

# We are IntechOpen, the world's leading publisher of Open Access books Built by scientists, for scientists

4,800

Open access books available

122,000

International authors and editors

135M

Downloads

Our authors are among the

154

Countries delivered to

TOP 1%

most cited scientists

12.2%

Contributors from top 500 universities



WEB OF SCIENCE™

Selection of our books indexed in the Book Citation Index  
in Web of Science™ Core Collection (BKCI)

Interested in publishing with us?  
Contact [book.department@intechopen.com](mailto:book.department@intechopen.com)

Numbers displayed above are based on latest data collected.  
For more information visit [www.intechopen.com](http://www.intechopen.com)



## Discriminative Cluster Analysis

Fernando De la Torre and Takeo Kanade  
*Robotics Institute, Carnegie Mellon University*  
 5000 Forbes Avenue Pittsburgh  
 USA

### 1. Introduction

Clustering is one of the most widely used statistical methods in data analysis (e.g. multimedia content-based retrieval, molecular biology, text mining, bioinformatics). Recently, with an increasing number of database applications that deal with very large high dimensional datasets, clustering has emerged as a very important research area in many disciplines. Unfortunately, many known algorithms tend to break down in high dimensional spaces because of the sparsity of the points. In such high dimensional spaces not all the dimensions might be relevant for clustering, outliers are difficult to detect, and the curse of dimensionality makes clustering a challenging problem. Also, when handling large amounts of data, time complexity becomes a limiting factor.

There are two types of clustering algorithms: partitional and hierarchical (Jain et al., 1999). Partitional methods (e.g.  $k$ -means, mixture of Gaussians, graph theoretic, mode seeking) only produce one partition of the data; whereas hierarchical ones (e.g. single link, complete link) produce several of them. In particular,  $k$ -means (MacQueen, 1967) is one of the simplest unsupervised learning algorithms that has been extensively studied and extended (Jain, 1988). Although  $k$ -means is a widely used technique due to its ease of programming and good performance, it suffers from several drawbacks. It is sensitive to initial conditions, it does not remove undesirable features for clustering, and it is optimal only for hyper-spherical clusters. Furthermore, its complexity in time is  $O(nkl)$  and in space is  $O(k)$ , where  $n$  is the number of samples,  $k$  is the number of clusters, and  $l$  the number of iterations. This degree of complexity can be impractical for large datasets.

To partially address some of these challenges, this paper proposes Discriminative Cluster Analysis (DCA). DCA jointly performs clustering and dimensionality reduction. In the first step, DCA finds a low dimensional projection of the data well suited for clustering by encouraging preservation of distances between neighboring data points belonging to the same class. Once the data is projected into a low dimensional space, DCA performs a "soft" clustering of the data. Later, this information is feedback into the dimensionality reduction step until convergence. Clustering in the DCA subspace is less prone to local minima, noisy dimensions that are irrelevant for clustering are removed, and clustering is faster to compute (especially for high dimensional data). Recently, other researchers (Ding & Li., 2007), (Ye et al., 2007) have further explored advantages of discriminative clustering methods versus generative approaches.

Source: Theory and Novel Applications of Machine Learning, Book edited by: Meng Joo Er and Yi Zhou,  
 ISBN 978-3-902613-55-4, pp. 376, February 2009, I-Tech, Vienna, Austria

## 2. Previous work

This section reviews previous work on k-means, spectral methods for clustering, and linear discriminant analysis in a unified framework.

