapped to codes which are infinitely spread out in the code space, so that the contribution of the finite channel noise to the bound (C.1) becomes negligible. In order to find the optimal bases of the subspaces spanned by the encoder and decoder weights \mathbf{W} and \mathbf{U} , we may constrain their singular values. It is straight-forward to see that in this case the right singular vectors of the optimal encoder weights \mathbf{W} span the principal subspace of \mathbf{S} (see e.g. Bishop (1995)). It is also clear that the obtained solutions are invariant under complementary rotations of \mathbf{W} and \mathbf{U} .

This demonstrates the simple result that for the considered variational problem, the optimal lower bound on $I(\mathbf{x}, \mathbf{y})$ is provided by PCA. Furthermore, by

comparing this result with Linsker’s *as-if* Gaussian bound (see Section 4.1.1), we see that for the considered channel nothing is gained by learning full covariances of linear Gaussian variational decoders. As before, this conclusion is reached without the need for a Gaussian assumption about the source distribution. It is intuitively clear that in order to go beyond the PCA solutions, more complex encoders/decoders are required. In the following section we consider the effects of increasing the complexity of the encoder, while still using a simple linear Gaussian variational decoder defined above.

C.2 Nonlinear Gaussian Channels: Linear Decoders

$$p(y|x) \sim \mathcal{N}_y(\mathbf{W}\phi(x), \Sigma_y), \quad q(x|y) \sim \mathcal{N}_x(\mathbf{U}y, \Sigma_x)$$

Here we consider the case of a nonlinear Gaussian channel with the encoder $p(y|x) \sim \mathcal{N}_y(\mathbf{W}\phi(x), \Sigma_y)$ and decoder $q(x|y) \sim \mathcal{N}_x(\mathbf{U}y, \Sigma_x)$. Note that for all the data points $\{\mathbf{x}^{(m)} | m = 1, \dots, M\}$, the set of encodings $\{\mathbf{y}^{(m)}\}$ is given by a noisy linear projection from the (potentially high-dimensional) feature space $\{\phi(\mathbf{x}^{(m)})\}$. In what follows we assume that $|\phi| > M$, i.e. dimensionality of the feature space exceeds the number of training patterns.

It is easy to see that for the considered case the lower bound on the mutual information $I(\mathbf{x}, \mathbf{y})$ is given by

$$\begin{aligned} \tilde{I}(\mathbf{x}, \mathbf{y}) &= -\frac{1}{2} \text{tr} \{ \Sigma_x^{-1} \mathbf{S} \} + \text{tr} \{ \Sigma_x^{-1} \mathbf{U} \mathbf{W} \langle \phi(\mathbf{x}) \mathbf{x}^T \rangle \} \\ &\quad - \frac{1}{2} \text{tr} \{ \mathbf{U}^T \Sigma_x^{-1} \mathbf{U} (\Sigma_y + \mathbf{W} \langle \phi(\mathbf{x}) \phi(\mathbf{x})^T \rangle \mathbf{W}^T) \}. \end{aligned} \quad (\text{C.9})$$

As before, we assume that $\mathbf{S} = \langle \mathbf{x} \mathbf{x}^T \rangle = \sum_m \mathbf{x}^{(m)} (\mathbf{x}^{(m)})^T / M$ is the sample covariance of the zero-mean data, and the averages are computed with respect to the empirical distribution $\tilde{p}(\mathbf{x}) = (1/M) \sum_{m=1}^M \delta(\mathbf{x} - \mathbf{x}^{(m)})$. Also, by analogy with Section 5.2.3 we assume that for high-dimensional feature spaces direct evaluation of the averages is implausible. It is therefore desirable to avoid explicit computations in $\{\phi\}$.

C.2.1 Kernelized Representation

Since each row $\tilde{\mathbf{w}}_i^T \in \mathbb{R}^{1 \times |\phi|}$ of the weight matrix $\mathbf{W} \in \mathbb{R}^{|\mathcal{Y}| \times |\phi|}$ has the same dimensionality as the feature vectors $\phi(\mathbf{x}^{(i)})^T$, it is representable as

$$\tilde{\mathbf{w}}_i = \sum_{m=1}^M \alpha_{im} \phi(\mathbf{x}^{(m)}) + \tilde{\mathbf{w}}_i^\perp, \quad (\text{C.10})$$

where $\tilde{\mathbf{w}}_i^\perp$ is orthogonal to the span of $\phi(\mathbf{x}^{(1)}), \dots, \phi(\mathbf{x}^{(M)})$. Then

$$\mathbf{W} = \mathbf{A} \mathbf{F}^T + \mathbf{W}^\perp, \quad \mathbf{F} \stackrel{\text{def}}{=} [\phi(\mathbf{x}^{(1)}), \dots, \phi(\mathbf{x}^{(M)})] \in \mathbb{R}^{|\phi| \times M}, \quad (\text{C.11})$$

where $\mathbf{A} = \{\alpha_{ij}\} \in \mathbb{R}^{|\mathcal{Y}| \times M}$ is the matrix of coefficients, and the transposed rows of \mathbf{W}^\perp are given by $\tilde{\mathbf{w}}_i^\perp$. In kernel literature \mathbf{F} is often referred to as the *design* matrix (e.g. Williams (1998), MacKay (1997), Cristianini and Shawe-Taylor (2000)).

From (C.11), we obtain expressions for the averages

$$\begin{aligned} \mathbf{W} \langle \phi(\mathbf{x}) \mathbf{x}^T \rangle &= \mathbf{A} \left\langle [\phi(\mathbf{x}^{(1)})^T \phi(\mathbf{x}), \dots, \phi(\mathbf{x}^{(M)})^T \phi(\mathbf{x})]^T \mathbf{x}^T \right\rangle \\ &= \frac{\mathbf{A}}{M} \left[\sum_m \mathbf{x}^{(m)} \phi(\mathbf{x}^{(m)})^T \phi(\mathbf{x}^{(1)}), \dots \right]^T = \frac{\mathbf{A} \mathbf{B}^T}{M} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{X}|}, \end{aligned} \quad (\text{C.12})$$

where we defined

$$\mathbf{B} \stackrel{\text{def}}{=} \sum_{m=1}^M \mathbf{x}^{(m)} \mathbf{k}(\mathbf{x}^{(m)})^T \in \mathbb{R}^{|\mathcal{X}| \times M}, \quad \mathbf{k}(\mathbf{x}^{(m)}) \stackrel{\text{def}}{=} \mathbf{F}^T \phi(\mathbf{x}^{(m)}) \in \mathbb{R}^M \quad (\text{C.13})$$

and used the fact that from the orthogonality assumption (C.11) we get $\mathbf{W}^\perp \mathbf{F} = \mathbf{0} \in \mathbb{R}^{|\mathcal{X}| \times M}$. Here $\mathbf{k}(\mathbf{x}^{(m)})$ corresponds to the m^{th} column of the Gram matrix $\mathbf{K} \stackrel{\text{def}}{=} \{K_{ij}\} \stackrel{\text{def}}{=} \{\phi(\mathbf{x}^{(i)})^T \phi(\mathbf{x}^{(j)})\} \in \mathbb{R}^{M \times M}$. Clearly, for a fixed $\mathbf{K} \in \mathbb{R}^{M \times M}$, the computed expectation $\mathbf{W} \langle \phi(\mathbf{x}) \mathbf{x}^T \rangle$ is a function of the coefficients $\mathbf{A} \in \mathbb{R}^{|\mathcal{Y}| \times |M|}$, which does not require explicit computations in the high-dimensional feature space.

Analogously, we can express

$$\begin{aligned} \mathbf{W} \langle \phi(\mathbf{x}) \phi(\mathbf{x})^T \rangle \mathbf{W}^T &= (\mathbf{A} \mathbf{F}^T + \mathbf{W}^\perp) \langle \phi(\mathbf{x}) \phi(\mathbf{x})^T \rangle (\mathbf{A} \mathbf{F}^T + \mathbf{W}^\perp)^T \\ &= \frac{1}{M} \mathbf{A} \sum_{m=1}^M \mathbf{k}(\mathbf{x}^{(m)}) \mathbf{k}(\mathbf{x}^{(m)})^T \mathbf{A}^T, \end{aligned} \quad (\text{C.14})$$

where we have used the fact that $\mathbf{W}^\perp \mathbf{F} = \mathbf{0}$, as follows from construction (C.11). Again, for the fixed Gram matrix the term is a quadratic function of coefficients \mathbf{A} alone.

By substitution, we may re-express the bound (C.9) as

$$\begin{aligned} \tilde{I}(\mathbf{x}, \mathbf{y}) &= \frac{1}{M} \text{tr} \{ \Sigma_x^{-1} \mathbf{U} \mathbf{A} \mathbf{B}^T \} - \frac{1}{2} \text{tr} \{ \mathbf{U}^T \Sigma_x^{-1} \mathbf{U} \Sigma_y \} \\ &\quad - \frac{1}{2M} \text{tr} \{ \mathbf{U}^T \Sigma_x^{-1} \mathbf{U} \mathbf{A} \mathbf{K}^2 \mathbf{A}^T \} - \frac{1}{2} \text{tr} \{ \Sigma_x^{-1} \mathbf{S} \} + c, \end{aligned} \quad (\text{C.15})$$

where $c = -|\mathcal{X}|/2 \log(2\pi\sigma^2)$ does not affect the optimization surface and be ignored in the rest of the discussion. In the simplest case when $\phi(\mathbf{x}) \equiv \mathbf{x} \in \mathbb{R}^{|\mathcal{X}|}$, we obtain $\mathbf{K}^2 \propto \mathbf{X}^T \mathbf{S} \mathbf{X} \in \mathbb{R}^{M \times M}$, $\mathbf{B} \propto \mathbf{S} \mathbf{X} \in \mathbb{R}^{|\mathcal{X}| \times M}$ where $\mathbf{X} \stackrel{\text{def}}{=} [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}] \in \mathbb{R}^{|\mathcal{X}| \times M}$ contains the training data. As expected, this transforms the bound (C.15) to the corresponding expression (C.7) for the linear Gaussian channel, thus resulting in PCA on the sample covariance \mathbf{S} for both the encoder and decoder weights \mathbf{U} , \mathbf{W}^T as the optimal choice.

The objective (C.15) may be used to learn optimal decoder weights $\mathbf{U} \in \mathbb{R}^{|\mathcal{X}| \times |\mathcal{Y}|}$ and optimal coordinates $\mathbf{A} \in \mathbb{R}^{|\mathcal{Y}| \times |M|}$ in the space spanned by the feature vectors $\{\phi(\mathbf{x}^{(i)}) | i \in [1, M] \cap \mathbb{N}\}$. Moreover, we may use the so-called kernel trick (e.g.

