



VARIATIONAL INFORMATION
MAXIMIZATION FOR POPULATION
CODING

Felix Agakov ^a David Barber ^b

IDIAP-RR 04-85

^a University of Edinburgh, EH12QL, Edinburgh, UK
^b IDIAP, Martigny, Switzerland

VARIATIONAL INFORMATION MAXIMIZATION FOR POPULATION CODING

Felix Agakov

David Barber

Abstract. The goal of neural processing assemblies is varied, and in many cases still rather unclear. However, a possibly reasonable *subgoal* is that sensory information may be encoded efficiently in a population of neurons. In this context, Mutual Information is a long studied measure of coding efficiency, and many attempts to apply this to *population coding* have been made. However, this is a numerically intractable task, and most previous studies redefine the criterion in forms of an approximation to Mutual Information, the Fisher Information being one such well-known approach. Here we describe a principled bound maximisation procedure for Mutual Information learning of population codes in a simple point neural model, and compare it with other approaches.

1 Introduction

Whilst the use of information theory cannot be justified as a central goal of neural processing – since arguably it is compression and decisions that are key to survival – nevertheless, much effort has been focussed on a possible *subgoal*, namely that of maximal information transmission from a stimulus to its representation as a set of spiking neurons [16, 3, 14].

The problem of encoding real-valued stimuli \mathbf{x} by a population of neural spikes \mathbf{y} may be addressed in many different ways. Essentially the framework is to, for a given set of patterns \mathbf{x} , to adapt the parameters of any mapping $p(\mathbf{y}|\mathbf{x})$ to make a desirable population code. There are many possible desiderata.

One could be that any reconstruction based on the population should be accurate. This is typically handled by appealing to the Fisher Information which, with care, can be used bound mean square reconstruction error. Another approach is to bound the *probability* of a correct reconstruction (to the best of our knowledge, surprisingly, this has not yet been studied). The approach we examine here, is to consider maximizing the amount of information which the spiking patterns contain about the stimuli [6, 13]. The fundamental information theoretic measure in this context is the mutual information

$$I(\mathbf{x}, \mathbf{y}) \equiv H(\mathbf{x}) - H(\mathbf{x}|\mathbf{y}), \quad (1)$$

which indicates the decrease of uncertainty in \mathbf{x} due to the knowledge of \mathbf{y} . Here $H(\mathbf{x}) \equiv -\langle \log p(\mathbf{x}) \rangle_{p(\mathbf{x})}$ 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 averages over all variables contained within the brackets.

It is sometimes argued that MI is desirable over and above other approaches which require the explicit definition of a reconstruction procedure [16]. However, we will argue that the specification of a reconstruction/decoding procedure is essentially unavoidable.

The principled information theoretic approach to learning neural codes involves maximization of the objective (1) with respect to parameters of the encoder $p(\mathbf{y}|\mathbf{x})$. However, it is easy to see that in large-scale systems exact evaluation of $I(\mathbf{x}, \mathbf{y})$ is in general computationally intractable. The key difficulty lies in the computation of the conditional entropy $H(\mathbf{x}|\mathbf{y})$ for the posterior distribution $p(\mathbf{x}|\mathbf{y})$, which is tractable only in a few special cases. Standard techniques address the problem of optimizing (1) by assuming that $p(\mathbf{x}, \mathbf{y})$ is jointly Gaussian [11], the output spaces are very low-dimensional [12], or the channels are deterministic and invertible [4]. Other popular methods suggest alternative objective functions (e.g. approximations based on the *Fisher Information* criterion [6]), which, however, do not retain proper bounds on $I(\mathbf{x}, \mathbf{y})$.

Recently we described a simple variational approach to information maximization which optimizes a proper lower bound on the mutual information [2]. In this paper we investigate applicability of the method in the context of neural coding. First, we briefly review the lower bound on $I(\mathbf{x}, \mathbf{y})$. Then we analyze the learning rule obtained by maximizing the bound for sigmoidal networks. Finally, we analyze the relation between our approach and standard techniques for approximate information maximization, focusing specifically on a comparison with the Fisher Information criterion.