### 2.1 k-means and spectral graph methods: a unified framework

k-means (MacQueen, 1967; Jain, 1988) is one of the simplest and most popular unsupervised learning algorithms used to solve the clustering problem. Clustering refers to the partition of  $n$  data points into  $c$  disjoint clusters. k-means clustering splits a set of  $n$  objects into  $c$  groups by maximizing the between-cluster variation relative to within-cluster variation. In other words, k-means clustering finds the partition of the data that is a local optimum of the following energy function:

$$J(\mathbf{m}_1, \dots, \mathbf{m}_c) = \sum_{i=1}^c \sum_{j \in C_i} \|\mathbf{d}_j - \mathbf{m}_i\|_2^2 \quad (1.1)$$

where  $\mathbf{d}_j$  (see notation<sup>1</sup>) is a vector representing the  $j^{\text{th}}$  data point and  $\mathbf{m}_i$  is the geometric centroid of the data points for class  $i$ . The optimization criterion in eq. (1.1) can be rewritten in matrix form as:

$$E_1(\mathbf{M}, \mathbf{G}) = \|\mathbf{D} - \mathbf{M}\mathbf{G}^T\|_F \text{ subject to } \mathbf{G}\mathbf{1}_c = \mathbf{1}_n \text{ and } g_{ij} \in \{0, 1\} \quad (1.2)$$

where  $\mathbf{G}$  is an indicator matrix, such that  $\sum_j g_{ij} = 1$ ,  $g_{ij} \in \{0, 1\}$  and  $g_{ij}$  is 1 if  $\mathbf{d}_i$  belongs to class  $j$ ,  $c$  denotes the number of classes and  $n$  is the number of samples.  $\mathbf{M} \in \mathfrak{R}^{d \times c}$  is the matrix containing all the means for each cluster. The columns of  $\mathbf{D} \in \mathfrak{R}^{d \times n}$  contain the original data points, and refers to the number of features. The equivalence between the k-means error function of eq. (1.1) and eq. (1.2) is only valid if  $\mathbf{G}$  strictly satisfies the constraints.

The k-means algorithm performs coordinate descent in  $E_1(\mathbf{M}, \mathbf{G})$ . Given the actual value of the means  $\mathbf{M}$ , the first step finds, for each data point  $\mathbf{d}_j$ , the value of  $\mathbf{g}^j$  minimizing eq. (1.2) subject to the constraints. The second step optimizes  $\mathbf{M} = \mathbf{D}\mathbf{G}(\mathbf{G}^T\mathbf{G})^{-1}$ , which effectively computes the mean of each cluster. Although it can be proven that alternating these two steps will always converge, the k-means algorithm does not necessarily find the optimal configuration of all possible assignments. The algorithm is significantly sensitive to the

---

<sup>1</sup> Bold capital letters denote matrices  $\mathbf{D}$ , and bold lower-case letters signify a column vector  $\mathbf{d}$ .  $\mathbf{d}_j$  represents the  $j^{\text{th}}$  column of the matrix  $\mathbf{D}$ .  $\mathbf{d}^j$  is a column vector that designates the  $j$ -th row of the matrix  $\mathbf{D}$ . All non-bold letters refer to scalar variables.  $d_{ij}$  corresponds to the scalar in the row  $i$  and column  $j$  of the matrix  $\mathbf{D}$ , as well as the  $i$ -th element of a column vector  $\mathbf{d}_j$ . *diag* is an operator that transforms a vector into a diagonal matrix or transforms the diagonal of a matrix into a vector. *vec* vectorizes a matrix into a vector.  $\mathbf{1}_k \in \mathfrak{R}^{k \times 1}$  is a vector of ones.  $\mathbf{I}_k \in \mathfrak{R}^{k \times k}$  denotes the identity matrix.  $\|\mathbf{d}\|_2$  denotes the norm of the vector  $\mathbf{d}$ .  $\text{tr}(\mathbf{A}) = \sum_i a_{ii}$  is the trace of the matrix  $\mathbf{A}$ , and  $|\mathbf{A}|$  denotes the determinant.  $\|\mathbf{A}\|_F = \text{tr}(\mathbf{A}^T\mathbf{A}) = \text{tr}(\mathbf{A}\mathbf{A}^T)$  designates the Frobenious norm of matrix  $\mathbf{A}$ .  $N_d(\mathbf{x}; \mu, \Sigma)$  indicates a  $d$ -dimensional Gaussian on the variable  $\mathbf{x}$  with mean  $\mu$  and covariance  $\Sigma$ .  $\circ$  denotes the Hadamard or point-wise product.

initial randomly selected cluster centers. It is typically run multiple times, and the solution with less error is chosen. Despite these limitations, the algorithm is used frequently as a result of its easiness of implementation and effectiveness.

