

---

# Supervised dimensionality reduction using mixture models

---

Sajama

Alon Orlitsky

University of California at San Diego

SAJAMA@UCSD.EDU

ALON@UCSD.EDU

## Abstract

Given a classification problem, our goal is to find a low-dimensional linear transformation of the feature vectors which retains information needed to predict the class labels. We present a method based on maximum conditional likelihood estimation of mixture models. Use of mixture models allows us to approximate the distributions to any desired accuracy while use of conditional likelihood as the contrast function ensures that the selected subspace retains maximum possible mutual information between feature vectors and class labels. Classification experiments using Gaussian mixture components show that this method compares favorably to related dimension reduction techniques. Other distributions belonging to the exponential family can be used to reduce dimensions when data is of a special type, for example binary or integer valued data. We provide an EM-like algorithm for model estimation and present visualization experiments using Gaussian and Bernoulli mixture models.

## 1. Introduction

Dimensionality reduction is a frequently used pre-processing step for supervised learning tasks. Reducing dimensions may improve classifier performance since it can suppress noise in the data and act as a form of regularization. Also, meaningful low dimensional representation can help in visualizing data sets and is an important tool in exploratory data analysis.

In this paper, we consider the problem of finding discriminative linear feature transformations. Given a collection of  $d$ -dimensional training samples and their

class labels, the goal is to find an  $L$ -dimensional hyperplane in  $\mathbb{R}^d$  such that the projected samples belonging to various classes are well separated. Our approach to this problem, termed supervised dimensionality reduction using mixture models (SDR-MM), is to model each class using a mixture model. The parameters of the model include affine parameters for a subspace to which the mixture means are constrained. Gaussian mixtures can approximate arbitrarily complex densities by lowering the minimum allowed variance and increasing the number of mixture components. Hence, this approach is *semi-parametric* - the subspace is determined by a set of affine parameters, while the distributions on the projected space are approximated non-parametrically. We use maximum conditional likelihood (MCL) estimation to determine the parameters of the lower dimensional subspace which ensures that the predictive information in the feature vectors is retained in the projected space. MCL has been widely used as a discriminative objective function for estimating hidden markov models in speech recognition and for Gaussian mixture models in the context of classification in (Jebara & Pentland, 1998).

Some dimension reduction methods make restrictive parametric assumptions about the distributions. For example, Fisher's linear discriminant analysis (LDA) can be obtained by maximum likelihood estimation assuming that the classes are Normally distributed with a common covariance matrix and different means, with the means constrained to lie in an  $L$  dimensional subspace. Other parametric methods include projection pursuit regression (Friedman & Stuetzle, 1981) and Generalized additive models (Hastie & Tibshirani, 1986). More recently, several semi-parametric methods have been proposed for supervised dimensionality reduction including sliced inverse regression (Li, 1991) and principal Hessian directions (pHd) (Li, 1992). Sufficient dimensionality reduction (Globerston & Tishby, 2003) is designed for the unsupervised case and uses maximum entropy principle for estimating the exponential models involved.

---

Appearing in *Proceedings of the 22<sup>nd</sup> International Conference on Machine Learning*, Bonn, Germany, 2005. Copyright 2005 by the author(s)/owner(s).

In terms of the density model used, the method most closely related to SDR-MM is Mixture discriminant analysis (MDA) (Hastie & Tibshirani, 1996) which generalizes LDA by approximating each of the classes by a mixture of Gaussians all of which have a common covariance matrix. SDR-MM differs from MDA in two important ways. Firstly, in SDR-MM, we use spherical Gaussian distributions while in MDA each Gaussian has the same full-covariance matrix. While this may mean that SDR-MM needs to use more mixture components for each class, the total number of parameters to be estimated is often reduced from not having to estimate the  $d^2$  parameters of the covariance matrix. Secondly, in MDA, parameters are estimated using maximum likelihood, while in SDR-MM, the parameters are estimated discriminatively by maximizing the conditional likelihood which also eliminates the need for subclass shrinkage used in MDA.

The other dimensionality reduction method closely related to SDR-MM is kernel dimensionality reduction (KDR) (Fukumizu et al., 2004) which also chooses the lower dimensional subspace based on maximum mutual information principle. SDR-MM differs from KDR in the way in which it measures the mutual information. While SDR-MM uses conditional likelihood, the KDR objective function is based on cross-covariance operators on reproducing kernel Hilbert spaces. A related method was proposed in (Torkkola & Campbell, 2000) in which instead of using the Shannon mutual information, a Renyi-entropy based expression for mutual information is estimated.

Recently, several methods have been proposed for probabilistic formulation of principal component analysis and its extension using the exponential family of distributions (see for e.g., (Sajama & Orlitsky, 2004) and the references therein). In SDR-MM also, we allow the mixture components to be drawn from the exponential family in order to allow the method to be suitable for the various data types. SDR-MM is an adaptation of the unsupervised method - semi-parametric principal component analysis (SP-PCA) (Sajama & Orlitsky, 2004) to the supervised scenario. We describe a simple and efficient EM-like algorithm for model estimation which uses iteratively re-weighted least squares in the maximization step. We present classification experiments which show that SDR-MM compares favorably to three related methods - pHd, MDA and KDR. We also show visualization examples for real-valued and binary data.

