

Most Informative Dimension Reduction

Amir Globerson and Naftali Tishby
School of Computer Science and Engineering and
The Interdisciplinary Center for Neural Computation
The Hebrew University, Jerusalem 91904, Israel

Abstract

Finding effective low dimensional features from empirical co-occurrence data is one of the most fundamental problems in machine learning and complex data analysis. One principled approach to this problem is to represent the data in low dimension with minimal loss of the information contained in the original data. In this paper we present a novel information theoretic principle and algorithm for extracting low dimensional representations, or feature-vectors, that capture as much as possible of the mutual information between the variables. Unlike previous work in this direction, here we do not cluster or quantize the variables, but rather extract continuous feature functions directly from the co-occurrence matrix, using a converging iterative projection algorithm. The obtained features serve, in a well defined way, as approximate sufficient statistics that capture the information in a joint sample of the variables. Our approach is both simpler and more general than clustering or mixture models and is applicable to a wide range of problems, from document categorization to bioinformatics and analysis of neural codes.

Introduction

The problem of complex data analysis can be understood as the search for a most compact representation, or model, which *explains* a given set of observations. One explicit view of data *explanation* stems from the notion of sufficient statistics in parameter estimation. Sufficient statistics are functions of samples that capture *all* the information about the parameters of a distribution. Every bit of information that can be extracted from the sample about the parameters is captured by such finite dimensional statistics. One may say that sufficient statistics 'explain', or encode, the *relevant* part of the sample with respect to the parameters. This analogy between sufficient statistics and features in pattern recognition and learning has been noticed before. Here we take this analogy much further and propose a general information theoretic approach for extracting features as approximate sufficient statistics. Exact sufficient statistics in the original statistical sense exist, in general, only for very special distributions - exponential parametric families (Degroot 1986). We first extend this narrow sense of sufficiency

Copyright © 2002, American Association for Artificial Intelligence (www.aaai.org). All rights reserved.

by considering general joint distributions of two variables. We then relax the requirement of *exact* sufficiency by introducing a variational approximation that enables us to calculate functions that approximately preserve the information, in any given dimension. For doing this we also consider the basic problem of "the information in an observation" and suggest a general and concrete answer to this question.

The following example can illustrate our motivation. Consider the familiar dice with unknown distribution of its six possible outcomes. What is the information provided by expected outcome of a dice roll about, say, the probability to obtain 6 in rolling this dice? One possible answer to this question was given by Jaynes using the "maximum entropy principle", which argues that the "most probable" (or "least informative") distribution of outcome is the one which maximizes the entropy, subject to the known observation as a constraint (Jaynes 1957). This notion, however, does not tell us *which observation* is most efficient, or most informative, about the unknown value. The answer to the latter question can be obtained directly by considering the mutual information between the observation - in this case the expected dice outcome - and the value of the probability to get 6 as an outcome. Interestingly, this value can be expressed as the mutual information of a very special joint distribution, for which the specific observations (expected outcome) and corresponding functions of the relevant unknowns (probability to obtain 6) capture all its mutual information. This (unique) distribution is the one of exponential form in those functions. This exponential distribution also solves an important variational problem - it is the joint distribution with minimal mutual information subject to the observations as constraints - given the marginal distributions. This provides us with the same exponential form given by the maximum entropy principle, but the rationale stems directly from the requirement to preserve mutual information. We now develop this new notion of *most informative observations*, or features, and introduce an algorithm for finding them from experimental data.

Consider a joint probability distribution, $p(x, y)$ over the variables X and Y . We say that the d -dimensional vector function $\vec{\psi}(Y) = (\psi_1(Y), \dots, \psi_d(Y))$ is *sufficient* for the variable X if $I(X; Y) = I(X; \vec{\psi}(Y))$, where $I(X; Y) = \sum_{x, y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$ is Shannon's mutual information

between X and Y (Cover & Thomas 1991).¹ In this case the information about X in a *sample* of Y , for a given value $X = x$, is completely captured by the empirical expectation of $\vec{\psi}(Y)$, namely, by $\langle \vec{\psi}(Y) \rangle_{p(Y|x)}$. Similarly, one can say that the functions $\vec{\phi}(X) = (\phi_1(X), \dots, \phi_d(X))$ are sufficient for Y if $I(X; Y) = I(\vec{\phi}(X); Y)$. In that case the dual problem of the information about a given $Y = y$ from a sample of X is captured by $\langle \vec{\phi}(X) \rangle_{p(X|y)}$. The interesting question we address in this work is if we can find such *dual* sets of functions, $\vec{\psi}(Y)$ and $\vec{\phi}(X)$, simultaneously.