After optimizing over  $\mathbf{M}$ , eq. (1.2), can be rewritten as:

$$E_2(\mathbf{G}) = \|\mathbf{D} - \mathbf{D}\mathbf{G}(\mathbf{G}^T\mathbf{G})^{-1}\mathbf{G}^T\|_F = tr(\mathbf{D}^T\mathbf{D}) - tr((\mathbf{G}^T\mathbf{G})^{-1}\mathbf{G}^T\mathbf{D}^T\mathbf{D}\mathbf{G}) \geq \sum_{i=c+1}^{\min(d,n)} \lambda_i \quad (1.3)$$

where  $\lambda_i$  are the eigenvalues of  $\mathbf{D}^T\mathbf{D}$ . Minimizing  $E_2(\mathbf{G})$ , eq. (1.3), is equivalent to maximizing  $tr((\mathbf{G}^T\mathbf{G})^{-1}\mathbf{G}^T\mathbf{D}^T\mathbf{D}\mathbf{G})$ . Ignoring the special structure of  $\mathbf{G}$  and considering  $g_{ij}$  in the continuous domain, the optimum  $\mathbf{G}$  value is given by the eigenvectors of the Gram matrix  $\mathbf{D}^T\mathbf{D}$ . The error of  $E_2$  with the optimal continuous  $\mathbf{G}$  is  $E_2 = \sum_{i=c+1}^{\min(d,n)} \lambda_i$ . A similar reasoning has been reported by (Ding & He, 2004; Zha et al., 2001), demonstrating that a lower bound of  $E_2(\mathbf{G})$ , eq. (1.3), is given by the sum of residual eigenvalues. The continuous solution of  $\mathbf{G}$  lies in the  $c-1$  subspace, spanned by the first  $c-1$  eigenvectors with highest eigenvalues (Ding & He, 2004) of  $\mathbf{D}^T\mathbf{D}$ .

Finally, it is worthwhile to point out the connections between  $k$ -means and standard spectral graph algorithms (Dhillon et al., 2004), such as Normalized Cuts (Shi & Malik, 2000), by means of kernel methods. The kernel trick is a standard method for lifting the points of a dataset to a higher dimensional space, where points are more likely to be linearly separable (assuming that the correct mapping is found). Consider a lifting of the original points to a higher dimensional space,  $\Gamma = [\phi(\mathbf{d}_1) \phi(\mathbf{d}_2) \dots \phi(\mathbf{d}_n)]$  where  $\phi$  represents a high dimensional mapping. The kernelized version of eq. (1.2) is:

$$E_3(\mathbf{M}, \mathbf{G}) = \|(\Gamma - \mathbf{M}\mathbf{G}^T)\mathbf{W}\|_F \quad (1.4)$$

in which we introduce a weighting matrix  $\mathbf{W}$  for normalization purposes. Eliminating  $\mathbf{M} = \Gamma\mathbf{W}\mathbf{W}^T\mathbf{G}(\mathbf{G}^T\mathbf{W}\mathbf{W}^T\mathbf{G})^{-1}$ , it can be shown that:

$$E_3 \propto tr((\mathbf{G}^T\mathbf{W}\mathbf{W}^T\mathbf{G})^{-1}\mathbf{G}^T\mathbf{W}\mathbf{W}^T\Gamma^T\Gamma\mathbf{W}\mathbf{W}^T\mathbf{G}) \quad (1.5)$$

where  $\Gamma^T\Gamma$  is the standard affinity matrix in Normalized Cuts (Shi & Malik, 2000). After a change of variable  $\mathbf{Z} = \mathbf{G}^T\mathbf{W}$ , the previous equation can be expressed as  $E_3(\mathbf{Z}) \propto tr((\mathbf{Z}\mathbf{Z}^T)^{-1}\mathbf{Z}\mathbf{W}^T\Gamma^T\Gamma\mathbf{W}\mathbf{Z}^T)$ . Choosing  $\mathbf{W} = \text{diag}(\Gamma^T\Gamma \mathbf{1}_n)^{-\frac{1}{2}}$  the problem is equivalent to solving the Normalized Cuts problem. This formulation is more general since it allows for arbitrary kernels and weights. In addition, the weight matrix can be used to reject the influence of pairs of data points with unknown similarity (i.e. missing data).