## 2. Model with Gaussian components

We are concerned with multi-class supervised problems where the feature vectors  $\mathbf{x}$  lie in  $\mathbb{R}^d$  and the class

labels  $y$  are drawn from the set  $\{1, \dots, M\}$ . We are given training data  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ , which are independent and identically distributed samples, drawn from a probability distribution  $P(y)P(\mathbf{x}|y)$ . Each class  $m$  is modelled by a mixture of  $c_m$  number of Gaussians  $\mathcal{N}(\mathbf{x}|\boldsymbol{\theta}, \sigma^2\mathbf{I})$  ( $\sigma$  common to all classes). Let  $c = \sum_{m=1}^M c_m$  be total number of mixture components over all classes,  $\Pi = \{\pi_1, \dots, \pi_c\}$  be the prior over these components and for each  $k \in \{1, \dots, c\}$ , let  $\psi_k(m)$  be given by

$$\psi_k(m) = \begin{cases} 1 & \text{if mixture component } k \in \text{class } m \\ 0 & \text{otherwise} \end{cases}$$

Let  $D(\mathbf{x}, \mathbf{w})$  denote the squared Euclidean distance between  $\mathbf{x}$  and  $\mathbf{w}$ . The distribution is given by

$$P(Y = m) = \sum_{k=1}^c \psi_k(m)\pi_k$$

$$P(\mathbf{x}, Y = m) = \sum_{k=1}^c \pi_k \psi_k(m) (2\pi)^{-d/2} e^{-D(\mathbf{x}, \boldsymbol{\theta}_k)/2\sigma^2}.$$

In order to obtain low dimensional representation and measure discriminative capability of feature transformations, we consider the *constrained* Gaussian mixture model. The means of Gaussians from all classes are restricted to lie in a lower ( $L$ ) dimensional hyperplane in  $\mathbb{R}^d$ . We represent this constraint on mixture parameters using  $L \times d$  rotation matrix  $V$  and  $d$ -dimensional displacement vector  $b$ . Each mean  $\boldsymbol{\theta}_k$  belonging to this hyperplane can be represented by the  $L$  dimensional vector  $\mathbf{a}_k$

$$\boldsymbol{\theta}_k = \mathbf{a}_k V + \mathbf{b}.$$

We use the matrix  $A$ , whose  $k$ 'th row is  $\mathbf{a}_k$ , to represent the mixture component parameters. Hence the SDR-MM model is parameterized by  $\Theta = \{\Pi, \psi, A, V, b\}$ . The assumption that the mixture components are spherical Gaussians with common variance ensures that we measure the discriminative capabilities of *linear* projection, since the direction perpendicular to the plane  $(V, b)$  is irrelevant in any metric involving relative values of likelihoods  $P(\mathbf{x}|\boldsymbol{\theta}_k)$ . To see why this is the case, consider  $\mathbf{x}_p$ , the point on the hyperplane  $(V, b)$  closest to  $\mathbf{x}$ . Now,  $P(\mathbf{x}|\boldsymbol{\theta}_k) \propto \exp(-\{D(\mathbf{x}, \mathbf{x}_p) + D(\mathbf{x}_p, \boldsymbol{\theta}_k)\}/2\sigma^2)$  and for a fixed  $\mathbf{x}$ , the factor involving  $D(\mathbf{x}, \mathbf{x}_p)$  is common to all  $\boldsymbol{\theta}_k$ 's on the hyperplane  $(V, b)$  and hence cancels out.

Like LDA and MDA, there is an inherent classifier associated with the SDR-MM model trained for reducing dimensions. Since each class is modelled by a mixture, the distribution  $P(y = m|\mathbf{x})$  can be obtained using Bayes rule and used to label any given test vector  $\mathbf{x}$ .

**Use of spherical Gaussians** We have already noted that use of fixed-variance spherical Gaussians

corresponds to measuring discriminative capability of a *linear* subspace when training samples are projected onto it. That sphericity is not a restrictive assumption follows from the universal approximation property of RBF networks with spherical gaussian kernels (Park & Sandberg, 1991). The idea is that spread of a given class along the subspace  $(V, b)$  can be approximated by spread of Gaussian means belonging to that class, assuming that a small enough variance is chosen. Use of full covariance matrices makes it necessary to regularize model estimation by penalizing the objective function. The assumption that all Gaussians have common spherical covariance reduces the number of parameters to be estimated by  $\mathcal{O}(d^2)$  and thereby improves model generalization. Experimental results in section 7 support these intuitive arguments.