A complete simultaneous dimensionality reduction, that preserves *all* the information in the original joint distribution, is possible in general only for special classes of $p(x, y)$, that are of an exponential form. However, these observations motivate a variational approach that finds such a reduction if one exists, but extracts functions $\vec{\psi}(Y)$ and $\vec{\phi}(X)$ that approximate such a reduction in a well defined information theoretic sense. The duality of these two function sets is a key component of our approach. We further show that this variational principle has other suggestive interpretations which make it a viable candidate for a general principled dimensionality reduction or feature extraction technique.

Problem formulation

We assume that we are given a joint distribution of two variables $p(x, y)$. In most practical cases we are in fact given only a finite sample from such a distribution, as a co-occurrence matrix. This empirical joint distribution enables us to estimate expectation values of functions of x and y . We also assume that the marginal distributions, $p(x)$ and $p(y)$, are known or can be well estimated. This latter assumption simplifies our analysis and algorithm, but can be relaxed. We first deal with just one of the two dual problems: finding the set $\vec{\psi}(Y)$. We later show that the resulting algorithm in fact solves the two problems simultaneously, and this apparent asymmetry will be removed.

We look for a set of d functions (or features) of the variable y , denoted by $\psi_i(y)$ for $i = 1 \dots d$, whose expectations can capture the information in the joint distribution in the following information theoretic sense.

The functions $\psi_i(y)$ should satisfy two complementary requirements. On the one hand their expectation should capture all the information about the values of x , thus no other information on y can be assumed. This means that we should *minimize* the information, subject to the known expected values, $\langle \psi_i(y) \rangle_{p(y|x)}$ for every x , as constraints. On the other hand, the functions $\psi_i(y)$ should be *selected* such that these expectations provide the maximum possible information on the variable X .

Given a candidate set of d functions ψ_i , we denote by $\tilde{p}(x, y)$ the distribution with *minimal* information under the expectation constraints. Namely

$$\tilde{p}(x, y) = \arg \min_{q(x, y) \in \mathcal{P}(\psi_i(y), p)} I[q(x, y)] \quad (1)$$

¹In the conventional notation of parameter estimation X is a d dimensional parameter vector Θ , and Y is an n dimensional i.i.d sample from $p(y|\Theta)$.

where $I[q(x, y)] = \sum_{x, y} q(x, y) \log \frac{q(x, y)}{p(x)p(y)}$ and the set $\mathcal{P}(\psi_i(y), p)$ is the set of distributions that satisfy the constraints, defined by

$$\mathcal{P}(\psi_i(y), p) = \left\{ \tilde{p}(x, y) : \begin{array}{l} \langle \psi_i \rangle_{\tilde{p}(y|x)} = \langle \psi_i \rangle_{p(y|x)} \quad \forall x \\ \tilde{p}(x) = p(x) \\ \tilde{p}(y) = p(y) \end{array} \right\}. \quad (2)$$

The second requirement should select the best possible candidate functions ψ_i , by *maximizing* this information over all possible d functions ψ .

Together this can be written as a single Max-Min problem for the optimal functions $\psi_i^*(y)$,

$$\psi_i^*(y) = \arg \max_{\psi_i(y)} \min_{\tilde{p}(x, y) \in \mathcal{P}(\psi_i(y), p)} I(\tilde{p}(x, y)). \quad (3)$$

Notice that this variational principle *does not* define a generative statistical model for the data and is in fact a model independent approach. As we show later, however, the resulting distribution $\tilde{p}(x, y)$ can be interpreted as a generative model in a class of exponential form. But there is no need to make any assumption about the validity of such a model for the data. The data distribution $p(x, y)$ is needed only to estimate the expectations $\langle \psi(y) \rangle$ for every x , given the candidate features.

In what follows we present an iterative projection algorithm, which finds a local solution of this variational problem, and prove its convergence using information geometrical tools. We also show that this asymmetric variational problem for finding the functions $\psi_i(y)$ is in fact equivalent to the dual problem of finding the functions $\phi_i(x)$ and in fact the two problems are solved simultaneously. Our iterative algorithm uses any MaxEnt algorithm, such as “iterative scaling” (Darroch & Ratcliff 1972), as its inner component. Using it we perform alternating projections between two (infinite) sets of linearly constrained distributions.