## 2.2 Linear discriminant analysis

The aim of LDA is to find a low dimensional projection, where the means of the classes are as far as possible from each other, and the intra-class variation is small. LDA can be computed in closed form using the following covariance matrices, conveniently expressed in matrix form (de la Torre & Kanade, 2005):

$$\begin{aligned}
 f\mathbf{S}_t &= \sum_{j=1}^n (\mathbf{d}_j - \mathbf{m})(\mathbf{d}_j - \mathbf{m})^T = \mathbf{D}\mathbf{P}_1\mathbf{D}^T \\
 f\mathbf{S}_w &= \sum_{i=1}^c \sum_{\mathbf{d}_j \in C_i} (\mathbf{d}_j - \mathbf{m}_i)(\mathbf{d}_j - \mathbf{m}_i)^T = \mathbf{D}\mathbf{P}_2\mathbf{D}^T \\
 f\mathbf{S}_b &= \sum_{i=1}^c n_i(\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^T = \mathbf{D}\mathbf{P}_3\mathbf{D}^T
 \end{aligned}$$

where  $f = n-1$ , and  $\mathbf{P}_i$ 's are projection matrices (i.e.  $\mathbf{P}_i^T = \mathbf{P}_i$  and  $\mathbf{P}_i^2 = \mathbf{P}_i$ ) with the following expressions:

$$\mathbf{P}_1 = \mathbf{I} - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^T \quad \mathbf{P}_2 = \mathbf{I} - \mathbf{G}(\mathbf{G}^T\mathbf{G})^{-1}\mathbf{G}^T \quad \mathbf{P}_3 = \mathbf{G}(\mathbf{G}^T\mathbf{G})^{-1}\mathbf{G}^T - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^T \quad (1.6)$$

$\mathbf{S}_b$  is the between-class covariance matrix and represents the average distance between the mean of the classes.  $\mathbf{S}_w$  is the within-class covariance matrix and it is a measure of the average compactness of each class. Finally,  $\mathbf{S}_t$  is the total covariance matrix. Through these matrix expressions, it can be easily verified that  $\mathbf{S}_t = \mathbf{S}_w + \mathbf{S}_b$ . The upper bounds on the ranks of the matrices are  $\min(c-1, d)$ ,  $\min(n-c, d)$ ,  $\min(n-1, d)$  for  $\mathbf{S}_b, \mathbf{S}_w$ , and  $\mathbf{S}_t$  respectively.

LDA computes a linear transformation of the data  $\mathbf{B} \in \mathfrak{R}^{d \times k}$  that maximizes the distance between class means and minimizes the variance within clusters. Rayleigh like quotients are among the most popular LDA optimization criterion (Fukunaga, 1990). For instance, LDA can be obtained by minimizing:

$$E_2(\mathbf{B}) = \text{tr}((\mathbf{B}^T\mathbf{S}_1\mathbf{B})^{-1}\mathbf{B}^T\mathbf{S}_2\mathbf{B}) \quad (1.7)$$