The SDR-MM method is a soft equivalent of prototype methods like LVQ and its probabilistic nature allows data to simultaneously influence multiple prototypes - attracting prototypes of the same class and repelling prototypes belonging to a different class during MCL estimation - thereby generating a large-margin like effect. This provides a simple alternative to *subclass shrinkage* used in MDA (Hastie & Tibshirani, 1996). There is a tradeoff between regularization and approximation capability - smaller variance is better for approximation and larger variance for the regularization effect described above.

### 3. The objective function

We propose using conditional likelihood of the training data as the objective function for choosing appropriate feature transformations, i.e., we pick the lower dimensional space specified by  $(V, b)$  using MCL estimation.

$$(V_{opt}, b_{opt}) = \arg \max_{(V, b)} \max_{A, \Pi} \prod_{i=1}^n P(y_i | \mathbf{x}_i, \Theta). \quad (1)$$

Use of this objective function can be motivated in several ways. In a classification problem, we are interested in finding a model which approximates the observed empirical conditional distribution  $P_{emp}(y|\mathbf{x})$ . Maximizing conditional likelihood is equivalent to minimizing the KL divergence between  $P_{emp}(y|\mathbf{x})$  and the model  $P_{(V, b)}(y|\mathbf{x})$ . Also, on a related note, MCL estimation is equivalent to maximum mutual information estimation (Jebara & Pentland, 1998; Klautau et al., 2003). Hence, this objective function is equivalent to picking transformations that preserve maximum amount of the relevant information (under the SDR-MM model) between distributions of  $\mathbf{x}$  and  $y$ .

We present simple examples of projecting two-dimensional samples onto a line to illustrate how MCL

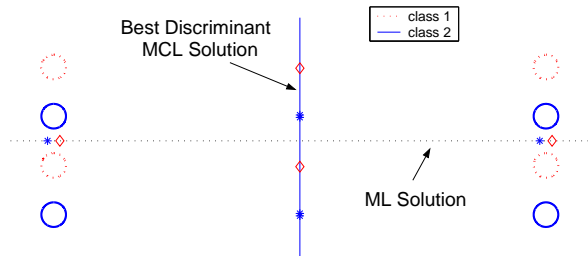


Figure 1. Advantage of MCL : Each class is a mixture of spherical Gaussians.  $\diamond$  and  $*$  denote means of gaussian components of classes 1 and 2 respectively. In this case the subspace MDA finds is the same as the ML solution.

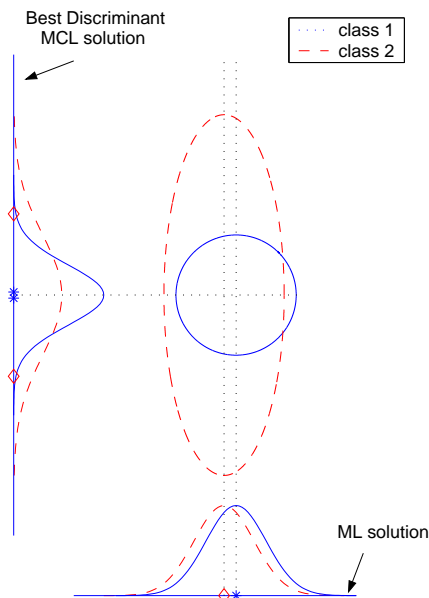


Figure 2. Advantage of MCL : Two classes with different covariance matrices.  $\diamond$  and  $*$  denote means of gaussian components of classes 1 and 2 respectively. In this case the subspace MDA finds is the same as the MCL solution.

estimation extends the applicability of previously studied methods that are also, like SDR-MM, based on constrained mixture of Gaussians. Figure 1 shows a two class example where each class is a mixture of four spherical Gaussians. Projection using low-rank ML estimation fully merges samples from the two classes while MCL estimated mixture model is able to find the best discriminant (see also (Jebara & Pentland, 1998)). Figure 2 shows an interesting example where each of the two classes are generated by a single Gaussian with almost the same mean, but they have very different variance in one direction. If we used ML estimation with *no* constraints on the covariance matrices to find a one-dimensional subspace, we would get the ML solution subspace shown in figure 2, *even*

if each class is allowed to be modelled by a mixture of several Gaussians. This is because no model can be better than the ‘true distribution’ in terms of likelihood of observed data (when data sample is large enough). However, since MDA imposes common covariance constraints on all mixture components of all classes, the MDA solution with three gaussian components for each class, coincides with the MCL solution in this case.

Simulation studies (Klautau et al., 2003) have found that MCL classifiers can compete with and sometimes outperform other discriminative and generative classifiers. For fixed  $(V, b)$ , picking the Gaussian means which maximize conditional likelihood is equivalent to estimating a discriminative mixture classifier based on data projected onto the subspace given by  $(V, b)$  (see also section 2). Hence optimizing the function (1) is equivalent to picking the best subspace for a discriminative Gaussian mixture classifier.