1.1 Variational Lower Bound on Mutual Information

A simple lower bound on the mutual information $I(\mathbf{x}, \mathbf{y})$ follows from non-negativity of the Kullback-Leibler divergence $KL(p(\mathbf{x}|\mathbf{y})||q(\mathbf{x}|\mathbf{y}))$ between the exact posterior $p(\mathbf{x}|\mathbf{y})$ and its variational approximation $q(\mathbf{x}|\mathbf{y})$. Clearly,

$$\langle \log p(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{x}|\mathbf{y})} - \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{x}|\mathbf{y})} \geq 0 \Rightarrow \langle \log p(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{x}|\mathbf{y})p(\mathbf{y})} \geq \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{x}|\mathbf{y})p(\mathbf{y})}. \quad (2)$$

This leads to the variational lower bound

$$I(\mathbf{x}, \mathbf{y}) \geq \tilde{I}(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} H(\mathbf{x}) + \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{x}, \mathbf{y})}, \quad (3)$$

where $q(x|y)$ is an arbitrary variational distribution saturating the bound for $q(x|y) \equiv p(x|y)$. Note that the objective (3) explicitly includes¹ both the encoder $p(y|x)$ (distribution of neural spikes for a given stimulus) and decoder $q(x|y)$ (reconstruction of the stimulus from a population of neural firings). It is possible to consider other lower bounds on the mutual information [9]. However, the flexibility of the choice of the decoder $q(x|y)$ makes the bound (3) particularly computationally convenient. Moreover, (3) corresponds to a moment matching approximation of $p(x|y)$ by $q(x|y)$, which is particularly beneficial in terms of decoding [15, 2].

2 Variational Learning of Population Codes

To learn optimal stochastic representations of the continuous training patterns x_1, \dots, x_M according to the lower bound $\tilde{I}(x, y)$, we need to choose a continuous density function for the decoder $q(x|y)$. Computationally, it is convenient to assume that the decoder is given by the isotropic Gaussian $q(x|y) \sim \mathcal{N}(Uy, \sigma^2 I)$, where $U \in \mathbb{R}^{|\mathcal{X}| \times |\mathcal{Y}|}$. In this case, exact evaluation of the bound $\tilde{I}(x, y)$ is straightforward, since it only involves computations of the second-order moments of y over the factorized distribution. Of course, other (correlated and non-linear) choices may be considered – however, for clarity we limit the discussion to isotropic linear Gaussian mappings. Effectively, this choice indicates that small changes in the post-synaptic firings do not significantly vary our guesses about the generating stimuli.

For the empirical distribution $p(x) = \sum_{m=1}^M \delta(x - x_m)/M$ the bound (3) is then given by

$$\tilde{I}(x, y) \propto \sum_{m=1}^M \text{tr} \left\{ U \langle y \rangle_{p(y|x_m)} x_m^T - \frac{1}{2} U^T U \langle yy^T \rangle_{p(y|x_m)} \right\} + c \quad (4)$$

where c is a constant. Expressing $\tilde{I}(x, y)$ as a function of the encoder $p(y|x)$ alone, we get

$$U = \langle xy^T \rangle \langle yy^T \rangle^{-1}, \quad \tilde{I}(x, y) \propto \text{tr} \left\{ \langle xy^T \rangle \langle yy^T \rangle^{-1} \langle yx^T \rangle \right\} + c. \quad (5)$$

Note that the objective (5) is a proper bound for any choice of the stochastic mapping $p(y|x)$. We may therefore² use it for optimizing a variety of channels with continuous source vectors.

Whilst making an objective without the explicit appearance of the decoding weights U is theoretically attractive, the resulting learning rules for the encoding weights will become biologically unrealistic. However, the reader should bear in mind that in the sequel, local rules can be formed by incremental optimisation of (4) jointly with respect to the encoding and decoding weights.

2.1 Sigmoidal Activations

Here we consider the case of high-dimensional continuous patterns $x \in \mathbb{R}^{|\mathcal{X}|}$ represented by stochastic firings of the post-synaptic neurons $y \in \{-1, +1\}^{|\mathcal{Y}|}$. For each neuron y_i we assume the logistic parameterization of the encoder $p(y_i|x)$, so that the probability of firing monotonically increases with an increase in the membrane potential. For conditionally independent activations, we obtain

$$p(y|x) = \prod_{i=1}^{|\mathcal{Y}|} p(y_i|x) \stackrel{\text{def}}{=} \prod_{i=1}^{|\mathcal{Y}|} \sigma(y_i(w_i^T x + b_i)) \quad (6)$$

where $w_i \in \mathbb{R}^{|\mathcal{X}|}$ is a vector of the synaptic weights for neuron y_i , b_i is the corresponding threshold, and $\sigma(a) \stackrel{\text{def}}{=} 1/(1 + e^{-a})$. If the decoder's weights U are unconstrained, we may optimize the bound

¹The bound (3) corresponds to the criteria used by Blahut-Arimoto algorithms (e.g. [7]); however, we optimize it for both encoder and decoder, which is constrained to lie in a tractable family.