where several combinations of  $\mathbf{S}_1$  and  $\mathbf{S}_2$  matrices lead to the same LDA solution (e.g.  $\mathbf{S}_1 = \{\mathbf{S}_b, \mathbf{S}_w, \mathbf{S}_t\}$  and  $\mathbf{S}_2 = \{\mathbf{S}_w, \mathbf{S}_t, \mathbf{S}_b\}$ ). The Rayleigh quotient of eq.(1.7) has a closed-form solution in terms of a Generalized Eigenvalue Problem (GEP),  $\mathbf{S}_2\mathbf{B} = \mathbf{S}_1\mathbf{B}\boldsymbol{\Lambda}$  (Fukunaga, 1990). In the case of high-dimensional data (e.g. images) the covariance matrices are not likely to be full rank due to the lack of training samples and alternative approaches to compute LDA are needed. This is the well-known small sample size (SSS) problem. There are many techniques to solve the GEP when  $\mathbf{S}_1$  and  $\mathbf{S}_2$  are rank deficient, see (Zhang & Sim, 2007; Ye, 2005) for a recent review. However, solving LDA with standard eigensolvers is not efficient (neither space or time) for large amounts of high dimensional data. Formulating LDA as a leastsquares problem suggests efficient methods to solve LDA techniques. Moreover, a least-squares formulation of LDA facilitates its analysis and generalization.

Consider the following weighted between-class covariance matrix  $\hat{\mathbf{S}}_b = \mathbf{D}\mathbf{G}\mathbf{G}^T\mathbf{D}^T = \sum_{i=1}^c \left(\frac{n_i}{n}\right)^2 \mathbf{m}_i\mathbf{m}_i^T$ , that favors classes with more samples.  $\mathbf{m}_i$  is the mean vector for class  $i$ , and we assume zero mean data (i.e.  $\mathbf{m} = \frac{1}{n}\mathbf{D}\mathbf{1}_n$ ). Previous work on neural networks (Gallinari et al., 1991; Lowe & Webb, 1991) have shown that maximizing  $J_4(\mathbf{B}) = \text{tr}((\mathbf{B}^T\hat{\mathbf{S}}_b\mathbf{B})(\mathbf{B}^T\mathbf{S}_t\mathbf{B})^{-1})$  is equivalent to minimizing:

$$E_4(\mathbf{B}, \mathbf{V}) = \|\mathbf{G}^T - \mathbf{V}\mathbf{B}^T\mathbf{D}\|_F \propto -\text{tr}((\mathbf{B}^T\mathbf{D}\mathbf{D}^T\mathbf{B})^{-1}\mathbf{B}^T\mathbf{D}\mathbf{G}\mathbf{G}^T\mathbf{D}^T\mathbf{B}) \quad (1.8)$$

This approach is attractive because (Baldi & Hornik, 1989) have shown that the surface of eq. (1.8) has a unique local minima, and several saddle points.

### 3. Discriminative cluster analysis

In the previous section, we have provided a least-squares framework for LDA (supervised dimensionality reduction) and k-means (unsupervised clustering). The aim of DCA is to combine clustering and dimensionality reduction in an unsupervised manner. In this section, we propose a least-squares formulation for DCA.

#### 3.1 Error function for LDA and DCA

The key aspect to simultaneously performing dimensionality reduction and clustering is the analysis of eq. (1.8). Ideally we would like to optimize eq. (1.8) w.r.t.  $\mathbf{B}$  and  $\mathbf{G}$ . However, directly optimizing eq. (1.8) has several drawbacks. First, eq. (1.8) biases the solution towards classes that have more samples because it maximizes  $\hat{\mathbf{S}}_b = \mathbf{D}\mathbf{G}\mathbf{G}^T\mathbf{D}^T = \sum_{i=1}^c (\frac{n_i}{n})^2 (\mathbf{m}_i)(\mathbf{m}_i)^T$ . Secondly, eq. (1.8) does not encourage sparseness in  $\mathbf{G}$  if  $g_{ij} > 0$ . That is, assuming that  $\mathbf{C} = \mathbf{B}^T\mathbf{D} \in \mathfrak{R}^{k \times n}$ , then eq. (1.8) is equivalent to  $E_4 = \text{tr}(\mathbf{G}^T\mathbf{G}) - \text{tr}(\mathbf{G}^T\mathbf{C}^T(\mathbf{C}\mathbf{C}^T)^{-1}\mathbf{C}\mathbf{G})$ . If  $g_{ij} \forall i, j$  is positive, minimizing the first term,  $\text{tr}(\mathbf{G}^T\mathbf{G})$ , does not encourage sparseness in  $\mathbf{g}^i \forall i$  ( $\mathbf{g}^i$  represents the  $i^{\text{th}}$  row of  $\mathbf{G}$ , see notation).