Vapnik (1998), Cristianini and Shawe-Taylor (2000), Scholkopf and Smola (2002)) and compute the entries of the Gram matrix $\mathbf{K} = \{K_{ij}\} \in \mathbb{R}^{M \times M}$ as $K_{ij} = \phi(\mathbf{x}^{(i)})^T \phi(\mathbf{x}^{(j)}) = \mathcal{K}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}; \Theta)$, where $\mathcal{K}_\Theta : |\mathbf{x}| \times |\mathbf{x}| \rightarrow \mathbb{R}$ defines a symmetric positive-definite *kernel* function. In principle, we may use the objective $\tilde{I}(\mathbf{x}, \mathbf{y})$ to learn the optimal parameters¹ of the kernel function. However, as we showed in Agakov and Barber (2004c) and will discuss later, such learning may be strongly influenced by the choice of constraints on channel parameters.

C.2.2 Nature of optimal solutions

In the following we assume for simplicity that $\Sigma_y = s^2 \mathbf{I}$ and $\Sigma_x = \sigma^2 \mathbf{I}$. We also assume that $|\mathbf{y}| \leq |\mathbf{x}| \leq |\phi|$ and $|\mathbf{x}| \leq M$, so that \mathbf{y} is a compressed representation of $\phi(\mathbf{x})$, and the number of training points is sufficient to ensure invertibility of the sample covariance.

C.2.2.1 Optimal Decoder

Similarly to Appendix C.2.2, we will express the bound on $I(\mathbf{x}, \mathbf{y})$ as a function of decoder parameters alone and analyze the optimal solutions. In order to do this, we will first find an analytical expression for optimal encoder coefficients, re-express the bound, and optimize it for the variational decoder.

Optimization of $\tilde{I}(\mathbf{x}, \mathbf{y})$ for the matrix of coefficients $\mathbf{A} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{M}|}$ leads to the fixed point condition

$$\partial \tilde{I}(\mathbf{x}, \mathbf{y}) / \partial \mathbf{A} = 0 \Rightarrow \mathbf{U}^T \Sigma_x^{-1} \mathbf{B} = \mathbf{U}^T \Sigma_x^{-1} \mathbf{U} \mathbf{A} \mathbf{K}^2 \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{M}|}, \quad (\text{C.16})$$

which for the considered isotropic case $\Sigma_x = \sigma^2 \mathbf{I}$ leads to

$$\mathbf{A} \mathbf{K}^2 = (\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T \mathbf{B} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{M}|}. \quad (\text{C.17})$$

Assuming that the Gram matrix \mathbf{K} is non-singular, we obtain $\mathbf{A} \mathbf{K} = (\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T \mathbf{B} \mathbf{K}^{-1}$. Then a substitution into (C.15) leads to

$$\begin{aligned} \tilde{I}(\mathbf{x}, \mathbf{y}) &= \frac{1}{2M\sigma^2} \text{tr} \{ \mathbf{U} \mathbf{A} \mathbf{K}^2 \mathbf{A}^T \mathbf{U}^T \} - \frac{s^2}{2\sigma^2} \text{tr} \{ \mathbf{U} \mathbf{U}^T \} - \frac{1}{2\sigma^2} \text{tr} \{ \mathbf{S} \} \\ &= \frac{1}{2M\sigma^2} \text{tr} \{ \mathbf{U} (\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T \mathbf{B} \mathbf{K}^{-2} \mathbf{B}^T \mathbf{U} (\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T \} - \frac{s^2}{2\sigma^2} \text{tr} \{ \mathbf{U} \mathbf{U}^T \} - \frac{1}{2\sigma^2} \text{tr} \{ \mathbf{S} \} \\ &= \frac{1}{2M\sigma^2} \text{tr} \{ \mathbf{U} (\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T \mathbf{B} \mathbf{K}^{-2} \mathbf{B}^T \} - \frac{s^2}{2\sigma^2} \text{tr} \{ \mathbf{U} \mathbf{U}^T \} - \frac{1}{2\sigma^2} \text{tr} \{ \mathbf{S} \}, \end{aligned} \quad (\text{C.18})$$

where we ignored the terms independent of the decoder parameters $\mathbf{U} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{y}|}$. Finally, by noticing that $\mathbf{B} = \mathbf{X} \mathbf{K} \in \mathbb{R}^{|\mathbf{x}| \times M}$ and

$$\mathbf{B} \mathbf{K}^{-2} \mathbf{B}^T = \mathbf{X} \mathbf{X}^T = \mathbf{M} \mathbf{S}, \quad (\text{C.19})$$

we may transform the bound (C.18) to

$$\tilde{I}(\mathbf{x}, \mathbf{y}) = \frac{1}{2\sigma^2} \text{tr} \{ \mathbf{U} (\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T \mathbf{X} \mathbf{X}^T \} - \frac{s^2}{2\sigma^2} \text{tr} \{ \mathbf{U} \mathbf{U}^T \} - \frac{1}{2\sigma^2} \text{tr} \{ \mathbf{S} \}, \quad (\text{C.20})$$

¹It is also possible to learn \mathbf{K} directly by constraining it to satisfy properties of inner products; this may be useful, for example, when the source alphabet is exhausted by M training patterns.

which is exactly the objective of the linear Gaussian channel (C.7).

Note that optimally $\|\mathbf{U}\|_F \rightarrow 0$, which (for a fixed \mathcal{K}_Θ) leads to the diverging encoder's coefficients $\mathbf{A} \in \mathbb{R}^{|y| \times |M|}$, resulting in the divergent weights $\|\mathbf{W}\|_F \rightarrow \infty$. As before, this may be easily explained from the information-theoretic perspective, as the resulting channel is approximately noiseless, and therefore it is characterized by a low value of the conditional entropy $H(\mathbf{y}|\mathbf{x})$ and a high value of the mutual information. Similarly to the case of a linear channel, the convergence of $\|\mathbf{W}\|_F$ for $s^2 > 0$ may be insured by imposing norm constraints on $\mathbf{U} \in \mathbb{R}^{|\mathbf{x}| \times |y|}$. Then, by the exact analogy with (C.7), we can see that under the specific assumption of isotropic Gaussian noise and a non-singular kernel matrix, optimal weights $\mathbf{U} \in \mathbb{R}^{|\mathbf{x}| \times |y|}$ of the linear Gaussian decoder correspond to principal components of the sample covariance \mathbf{S} (and their rotations). Fundamentally, we may note that the considered linear Gaussian variational decoder restricts the power of the approach, as for the optimal settings of encoder and decoder parameters we cannot improve on the PCA bounds (C.7).

C.2.2.2 Optimal Encoder

For completeness, we may express the optimal encoder weights. From (C.16) it is clear that optimal solutions for the encoder are given by

$$\mathbf{W}\mathbf{F} = \mathbf{A}\mathbf{K} \propto \mathbf{U}^+\mathbf{X} \in \mathbb{R}^{|y| \times M}, \quad (\text{C.21})$$

where $\mathbf{U}^+ \in \mathbb{R}^{|y| \times |\mathbf{x}|}$ denotes the pseudo-inverse, and left singular values of $\mathbf{U} \in \mathbb{R}^{|\mathbf{x}| \times |y|}$ correspond to principal eigenvectors of $\mathbf{S} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|}$. In the case when $\phi(\mathbf{x}) \equiv \mathbf{x} \in \mathbb{R}^{|\mathbf{x}|}$ and the sample covariance $\mathbf{X}\mathbf{X}^T/M$ is non-singular, condition (C.21) results in $\mathbf{W} = \mathbf{U}^+ \in \mathbb{R}^{|y| \times |\mathbf{x}|}$, which is the PCA solution of the linear Gaussian channel. However, for general nonlinear mappings, optimal encoder weights $\mathbf{W}^T \in \mathbb{R}^{|\mathbf{x}| \times |y|}$ do not necessarily give rise to the nonlinear PCA solution.

Finally, note that if the channel noise is isotropic and there are no constraints preventing the weights from taking optimal solutions according to (C.17) and (C.21), then the bound is given by the summation of $|y|$ principal eigenvalues of the sample covariance \mathbf{S} . It is also important to note that for both of the considered channels (with linear and kernelized Gaussian encoders), the bound on the mutual information is maximized when the projection weights are unconstrained, since this situation leads to insignificant channel noise contributions. By imposing norm constraints on the weights of a linear Gaussian decoder, both channels are shown to result in the same optimization surface for \mathbf{U} independently of the choice of nonlinearity. Hence, we reach an important conclusion: *for isotropic channel noise, if the decoder is linear, nothing is gained by using a nonlinear encoder* in the proposed variational settings with the considered norm constraints. This agrees with the related result of Bourlard and Kamp (1988) and Bourlard (2000) for noiseless autoencoders, but is derived in the context of variational information maximization for a stochastic channel $\mathbf{x} \rightarrow \mathbf{y}$. Moreover, in our derivation we did not have to assume approximate linearity of $\phi(x)$ around $x = 0$ (*cf* Bourlard (2000)).

These results are somewhat disappointing. In order to improve the power of the method, we need to consider both nonlinear encoders and decoders. However,

from (2.2) it is clear that in the stochastic context, the naive approach of using a nonlinear decoder will typically result in intractable averages over \mathbf{y} in the expression for the variational bound $\tilde{I}(\mathbf{x}, \mathbf{y})$. In order to avoid this computational difficulty, we derive a modified bound on the mutual information by considering further relaxations of the generic bound and performing decoding in the feature space.

C.3 Nonlinear Gaussian Channels: Nonlinear Decoders and KPCA

Since projections to the feature space are deterministic, one may derive an alternative form of the bound on the mutual information. For any choice of a functional nonlinearity $\phi : \mathbf{x} \mapsto \mathbf{f}$, the codes \mathbf{y} are as predictable from the feature vectors $\mathbf{f} \stackrel{\text{def}}{=} \phi(\mathbf{x}) \in \mathbb{R}^{|\phi|}$, as they are from the source variables \mathbf{x} themselves. This may be used to modify the bound on $I(\mathbf{x}, \mathbf{y})$ in such a way that the reconstruction is performed in the feature space, which results in a simple nonlinear generalization of the PCA solutions for optimal parameters.

The nonlinear Gaussian channel discussed in Section C.2 may be represented by the Markov chain $\mathbf{x} \rightarrow \mathbf{f} \rightarrow \mathbf{y}$, where $\mathbf{f} \in \mathbb{R}^{|\phi|}$ and $p(\mathbf{f}|\mathbf{x}) \sim \delta(\mathbf{f} - \phi(\mathbf{x}))$, $p(\mathbf{y}|\mathbf{f}) \sim \mathcal{N}_y(\mathbf{W}\mathbf{f}, \Sigma_y)$. Indeed, by marginalizing the feature variables \mathbf{f} it is clear that the encoder is given² by $p(\mathbf{y}|\mathbf{x}) = \int_{\mathbf{f}} \delta(\mathbf{f} - \phi(\mathbf{x})) \mathcal{N}_y(\mathbf{W}\mathbf{f}, \Sigma_y) = \mathcal{N}_y(\mathbf{W}\phi(\mathbf{x}), \Sigma_y)$.