²From (5) it is clear that if $\langle yy^T \rangle$ is near-singular, the varying part of the objective $\tilde{I}(x, y)$ may be infinitely large. However, if the mapping $x \mapsto y$ is probabilistic and the number of training stimuli M exceeds the dimensionality of the neural codes $|\mathcal{Y}|$, the optimized criterion is typically positive and finite.

(5), which in this case is given by

$$\tilde{I}(\mathbf{x}, \mathbf{y}) \propto \text{tr} \left\{ \langle \mathbf{x} \boldsymbol{\lambda}_x^T \rangle \langle \boldsymbol{\lambda}_x \boldsymbol{\lambda}_x^T - \mathbf{D}_{\lambda_x} + \mathbf{I} \rangle^{-1} \langle \boldsymbol{\lambda}_x \mathbf{x}^T \rangle \right\}. \quad (7)$$

Here $\boldsymbol{\lambda}_x \in [-1, 1]^{|y|}$ is a vector whose elements $\lambda_i(\mathbf{x}) \stackrel{\text{def}}{=} \langle y_i \rangle_{p(y_i|\mathbf{x})} = 2\sigma(\mathbf{w}_i^T \mathbf{x} + b_i) - 1$ correspond to expected firings of the i^{th} unit for a fixed stimulus \mathbf{x} , and

$$\mathbf{D}_{\lambda_x} \stackrel{\text{def}}{=} \text{diag} \left(\lambda_1^2(\mathbf{x}), \dots, \lambda_{|y|}^2(\mathbf{x}) \right) = \mathbf{I} - \text{cov}(\mathbf{y}|\mathbf{x}) \in [0, 1]^{|y| \times |y|} \quad (8)$$

is a measure of consistency of neural firings.

Since the lower bound (7) depends only on the thresholds and synaptic weights, the learning rule is easily obtained by differentiating (7) with respect to $\mathbf{b} \in \mathbb{R}^{|y|}$ and $\mathbf{W} \in \mathbb{R}^{|y| \times |x|}$ (where rows of \mathbf{W} are given by $\mathbf{w}_i^T \in \mathbb{R}^{1 \times |x|}$). This leads to

$$\Delta \mathbf{W} \propto \sum_{m=1}^M (\mathbf{I} - \mathbf{D}_{\lambda_{x_m}}) \left(\tilde{\mathbf{D}} \boldsymbol{\lambda}_{x_m} + \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{yx} (\mathbf{x}_m - \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\lambda}_{x_m}) \right) \mathbf{x}_m^T, \quad (9)$$

where $\boldsymbol{\Sigma}_{yy} \stackrel{\text{def}}{=} \langle \mathbf{y} \mathbf{y}^T \rangle$, $\boldsymbol{\Sigma}_{yx} \equiv \boldsymbol{\Sigma}_{xy}^T \stackrel{\text{def}}{=} \langle \mathbf{y} \mathbf{x}^T \rangle$ are the second-order moments, and $\tilde{\mathbf{D}}$ corresponds to the diagonal of $\boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{yx} (\boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{yx})^T$. The update for the threshold $\Delta \mathbf{b}$ has the same form as (9) without the post-multiplication of each term by the training stimulus \mathbf{x}_m^T .

From (9) it is clear that the magnitude of each weight update $\Delta \mathbf{w}_i \in \mathbb{R}^{|x|}$ decreases with a decrease in the corresponding conditional variance $\text{var}(y_i|\mathbf{x}_m)$. Effectively, this corresponds to a variable learning rate – as training continues and magnitudes of the synaptic weights increase, the firings become more deterministic, and learning slows down.

Although this update rule is biologically difficult to interpret, as we stated before, interpretable rules are available by optimising (4) with respect to both \mathbf{U} and \mathbf{W} .

3 Fisher Information and Mutual Information

Here we briefly review a popular class of approximations to the mutual information based on the Fisher Information criterion [6]. We also outline the corresponding learning rule for sigmoidal channels.