The problem as formulated provides a tool for data dimensionality reduction, and as such is applicable for a wide range of problems, from natural language processing to neural code analysis and bioinformatics. An illustrative document classification application will be presented. As we discuss further, this variational problem is similar in structure to the fundamental problem of the capacity of an unknown channel which suggests other interesting interpretations for our procedure. It is also related to the recently proposed *Information Bottleneck Method* (Tishby, Pereira, & Bialek 1999) as well as to other dimensionality reduction algorithms, such as LSA (Deerwester *et al.* 1990), LLE (Roweis & Saul 2000) and non-negative matrix factorization (Lee & Seung 1999).

The nature of the solution

We first show that the problem as formulated in Eq. 3 is equivalent to the problem of minimizing the KL divergence between the empirical distribution $p(x, y)$ and a large (parametric) family of distributions of an exponential form.

To simplify notation, we sometimes omit the suffix of (x, y) from the distributions. Thus \tilde{p}_t stands for $\tilde{p}_t(x, y)$ and p for $p(x, y)$

Minimizing the mutual information in Eq. 1 under the linear constraints on the expectations $\langle \vec{\psi}(y) \rangle$ is equivalent to maximizing the joint entropy, $H[\tilde{p}(x, y)] = -\sum_{x, y} \tilde{p}(x, y) \log \tilde{p}(x, y)$, under these constraints, with the additional requirement on the marginals, $\tilde{p}(x) = p(x)$ and $\tilde{p}(y) = p(y)$. Due to the concavity of the entropy and the convexity of the linear constraints, there exists a unique maximum entropy distribution (for compact domains of x and y) which has the exponential form,

$$\tilde{p}_{\vec{\psi}(y)}^*(x, y) = \frac{1}{Z} \exp \left(\sum_{i=1}^d \phi_i(x) \psi_i(y) + A(x) + B(y) \right), \quad (4)$$

where $Z = \sum_{x, y} \exp \left(\sum_{i=1}^d \phi_i(x) \psi_i(y) + A(x) + B(y) \right)$ is the normalization (partition) function².

The functions $\phi_i(x)$ are uniquely determined as Lagrange multipliers from the expectation values $\langle \psi_i(y) \rangle$. This exponential form is also directly linked to our interpretation of the functions ψ_i and ϕ_i as conjugate sufficient statistics.

One key property of the exponential form, Eq. 4, is that it is the only joint distribution for which the mutual information between x and y is completely captured by the dual function vectors, $\psi_i(y), \phi_i(x)$. Namely, it is the only joint distribution $\tilde{p}(x, y)$ for which the following equality holds,

$$I[\tilde{p}(x, y)] = I(\psi; \phi). \quad (5)$$

This property manifests the fact that the reduced description of the x and y by ϕ and ψ indeed captures the relevant structure of the joint distribution. It makes it clear why one would like to find a good approximation to the given joint distribution, which is of this form.

The set of distributions of the exponential form in Eq. 4 can also be considered as a parametric family, parametrized by the infinite family of functions $\Theta = [\psi_i(y), \phi_i(x), A(x), B(y)]$ (note we treat ψ and ϕ symmetrically).

We denote this family by P_Θ and a member of the family by p_Θ . Naturally, $\tilde{p}_{\vec{\psi}(y)}^* \in P_\Theta$.

We define the set of distributions $\mathcal{P}_\Psi \subset P_\Theta$:

$$\mathcal{P}_\Psi = \left\{ \tilde{p} \in P_\Theta : \begin{array}{l} \langle \psi_i \rangle_{\tilde{p}(y|x)} = \langle \psi_i \rangle_{p(y|x)} \quad \forall x \\ \tilde{p}(x) = p(x) \\ \tilde{p}(y) = p(y) \end{array} \right\}. \quad (6)$$

The problem presented above (Eq. 3) is then equivalent to finding the information maximizing distribution in \mathcal{P}_Ψ ,

$$\tilde{p}^* = \arg \max_{\tilde{p} \in \mathcal{P}_\Psi} I[\tilde{p}]. \quad (7)$$

We now return to the information maximization problem, i.e. the selection of the feature candidates $\vec{\psi}(y)$. For every $\tilde{p} \in \mathcal{P}_\Psi$ one can easily show that

$$I[\tilde{p}] = I[p] - D_{KL}[p|\tilde{p}], \quad (8)$$

where $D_{KL}[p|q] = \sum_{x, y} p \log \frac{p}{q}$. The above expression has two important consequences. The first is that maximizing

²Note that the unique distribution can actually be on the closure of such exponential forms. We do not address this detail here.