Proposition C.1. *Let $\mathbf{s} \rightarrow \mathbf{t} \rightarrow \mathbf{r}$ define a Markov chain, such that $p(\mathbf{t}|\mathbf{s}) = \delta(\mathbf{t} - \mathbf{f}(\mathbf{s}))$, and $p(\mathbf{r}|\mathbf{t})$ is a continuous differentiable density function satisfying $\forall \mathbf{r}. \forall \mathbf{t}. p(\mathbf{r}|\mathbf{t}) \neq 0$. Then $I(\mathbf{s}, \mathbf{r}) = I(\mathbf{t}, \mathbf{r})$.*

Proof. From basic properties of the mutual information (see e.g. Cover and Thomas (1991)) it is easy to see that

$$I(\mathbf{s}, \mathbf{t}; \mathbf{r}) = H(\mathbf{r}) - H(\mathbf{r}|\mathbf{t}, \mathbf{s}) = I(\mathbf{t}, \mathbf{r}), \quad (\text{C.22})$$

$$\begin{aligned} I(\mathbf{s}, \mathbf{t}; \mathbf{r}) &= H(\mathbf{s}) + H(\mathbf{t}|\mathbf{s}) - H(\mathbf{s}|\mathbf{r}) - H(\mathbf{t}|\mathbf{s}, \mathbf{r}) \\ &= I(\mathbf{s}, \mathbf{r}) + H(\mathbf{t}|\mathbf{s}) - H(\mathbf{t}|\mathbf{s}, \mathbf{r}). \end{aligned} \quad (\text{C.23})$$

Utilizing the chain structure and the deterministic mapping $p(\mathbf{t}|\mathbf{s})$, we obtain

$$p(\mathbf{t}|\mathbf{s}, \mathbf{r}) = \frac{\delta(\mathbf{t} - \mathbf{f}(\mathbf{s}))p(\mathbf{r}|\mathbf{t})}{\int_{\mathbf{t}} \delta(\mathbf{t} - \mathbf{f}(\mathbf{s}))p(\mathbf{r}|\mathbf{t})} = \frac{\delta(\mathbf{t} - \mathbf{f}(\mathbf{s}))p(\mathbf{r}|\mathbf{f}(\mathbf{s}))}{p(\mathbf{r}|\mathbf{f}(\mathbf{s}))}, \quad (\text{C.24})$$

i.e. $p(\mathbf{t}|\mathbf{s}, \mathbf{r}) = p(\mathbf{t}|\mathbf{s})$. Here we used $f(x)\delta(x-a) = \lim_{\epsilon \rightarrow 0} [f(a-\epsilon) + f(a+\epsilon)]\delta(x-a)/2$ (see e.g. Korn and Korn (1968)). Then $H(\mathbf{t}|\mathbf{s}, \mathbf{r}) = H(\mathbf{t}|\mathbf{s})$, and from (C.22), (C.23) we obtain $I(\mathbf{s}, \mathbf{r}) = I(\mathbf{t}, \mathbf{r})$. \square

²We assume Cartesian coordinates, i.e. $\delta(\mathbf{x} - \mathbf{a}) = \prod_i \delta(x^{(i)} - a_i)$, see e.g. Korn and Korn (1968).

From proposition C.1, the mutual information $I(\mathbf{x}, \mathbf{y})$ may be bounded as

$$I(\mathbf{x}, \mathbf{y}) = I(\mathbf{f}, \mathbf{y}) \geq \tilde{I}(\mathbf{f}, \mathbf{y}), \text{ where } \tilde{I}(\mathbf{f}, \mathbf{y}) \stackrel{\text{def}}{=} \langle \log q(\mathbf{f}|\mathbf{y}) \rangle_{\tilde{p}(\mathbf{x})p(\mathbf{f}|\mathbf{x})p(\mathbf{y}|\mathbf{f})} + H(\mathbf{f}). \quad (\text{C.25})$$

We make the simple assumption that the feature decoder is Gaussian, $q(\mathbf{f}|\mathbf{y}) \sim \mathcal{N}_{\mathbf{f}}(\mathbf{U}\mathbf{y}, \Sigma_{\mathbf{f}})$, which leads to

$$\tilde{I}(\mathbf{x}, \mathbf{y}) = -\frac{1}{2} \text{tr} \{ \mathbf{U}^T \Sigma_{\mathbf{f}}^{-1} \mathbf{U} (\Sigma_{\mathbf{y}} + \mathbf{W} \mathbf{S}_F \mathbf{W}^T) \} + \text{tr} \{ \Sigma_{\mathbf{f}}^{-1} \mathbf{U} \mathbf{W} \mathbf{S}_F \} + H(\mathbf{f}) + c \quad (\text{C.26})$$

where $\mathbf{S}_F \stackrel{\text{def}}{=} \langle \mathbf{f} \mathbf{f}^T \rangle_{p(\mathbf{f})} \in \mathbb{R}^{|\mathbf{f}| \times |\mathbf{f}|}$ corresponds to the matrix of second-order moments in the feature space. By centering the data in the feature space (see e.g. Schoelkopf et al. (1998)), we may also express the covariance of the feature vectors; we will omit this discussion, as it proves to be secondary for the consequent analysis.

The matrix of the second order moments is easy to compute from the training set as

$$\mathbf{S}_F = \frac{1}{M} \sum_{i=1}^M \phi(\mathbf{x}^{(i)}) \phi(\mathbf{x}^{(i)})^T, \quad (\text{C.27})$$

where $\mathbf{f} = \phi(\mathbf{x})$. As expected, in the special case of linear mappings $\phi(\mathbf{x}) \equiv \mathbf{x}$ we get $\mathbf{S}_F = \langle \mathbf{x} \mathbf{x}^T \rangle \equiv \mathbf{S}$, which transforms the objective (C.26) to the simpler bound (C.1).

C.3.1 Constraints on the Feature Mappings $p(\mathbf{f}|\mathbf{x})$

Evaluation of the bound (C.26) is complicated by the need of computing the entropic term $H(\mathbf{f})$. Despite the fact that the mapping to the feature space is deterministic, generally we do not know explicit feature space representations of the training patterns, i.e. numeric approximations due to Brunel and Nadal (1998), Shriki et al. (2002), Corduneanu and Jaakkola (2003) are not directly applicable (moreover, such approximations will not generally retain a proper bound on $I(\mathbf{x}, \mathbf{y})$). To simplify the problem of computing $H(\mathbf{f})$, we may note that if $\phi : \mathbf{x} \rightarrow \mathbf{f}$ is deterministic and $\tilde{p}(\mathbf{x}) = \sum_{i=1}^M \delta(\mathbf{x} - \mathbf{x}^{(i)})/M$ then $\phi(\mathbf{x}^{(i)})$ corresponds to a re-labeled source pattern. If no two distinct points $\mathbf{x}^{(i)}$ and $\mathbf{x}^{(j)}$ in the data space are mapped to the same point in the feature space, then $H(\mathbf{f}) = H(\mathbf{x}) = \log M/\delta(0) = \text{const}$, i.e. it may be ignored during the optimization. One way to ensure that this one-to-one condition is satisfied is by imposing inequality constraints on the kernel function, so that

$$\forall i, j \in \{1, \dots, M\}. \mathcal{K}_{i,i} - \mathcal{K}_{i,j} > \epsilon > 0. \quad (\text{C.28})$$

Clearly, this is a sufficient (but not necessary) condition for the distinct sources $\mathbf{x}^{(i)}$ and $\mathbf{x}^{(j)}$ to have different images in the feature space.

By enforcing the constraints, one may simply ignore the (constant) entropic term $H(\mathbf{f})$ and optimize the objective (C.26) by analogy with the simple bound (C.1) for isotropic linear Gaussians. In conjunction with the set of constraints

(C.28), the objective may be further optimized with respect to parameters of the kernel function. Note that (C.26) remains a proper bound on $I(\mathbf{x}, \mathbf{y})$ for the considered constrained nonlinear Gaussian channel.

For the considered choice of the variational feature-space decoder $q(\mathbf{f}|\mathbf{y}) \sim \mathcal{N}_f(\mathbf{U}\mathbf{y}, \sigma_f^2 \mathbf{I})$, we may express the corresponding stochastic mapping from the code to the data space as $q(\mathbf{x}|\mathbf{y}) = \int_{\mathbf{f}} p(\mathbf{x}|\mathbf{f}) \mathcal{N}_f(\mathbf{U}\mathbf{y}, \sigma_f^2 \mathbf{I})$. Unfortunately, these variational posteriors may in principle be difficult to compute. This is a possible limitation of using the variational decoders for reconstruction, since for a general feature mapping $\phi(\mathbf{x})$ it may be difficult to reconstruct a source vector \mathbf{x} from the corresponding feature vector \mathbf{f} . However, as we mentioned in Chapter 2.1.3, there could be multiple ways of reconstructing source vectors from their noisy representations. In fact, one may view optimization of the bound (C.26) as a way of learning optimal encoders. Once the encoder $p(\mathbf{y}|\mathbf{x})$ is learned, it may be used by any method which can perform exact or approximate inference. For example, if the empirical distribution is available, and if it exhausts all the possible training instances, one may retrieve the transmitted source \mathbf{x} from the noisy encoding \mathbf{y} by reconstructing with the exact posterior $p(\mathbf{x}|\mathbf{y}) \propto \sum_{i=1}^M \delta(\mathbf{x} - \mathbf{x}^{(i)}) p(\mathbf{y}|\mathbf{x})$.

C.3.2 Variational Bounds on $H(\mathbf{f})$

As we have discussed above, one way to avoid computation and optimization of the generally intractable entropy of the feature variables $H(\mathbf{f})$ is to impose constraints on the kernel function. This ensures that $H(\mathbf{f}) = H(\mathbf{x}) = \text{const}$, which may significantly simplify optimization of (C.26). Moreover, constraints on the off-diagonal elements of the kernel matrices may help to avoid possible singularities, as they will intuitively favor more uniform eigenspectra.