3.1 Fisher Approximation

Let $\hat{\mathbf{x}} \in \mathbb{R}^{|x|}$ be a statistical estimator of the input stimulus \mathbf{x} obtained from the stochastic neural firings \mathbf{y} . It is easy to see that $\mathbf{x} \rightarrow \mathbf{y} \mapsto \hat{\mathbf{x}}$ forms a Markov chain with $p(\hat{\mathbf{x}}|\mathbf{y}) \sim \delta(\hat{\mathbf{x}} - \hat{\mathbf{x}}(\mathbf{y}))$. If $\hat{\mathbf{x}}$ is *efficient*, its covariance saturates the Cramer-Rao bound (see e.g. [7]), which results in an upper bound on the entropy of the conditional distribution $H(p(\hat{\mathbf{x}}|\mathbf{x}))$. From the data processing inequality, one may obtain a lower bound on the mutual information

$$I(\mathbf{x}, \mathbf{y}) \geq H(\hat{\mathbf{x}}) + \frac{1}{2} \langle \log |\mathbf{F}_x| \rangle_{p(\mathbf{x})} + c, \quad (10)$$

where $\mathbf{F}_x = \{F_{ij}(\mathbf{x})\} \stackrel{\text{def}}{=} -\langle \partial^2 \log p(\mathbf{y}|\mathbf{x}) / \partial x_i \partial x_j \rangle_{p(\mathbf{y}|\mathbf{x})}$ is the Fisher Information matrix and c is an irrelevant constant. 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 (10) is in general computationally intractable. [6] show that under some assumptions $H(\hat{\mathbf{x}}) \approx H(\mathbf{x})$, leading to the approximation

$$I(\mathbf{x}, \mathbf{y}) \gtrsim \tilde{I}_F(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} H(\mathbf{x}) + \frac{1}{2} \langle \log |\mathbf{F}_x| \rangle_{p(\mathbf{x})} + c. \quad (11)$$

Note that since $H(\mathbf{x})$ is independent of the parameters of $p(\mathbf{y}|\mathbf{x})$, maximization of (11) is equivalent to maximization of (10) where the intractable entropic term is ignored.

3.2 Sigmoidal Activations

For sigmoidal activations (6), the criterion (11) is given by

$$\tilde{I}_F(x, y) \propto \sum_{m=1}^M \log |\mathbf{W}^T (\mathbf{I} - \mathbf{D}_{\lambda_{x_m}}) \mathbf{W}| + c. \quad (12)$$

Again, $\mathbf{W} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{X}|}$ is the matrix of synaptic weights, $\mathbf{I} - \mathbf{D}_{\lambda_{x_m}}$ is the conditional covariance of the stochastic spikes (see (8)), and c corresponds to the remaining terms which do not affect the optimization.

It is interesting to note that for the square model with $|\mathcal{X}| = |\mathcal{Y}|$, optimization of (12) leads to

$$\Delta \mathbf{W} = 2\mathbf{W}^{-T} - \langle \boldsymbol{\lambda}_x \mathbf{x}^T \rangle, \quad (13)$$

which (apart from the coefficient at the inverse weight – *redundancy* term) has the same form as the learning rule of [4] derived for the invertible channel with $p(y|x) \sim \delta(y - \sigma(x))$. If $|\mathcal{X}| < |\mathcal{Y}|$ and the conditional variances of the responses are invariant under the training stimuli, the redundancy term in (13) is replaced by the transposed pseudo-inverse $\mathbf{W}(\mathbf{W}^T \mathbf{W})^{-1} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{X}|}$. Notably, the weight update (13) has no Hebbian terms.

More importantly, from (12) it is clear that as the variance of the stochastic firings decreases, the objective $\tilde{I}_F(x, y)$ may become infinitely loose. Since directions of low variation swamp the volume of the manifold, neural spikes generated by a fixed stimulus may often be inconsistent. It is also clear that optimization of the Fisher information-based objective (12) is limited to the cases when $\mathbf{W}^T \mathbf{W} \in \mathbb{R}^{|\mathcal{X}| \times |\mathcal{X}|}$ is full-rank, which complicates applicability of the method for a variety of tasks involving relatively low-dimensional encodings of very high-dimensional stimuli.