#### 4. Exponential family components

Using Gaussian means and constraining them to a lower dimensional subspace of data space is equivalent to using a ‘soft’ prototype method where the prototypes are real valued and  $D(\mathbf{x}, \boldsymbol{\theta})$ , the distance between a point  $\mathbf{x}$  and prototype  $\boldsymbol{\theta}$ , is Euclidean. This Gaussian model may not be appropriate for other data types, for instance binary or integer data. The Bernoulli distribution may be better for binary data and Poisson for integer data. These three distributions, along with several others, belong to a family of distributions known as the *exponential family* (McCullagh & Nelder, 1983) and can be written in the form

$$\log P(x|\theta) = \log P_0(x) + x\theta - G(\theta).$$

Here,  $\theta$  is called the *natural parameter* and  $G(\theta)$  is a function that ensures that the probabilities sum to one. Studies in the area of unsupervised dimensionality reduction of special data types, have found that use of exponential family models yields better low dimensional representations (e.g., (Sajama & Orlitsky, 2004) and the references therein). Hence we extend the model described in section 2 by using multivariate exponential family distributions for mixture components in the place of fixed variance Gaussians,

$$\log P(\mathbf{x}|\boldsymbol{\theta}) = \sum_{j=1}^d \{\log P_{0j}(x_j) + x_j\theta_j - G_j(\theta_j)\}, \quad (2)$$

where  $x_j$  and  $\theta_j$  are the  $j$ ’th components of  $\mathbf{x}$  and  $\boldsymbol{\theta}$ . Note that by using different distributions for different components of the feature vector  $\mathbf{x}$ , we can model mixed data types.

#### 5. Low dimensional representation

We discuss two of the several ways in which low dimensional representations can be obtained using the model  $\Theta$ . The first method is to represent  $\mathbf{x}$  by that point  $\boldsymbol{\theta}$  on  $(V, b)$  that is closest according to the appropriate Bregman (exponential family-based) distance. It can be shown that there is a unique such  $\boldsymbol{\theta}_{opt}$  on the plane. This representation is a generalization of the standard Euclidean projection. The second method of low dimensional representation is based on Bayes rule. Each feature vector  $\mathbf{x}$  induces a posterior distribution over the latent domain  $P(\boldsymbol{\theta}_i|\mathbf{x}) = \pi_i P(\mathbf{x}|\boldsymbol{\theta}_i)/P(\mathbf{x})$ . Under the SDR-MM model, all the information in  $\mathbf{x}$  about  $y$  is contained in this posterior distribution since  $y$  and  $\mathbf{x}$  are independent when conditioned upon the latent variable  $\boldsymbol{\theta}$ . Hence  $\mathbf{x}$  can be represented by a suitable function of this posterior and we choose to use the mean. This representation has been used successfully by several probabilistic methods in the unsupervised case, to get meaningful low dimensional views.

#### 6. Algorithm

Several iterative algorithms have been proposed for MCL estimation of mixture models, see for example (Jebara & Pentland, 1998; Klautau et al., 2003). The common thread in these algorithms is that each iteration involves evaluating a tight lower bound which touches the objective function at the current parameter value. Model parameters are then updated by maximizing this lower bound. This technique was called bound maximization in (Jebara & Pentland, 1998) and is the basis of many iterative algorithms including the expectation maximization (EM) algorithm.

We use the idea of bound maximization and derive an algorithm for MCL estimation under low rank constraint on mixture component parameters  $\Theta$ . Let  $\Theta^t$  and  $\Theta^{t+1}$  denote the current and updated parameter values at iteration  $t$ . The change in conditional log-likelihood at iteration  $t$  can be written as

$$\begin{aligned} \Delta l &= \sum_{i=1}^n \{\log P(y_i|\mathbf{x}_i, \Theta^{t+1}) - \log P(y_i|\mathbf{x}_i, \Theta^t)\} \\ &\geq \sum_{i=1}^n \sum_{k=1}^c \hat{z}_{ik} \log P(\boldsymbol{\theta}_k, \mathbf{x}_i, y_i|\Theta^{t+1}) \\ &\quad - \sum_{i=1}^n \rho_i P(\mathbf{x}_i|\Theta^{t+1}) + \text{constant}, \end{aligned}$$

$$\text{where } \hat{z}_{ik} = \frac{P(\boldsymbol{\theta}_k, \mathbf{x}_i, y_i|\Theta^t)}{\sum_{k'=1}^c P(\boldsymbol{\theta}_{k'}, \mathbf{x}_i, y_i|\Theta^t)} \text{ \& } \rho_i = \frac{1}{P(\mathbf{x}_i|\Theta^t)}.$$

Here the first term was lower bounded using Jensen’s

inequality (similar to the EM algorithm) and the second term using  $\log w \leq w - 1$ . At each iteration, we compute the lower bound by computing  $\hat{z}_{ik}$  and  $\rho_i$  for  $i = 1, \dots, n$  and  $k = 1, \dots, c$ . The lower bound is then optimized by alternately maximizing over each of  $\Pi$ ,  $A$ ,  $V$  and  $b$  while holding the rest of the parameters constant.