An alternative way to address intractability of $H(\mathbf{f})$ is by relaxing the bound on $I(\mathbf{x}, \mathbf{y})$. Clearly, from the joint rule for entropies we get

$$H(\mathbf{f}) = H(\mathbf{f}|\mathbf{x}) + H(\mathbf{x}) - H(\mathbf{x}|\mathbf{f}) \geq H(\mathbf{x}) + \langle \log q(\mathbf{x}|\mathbf{f}) \rangle_{p(\mathbf{x}, \mathbf{f})} + c, \quad (\text{C.29})$$

where we used the fact that the mapping $\mathbf{x} \mapsto \mathbf{f}$ is deterministic (but of course not in general one-to-one). By substituting (C.29) into the bound (C.25), we obtain

$$\tilde{I}(\mathbf{x}, \mathbf{y}) \geq \langle \log q(\mathbf{f}|\mathbf{y}) \rangle_{\tilde{p}(\mathbf{x})p(\mathbf{f}|\mathbf{x})p(\mathbf{y}|\mathbf{f})} + \langle \log q(\mathbf{x}|\mathbf{f}) \rangle_{p(\mathbf{x}, \mathbf{f})} + c \quad (\text{C.30})$$

where c are irrelevant constants. As before, we will be optimizing (C.30) with respect to parameters of the encoder $p(\mathbf{y}|\mathbf{f})$, the feature decoder $q(\mathbf{f}|\mathbf{y})$, and the data decoder $q(\mathbf{x}|\mathbf{f})$.

One tractable way to constrain $q(\mathbf{x}|\mathbf{f})$ is again to use a linear Gaussian $q(\mathbf{x}|\mathbf{f}) \sim \mathcal{N}_{\mathbf{x}}(\mathbf{V}\mathbf{f}, \Sigma_{\mathbf{x}|\mathbf{f}})$. This parameterization results in a simple linear Gaussian decoder $q(\mathbf{x}|\mathbf{y})$ with a structured covariance matrix, which may be easily used for variational decoding. Other tractable parameterizations may also be possible. Indeed, since integration over \mathbf{x} and \mathbf{f} reduces to evaluations of the empirical averages, computing the average $\langle \log q(\mathbf{x}|\mathbf{f}) \rangle_{p(\mathbf{x}, \mathbf{f})}$ is easy for any data decoder $q(\mathbf{x}|\mathbf{f})$ (provided that it is kernelized, so that there are no explicit computations in the feature space). Fundamentally, however, the choice of the data decoder $q(\mathbf{x}|\mathbf{f})$ does not affect the nature of optimal solutions for parameters of the encoder and feature decoder, provided that at each iteration of the IM the kernel function \mathcal{K}_{Θ} is fixed.

C.3.3 Kernelized Representation

Here we show that for a fixed kernel function \mathcal{K}_Θ , optimal parameters of the encoder $p(y|\mathbf{f}) \sim \mathcal{N}_y(\mathbf{W}\mathbf{f}, \Sigma_y)$ and feature decoder $q(\mathbf{f}|y) \sim \mathcal{N}_f(\mathbf{U}y, \Sigma_f)$ result in a simple nonlinear generalization of the PCA solutions, which justifies nonlinear PCA as a lower bound on mutual information $I(\mathbf{x}, y)$ between the sources and the codes³.

In what follows we assume that $\Sigma_y = s^2\mathbf{I} \in \mathbb{R}^{|\mathbf{y}|}$, $\Sigma_f = \sigma_f^2\mathbf{I} \in \mathbb{R}^{|\phi|}$, and $|\mathbf{y}| < M < |\phi|$. By analogy with Section C.2 we notice that rows of \mathbf{W} and columns of \mathbf{U} have dimension $|\phi|$. Then they may be represented in the basis defined by the span of $\{\phi(\mathbf{x}^{(m)}) | m = 1, \dots, M\}$ and its orthogonal compliment as

$$\mathbf{W} = \mathbf{A}\mathbf{F}^T + \mathbf{W}^\perp \in \mathbb{R}^{|\mathbf{y}| \times |\phi|}, \quad (\text{C.31})$$

$$\mathbf{U} = \mathbf{F}\mathbf{C} + \mathbf{U}^\perp \in \mathbb{R}^{|\phi| \times |\mathbf{y}|}. \quad (\text{C.32})$$

Here $\mathbf{F} \in \mathbb{R}^{|\phi| \times M}$ is the design matrix, $\mathbf{U}^\perp, \mathbf{W}^\perp$ are orthogonal to \mathbf{F} , and $\mathbf{A} \in \mathbb{R}^{|\mathbf{y}| \times M}$, $\mathbf{C} \in \mathbb{R}^{M \times |\mathbf{y}|}$ are matrices of coefficients to be learned. A substitution into expression (C.30) results in

$$\begin{aligned} \tilde{I}(\mathbf{x}, y) \propto & 2\text{tr}\{\mathbf{C}\mathbf{A}\mathbf{K}^2\} - s^2M\text{tr}\{\mathbf{C}^T\mathbf{K}\mathbf{C}\} - \text{tr}\{\mathbf{A}\mathbf{K}^2\mathbf{A}^T[\mathbf{C}^T\mathbf{K}\mathbf{C} + (\mathbf{U}^\perp)^T\mathbf{U}^\perp]\} \\ & - Ms^2\text{tr}\{(\mathbf{U}^\perp)^T\mathbf{U}^\perp\} - \text{tr}\{\mathbf{K}\} + \frac{1}{2M\sigma_f^2}\langle \log q(\mathbf{x}|\mathbf{f}) \rangle_{p(\mathbf{x}, \mathbf{f})}, \end{aligned} \quad (\text{C.33})$$

where we have used the orthogonality conditions $\mathbf{W}^\perp\mathbf{F} = \mathbf{0} \in \mathbb{R}^{|\mathbf{y}| \times M}$ and $\mathbf{F}^T\mathbf{U}^\perp \in \mathbb{R}^{M \times |\mathbf{y}|}$. It is clear that unconstrained optimization of the objective (C.33) for the complimentary basis \mathbf{U}^\perp leads to $\mathbf{U}^\perp = \mathbf{0} \in \mathbb{R}^{|\phi| \times |\mathbf{y}|}$.

The objective (C.33) may be readily used for learning coefficients \mathbf{A} , \mathbf{C} and parameters of $q(\mathbf{x}|\mathbf{f})$. In the considered case it may also be applied to learning parameters Θ of the kernel function $\mathcal{K}_\Theta(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}; \Theta)$ which gives rise to the Gram matrix \mathbf{K} .

C.3.4 Nature of Optimal Solutions

Optimization of the bound (C.33) for the coefficients \mathbf{A} leads to

$$\partial\tilde{I}(\mathbf{x}, y)/\partial\mathbf{A} \propto \mathbf{C}^T\mathbf{K}^2 - \mathbf{C}^T\mathbf{K}\mathbf{C}\mathbf{A}\mathbf{K}^2, \quad (\text{C.34})$$

resulting in

$$\mathbf{A} = (\mathbf{C}^T\mathbf{K}\mathbf{C})^{-1}\mathbf{C}^T \quad (\text{C.35})$$

for the case when the Gram matrix $\mathbf{K} \in \mathbb{R}^{M \times M}$ is non-singular. This transforms the objective (C.33) to

$$\tilde{I}(\mathbf{x}, y) = \text{tr}\{\mathbf{K}^2\mathbf{C}(\mathbf{C}^T\mathbf{K}\mathbf{C})^{-1}\mathbf{C}^T\} - s^2M\text{tr}\{\mathbf{C}^T\mathbf{K}\mathbf{C}\} + \text{const}. \quad (\text{C.36})$$

³For clarity, we assume that the feature vectors $\mathbf{f} \in \mathbb{R}^{|\phi|}$ are centered (see e.g. Schoelkopf et al. (1998)); otherwise we would consider $p(y|\mathbf{f}) \sim \mathcal{N}_y(\mathbf{W}(\mathbf{f} - \langle \mathbf{f} \rangle), \Sigma_y)$, $q(\mathbf{f}|y) \sim \mathcal{N}_f(\mathbf{U}y + \langle \mathbf{f} \rangle, \Sigma_f)$. Note that the explicit computations in the feature space will not be required, as evaluation of $\langle \log q(\mathbf{f}|y) \rangle$ will only involve computations of scalar products of the feature vectors.

Here we ignored the terms independent of the coefficients \mathbf{A} and \mathbf{C} , and used parameterization (C.31 – C.32) with the optimal settings $\mathbf{U}^\perp = \mathbf{0}$. Moreover, for the fixed kernel matrix \mathbf{K} we obtain

$$\begin{aligned}\tilde{I}(\mathbf{x}, \mathbf{y}) &= M \text{tr} \{ \mathbf{F}^T \mathbf{S}_F \mathbf{F} \mathbf{C} (\mathbf{C}^T \mathbf{F}^T \mathbf{F} \mathbf{C})^{-1} \mathbf{C}^T \} - s^2 M \text{tr} \{ \mathbf{C}^T \mathbf{F}^T \mathbf{F} \mathbf{C} \} + \text{const} \\ &\propto \text{tr} \{ \mathbf{U} (\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T \mathbf{S}_F \} - s^2 \text{tr} \{ \mathbf{U} \mathbf{U}^T \} + \text{const}.\end{aligned}\tag{C.37}$$

By analogy with Section C.2, maximization of (C.37) gives rise to the nonlinear PCA solution (and rotations) for the left singular vectors of \mathbf{U} and \mathbf{W}^T . (As before, we have ignored \mathbf{W}^\perp and \mathbf{U}^\perp in the definitions (C.31), (C.32) since they should optimally cancel).

Just as in the linear case, in order to prevent divergence of $\|\mathbf{W}\|_F$ for $s^2 \neq 0$, it is useful to constrain the singular values of $\mathbf{U} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{y}|}$. In the special case when $\mathbf{U}^T \mathbf{U} = \mathbf{C}^T \mathbf{K} \mathbf{C} = \mathbf{I}_{|\mathbf{y}|}$, expressions (C.32) and (C.37) lead to the optimal settings

$$\mathbf{K}^2 \mathbf{C} \mathbf{R} = \mathbf{K} \mathbf{C} \mathbf{R} \boldsymbol{\lambda}_{\mathbf{S}_F} \in \mathbb{R}^{M \times |\mathbf{y}|}.\tag{C.38}$$

Here $\boldsymbol{\lambda}_{\mathbf{S}_F} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{y}|}$ is a diagonal matrix of $|\mathbf{y}|$ eigenvalues of \mathbf{S}_F (and \mathbf{K}), and $\mathbf{R} \in \mathbb{R}^{|\mathbf{y}| \times |\mathbf{y}|}$ is a rotation matrix. From (C.35) and (C.38) it is clear that optimal \mathbf{C} and \mathbf{A}^T correspond to rotations of principal eigenvectors of the Gram matrix \mathbf{K} , which is the kernel PCA solution. Hence, *for a fixed kernel function \mathcal{K}_Θ and data decoder $q(\mathbf{x}|\mathbf{f})$ of a nonlinear Gaussian channel, the variational lower bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ is maximized by the kernel PCA solutions for encoder and feature decoder weights.*

C.3.5 Optimal Kernel Functions

The bound (C.30) enables us, in a principled way, to choose between different parameters of the kernel function, or to choose between competing kernels. Indeed, for any kernelized representation of the data decoder $q(\mathbf{x}|\mathbf{f})$, optimal parameters Θ of the kernel function $\mathcal{K}_\Theta : |\mathbf{x}| \times |\mathbf{x}| \rightarrow \mathbb{R}$ may be obtained by maximizing the general objective (C.30). Another alternative which we have discussed is optimization of (C.26) subject to the distance constraints on the Gram matrix (C.28). The optimization procedure may be viewed as a special case of the IM algorithm, which for this case may be formulated as follows:

1. For the fixed \mathcal{K}_Θ , optimize the bound (C.26) with respect to \mathbf{U} , \mathbf{W} (or the dual parameters \mathbf{C} , \mathbf{A}), and parameters of the data decoder $q(\mathbf{x}|\mathbf{f})$.
2. For the fixed \mathbf{C} , \mathbf{A} , and $q(\mathbf{x}|\mathbf{f})$, optimize the bound (C.26) with respect to the kernel parameters Θ of \mathcal{K}_Θ .