4 Variational Lower Bound vs. Fisher Approximation

Since $\tilde{I}_F(x, y)$ is in general not a proper lower bound on the mutual information, it is difficult to analyze its tightness or compare it with the variational bound (3). To illustrate a relation between the approaches, we may consider a Gaussian decoder $q(x|y) \sim \mathcal{N}_x(\boldsymbol{\mu}_y; \boldsymbol{\Sigma})$, which transforms the variational bound into

$$\tilde{I}(x, y) = -\frac{1}{2} \langle \text{tr} \{ \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_y) (\mathbf{x} - \boldsymbol{\mu}_y)^T \} \rangle_{p(x, y)} + \frac{1}{2} \log |\boldsymbol{\Sigma}^{-1}| + \acute{c}. \quad (14)$$

Here \acute{c} incorporates $H(x)$ and other irrelevant constants, and $\boldsymbol{\Sigma} \in \mathbb{R}^{|\mathcal{X}| \times |\mathcal{X}|}$ is a function of parameters of the conditional $p(y|x)$. Clearly, if the log eigenspectrum of the inverse covariance of the decoder is constrained to satisfy

$$\sum_{i=1}^{|\mathcal{X}|} \log l_i(\boldsymbol{\Sigma}^{-1}) = \sum_{i=1}^{|\mathcal{X}|} \langle \log l_i(\mathbf{F}_x) \rangle_{p(x)}, \quad (15)$$

where $\{l_i(\boldsymbol{\Sigma}^{-1})\}$ and $\{l_i(\mathbf{F}_x)\}$ are eigenvalues of $\boldsymbol{\Sigma}^{-1}$ and \mathbf{F}_x respectively, then the lower bound (14) reduces to the objective (11) amended with the average quadratic reconstruction error

$$\tilde{I}(x, y) = -\frac{1}{2} \underbrace{\langle \text{tr} \{ \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_y) (\mathbf{x} - \boldsymbol{\mu}_y)^T \} \rangle_{p(x, y)}}_{\text{reconstruction error}} + \frac{1}{2} \underbrace{\langle \log |\mathbf{F}_x| \rangle_{p(x)}}_{\text{Fisher criterion}} + \acute{c}. \quad (16)$$

Arguably, it is due to the subtraction of the non-negative quadratic term that (14) remains a general lower bound independently of the parameterization of the model and spectral properties of \mathbf{F}_x .

Another principal advantage of the variational approach to information maximization is the flexibility in the choice of the decoder. Clearly, if the Fisher Information matrices are small or nearly

Table 1: Objective functions for approximate information maximization

1. *Invertible channels:* $I(\mathbf{x}, \mathbf{y}) = \langle \log |\mathbf{J}_x| \rangle_{p(\mathbf{x})}$
2. *Fisher approximation:* $\tilde{I}_F(\mathbf{x}, \mathbf{y}) = \langle \log |\mathbf{F}_x| \rangle_{p(\mathbf{x})}$
3. *Variational lower bound:* $\tilde{I}(\mathbf{x}, \mathbf{y}) = \langle \log q(\mathbf{x}|\mathbf{y}) \rangle_{p(\mathbf{x}, \mathbf{y})}$

singular, both (11) and (16) are quite weak. However, by relaxing (15) and choosing other non-singular covariances (or completely different decoder types), the variational bound may be significantly strengthened.

Table 1 summarizes effective criteria optimized by [4], [6], and the variational approach of [2]. Here $\mathbf{J}_x = \{J_{ij}(\mathbf{x})\} \stackrel{\text{def}}{=} \{\partial y_i(\mathbf{x})/\partial x_j\} \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|}$ is the Jacobian of the deterministic invertible mapping $\mathbf{x} \mapsto \mathbf{y}$ (with $|\mathbf{y}| = |\mathbf{x}|$), $\mathbf{F}_x \in \mathbb{R}^{|\mathbf{x}| \times |\mathbf{x}|}$ is the Fisher Information matrix, and $q(\mathbf{x}|\mathbf{y})$ is an approximate decoder lying in a tractable family. In contrast to the first two 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 generalizes applications to population coding.