The lower bound can be written as (ignoring constants since they do not affect the optimization steps)

$$\Delta l = \sum_i \sum_k \hat{z}_{ik} \log \pi_k + \sum_i \sum_k \hat{z}_{ik} \log \psi_k(y_i) \quad (3)$$

$$+ \sum_i \sum_k \hat{z}_{ik} \log P(\mathbf{x}_i | \boldsymbol{\theta}_k) - \sum_i \sum_k \rho_i \pi_k P(\mathbf{x}_i | \boldsymbol{\theta}_k).$$

UPDATING  $\Pi$  :  $\Pi^{t+1}$  is obtained by maximizing the Lagrangian (formed using terms in  $\Delta l$  involving  $\pi_k$ )

$$L = \sum_{k=1}^c \{c_{1k} \log \pi_k - c_{2k} \pi_k\} + \lambda \left( \sum_{k=1}^c \pi_k - 1 \right),$$

where  $c_{1k} = \sum_{i=1}^n \hat{z}_{ik}$ ,  $c_{2k} = \sum_{i=1}^n \rho_i P(\mathbf{x}_i | \boldsymbol{\theta}_k)$  and  $\lambda$  is a lagrange multiplier used to impose the constraint that the latent distribution sums to one. This optimization is a little more complicated than its counterpart in the EM algorithm for ML estimation since we have both linear and logarithmic terms. Differentiating  $L$  and setting the derivative to zero, we get  $\pi_k = c_{1k} / (c_{2k} - \lambda)$ . We need to find  $\lambda$  that satisfies  $f(\lambda) = \sum_{k=1}^c c_{1k} / (c_{2k} - \lambda) = 1$ . There is no explicit solution for this equation, but it is easy to verify that at  $\lambda_0 = \min_k (c_{2k} - c_{1k})$ ,  $f(\lambda_0) > 1$  and that as  $\lambda \rightarrow -\infty$ ,  $f(\lambda) \rightarrow 0$ . Moreover,  $f(\lambda)$  is continuous and monotone in the region  $[-\infty, \lambda_0]$  implying that there is a unique  $\lambda_{opt}$  such that  $f(\lambda_{opt}) = 1$ , which can be found using bisection line search.

OPTIMIZING  $A$ ,  $V$  AND  $b$  : For optimizing  $A$  and  $V$ , we use an iterative weighted least squares method similar to that used in fitting generalized linear models (McCullagh & Nelder, 1983), i.e., we apply the Newton-Raphson procedure to the equations obtained by setting the derivative of  $\Delta l$  to zero. Upon taking the first and second derivatives with respect to the components of the matrix  $A$ , it turns out that each row can be updated independently of the others in a given iteration. This decoupling is convenient since it means that updating the parameters involves smaller matrix operations. Similarly, we find that each column of  $V$  and each component of  $b$  can be updated independently. Update equations for  $A$  and  $V$  are given here, and can be derived similarly for  $b$  (not included here because of space constraints).  $\Delta l$  depends on  $A$ ,  $V$  and  $b$  only through the last two terms in equation

3. Hence, ignoring constants, we want to maximize

$$\sum_{k=1}^c \sum_{j=1}^d (\boldsymbol{\theta}_{kj} \tilde{x}_{kj} - G(\boldsymbol{\theta}_{kj}) \tilde{z}_k) - \sum_{i=1}^n \sum_{k=1}^c \rho_i \pi_k P(\mathbf{x}_i | \boldsymbol{\theta}_k), \quad (4)$$

where,  $\tilde{x}_{kj} = \sum_{i=1}^n \hat{z}_{ik} x_{ij}$  and  $\tilde{z}_k = \sum_{i=1}^n \hat{z}_{ik}$  and  $P(\mathbf{x}_i | \boldsymbol{\theta}_k)$  is as defined before in equation 2.

Each row of  $A$ ,  $\mathbf{a}_r$  is updated by adding  $\delta \mathbf{a}_r$  which is calculated using  $(V \Omega_r V^t) \delta \mathbf{a}_r = GR_r$ , where the  $d \times d$  matrix  $\Omega_r$  and the  $L \times 1$  matrix  $GR_r$  are given by

$$[\Omega_r]_{jj'} = \left\{ \tilde{z}_r - \sum_{i=1}^n \rho_i \pi_r P(\mathbf{x}_i | \boldsymbol{\theta}_r) \right\} \frac{\partial g(\theta_{rj})}{\partial \theta_{rj}} \delta(j = j')$$

$$+ \sum_{i=1}^n \rho_i \pi_r P(\mathbf{x}_i | \boldsymbol{\theta}_r) (x_{ij'} - g(\theta_{rj'})) (x_{ij} - g(\theta_{rj}))$$

and

$$[GR_r]_s = \sum_{j=1}^d v_{sj} \tilde{x}_{rj} - \tilde{z}_r g(\theta_{rj})$$

$$- \sum_{i=1}^n \rho_i \pi_r P(\mathbf{x}_i | \boldsymbol{\theta}_r) (x_{ij} - g(\theta_{rj})).$$