Depending on parameterization of the data decoder $q(\mathbf{x}|\mathbf{f})$ or the choice of constraints (C.28), this procedure generally results in non-trivial settings of the kernel parameters Θ .

C.3.5.1 Learning Kernel Functions for KPCA Channels

Before discussing a general way of optimizing the bound (C.26) for kernel parameters Θ , we will mention a simple pitfall of a careless interpretation of projections into the feature space. Effectively, this case corresponds to learning optimal kernels for KPCA channels by maximizing $I(\mathbf{f}, \mathbf{y})$.

As discussed in Section C.3.1, by presuming that no two data patterns $\mathbf{x}^{(i)}, \mathbf{x}^{(j)}$ have identical feature space representations $\forall i \neq j$, one could ignore the entropic term $H(\mathbf{f})$ and optimize the simpler bound (C.37) instead. In our variational formulation, this would correspond to a simple special case when the contribution of the data decoder $\langle \log q(\mathbf{x}|\mathbf{f}) \rangle_{p(\mathbf{x}, \mathbf{f})}$ to the bound (C.26) is independent of parameters of the kernel function \mathcal{K}_Θ ; for example, this is the case for the trivial feature-independent setting of the variational data decoder $q(\mathbf{x}|\mathbf{f}) \equiv q(\mathbf{x})$. In our framework, this special case would effectively be identical to the unconstrained learning of optimal kernel transformations for nonlinear PCA channels.

For this case, optimization of the objective (C.37) for Θ would reduce to maximizing

$$\tilde{I}(\mathbf{x}, \mathbf{y}) = \langle \log q(\mathbf{f}|\mathbf{y}) \rangle_{p(\mathbf{f}, \mathbf{y})} = -\frac{1}{2} \log (\sigma_y^2 + \text{tr} \{\mathbf{K}\} - \text{tr} \{\mathbf{C}^T \mathbf{K} \mathbf{C}\}), \quad (\text{C.39})$$

where the matrix of the decoder coefficients $\mathbf{C} \in \mathbb{R}^{M \times |\mathbf{y}|}$ performs PCA on $\mathbf{K} \in \mathbb{R}^{M \times M}$ (see expressions (C.37), (C.38) and the discussion in section C.3.2). Equivalently, it may be written in terms of the eigenvalues of the Gram matrix as

$$\tilde{I}(\mathbf{x}, \mathbf{y}) = -\frac{1}{2} \log \left(\sigma_y^2 + \sum_{i=|\mathbf{y}|+1}^M \lambda_i(\mathbf{K}(\Theta)) \right), \quad (\text{C.40})$$

where $\lambda_i(\mathbf{K}(\Theta))$ is the i^{th} principal component of the Gram matrix $\mathbf{K} \in \mathbb{R}^{M \times M}$ corresponding to the kernel function \mathcal{K}_Θ (everywhere in our discussion, we implied that $\mathbf{K} = \mathbf{K}(\Theta)$ is a function of the kernel parameters Θ).

Clearly, an alternative formulation of the optimization problem for Θ in this case is given by

$$\Theta^* = \arg \max_{\Theta} \sum_{i=1}^{|\mathbf{y}|} \lambda_i(\mathbf{K}(\Theta)), \quad (\text{C.41})$$

where Θ^* indicates the optimal kernel parameter settings. It is easy to see that if the achievable ranks and norms of $\mathbf{K}(\Theta)$ are unconstrained, nothing would prevent the method optimizing (C.41) from generating degenerate Gram matrices (Agakov and Barber (2004c)). In fact, there are at least two possible sources of degeneracy of \mathbf{K} : due to the norm divergence and due to \mathbf{K} 's singularity. The first case is somewhat analogous to what has been described for the simple linear case (see Section C.1). It arises, for example, when $\Theta > 0$ is a scaling factor of some fixed positive semi-definite function. In this case it is intuitive that the optimal channel will be characterized by the divergent norm $\|\mathbf{K}\|_F$, resulting in the diminishing noise effects and the optimality of the divergent settings of $\Theta \rightarrow \infty$.

Effects of near-singular Gram matrices on parameters of kernel functions may in general be less trivial. From (C.41) it is clear that the worst kernel has a flat spectrum (which is the case for $\mathbf{K} = \text{cl}$), while the optimal kernel function results in the Gram matrix \mathbf{K} with the eigenspectrum concentrated at $|\mathbf{y}|$ principal components. Intuitively, by allowing changes in eigenspectra of Gram matrices, we effectively choose an M -dimensional subspace of the feature space which could be well modeled by $|\mathbf{y}|$ -KPCA. For example, for trace-constrained matrices, we expect to reach the maximum of the objective (C.41) by choosing the parameters Θ in such a way that $\text{rank}(\mathbf{K}(\Theta)) \leq |\mathbf{y}|$. A degenerate solution which satisfies the optimality condition is when \mathbf{K} is approximately rank-1, i.e. $\text{tr}\{\mathbf{K}\} \approx \lambda_1(\mathcal{K}_\Theta)$ and $K_{ij} \approx K_{ii}$ for all $i, j \in \{1, \dots, M\}$.

It may be noticed that the rank degeneracy is largely an artifact of the trivial data decoder $q(\mathbf{x}|\mathbf{f})$. Indeed, for noiseless channels and rank-1 kernels, it is possible to achieve a perfect reconstruction of *feature* vectors \mathbf{f} from their encoded representations \mathbf{y} . However, the objective (C.39) does not explicitly favour a good reconstruction of the source vectors \mathbf{x} from features \mathbf{f} . As a result, all the source vectors may potentially be mapped to a single vector in the feature space, which would lead to a degenerate optimum of the bound. For many interesting kernels, this solution would be characterized by degenerate values of Θ . For example, for trace-constrained radial basis and mixture kernels, the objective (C.39) would be a monotonic function of Θ , which leads to trivial settings of the kernel parameters (Agakov and Barber (2004c)).

C.3.5.2 Learning Kernel Functions by Variational Information Maximization

As we mentioned earlier, one simple way to handle a possible degeneracy of the kernel parameters is by constraining the off-diagonal values of the Gram matrix according to (C.28). It is easy to see that this constraint helps to avoid rank-deficiency of the trivial projections to the feature space. Moreover, numerically it ensures that $H(\mathbf{f}) = \text{const}$, i.e. the entropy of the features may be safely ignored during the optimization of the general objective (C.25). Of course, the choice of the threshold ϵ may in practice require similar kinds of heuristics as a choice of kernel parameters themselves. Nevertheless, the thresholds may be easier to interpret geometrically, as they define minimal angles between distinct vectors in the feature space.

Alternatively, more complex *data decoders* $q(\mathbf{x}|\mathbf{f}, \Theta_q)$ may be considered, though by analogy with Section C.2, a proper care should be taken for $q(\mathbf{x}|\mathbf{f}, \Theta_q)$ to be kernelized and Θ_q to be properly constrained. For both of these cases, a formal analysis of the optimal solutions for the kernel parameters Θ depends on the constraints and the specifics of \mathcal{K}_Θ , and is difficult in general.

Bibliography

- Agakov, F. V. and Barber, D. (2003). Approximate Learning in Temporal Hidden Hopfield Models. In *13th International Conference on Artificial Neural Networks*. Springer.
- Agakov, F. V. and Barber, D. (2004a). An Auxiliary Variational Method. In *International Conference on Neural Information Processing*. Springer.
- Agakov, F. V. and Barber, D. (2004b). Variational Information Maximization for Neural Coding. In *International Conference on Neural Information Processing*. Springer.
- Agakov, F. V. and Barber, D. (2004c). Variational Information Maximization in Gaussian Channels. Technical Report EDI-INF-RR-0206, School of Informatics, University of Edinburgh.
- Agakov, F. V. and Barber, D. (2005a). Auxiliary Variational Information Maximization for Dimensionality Reduction. In *PASCAL: Subspace, Latent Structure and Feature Selection techniques: Statistical and Optimisation perspectives Workshop*. Available at http://homepages.inf.ed.ac.uk/felixa/Papers/aux_pasc2.pdf.
- Agakov, F. V. and Barber, D. (2005b). Kernelized Infomax Clustering. In *Neural Information Processing Systems*. MIT Press.
- Agakov, F. V. and Barber, D. (2005c). Nonlinear Encoder Models for Information-Theoretic Clustering. In *PASCAL: Statistics and Optimization of Clustering Workshop*.
- Aizerman, M., Braverman, E., and Rozonoer, L. (1964). Theoretical foundations of the potential function method in pattern recognition learning. *Automation and Remote Control*, 25.
- Arfken, G. (1985). *Mathematical Methods for Physicists*. Academic Press.
- Arimoto, S. (1972). An Algorithm for computing the capacity of arbitrary discrete memoryless channels. *IT*, 18.
- Atick, J. J. (1992). Entropy minimization: A design principle for sensory perception? *International Journal of Neural Systems*, 3:81–90.