In this section, we correct eq. (1.8) to obtain the unbiased LDA criterion by normalizing  $E_4$  as follows:

$$E_5(\mathbf{B}, \mathbf{V}, \mathbf{G}) = \|(\mathbf{G}^T\mathbf{G})^{-\frac{1}{2}}(\mathbf{G}^T - \mathbf{V}\mathbf{B}^T\mathbf{D})\|_F \quad (1.9)$$

where  $(\mathbf{G}^T\mathbf{G})^{-\frac{1}{2}}$  is the normalization factor. After eliminating  $\mathbf{V}$ , eq. (1.9) can be written as:

$$E_5(\mathbf{B}, \mathbf{G}) = \|(\mathbf{G}^T\mathbf{G})^{-\frac{1}{2}}\mathbf{G}^T(\mathbf{I}_n - \mathbf{C}^T(\mathbf{C}\mathbf{C}^T)^{-1}\mathbf{C})\|_F \\ \propto \text{tr}(\underbrace{(\mathbf{B}^T\mathbf{D}\mathbf{D}^T\mathbf{B})^{-1}}_{f\mathbf{S}_t} \mathbf{B}^T \underbrace{\mathbf{D}\mathbf{G}(\mathbf{G}^T\mathbf{G})^{-1}\mathbf{G}^T\mathbf{D}^T\mathbf{B}}_{f\mathbf{S}_b}) \quad (1.10)$$

If  $\mathbf{G}$  is known, eq. (1.10) is the exact expression for LDA.

Eq. (1.10) is also the basis for DCA. Unlike LDA, DCA is an unsupervised technique and  $\mathbf{G}$  will be a variable to optimize, subject to the constraints that  $g_{ij} \in \{0,1\}$ , and  $\mathbf{G}\mathbf{1}_c = \mathbf{1}_n$ . DCA jointly optimizes the data projection matrix  $\mathbf{B}$  and the indicator matrix  $\mathbf{G}$ .

#### 3.2 Updating $\mathbf{B}$

The optimal  $\mathbf{B}$  given  $\mathbf{G}$  can be computed in closed form by solving the following GEP:

$$\mathbf{R}\mathbf{D}\mathbf{D}^T\mathbf{B} = \mathbf{D}\mathbf{D}^T\mathbf{B}\Lambda_1 \quad \text{where} \quad \mathbf{R} = \mathbf{G}(\mathbf{G}^T\mathbf{G})^{-1}\mathbf{G}^T \quad (1.11)$$

There are many methods for efficiently solving the GEP in the case of highdimensional data when  $(d \gg n)$  (de la Torre et al., 2005; Zhang & Sim, 2007; Ye, 2005). In this section, we propose a regularized stable closed form solution. Assuming  $\mathbf{D}^T\mathbf{D}$  is full rank, computing  $(\mathbf{D}^T\mathbf{D})^{-1}$  can be a numerically unstable process, especially if  $\mathbf{D}^T\mathbf{D}$  has eigenvalues close to zero. A common method to solve ill-conditioning is to regularize the solution by factorizing  $\Sigma = \mathbf{D}^T\mathbf{D}$  as the sum of the outer products plus a scaled identity matrix, i.e.  $\Sigma \approx \mathbf{V}\Lambda\mathbf{V}^T + \sigma^2\mathbf{I}_d$ .  $\mathbf{V} \in \mathfrak{R}^{n \times k}$ ,  $\Lambda \in \mathfrak{R}^{k \times k}$  is a diagonal matrix. The parameters  $\sigma^2$ ,  $\mathbf{V}$  and  $\Lambda$  are estimated by minimizing:

$$E_c(\mathbf{V}, \Lambda, \sigma^2) = \|\Sigma - \mathbf{V}\Lambda\mathbf{V}^T - \sigma^2\mathbf{I}_n\|_F \quad (1.12)$$