Each column of the matrix  $V$ ,  $\mathbf{v}_s$  is updated by adding  $\delta \mathbf{v}_s$  obtained by solving  $(A^t \Omega_s A) \delta \mathbf{v}_s = GR_s$ , where the  $c \times c$  diagonal matrix  $\Omega_s$ , and the  $L \times 1$  matrix  $GR_s$  are given by,

$$[\Omega_s]_{kk} = \left\{ \tilde{z}_k - \sum_{i=1}^n \rho_i \pi_k P(\mathbf{x}_i | \boldsymbol{\theta}_k) \right\} \frac{\partial g(\theta_{ks})}{\partial \theta_{ks}}$$

$$+ \sum_{i=1}^n \rho_i \pi_k P(\mathbf{x}_i | \boldsymbol{\theta}_k) (x_{is} - g(\theta_{ks}))^2$$

and

$$[GR_s]_r = \sum_{k'=1}^c a_{k'r} \{ \tilde{x}_{k's} - \tilde{z}_{k'} g(\theta_{k's})$$

$$+ \sum_{i=1}^n \rho_i \pi_{k'} P(\mathbf{x}_i | \boldsymbol{\theta}_{k'}) (x_{is} - g(\theta_{k's})) \}$$

Note that using the Newton-Raphson method does not guarantee monotone increase in the value of  $\tilde{L}$ . Monotonicity can be enforced using standard optimization procedures like line search or the trust regions method.

COMPUTATIONAL COMPLEXITY : Time taken for each iteration of this algorithm is  $\mathcal{O}(cdnL^2)$ . Computing  $\hat{z}_{ik}$  and  $\rho_i$  involve computing  $P(\mathbf{x}_i | \boldsymbol{\theta}_k)$  which is expensive and is a common problem faced in maximum likelihood

estimation and in training of RBF networks. (Omhundro, 1987) gives a procedure for speeding up this procedure using the k-d tree data structure by identifying relevant prototypes (for each  $\mathbf{x}$ ) thereby avoiding unnecessary computation.

## 7. Experiments

We experimented with the Gaussian mixture model on four real-valued datasets and with the Bernoulli mixture model on a binary set. As noted in section 2, for the Gaussian mixture model, an appropriate variance should be chosen to achieve the right tradeoff between regularization and approximation capability. Also, the value of  $P(\mathbf{x}_i|\boldsymbol{\theta}_k)$  can become very small and lead to computational difficulties if the variance is chosen to be too small. In the experiments reported here, we used fixed variance Gaussians and the data was sphered. The variance was selected by trying a few values ranging between 0.5 and 2 and choosing the variance that maximized conditional log-likelihood (a part of the training set was used for validation). As with most iterative optimization methods, the model estimated by the SDR-MM algorithm depends on parameter initialization. We tried a few different random starts and chose the model which gives highest conditional log-likelihood on training data (validation was not used for this purpose).

### 7.1. Classification results

Table 1. Description of data sets for the classification problem.

DATA SET	DATA DIMENSION	TRAINING SET SIZE	TEST SET SIZE
HEART DISEASE	13	149	148
IONOSPHERE	34	151	200
BREAST CANCER	30	200	369
WAVEFORM	21	300	500

We give classification results comparing SDR-MM with KDR, MDA and pHd. We modified the matlab package of Kernel ICA (Bach, 2002) to obtain the KDR results. The variance parameter for KDR was gradually decreased (between iterations) to two as suggested in (Fukumizu et al., 2004). For the experiments with MDA and pHd, we used the *mda* and *dr* packages in the R language. We used four data sets from the UCI machine learning repository, viz. Heart disease, Ionosphere, Breast cancer and waveform data sets (summarized in Table 1).

Table 2 shows classification results obtained by first projecting data using the various methods and then

Table 2. Accuracies for best SVM classifiers associated with projection onto various lower dimensions.

DATA SET	L	PHD	KDR	MDA	SDR-MM
HEART	1	52.37	80.68	77.84	80.81
	3	68.92	77.43	77.97	80.95
	5	73.31	76.82	80.74	81.49
IONOSPHERE	1	68.80	90.28	75.75	87.14
	3	82.75	95.28	86.9	89.71
	5	87.65	94.88	88.85	91.14
BREAST	1	73.88	93.82	92.55	95.50
	3	84.23	90.92	93.36	95.83
	5	90.41	88.59	93.88	95.85
WAVEFORM	1	-	59.32	60.58	60.98
	2	-	82.80	84.40	85.16
	4	61.6	79.08	83.78	84.36

Table 3. Calculated t-values for comparison between various dimension reduction methods followed by SVM classifier. Paired samples test of significance for 10-fold cross validation is significant with probability 0.05/0.01/0.001 if t-value is higher than 2.23/3.17/4.59, respectively. Positive/negative t-value means that the first/second classifier, respectively, is better than the other.