- Bach, F. R. and Jordan, M. I. (2003). Learning spectral clustering. In *Proceedings of Neural Information Processing Systems*, volume 16.
- Baldi, P. and Hornik, K. (1989). Neural networks and principal components analysis: Learning from examples without local minima. *Neural Networks*, 2.
- Barber, D. and Agakov, F. V. (2003). The IM Algorithm: A Variational Approach to Information Maximization. In *NIPS*. MIT Press.
- Barber, D. and Sollich, P. (2000). Gaussian Fields for Approximate Inference. In Solla, S. A., Leen, T., and Müller, K.-R., editors, *Advances in Neural Information Processing Systems 12*, pages 393–399. MIT Press, Cambridge, MA.
- Barber, D. and Wiering, W. (1998). Tractable variational structures for approximating graphical models. *Neural Information Processing*.
- Barber, D. and Williams, C. K. I. (1997). Gaussian processes for bayesian classification via hybrid monte carlo. In Mozer, M. C., Jordan, M. I., and Petsche, T., editors, *Advances in Neural Information Processing Systems*, volume 9. The MIT Press.
- Barlow, H. (1989). Unsupervised Learning. *Neural Computation*, 1:295–311.
- Bartholomew, D. (1987). *Latent Variable Models and Factor Analysis*. Charles Griffin and Co. Ltd., London.
- Beal, M. J. (2003). *Variational Algorithms for Approximate Bayesian Inference*. PhD thesis, University College London.
- Becker, S. (1992). *An Information-theoretic unsupervised learning algorithm for neural networks*. PhD thesis, University of Toronto.
- Becker, S. and Hinton, G. (1992). A self-organized neural network that discovers surfaces in random-dot stereograms. *Nature*, 355:161–163.
- Bell, A. J. and Sejnowski, T. J. (1994). A non-linear information maximization algorithm that performs blind separation. In *Advances in Neural Information Processing Systems*, volume 7. MIT Press.
- Bell, A. J. and Sejnowski, T. J. (1995). An information-maximization approach to blind separation and blind deconvolution. *Neural Computation*, 7(6):1129–1159.
- Berger, A., Pietra, S. D., and Pietra, V. D. (1996). A maximum entropy approach to natural language processing. *Computational Linguistics*, 22(1).
- Bertsekas, D. P. (1996). *Constrained Optimization and Lagrange Multiplier Methods*. Athena Scientific.
- Bertsekas, D. P. (1999). *Nonlinear Programming, 2nd ed.* Athena Scientific.

- Bethge, M., Rotermund, D., and Pawelzik, K. (2002). Optimal Short-Term Population Coding: When Fisher Information Fails. *Neural Computation*, 14:2317–2351.
- Bishop, C. (1994). Mixture Density Networks. Technical Report NCRG 4288, Aston University.
- Bishop, C. (1995). *Neural Networks for Pattern Recognition*. Clarendon Press, Oxford.
- Bishop, C. M., Lawrence, N., Jaakkola, T., and Jordan, M. I. (1998). Approximating posterior distributions in belief networks using mixtures. In Jordan, M. I., Kearns, M. J., and Solla, S. A., editors, *Advances in Neural Information Processing Systems*, volume 10. The MIT Press.
- Bishop, Svensen and Williams (1998a). Developments of the Generative Topographic Mapping. *Neurocomputing*, 21.
- Bishop, Svensen and Williams (1998b). GTM: The Generative Topographic Mapping. *Neural Computation*, 10(1):215–234.
- Blahut, R. (1972). Computation of channel capacity and rate-distortion functions. *IT*, 18.
- Boser, B. E., Guyon, I. M., and Vapnik, V. N. (1992). A training algorithm for optimal margin classifiers. In *Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory*. ACM Press.
- Bourlard, H. (2000). Auto-association by multilayer perceptrons and singular value decomposition. Technical Report IDIAP RR 00-16, IDIAP.
- Bourlard, H. and Kamp, Y. (1988). Auto-association by multilayer perceptrons and singular value decomposition. *Biological Cybernetics*, 59.
- Boyd, S. and Vandenberghe, L. (2004). *Convex Optimization*. Cambridge University Press.
- Brunel, N. and Nadal, J.-P. (1998). Mutual Information, Fisher Information and Population Coding. *Neural Computation*, 10:1731–1757.
- Burshtein, D. and Miller, G. (2002). Bounds on the performance of belief propagation decoding. *IEEE Transactions on Information Theory*, 48(1):112–122.
- Cardoso, J. F. (1997). Infomax and maximum likelihood for blind source separation. In *IEEE Signal Processing Letters*, volume 4.
- Chechik, G., Globerson, A., Tishby, N., and Weiss, Y. (2003). Information bottleneck for gaussian variables. In *Advances in Neural Information Processing Systems*, volume 16. The MIT Press.

- Chechik, G. and Tishby, N. (2002). Extracting relevant structures with side information. In Becker, S., Thrun, S., and Obermayer, K., editors, *Advances in Neural Information Processing Systems*, volume 15. The MIT Press.
- Collins, M., Schapire, R. E., and Singer, Y. (2002). Logistic regression, AdaBoost and Bregman distances. *Machine Learning*, 48.
- Corduneanu, A. and Jaakkola, T. (2003). On Information Regularization. In *Uncertainty in Artificial Intelligence: Proceedings of the Nineteenth Conference (UAI-2003)*. Morgan Kaufmann Publishers.
- Cottrell, G., Munro, P., and Zipser, D. (1987). Learning internal representations from gray-scale images: An example of extensional programming. In *Proc. 9th Ann. Conf. of the Cognitive Science Society*.
- Courant, R. and Hilbert, D. (1953). *Methods of Mathematical Physics*. Interscience.
- Cover, T. M. and Thomas, J. A. (1991). *Elements of Information Theory*. Wiley, New York.
- Cowell, R. G., Dawid, A. P., Lauritzen, S. L., and Spiegelhalter, D. J. (1999). *Probabilistic Networks and Expert Systems (Information Science and Statistics S.)*. Springer-Verlag.
- Cramer, H. (1946). *Mathematical Methods of Statistics*. Princeton University Press.
- Cristianini, N. and Shawe-Taylor, J. (2000). *An Introduction to Support Vector Machines and Other Kernel-based Learning Methods*. Cambridge University Press.
- Dahlquist, G. (2003). *Numerical Methods*. Dover Publications.
- Darroch, J. N. and Ratcliff, D. (1972). Generalized Iterative Scaling for Log-Linear Models. *The Annals of Mathematical Statistics*, 43(1).
- Dayan, P. (2001). Unsupervised Learning. In Wilson, R. A. and Keil, F., editors, *The MIT Encyclopedia of the Cognitive Sciences*. MIT Press.
- Dayan, P., Hinton, G., Neal, R., and Zemel, R. (1995). The helmholtz machine. *Neural Computation*, 7.
- Dempster, A. P., Laird, M., and Rubin, D. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society*, 39(1).
- Dennis, J. E. and Schnabel, R. B. (1996). *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*. Dover.

- Dhillon, I. S. and Guan, Y. (2003). Information Theoretic Clustering of Sparse Co-Occurrence Data. In *Proceedings of the Third IEEE International Conference on Data Mining*.
- Dhillon, I. S., Guan, Y., and Kogan, J. (2002). Iterative clustering of high dimensional text data augmented by local search. In *Proceedings of the Second IEEE International Conference on Data Mining*.
- Dhillon, I. S., Guan, Y., and Kulis, B. (2004). Kernel k-means, Spectral Clustering and Normalized Cuts. In *Proceedings of The Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining(KDD)*.
- Diamantaras, K. I. and Kung, S. Y. (1996). *Principal Component Neural Networks: Theory and Applications*. John Wiley, New York.
- Doi, E. and Lewicki, M. (2004). Sparse coding of natural images using an over-complete set of limited capacity units. In *NIPS*.
- Eckmiller, R., Neumann, D., and Baruth, O. (2005). Tunable retina encoders for retina implants: why and how. *J. Neural Eng.*, 2:91–104.
- El-Hay, T. and Friedman, N. (2002). Incorporating Expressive Graphical Models in Variational Approximations: Chain-Graphs and Hidden Variables. In *UAI: Proceedings of the 18th Conference*.
- Everitt, B. S. (1984). *An Introduction to Latent Variable Models*. Chapman and Hall.
- Ewig, G. M. (1985). *Calculus of Variations with Applications*. Dover.
- Fano, R. M. (1961). *Transmission of Information: A Statistical theory of Communications*. Wiley: New York.
- Feller, W. (1971). *An Introduction to Probability Theory and Its Applications, Vol.2, 2nd edition*. John Wiley.
- Field, D. J. (1994). What is the goal of sensory coding? *Neural Computation*, 6:559–601.
- Fisher, J. W., Darrell, T., Freeman, W. T., and Viola, P. (2000). Learning joint statistical models for audio-visual fusion and segregation. In *Proc. Conf. Advances in Neural Information Processing Systems, NIPS*, volume 13. MIT Press.
- Fisher, J. W. and Principe, J. C. (1998). A methodology for information theoretic feature extraction. In *Proc. of the IEEE International Joint Conference on Neural Networks*.
- Fisher, R. A. (1922). On the Mathematical Foundations of the Theoretical Statistics. *Philosophical Transactions of the Royal Society, A*, 222:309–368.

- Fisher, R. A. (1925). Theory of statistical estimation. In *Proceedings of the Cambridge Philosophical Society*, volume 22. Cambridge University Press.
- Fisher, R. A. (1950). *Contributions to Mathematical Statistics*. John Wiley and Sons, New York, London.
- Fox, C. (1987). *An Introduction to the Calculus of Variations*. Dover.
- Friedman, N., Mosenzon, O., Slonim, N., and Tishby, N. (2001). Multivariate information bottleneck. In *Uncertainty in Artificial Intelligence: Proceedings of the Seventeenth Conference (UAI-2001)*, pages 152–161, San Francisco, CA. Morgan Kaufmann Publishers.
- Galeev, E. M. and Tihomirov, V. M. (2000). *Optimization*. Editorial URSS, Moscow.
- Gallager, R. G. (1963). *Low Density Parity Check Codes*. MIT Press, Cambridge, Massachusetts.
- Gamerman, D. (1997). *Markov Chain Monte Carlo: Stochastic Simulation of Bayesian Inference (Chapman and Hall Texts in Statistical Science)*. Taylor and Francis.
- Gelfand, I. M. and Fomin, S. V. (1963). *Calculus of Variations*. Prentice-Hall.
- Ghahramani, Z. and Hinton, G. E. (1996). The EM Algorithm for Mixtures of Factor Analyzers. Technical Report CRG-TR-96-1, University of Toronto.
- Ghahramani, Z. and Hinton, G. E. (1998). Switching state-space models. Technical report, Dept. of Computer Science, University of Toronto, 6 King’s College Road, Toronto M5S 3H5, Canada.
- Ghahramani, Z. and Hinton, G. E. (2000). Variational learning for switching state-space models. *Neural Computation*, 12(4):831–864.
- Ghahramani, Z. and Jordan, M. (1995). Factorial hidden Markov models. In Touretzky, D. S., Mozer, M. C., and Hasselmo, M. E., editors, *Proc. Conf. Advances in Neural Information Processing Systems, NIPS*, volume 8, pages 472–478. MIT Press.
- Gibbs, M. N. and MacKay, D. J. C. (2000). Variational gaussian process classifiers. *IEEE-NN*, 11(6).
- Glover, I. and Grant, P. (2003). *Digital Communications*. Prentice Hall.
- Gokcay, E. and Principe, J. C. (2002). Information theoretic clustering. In *IEEE Trans. Pattern Analysis and Machine Intelligence*, volume 24, pages 158–171.
- Golub, G. H. and Loan, C. F. V. (1996). *Matrix Computations (3rd ed.)*. Johns Hopkins University Press.