$I[\tilde{p}]$ for $\tilde{p} \in \mathcal{P}_\Psi$ is equivalent to minimizing $D_{KL}[p|\tilde{p}]$ for $\tilde{p} \in \mathcal{P}_\Psi$:

$$\tilde{p}^* = \arg \min_{\tilde{p} \in \mathcal{P}_\Psi} D_{KL}[p|\tilde{p}]. \quad (9)$$

In addition, Eq. 8 shows that the information in $I[\tilde{p}^*]$ can not be larger than the information in the original data. This supports the intuition that the model \tilde{p}^* maintains only properties present in the original distribution that are captured by the selected features $\psi_i^*(y)$, for any value of x .

The problem in Eq. 9 is a maximization of a function over a subset of P_Θ , namely \mathcal{P}_Ψ . The following proposition shows that this is in fact equivalent to maximizing the same function over all of P_Θ .

Proposition 1 $\arg \min_{\tilde{p} \in \mathcal{P}_\Psi} D_{KL}[p|\tilde{p}] = \arg \min_{\tilde{p} \in P_\Theta} D_{KL}[p|\tilde{p}]$

Proof: We need to show that the distribution which minimizes the right hand side is in \mathcal{P}_Ψ . Indeed, by taking the (generally functional) derivative of $D_{KL}[p(x, y)|\tilde{p}]$ w.r.t. the parameters Θ in \tilde{p} , one obtains the following conditions:

$$\begin{array}{ll} \forall x, i & \langle \psi_i(y) \rangle_{\tilde{p}(y|x)} = \langle \psi_i(y) \rangle_{p(y|x)} \\ \forall y, i & \langle \phi_i(x) \rangle_{\tilde{p}(x|y)} = \langle \phi_i(x) \rangle_{p(x|y)} \\ \forall x & \tilde{p}(x) = p(x) \\ \forall y & \tilde{p}(y) = p(y). \end{array} \quad (10)$$

Clearly, this distribution satisfies the constraints in $\mathcal{P}(\psi_i(y), p)$ and is therefore in \mathcal{P}_Ψ . \square

Our problem is thus equivalent to the minimization problem:

$$p^* = \arg \min_{\tilde{p} \in P_\Theta} D_{KL}[p|\tilde{p}]. \quad (11)$$

Equation 11 is symmetric with respect to ϕ and ψ , thus removing the asymmetry between X and Y in the formulation of Eq. 3.

Notice that this minimization problem can be viewed as a Maximum Likelihood fit to the given $p(x, y)$ in the class P_Θ , but since nothing on its own guarantees the quality of this fit, nor justifies the class P_Θ , we prefer the information theoretic, model independent, interpretation of our approach.

An Iterative Projection Algorithm

Our main result is an iterative algorithm for solving the Max-Min variational problem, Eq. 3, which provably converges to a minimum of Eq. 11. We describe the algorithm using the information geometric notion of *I-projections* (Csiszar 1975).

The *I-projection* of a distribution $q(x)$ on a set of distributions \mathcal{P} is defined as the distribution $p^*(x)$ in \mathcal{P} which minimizes the KL-divergence $D_{KL}[p|q]$.

$$p^*(x) = \arg \min_{p \in \mathcal{P}} D_{KL}[p|q]. \quad (12)$$

We shall use the notation $p^*(x) = IPR(q, \mathcal{P})$.

An important property of the *I-projection* is the, so called, Pythagorean property (Cover & Thomas 1991). For every distribution p in \mathcal{P} the following holds:

$$D_{KL}[p|q] \leq D_{KL}[p|p^*] + D_{KL}[p^*|q]. \quad (13)$$

We now focus on the case where the set \mathcal{P} is a convex set determined by expectation values. Given a set of d functions $f_i(x)$ and a distribution $p(x)$, we denote the set of distributions which agree with $p(x)$ on the expectation values of $f_i(x)$, by $\mathcal{P}(f_i(x), p(x))$. Namely,

$$\mathcal{P}(f_i(x), p(x)) = \left\{ \tilde{p}(x) : \langle f_i(x) \rangle_{\tilde{p}(x)} = \langle f_i(x) \rangle_{p(x)} \right\} \quad (14)$$

The *I-projection* in this case has the exponential form

$$IPR(q(x), \mathcal{P}(f_i(x), p(x))) = \frac{1}{Z^*} q(x) \exp \sum_i \lambda_i^* f_i(x), \quad (15)$$

where λ_i^* is the Lagrange multiplier corresponding to $\langle f_i(x) \rangle_{p(x)}$. In addition, for this special set, the Pythagorean inequality actually becomes an equality (Csiszar 1975).