DATA SET	L	SDR-MM vs KDR	SDR-MM vs MDA	KDR vs MDA
HEART	1	0.13	0.90	0.70
	3	2.16	0.94	-0.17
	5	4.60	0.91	-2.82
IONOSPHERE	1	-1.62	3.44	6.06
	3	-3.34	1.94	7.37
	5	-2.78	1.18	7.06
BREAST	1	2.50	4.52	1.69
	3	4.12	4.00	-1.68
	5	5.23	2.44	-3.67
WAVEFORM	1	2.11	0.47	-1.40
	2	3.58	1.69	-4.18
	4	6.53	1.08	-6.06

using SVM to classify the projected data. For MDA and SDR-MM, we obtained results similar to SVM using the inherent classifier, that uses the probability densities estimated in the process of finding the lower dimensional space (not shown here for lack of space). The classification rates shown in the table are averaged 10-fold cross validation results. The t-values of the paired significance tests comparing SDR-MM, MDA and KDR are given in Table 3. We found that SDR-MM performs significantly better than KDR on all of the data sets except one - the Ionosphere data. SDR-MM also did better than MDA consistently, but the significance t-values were not (on an average) as high as the comparison with KDR.

## 7.2. Visualization - Gaussian case

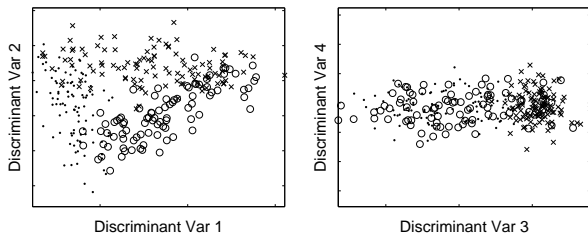


Figure 3. Some two dimensional views of waveform dataset projected onto the four basis vectors obtained using SDR-MM

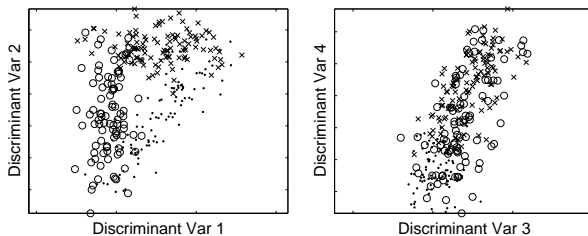


Figure 4. Some two dimensional views of waveform dataset projected onto the four basis vectors obtained using KDR

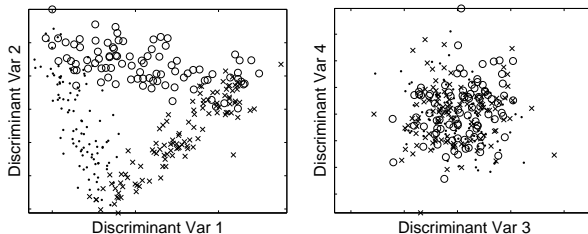


Figure 5. Some two dimensional views of waveform dataset projected onto the four basis vectors obtained using MDA

For the visualization experiment we used the Waveform data set. We trained a model with 30 Gaussian components (10 for each class) and with mean parameters constrained to a four-dimensional subspace. The estimated matrix  $V$  was processed using the Gram-Schmidt procedure to obtain orthogonal basis for the lower dimensional subspace and the training data was projected onto this subspace. Figure 3 shows two views of this four-dimensional projected set. The first two coordinates were sufficient to discriminate between the three classes since the two-dimensional model achieves an error rate close to the minimum possible (Bayes) error (see Table 2). However, we see that the third coordinate distinguishes one class from the other two, indicating that maximum mutual information based methods may be able to discover

more discriminating information than what is needed for classification. KDR projection gave similar lower dimensional views, but with greater overlap among the three classes (figure 4). In the corresponding projections obtained using MDA, shown in figure 5, the third and fourth discriminants do not significantly discriminate between the classes.

## 7.3. Visualization - Binary case

We demonstrate the binary data visualization capability of SDR-MM with Bernoulli conditional distribution. While performing the experiments we found that the algorithm was much more likely to get stuck in local minima when the Bernoulli mixture components are used than in the Gaussian case. The visualization shown in this section was obtained by running the SDR-MM algorithm several times and picking the best view. For this purpose, we use the ICU data set (Lemeshow et al., 1988) which consists of a sample of 200 subjects who were part of a study on survival of patients following admission to an adult intensive care unit (ICU). We picked 190 patients and 16 binary features from this data-set.

The goal is to extract and understand features that predict whether a patient will leave the ICU alive. The features considered include presence of coma, cancer, fracture and infection, the patient’s gender and race and whether the admission to ICU was elective or due to an emergency. The two dimensional projection obtained using MCL estimation of constrained Bernoulli mixture model is shown in Fig. 6. We examined the basis vectors of the lower-dimensional parameter space obtained using SDR-MM, and found that the features that change most significantly along the horizontal direction are the type of admission (elective versus emergency) and whether a fracture was involved. Along the vertical direction, the feature with maximum change is presence of cancer.