- Hardy, G. H., Littlewood, J. E., and Polya, G. (1988). *Inequalities, 2nd ed.* Cambridge University Press.
- Hartigan, J. A. and Wong, M. A. (1979). A K-means Clustering Algorithm. *Applied Statistics*, 28:100–108.
- Haussler, D. and Opper, M. (1997). Mutual Information, Metric Entropy, and Cumulative Relative Entropy Risk. *The Annals of Statistics*, 25(6).
- Hertz, J., Krogh, A., and Palmer, R. G. (1991). *Introduction to the Theory of Neural Computation*. MA: Addison-Wesley Publishing Company.
- Heskes, T. (2002). Stable fixed points of loopy belief propagation are minima of the Bethe free energy. In *NIPS*.
- Higdon, D. M. (1998). Auxiliary variable methods for Markov chain Monte Carlo with applications. *Journal of the American Statistical Association*, 93(442):585–595.
- Hinton, G. (1989). Connectionist learning procedures. *Artificial Intelligence*, 40.
- Hinton, G. E. and Sejnowski, T. J., editors (1999). *Unsupervised Learning*. MIT Press.
- Hoch, T., Wenning, G., and Obermayer, K. (2003). Optimal noise-aided signal transmission through populations of neurons. *Physical Review E*, 68.
- Ihler, A. T., III, J. W. F., and Willsky, A. S. (2005). Loopy Belief Propagation: Convergence and Effects of Message Errors. *JMLR*, 6.
- Jaakkola, T. (1997). *Variational Methods for Inference and Estimation in Graphical Models*. PhD thesis, Department of Cognitive Sciences, MIT.
- Jaakkola, T. S. and Jordan, M. I. (1996). Computing upper and lower bounds on likelihoods in intractable networks. Technical Report AIM-1571, MIT.
- Jaakkola, T. S. and Jordan, M. I. (1998). Improving the Mean Field Approximation via the Use of Mixture Distributions. In Jordan, M. I., editor, *Learning in Graphical Models*. Kluwer Academic Publishers.
- Jacobs, R. A., Jordan, M., Nowlan, S. J., and Hinton, G. E. (1991). Adaptive Mixtures of Local Experts. *Neural Computation*, 3.
- Jaynes, E. T. (1996). *Probability Theory: The Logic of Science*. Washington University, St. Louis. Fragmentary edition.
- Jensen, F. (1996). *An introduction to Bayesian networks*. UCL Press, London.
- Jensen, J. L. W. V. (1906). Sur les fonctions convexes et les inegalites entre les valeurs moyennes. *Acta Math*, 30.

- Jenssen, R., Hild, K. E., Erdogmuz, D., Principe, J. C., and Eltoft, T. (2003). Clustering using Renyi's Entropy. In *International Joint Conference on Artificial Neural Networks*.
- Johnson, D. (2003). Cramer-Rao Bound. The Connexions Project: m11266.
- Jordan, M. (2005). Introduction to graphical models. Unfinished.
- Jordan, M., Ghahramani, Z., Jaakkola, T., and Saul, L. (1998). An Introduction to Variational Methods for Graphical Models. In Jordan, M. I., editor, *Learning in Graphical Models*, chapter 1. Kluwer Academic Publishers.
- Jordan, M. I., editor (1998). *Learning in Graphical Models*. Kluwer Academic Publishers.
- Kang, K. and Sompolinsky, H. (2001). Mutual information of population codes and distance measures in probability space. *Physical Review Letters*, 86(21).
- Kannan, R., Vempala, S., and Vetta, A. (2000). On clusterings – good, bad, and spectral. In *Proceedings of the 41st Annual Symposium on Foundations of Computer Science*, volume 41.
- Kikuchi, R. (1951). A theory of cooperative phenomena. *Phys. Rev.*, 81(6).
- Korn, G. A. and Korn, T. M. (1968). *Mathematical Handbook for Scientists and Engineers: Definitions, Theorems, and Formulas for Reference and Review*. McGraw-Hill, New York.
- Kullback, S. (1959). *Information Theory and Statistics*. New York: Dover Publications.
- Kullback, S. and Leibler, R. A. (1951). On information and sufficiency. *Annals of Mathematical Statistics*, 22.
- Landau, L. D. and Lifshitz, E. M. (1996). *Statistical Physics: Course of Theoretical Physics*. Butterworth-Heinemann.
- Lauritzen, S. L. and Spiegelhalter, D. J. (1988). Local computations with probabilities on graphical structures and their applications to expert systems. *Journal of Royal Statistical Society B*, 50(2):157 – 224.
- Lawrence, N. (2000). *Variational Inference in Probabilistic Models*. PhD thesis, Computer Laboratory, University of Cambridge.
- Lawrence, N. D. (2003). Gaussian Process Latent Variable Models for Visualization of High Dimensional Data. In *NIPS*.
- Lawrence, N. D., Bishop, C. M., and Jordan, M. I. (1998). Mixture Representations for Inference and Learning in Boltzmann Machines. In *UAI: Proceedings of the 14th Conference*.

- LeCun, Y. and Cortes, C. (1998). The MNIST Database. Available at <http://yann.lecun.com/exdb/mnist/>.
- Lee, B. B., Kremer, J., and Yeh, T. (1998). Receptive fields of primate retinal ganglion cells studied with a novel technique. *Visual Neuroscience*, 15.
- Linsker, R. (1988). Towards an Organizing Principle for a Layered Perceptual Network. In *Advances in Neural Information Processing Systems*. American Institute of Physics.
- Linsker, R. (1989a). An Application of the Principle of Maximum Information Preservation to Linear Systems. In Touretzky, D., editor, *Advances in Neural Information Processing Systems*, volume 1. Morgan-Kaufmann.
- Linsker, R. (1989b). How to generate ordered maps by maximizing the mutual information between input and output signals. *Neural Computation*, 1.
- Linsker, R. (1992). Deriving Receptive Fields Using an Optimal Encoding Criterion. In Steven Hanson, Jack Cowan, L. G. e., editor, *Advances in Neural Information Processing Systems*, volume 5. Morgan-Kaufmann.
- Linsker, R. (1997). A local learning rule that enables information maximization for arbitrary input distributions. *Neural Computation*, 9.
- Luby, M. G., Mitzenmacher, M., Shokrollahi, M. A., and Spielman, D. A. (2001). Improved low-density parity-check codes using irregular graphs and belief propagation. *IEEE Trans. Info Theory*, 47(2).
- Luenberger, D. (1973). *Introduction to Linear and Nonlinear Programming*. Addison-Wesley.
- Luenberger, D. G. (1998). *Optimization by Vector Space Methods*. Wiley Professional Paperbacks.
- MacKay, D. (1998). Introduction to Monte Carlo methods. In Jordan, M., editor, *Learning in Graphical Models*, chapter 1. MIT Press.
- MacKay, D. (1999a). Good error correcting codes based on very sparse matrices. *IEEE Trans. Info. Theory*, 45(2).
- MacKay, D. (1999b). Maximum likelihood and covariant algorithms for independent components analysis. Technical report, University of Cambridge.
- MacKay, D. and Neal, R. (1999). Good Error-Correcting Codes based on Very Sparse Matrices. In *IEEE-IT*.
- MacKay, D. J. C. (1997). Gaussian processes - a replacement for supervised neural networks? Lecture notes for a tutorial at NIPS.
- MacKay, D. J. C. (2003). *Information Theory, Inference and Learning Algorithms*. Cambridge University Press.

- Magnus, J. R. and Neudecker, H. (1999). *Matrix differential calculus with applications in statistics and econometrics*. Wiley, New York, second edition.
- Markram, H., Wang, Y., and Tsodyksf, M. (1998). Differential signalling via the same axon of neocortical pyramidal neurons. *Neurobiology*, 95(9).
- McEliece, R. J. (1977). *The Theory of Information and Coding*. Addison-Wesley.
- Mercer, J. (1909). Functions of positive and negative type and their connection with the theory of integral equations. *Philos. Trans. Roy. Soc.*, 209.
- Mika, S., Schoelkopf, B., Smola, A., Muller, K.-R., Scholz, M., and Ratsch, G. (1999). Kernel PCA and De-Noising in Feature Spaces. *Advances in Neural Information Processing Systems*, 11.
- Minka, T. (2000). Old and new matrix algebra useful for statistics. Technical report, MIT, Available at <ftp://vismod.www.media.mit.edu/pub/tpminka/papers/minka-matrix.ps.gz>.
- Minka, T. (2003). A comparison of numerical optimizers for logistic regression. Technical report, CMU.
- Mooij, J. M. and Kappen, H. J. (2005). Sufficient Conditions for Convergence of Loopy Belief Propagation. In *UAI: Proceedings of the 21st Conference*.
- Morris, R. (1999). Auxiliary variables for markov random fields with higher order interactions. In Hancock, E. R. and Pelillo, M., editors, *Energy Minimization Methods in Computer Vision and Pattern Recognition*, volume 1654 of *Lecture Notes in Computer Science*. Springer.
- Murphy, K. P., Weiss, Y., and Jordan, M. I. (1999). Loopy belief propagation for approximate inference: An empirical study. In *Uncertainty in Artificial Intelligence*, volume 9.
- Nadal, J.-P., Brunel, N., and Parga, N. (1998). Nonlinear feedforward networks with stochastic outputs: infomax implies redundancy reduction. *Network: Computation in Neural Systems*, 9(2):207–217.
- Nadal, J.-P. and Parga, N. (1994). Nonlinear neurons in the low-noise limit: a factorial code maximizes information transfer. *NETWORK*, 5:565 – 581.
- Neal, R. (1993). Probabilistic inference using Markov Chain Monte Carlo methods. CRG-TR-93-1, Dept. of Computer Science, University of Toronto.
- Neal, R. M. (1992). Connectionist learning of belief networks. *Artificial Intelligence*, 56:71 – 113.
- Neal, R. M. and Hinton, G. E. (1998). A View of the EM Algorithm That Justifies Incremental, Sparse, and Other Variants. In Jordan, M., editor, *Learning in Graphical Models*, chapter 1. MIT Press.