5 Experiments

Variational IM vs Fisher criterion

In the first set of experiments we were interested to see if, by maximising our bound $\tilde{I}(\mathbf{x}, \mathbf{y})$ on the MI, the parameters found at each iteration indeed increased the true MI. We compared this with how the value of the true MI changed as the parameters were updated by maximising the Fisher criterion $\tilde{I}_F(\mathbf{x}, \mathbf{y})$. The dimensionality of the response variables $|\mathbf{y}|$ was set to be sufficiently small, so that the true mutual information $I(\mathbf{x}, \mathbf{y})$ could be computed. Figure 1 illustrates changes in $I(\mathbf{x}, \mathbf{y})$ with iterations of the variational and Fisher-based learning, where the variational decoder was chosen to be an isotropic linear Gaussian with the optimal weights (5). We found that for $|\mathbf{x}| \leq |\mathbf{y}|$ (Figure 1 (*left*)), both approaches tend to lead to consistent improvements in the true mutual information (however, the variational approach typically resulted in higher values of $I(\mathbf{x}, \mathbf{y})$ after just a few iterations). For $|\mathbf{x}| > |\mathbf{y}|$ (Figure 1 (*right*)), optimization of the Fisher criterion was numerically unstable and lead to no visible improvements of $I(\mathbf{x}, \mathbf{y})$ over its starting value at initialization. Further approximations aimed at handling singularities of the gradients could lead to slight improvements, though their dynamics was rather inconsistent.

Variational IM: stochastic representations of the digit data

Here we apply the variational learning to stochastic coding and reconstruction of visual patterns. In our experiments we used the simplest form of the linear Gaussian decoder discussed in Section 2.1. After numerical optimization with an explicit constraint on the channel noise, we performed reconstruction of 196-dimensional continuous visual stimuli from 10 spiking neurons. The training stimuli consisted of 30 instances of digits 1 and 8 (15 of each class). The source variables were reconstructed from 50 stochastic spikes at the mean of the optimal approximate decoder $q(\mathbf{x}|\mathbf{y})$. Note that since $|\mathbf{x}| > |\mathbf{y}|$, the problem cannot be efficiently addressed by optimization of the Fisher Information-based criterion (12). Clearly, the approach of [4] is not applicable either, due to its fundamental assumption of invertible mappings between the spikes and the visual stimuli. Figure 2 illustrates a subset of the original source signals, samples of the corresponding binary responses, and reconstructions of the source data.

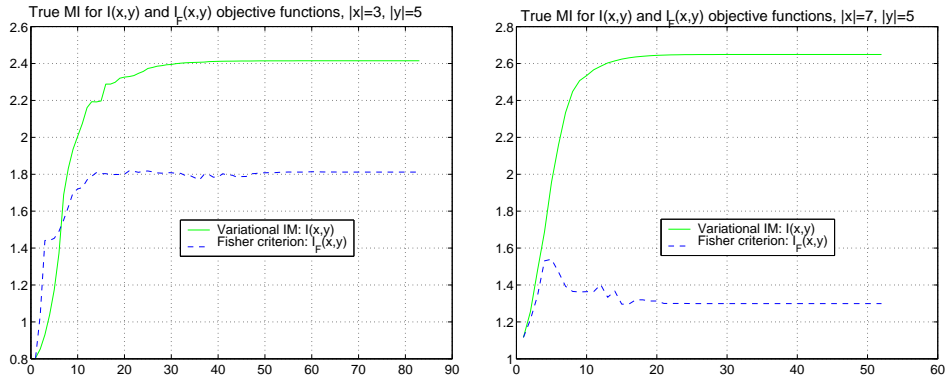


Figure 1: Changes in the exact mutual information $I(x,y)$ for parameters of the coder $p(y|x)$ obtained by maximizing the variational lower bound and the Fisher information criterion for $M = 20$ training stimuli. *Left:* $|x| = 3$, $|y| = 5$ *Right:* $|x| = 7$, $|y| = 5$. Generally, using the parameters given at each iteration from the bounding procedure, *does* increase the MI. However, those parameters given by the Fisher criterion, do not generally increase the MI.

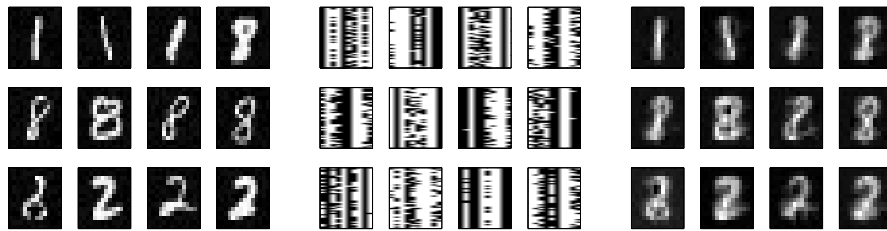


Figure 2: *Left:* a subset of the original visual stimuli. *Middle:* 20 samples of the corresponding spikes generated by each of the 10 neurons. *Right:* Reconstructions from 50 samples of neural spikes. Note that we numerically constrained slightly the weights so as to a solution being reached which would constitute deterministic firings.