After optimizing over  $\mathbf{V}, \Lambda, \sigma^2$ , it can be shown (de la Torre & Kanade, 2005) that:  $\sigma^2 = \text{tr}(\Sigma - \mathbf{V}\hat{\Lambda}\mathbf{V}^T)/d - k$ ,  $\Lambda = \hat{\Lambda} - \sigma^2\mathbf{I}_d$ , where  $\hat{\Lambda}$  is a matrix containing the eigenvalues of the covariance matrix  $\Sigma$  and  $\mathbf{V}$  the eigenvectors. This expression is equivalent to probabilistic PCA (Moghaddam & Pentland, 1997; Roweis & Ghahramani, 1999; Tipping & Bishop, 1999). After the factorization, the matrix inversion lemma (Golub & Loan, 1989)  $(\mathbf{A}^{-1} + \mathbf{V}\mathbf{C}^{-1}\mathbf{V}^T)^{-1} = \mathbf{A} - \mathbf{A}\mathbf{V}(\mathbf{C} + \mathbf{V}^T\mathbf{A}\mathbf{V})^{-1}\mathbf{V}^T\mathbf{A}$  is applied to invert  $(\mathbf{V}\Lambda\mathbf{V}^T + \sigma^2\mathbf{I}_n)^{-1}$ , which results in:

$$(\mathbf{V}\Lambda\mathbf{V}^T + \sigma^2\mathbf{I}_n)^{-1} = \frac{1}{\sigma^2}(\mathbf{I}_n - \frac{1}{\sigma^2}\mathbf{V}(\Lambda^{-1} + \frac{\mathbf{I}_n}{\sigma^2})^{-1}\mathbf{V}^T)$$

Now, solving  $(\mathbf{I}_n - \frac{1}{\sigma^2}\mathbf{V}(\Lambda^{-1} + \frac{\mathbf{I}_n}{\sigma^2})^{-1}\mathbf{V}^T)\mathbf{R}\mathbf{D}^T\mathbf{D}\alpha = \alpha\Lambda$  becomes a better conditioned problem.

The number of bases ( $k$ ) are bounded by the number of classes ( $c$ ), because the  $\text{rank}(\mathbf{D}\mathbf{R}\mathbf{D}^T) = c$ . We typically choose  $c-1$  to be consistent with LDA. Moreover, the best clustering results are achieved by projecting the data into a space of  $c-1$  dimensions. Also, observe that there is an ambiguity in the result, because for any invertible matrix  $\mathbf{T}_1 \in \mathbf{R}^{k \times k}$ ,  $E_5(\mathbf{B}) = E_5(\mathbf{B}\mathbf{T}_1)$ .

### 3.3 Optimizing $\mathbf{G}$

Let  $\mathbf{A} = \mathbf{C}^T(\mathbf{C}\mathbf{C}^T)^{-1}\mathbf{C} \in \mathfrak{R}^{n \times n}$ , where  $\mathbf{C} = \mathbf{B}^T\mathbf{D}$ , then eq. (1.10) can be rewritten as:

$$E_5(\mathbf{G}) \propto \text{tr}((\mathbf{G}^T\mathbf{G})^{-1}\mathbf{G}^T\mathbf{A}\mathbf{G}) \quad (1.13)$$

Optimizing eq. (1.13) subject to  $g_{ij} \in \{0,1\}$  and  $\mathbf{G}\mathbf{1}_c = \mathbf{1}_n$  is an *NP* complete problem. To make it tractable, we relax the discrete constraint on  $g_{ij}$  allowing to take values in the range  $(0,1)$ . To use a gradient descent search mechanism, we parameterize  $\mathbf{G}$  as the Hadamard (pointwise) product of two matrices  $\mathbf{G} = \mathbf{V} \circ \mathbf{V}$  (Liu & Yi, 2003), and use the following updating scheme:

$$\begin{aligned} \mathbf{V}^{n+1} &= \mathbf{V}^n - \eta \frac{\partial E_5(\mathbf{G})}{\partial \mathbf{V}} \\ \frac{\partial E_5(\mathbf{G})}{\partial \mathbf{V}} &= (\mathbf{I}_c - \mathbf{G}(\mathbf{G}^T\mathbf{G})^{-1}\mathbf{G}^T)\mathbf{A}\mathbf{G}(\mathbf{G}^T\mathbf{G})^{-1} \circ \mathbf{V} \end{aligned} \quad (1.14)$$