Before describing the algorithm, we need some additional notations:

- $\tilde{p}_t(x, y)$ - the distribution after t iterations.
- $\psi_{i,t}(y)$ - the $\psi_i(y)$ functions for $\tilde{p}_t(x, y)$, Lagrange multipliers for $\langle \phi_{i,t-1}(x) \rangle$.
- $\phi_{i,t+1}(x)$ - the $\phi_i(x)$ functions for $\tilde{p}_{t+1}(x, y)$, Lagrange multipliers for $\langle \psi_{i,t}(y) \rangle$.
- Θ_t - the full parameter set for $\tilde{p}_t(x, y)$.

The iterative projection algorithm is outlined in figure 1. In this figure the algorithm is described through iterated *I-projections* of (exponential) distributions, once for fixed $\psi_i(y)$ and their expectations, and then for fixed $\phi_i(x)$ and their expectations. Interestingly, during the first projection, the functions $\phi_i(x)$ are modified as Lagrange multipliers for $\langle \psi_i(y) \rangle$, and vice-versa in the second projection. The iteration can thus be viewed as alternating mapping between the two sets of d -dimensional functions, ψ_i and ϕ_i . This is also the direct goal of the variational problem.

We proceed to prove the convergence of the algorithm. We first show that every step reduces $D_{KL}[p(x, y)|\tilde{p}_t(x, y)]$.

$$\text{Proposition 2 } D_{KL}(p(x, y), \tilde{p}_{t+1}(x, y)) \leq D_{KL}(p(x, y), \tilde{p}_t(x, y)).$$

Proof: For each x , the following is true:

- $\tilde{p}_{t+1}(y|x)$ is the *I-projection* of $\tilde{p}_t(y|x)$ on the set $\mathcal{P}(\psi_{i,t}(y), p(y|x))$.
- $p(y|x)$ is also in $\mathcal{P}(\psi_{i,t}(y), p(y|x))$.

Using the Pythagorean property, (which is an equality here) we have that $D_{KL}[p(y|x)|\tilde{p}_t(y|x)]$ is equal to:

$$D_{KL}[p(y|x)|\tilde{p}_{t+1}(y|x)] + D_{KL}[\tilde{p}_{t+1}(y|x)|\tilde{p}_t(y|x)]. \quad (16)$$

Multiplying by $p(x)$ and summing over all x values, we obtain:

$$D_{KL}[p|\tilde{p}_t] = D_{KL}[p|\tilde{p}_{t+1}] + D_{KL}[\tilde{p}_{t+1}|\tilde{p}_t]. \quad (17)$$

where we have used $\tilde{p}_t(x) = p(x)$. Note that the summation resulted in the term $D_{KL}[p(x)|\tilde{p}_{t+1}(x)]$ on both sides of the equation.

Using the non-negativity of $D_{KL}[\tilde{p}_{t+1}|\tilde{p}_t]$ we have:

$$D_{KL}[p|\tilde{p}_t] \geq D_{KL}[p|\tilde{p}_{t+1}]. \quad (18)$$

Note that equality is obtained iff $D_{KL}[\tilde{p}_{t+1}|\tilde{p}_t] = 0$. \square

Input: Joint (empirical) distribution $p(x, y)$

Output: $2d$ feature functions: $\psi_i(y)$ $\phi_i(x)$ that result from \tilde{p}^* , a solution of the variational problem Eq. 3 and a (local) minimum of Eq. 9.