6 Discussion

We described a variational approach to information maximization for the case when continuous source stimuli are represented by stochastic binary responses. Our results indicate that other approximate methods for information maximization [12, 6] may be viewed as approximations of our approach, which, however, do not always preserve a proper bound on the mutual information. We do not wish here to discredit generally the use of the Fisher Criterion, since this can be relevant for bounding reconstruction error. However, for the case considered here as a method for maximising information, we believe that our method is more attractive. An obvious extension of our work would be to more realistic spiking networks[1]. However, this case is difficult in the variational information maximisation framework, which cannot easily deal with correlations in the neural firings. Surrogate Mutual Information methods have been applied in this case[8], although we believe that an approach similar to a standard variational treatment of input-output Hidden Markov Models[5, 10], under a maximum probability of reconstruction measure would be relatively straightforward and arguably more desirable surrogate.

Whilst this paper is largely theoretical, we hope to have sustained the reader's interest at least to demonstrate that there is indeed a *principled* way to maximise information transmission to form efficient population codes. Its possible biological interpretation may interest others more expert than

ourselves in that arena.

- [1] D. Barber. Learning in Spiking Neural Assemblies. *NIPS*, 2002.
- [2] D. Barber and F. V. Agakov. The IM Algorithm: A Variational Approach to Information Maximization. In *NIPS*. MIT Press, 2003.
- [3] H. Barlow. Unsupervised Learning. *Neural Computation*, 1:295–311, 1989.
- [4] A. J. Bell and T. J. Sejnowski. An information-maximization approach to blind separation and blind deconvolution. *Neural Computation*, 7(6):1129–1159, 1995.
- [5] Y. Bengio and P. Frasconi. An input-output HMM architecture. *NIPS 7*, 1995.
- [6] N. Brunel and J.-P. Nadal. Mutual Information, Fisher Information and Population Coding. *Neural Computation*, 10:1731–1757, 1998.
- [7] T. M. Cover and J. A. Thomas. *Elements of Information Theory*. Wiley, New York, 1991.
- [8] G. Deco and B. Schürmann. Spatiotemporal Coding in the Cortex: Information Flow-Based Learning in Spiking Neural Networks. *Neural Computation*, 11(4):919–934, 1999.
- [9] T. S. Jaakkola and M. I. Jordan. Improving the Mean Field Approximation via the Use of Mixture Distributions. In M. I. Jordan, editor, *Learning in Graphical Models*. Kluwer Academic Publishers, 1998.
- [10] M.I. Jordan, Z. Ghahramani, T.S. Jaakkola, and L.K. Saul. An Introduction to Variational Methods for Graphical Models. In M.J. Jordan, editor, *Learning in Graphical Models*, chapter 1. MIT Press, 1998.
- [11] R. Linsker. An Application of the Principle of Maximum Information Preservation to Linear Systems. In David Touretzky, editor, *Advances in Neural Information Processing Systems*, volume 1. Morgan-Kaufmann, 1989.
- [12] R. Linsker. How to generate ordered maps by maximizing the mutual information between input and output signals. *Neural Computation*, 1, 1989.
- [13] J.-P. Nadal, N. Brunel, and N. Parga. Nonlinear feedforward networks with stochastic outputs: infomax implies redundancy reduction. *Network: Computation in Neural Systems*, 9(2):207–217, 1998.
- [14] F. Rieke, D. Warland, R. R. de Ruyter van Steveninck, and W. Bialek. *Spikes – Exploring the Neural Code*. MIT Press, Cambridge, MA, 1996.
- [15] D. Saad and M. Opper. *Advanced Mean Field Methods Theory and Practice*. MIT Press, 2001.
- [16] E. Schneidman, W. Bialek, and M.J. Berry II. Synergy, Redundancy, and Independence in Population Codes. *the Journal of Neuroscience*, 23(37):11539–11553, 2003.