- Neuneier, R., Hergert, F., Finnoff, W., and Ormoneit, D. (1994). Estimation of Conditional Densities: A Comparison of Neural Network Approaches. In *ICANN*.
- Ng, A. Y., Jordan, M., and Weiss, Y. (2001). On spectral clustering: Analysis and an algorithm. In *Proceedings of Neural Information Processing Systems*, volume 14.
- Oja, E. (1989). Neural networks, principal components, and subspaces. *International Journal on Neural Systems*, 1.
- Opper, M. and Haussler, D. (1995). Bounds for Predictive Errors in the Statistical Mechanics of Supervised Learning. *Physical Review Letters*, 75(20).
- Papoulis, A. (1984). *Probability, Random Variables, and Stochastic Processes*, 2nd ed. McGraw-Hill.
- Pearl, J. (1988). *Probabilistic Reasoning in Intelligent Systems*. Morgan Kaufman, San Mateo.
- Pearlmutter, B. A. and Parra, L. C. (1996). A context-sensitive generalization of ica. In *Proceedings of the International Conference on Neural Information Processing*.
- Penney, R. W. and Sherrington, D. (1993). The weight-space of the binary perceptron. *J. Phys. A: Math. Gen.*, 26:6173–6185.
- Perona, P. and Freeman, W. T. (1998). A factorization approach to grouping. In Burkardt, H. and Neumann, B., editors, *Proceedings of European Conference on Computer Vision*.
- Pinsker, M. S. (1964). *Information and information stability of random variables and processes*. San Francisco: Holden-Day.
- Pouget, A., Zhang, K. C., Deneve, S., and Latham, P. E. (1998). Statistically efficient estimation using population code. *Neural Computation*, 10:373–401.
- Press, W. H., Teukolsky, S. A., Vetterling, W. T., and Flannery, B. P. (1992). *Numerical Recipes in C*. Cambridge University Press, Second edition.
- Pretzel, O. (1996). *Error-Correcting Codes and Finite Fields*. Clarendon Press, Oxford.
- Principe, J., Fisher, J., and Xu, D. (2000). Information theoretic learning. In Haykin, S., editor, *Unsupervised Adaptive Filtering*. Wiley.
- Principe, J., Xu, D., and Fisher, J. (1998). Pose estimation in sar using an information-theoretic criterion. In *Proceedings of SPIE98*.
- Redlich, A. N. (1993). Redundancy reduction as a strategy for unsupervised learning. *Neural Computation*, 5:289–304.

- Riley, K. F., Hobson, M. P., and Bence, S. J. (2002). *Mathematical Methods for Physics and Engineering: A Comprehensive Guide*. Cambridge U Press.
- Rosen-Zvi, M. and Kanter, I. (2001). Training a perceptron in a discrete weight space. *Physical Review E*, 64.
- Roweis, S. and Ghahramani, Z. (1999). A unifying review of linear gaussian models. *Neural Computation*, 11(2).
- Rubinov, A. and Yang, X. (2003). *Lagrange-type Functions in Constrained Non-convex Optimization (Applied Optimization S.)*. Kluwer Academic Publishers.
- Saad, D. and Opper, M. (2001). *Advanced Mean Field Methods Theory and Practice*. MIT Press.
- Samengo, I. and Treves, A. (2001). Representational capacity of a set of independent neurons. *Physical Review E*, 63.
- Saul, L., Jaakkola, T., and Jordan, M. (1996). Mean field theory for sigmoid belief networks. *Journal of Artificial Intelligence Research*, 4.
- Saul, L. and Jordan, M. (1998). A mean field learning algorithm for unsupervised neural networks. In Jordan, M. I., editor, *Learning in Graphical Models*, chapter 3. MIT Press.
- Schoelkopf, B., Smola, A., and Mueller, K. (1998). Nonlinear Component Analysis as a Kernel Eigenvalue Problem. *Neural Computation*, 10.
- Scholkopf, B. and Smola, A. (2002). *Learning with Kernels*. MIT Press.
- Shannon, C. E. (1948). A Mathematical Theory of Communication. *The Bell Systems Technical Journal*, 27.
- Shi, J. and Malik, J. (2000). Normalized Cuts and Image Segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):888–905.
- Shriki, O., Sompolinsky, H., and Lee, D. D. (2002). An information maximization approach to overcomplete and recurrent representations. In *NIPS*, pages 612 – 618. MIT Press.
- Slonim, N. (2002). *The Information Bottleneck: Theory and Applications N. Slonim*. PhD thesis, The Hebrew University.
- Slonim, N., Friedman, N., and Tishby, N. (2001). Agglomerative multivariate information bottleneck. In *Advances in Neural Information Processing Systems (NIPS)*.
- Slonim, N. and Weiss, Y. (2002). Maximum Likelihood and the Information Bottleneck. In *NIPS*. MIT Press.
- Smith, D. R. (1998). *Variational Methods in Optimization*. Dover.

- Smola, A. J. (1998). *Learning with Kernels*. PhD thesis, Technische Universität Berlin.
- Still, S. and Bialek, W. (2004). How many clusters? an information theoretic perspective. *Neural Computation*, 16(12).
- Still, S., Bialek, W., and Bottou, L. (2004). Geometric Clustering using the Information Bottleneck method. In *NIPS*. MIT Press.
- Stocks, N. G. and Mannella, R. (2001). Generic noise-enhanced coding in neuronal arrays. *Physical Review E*, 64.
- Swendsen, R. and Wang, J.-S. (1987). Nonuniversal critical dynamics in monte carlo simulations. *Physical Review Letters*, 58:86–88.
- Szummer, M. and Jaakkola, T. (2002). Information regularization with partially labeled data. In *Advances in Neural Information Processing Systems (NIPS)*.
- Tipping, M. E. and Bishop, C. M. (1999a). Mixtures of Probabilistic Principal Component Analyzers. *Neural Computation*, 11(2):443–482.
- Tipping, M. E. and Bishop, C. M. (1999b). Probabilistic principal components analysis. *J. Roy. Statistical Society B*, 61(3):611–622.
- Tishby, N., Pereira, F. C., and Bialek, W. (1999). The information bottleneck method. In *Proc. of the 37-th Annual Allerton Conference on Communication, Control and Computing*, pages 368–377.
- Torkkola, K. (2000). Visualizing class structure in data using mutual information. In *NNSP*.
- Torkkola, K. (2001). Learning discriminative feature transforms to low dimensions in low dimensions. In *NIPS*.
- Torkkola, K. and Campbell, W. M. (2000). Mutual Information in Learning Feature Transformations. *Proc. 17th International Conf. on Machine Learning*.
- Vapnik, V. (1998). *Statistical Learning Theory*. Wiley.
- Viterbi, A. J. (1995). *CDMA, Principles of Spread Spectrum*. Addison-Wesley.
- von Mises, R. (1964). *Mathematical Theory of Probability and Statistics*. Academic Press, New York.
- Wainwright, M. (2002). *Stochastic processes on graphs with cycles: geometric and variational approaches*. PhD thesis, Electrical Engineering and Computer Science, MIT.
- Wainwright, M., Jaakkola, T., and Willsky, A. (2001). Tree-based reparameterization framework for approximate estimation on graphs with cycles. In *IEEE Trans. on Information Theory*.

- Wainwright, M., Jaakkola, T., and Willsky, A. (2002). A new class of upper bounds on the log partition function. In *Uncertainty in Artificial Intelligence*.
- Wang, Y., Jiar, Y., Hu, C., and Turk, M. (2004). Face recognition based on kernel radial basis function networks. In *Asian Conference on Computer Vision*.
- Watanabe, M. and Rodieck, R. W. (1989). Parasol and midget ganglion cells of the primate retina. *J. Comp. Neurol.*, 289:333–353.
- Weinstock, R. (1974). *Calculus of Variations, With Applications to Physics and Engineering: With Applications to Physics and Engineering*. Dover.
- Weiss, Y. (1999). Segmentation using eigenvectors: A unifying view. In *International Conference on Computer Vision*.
- Weiss, Y. and Freeman, W. T. (2001). Correctness of belief propagation in Gaussian graphical models of arbitrary topology. *Neural Computation*, 13(10).
- Williams, C. K. I. (1997). Regression with Gaussian Processes . In S. W. Ellacott, J. C. M. and Anderson, I. J., editors, *Mathematics of Neural Networks: Models, Algorithms and Applications*. Kluwer.
- Williams, C. K. I. (1998). Prediction with Gaussian Processes: From Linear Regression to Linear Prediction and Beyond. In Jordan, M. I., editor, *Learning in Graphical Models*. Kluwer Academic Publishers.
- Williams, C. K. I. and Barber, D. (1998). Bayesian Classification with Gaussian Processes. *IEEE Trans Pattern Analysis and Machine Intelligence*, 20(12).
- Williams, C. K. I. and Rasmussen, C. E. (1996). Gaussian processes for regression. In Touretzky, D. S., Mozer, M. C., and Hasselmo, M. E., editors, *Advances in Neural Information Processing Systems*, volume 8. The MIT Press.
- Yedidia, J., Freeman, W., and Weiss, Y. (2000a). Bethe Free Energy, Kikuchi Approximations, and Belief Propagation Algorithms. Technical report, MERL.
- Yedidia, J. S., Freeman, W. T., and Weiss, Y. (2000b). Generalized belief propagation. In *NIPS*, pages 689–695.
- Yedidia, J. S., Freeman, W. T., and Weiss, Y. (2004). Constructing Free Energy Approximations and Generalized Belief Propagation Algorithms. Technical Report TR-2004-040, MERL.
- Yu, S. X. and Shi, J. (2003). Multiclass spectral clustering. In *International Conference on Computer Vision*.
- Yuille, A. L. (2002). CCCP algorithms to minimize the Bethe and Kikuchi free energies: Convergent alternatives to belief propagation. *Neural Computation*, 14(7).
- Zemel, R. (1993). *A Minimum Description Length Framework for Unsupervised Learning*. PhD thesis, University of Toronto.

- Zemel, R. and Hinton, G. (1994). Autoencoders, minimum description length and helmholtz free energy. In *NIPS*.
- Zemel, R. and Hinton, G. (1995). Developing population codes by minimizing description length. *Neural Computation*, 7(3).
- Zha, H., Ding, C., Gu, M., He, X., and Simon, H. (2001). Spectral relaxation for k-means clustering. In *Proceedings of Neural Information Processing Systems*, volume 14.
- Zhang, K. and Sejnowski, T. J. (1999). Neuronal tuning: To sharpen or broaden? *Neural Computation*, 11:75–84.
- Zhang, R. and Rudnicky, A. I. (2002). A large scale clustering scheme for kernel k-means. In *Proceedings of International Conference on Pattern Recognition*.