Initialization:

- Initialize $\tilde{p}_0(x, y) \in P_\Theta$ randomly

Iterate:

- For all x , set:

$$\begin{aligned} \tilde{p}_{t+1}(y|x) &= IPR(\tilde{p}_t(y|x), \mathcal{P}(\psi_{i,t}(y), p(y|x))) \\ \tilde{p}_{t+1}(x, y) &= \tilde{p}_{t+1}(y|x)p(x). \end{aligned}$$

The functions $\phi_{i,t+1}(x)$ are determined as the Lagrange multipliers

- For all y , set:

$$\begin{aligned} \tilde{p}_{t+2}(x|y) &= IPR(\tilde{p}_{t+1}(x|y), \mathcal{P}(\phi_{i,t+1}(x), p(x|y))) \\ \tilde{p}_{t+2}(x, y) &= \tilde{p}_{t+2}(x|y)p(y). \end{aligned}$$

The functions $\psi_{i,t+1}(y)$ are determined as Lagrange multipliers

- Halt on convergence (when small enough change in $\tilde{p}_t(x, y)$).

Figure 1: The iterative projection algorithm.

An analogous argument proves that:

$$D_{KL}[p|\tilde{p}_{t+2}] \leq D_{KL}[p|\tilde{p}_{t+1}]. \quad (19)$$

The following easily provable proposition states that the stationary points of the algorithm coincide with extremum points of the target function $D_{KL}[p|p_\Theta]$. Its proof uses the properties of the projections in the algorithm, and the characterization of the extremum point in Eq. 10.

Proposition 3 If $\tilde{p}_t = \tilde{p}_{t+2}$ then the corresponding Θ_t satisfies $\frac{\partial}{\partial \Theta} D_{KL}[p|p_\Theta] = 0$.

It is now easy to prove convergence of the algorithm to a local minimum of $D_{KL}[p|p_\Theta]$ using continuity arguments and the improvement at each iteration given in Eq. 17. We will give a detailed proof in a separate paper.

Implementation - Partial I-projections

The description of the iterative algorithm assumes the existence of a module which calculates *I-projections* on linear constraints. Because no closed form solution for such a projection is known, it is found by successive iterations which asymptotically converge to the solution. It is straightforward to show that even if our projection algorithm uses such a "Partial I projection" algorithm as its IPR module, it still converges to a minimum.

In this work we use as an *IPR* algorithm the Generalized Iterative Scaling (GIS) procedure (Darroch & Ratcliff 1972), described in figure 2. Alternative algorithms such as

Improved Iterative Scaling (Pietra & Lafferty 1997) or conjugate gradient methods, can also be used and may improve convergence rate.

Input: Distributions $q(x), p(x)$, d functions $f_i(x)$
Output: $IPR(q(x), \mathcal{P}(f_i(x), p(x)))$
Initialize: $p_0(x) = q(x)$
Iterate:
<ul style="list-style-type: none"> • $p_{t+1}(x) = \frac{1}{Z_t} p_t(x) \exp \sum_{i=1}^d f_i(x) \log \frac{\langle f_i(x) \rangle_{p(x)}}{\langle f_i(x) \rangle_{p_t(x)}}$ • Z_t is a normalization constant

Figure 2: The Generalized Iterative Scaling Algorithm.

Discussion

Our proposed method is a new dimensionality reduction technique. It is nonlinear, unlike PCA or ICA, and it aims directly at preserving mutual information in a given empirical co-occurrence matrix. We achieved that through an information variation principle that enables us to calculate simultaneously informative feature functions for *both* random variables. In addition we obtain an exponential model approximation to the given data which has precisely these features as *dual sets* of sufficient statistics. We described an alternating projection algorithm for finding these features and proved its convergence to a (local) optimum. This is in fact an algorithm for extracting optimal sets of constraints from statistical data. We briefly address now several other important issues that our procedure raises.

Finite samples The basic assumption behind our problem formulation is that we have the joint distribution of the variables x and y , $p(x, y)$. For the machine learning community this assumption may look strange, but one should remember that the goal is *not* to learn the mapping from X to Y , but rather to extract good features of X w.r.t. Y . In fact, $p(x, y)$ is not needed explicitly, but only to estimate the expectation values $\langle \vec{\psi}(y) \rangle$ and $\langle \vec{\phi}(x) \rangle$. These can be estimated *uniformly well* from a finite sample under the standard uniform convergence conditions. In other words, standard learning theoretical techniques can give us the sample complexity bounds, given the dimensions of X and Y and the reduced dimension d . For continuous x and y further assumptions must be made, such as the VC dimension of the features, etc.