The projected data can be visually divided into five clusters (figure 6). Four of the clusters, numbered 1, 2, 4 and 5, were relatively ‘pure’, i.e., consist of either people who left the ICU alive or those who did not, while cluster 3 consists of both types of people. Some conclusions that can be readily drawn from this are that people who elected to join ICU to receive medical attention survived with high probability. Among those who joined the ICU because of an emergency, those who joined because of a fracture survived with high probability (cluster 1), though some of these (presumably with severe damage) did not survive. The type of service at admission and type of admission are highly correlated for this cluster.

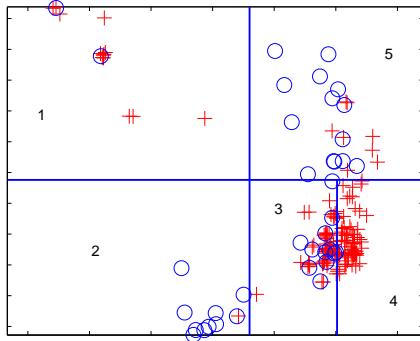


Figure 6. Two dimensional representation of binary data from the ICU data set : patients who left the ICU alive are shown by '+' and the patients who did not by 'o'.

## 8. Conclusion

Semi-parametric PCA is a recently proposed probabilistic alternative to principal component analysis based on maximum likelihood estimation of latent variable models. In this paper, we argued that use of maximum conditional likelihood estimation is a natural way to extend this method to the supervised case. Experiments demonstrate the potential of this method to learn discriminating transformations and for supervised visualization of high dimensional data.

There are many promising directions for future work. Typically, supervised multi-class dimension reduction experiments involve learning directions which discriminate among all classes simultaneously. Finding projections suitable for separating pairs (or more generally subsets) of classes can give better discriminative directions. Outputs from these low-complexity classifiers can then be combined to obtain full classifiers with good performance. Another interesting extension would be to use mixture modelling approach with a suitable objective function for semi-supervised dimensionality reduction.

## Acknowledgments

We thank Sanjoy Dasgupta and Thomas John for helpful discussions.

## References

- Bach, F. (2002). The kernel-ica package, <http://www.cs.berkeley.edu/~fbach/kernel-ica/index.htm>.
- Friedman, J. H., & Stuetzle, W. (1981). Projection pursuit regression. *Journal of the american statistical association*, 76, 817–823.

- Fukumizu, K., Bach, F. R., & Jordan, M. I. (2004). Dimensionality reduction for supervised learning with reproducing kernel hilbert spaces. *Journal of Machine Learning Research*, 5, 73–99.
- Globerson, A., & Tishby, N. (2003). Sufficient dimensionality reduction. *Journal of Machine Learning Research*, 3, 1307–1331.
- Hastie, T., & Tibshirani, R. (1986). Generalized additive models. *Statistical Science*, 1, 297–318.
- Hastie, T., & Tibshirani, R. (1996). Discriminant analysis by Gaussian mixtures. *Journal of the Royal Statistical Society series B*, 58, 158–176.
- Jebara, T., & Pentland, A. (1998). Maximum conditional likelihood via bound maximization and the CEM algorithm. *Neural Information Processing Systems 11*.
- Klautau, A., Jevtic, N., & Orlitsky, A. (2003). Discriminative gaussian mixture models: A comparison with kernel classifiers. *Proc. 20th International Conf. on Machine Learning* (pp. 353–360).
- Lemeshow, S., Teres, D., Avrunin, J. S., & Pastides, H. (1988). Predicting the outcome of intensive care unit patients. *Journal of the American Statistical Association*, 83, 348–356.
- Li, K.-C. (1991). Sliced inverse regression for dimension reduction (with discussion). *Journal of american statistical association*, 86, 316–342.
- Li, K.-C. (1992). On principal hessian directions for data visualization and dimension reduction: Another application of stein's lemma. *Journal of american statistical association*, 87, 1026–1039.
- McCullagh, P., & Nelder, J. A. (1983). *Generalized linear models*. Monographs on Statistics and Applied Probability. Chapman and Hall.
- Omohundro, S. M. (1987). Efficient algorithms with neural networks behaviour. *Complex Systems*, 1, 273–347.
- Park, J., & Sandberg, L. W. (1991). Universal approximation using radial basis function networks. *Neural Computation*, 3, 246–257.
- Sajama, & Orlitsky, A. (2004). Semi-parametric exponential family PCA. *Advances in Neural Information Processing 17 (NIPS)*.
- Torkkola, K., & Campbell, W. M. (2000). Mutual information in learning feature transformations. *Proc. 17th International Conf. on Machine Learning* (pp. 1015–1022).