Uniqueness of the solution Interestingly, while the goal of the algorithm are features that preserve information, the information in the functions $\vec{\psi}(y)$ and $\vec{\phi}(x)$ is not estimated directly at any point. Furthermore, there is a freedom in selecting the features that stems from the fact that only the dot-product $\vec{\phi}(x) \cdot \vec{\psi}(y)$ appears in the distribution $\tilde{p}(x, y)$. Any invertible matrix R can be applied such that $\vec{\phi}(x)R^{-1}$ and $R\vec{\psi}(y)$ are equally good features. One can remove this ambiguity by orthogonalization and scaling of the feature

functions, for example by applying SVD to $\log \tilde{p}(x, y)$. Notice, however, that our procedure is very different from direct application of SVD to $\log p(x, y)$. These two coincide *only* when the original joint distribution is already of the exponential form of Eq.(4). In all other cases SVD based approximations (LSA included) will not preserve information as well as our features at the same dimension reduction. The resulting functions $\vec{\psi}(y)$ and $\vec{\phi}(x)$ thus depend on the initial point of the iterations, but the information extracted does not (for the same optimum).

Information theoretic interpretation Our information MaxMin principle is close in its formal structure to the problem of channel capacity with some channel uncertainty (see e.g. (Lapidoth & Narayan 1998)). This suggests the interesting interpretation for the features as channel characteristics. If the channel only enables the reliable transmission of d expected values, then our $\vec{\psi}(y)$ exploit this channel in an optimal way. The channel decoder of this case is provided by the dual vector $\vec{\phi}(x)$ and the decoding is performed through a dot-product in of these two vectors. This intriguing interpretation of our algorithm obviously requires further analysis.

Relations to other algorithms

Dimension reduction algorithms have become a fundamental component in unsupervised large scale data analysis, together with clustering. While linear methods, as PCA and ICA, provide very useful first steps, they do not provide a principled method for statistical co-occurrence data, where such linear assumptions about the matrix are unjustified. Several interesting non-linear methods have been proposed in the past years. Of particular interest are

LLE (Roweis & Saul 2000) and non-negative matrix factorization (Lee & Seung 1999). We feel that none of these new algorithms directly address the question of information preserving as we suggest here. Furthermore, preliminary comparisons of our algorithm with LSA (Deerwester *et al.* 1990) for document indexing are very encouraging and justify our non-linear approach. A closely related idea is the *Information Bottleneck Method* (Tishby, Pereira, & Bialek 1999) which aims at a clustering that preserves information. However, clustering may not be the correct answer for many problems where the relationship between the variables comes from some hidden low dimensional continuous structures. In such cases clustering tends to quantize the data in a rather arbitrary way, while low dimensional features are simpler and easier for interpretation. The resulting algorithm is also computationally simpler, with no need for complicated splitting or merging of clusters.

Acknowledgement

We thank Noam Slonim and Gal Chechik for helpful discussions and comments and for the help with experimental data. This work is partly supported by a grant from the Israeli Academy of Science. A.G. is supported by the Eshkol Foundation.

References

- Cover, T., and Thomas, J. 1991. *Elements of information theory*. Wiley.
- Csiszar, I. 1975. I-divergence geometry of probability distributions and minimization problems. *Annals of Probability* 3(1):146–158.
- Darroch, J., and Ratcliff, D. 1972. Generalized iterative scaling for log-linear models. *Ann. Math. Statist.* 43:1470–1480.
- Deerwester, S.; Dumais, S.; Landauer, T.; Furnas, G.; and Harshman, R. 1990. Indexing by latent semantic analysis. *Journal of the American Society of Information Science* 41(6):391–407.
- Degroot, M. 1986. *Probability and Statistics*. Addison-Wesley.
- Jaynes, E. 1957. Information theory and statistical mechanics. *Physical Review* 106:620.
- Lapidoth, A., and Narayan, P. 1998. Reliable communication under channel uncertainty. *IEEE Transactions on Information Theory* 44(6):2148–2177.
- Lee, D., and Seung, H. 1999. Learning the parts of objects by non-negative matrix factorization. *Nature* 401(6755):788–791.
- Pietra, S. D., and Lafferty, V. D. P. J. 1997. Inducing features of random fields. *IEEE Transactions on PAMI* 19(4):380–393.
- Roweis, S., and Saul, L. 2000. Nonlinear dimensionality reduction by locally linear embedding. *Science* 290(5500):2323–6.
- Tishby, N.; Pereira, F.; and Bialek, W. 1999. The information bottleneck method. In *Proc. of the 37-th Annual Allerton Conference on Communication, Control and Computing*, 368–377.