The increment of the gradient,  $\eta$ , in eq. (1.14) is determined with a line search strategy (Fletcher, 1987). To impose  $\mathbf{G}\mathbf{1}_c = \mathbf{1}_n$  in each iteration,  $\mathbf{V}$  is normalized to satisfy the constraint. Because eq. (1.14) is prone to local minima, this method starts from several random initial points and selects the solution with smallest error.

This optimization problem is similar in spirit to recent work on clustering with non-negative matrix factorization (Zass & Shashua, 2005; Ding et al., 2005; Lee & Seung, 2000). However, we optimize a discriminative criterion rather than a generative one. Moreover, we simultaneously compute dimensionality reduction and clustering, using a different optimization technique.

### 3.4 Initialization

At the beginning, neither  $\mathbf{G}$  nor  $\mathbf{B}$  are known, but the matrix  $\mathbf{G}(\mathbf{G}^T\mathbf{G})^{-1}\mathbf{G}^T$  can be estimated from the available data. Similar to previous work (He & Niyogi, 2003), we compute an

estimate of a local similarity matrix,  $\mathbf{G}(\mathbf{G}^T\mathbf{G})^{-1}\mathbf{G}^T \in \mathfrak{R}^{n \times n}$ , from data. We assume that  $(\mathbf{G}^T\mathbf{G}) \approx s\mathbf{I}_c$ , so that all classes are equally distributed and  $s$  is the number of samples per class.  $\mathbf{R} = \frac{1}{s}\mathbf{G}\mathbf{G}^T$  is a hard-affinity matrix, where  $r_{ij}$  will be 1 if  $\mathbf{d}_i$  and  $\mathbf{d}_j$  are considered neighbors (i.e. belong to the same class).  $\mathbf{R}$  can be estimated by simply computing the  $k$  nearest neighbors for each data point using the Euclidian distance. To make  $\mathbf{R}$  symmetric, if  $\mathbf{d}_i$  is within the  $k$ -neighborhood of  $\mathbf{d}_j$ , but not the contrary, then its similarity is set to zero. Figure 1.5.b shows an estimate of  $\mathbf{R}$  for 15 subjects in the ORL database. Each subject (class) has ten samples and for each sample the nearest nine neighbors are selected. The samples are ordered by class. After factorizing  $\mathbf{R} = \mathbf{U}\Sigma\mathbf{U}^T$ , we normalize  $\mathbf{R}$  as  $\hat{\mathbf{R}} \approx \mathbf{U}_c\mathbf{U}_c^T$ , where  $\mathbf{U}_c \in \mathfrak{R}^{n \times c}$  are the first  $c$  eigenvectors of  $\mathbf{R}$ .  $\hat{\mathbf{R}}$  is the initial neighbor matrix.

### 3.5 Interpreting the weighted covariance matrix

A key aspect to understand DCA is the interpretation of the weighted covariance matrix  $\mathbf{DRD}^T = \sum_{i=1}^n \sum_{j=1}^n r_{ij}\mathbf{d}_i\mathbf{d}_j^T$ . Principal Component Analysis (PCA) (Jolliffe, 1986) computes a basis  $\mathbf{B}$  that maximizes the variance of the projected samples, i.e. PCA finds an orthonormal basis that maximizes  $\text{tr}(\mathbf{B}^T\mathbf{D}\mathbf{D}^T\mathbf{B}) = \sum_{i=1}^n \|\mathbf{B}^T\mathbf{d}_i\|_2^2$ . The PCA solution  $\mathbf{B}$  is given by the eigenvectors of  $\mathbf{D}\mathbf{D}^T$ . Finding the leading eigenvectors of  $\mathbf{DRD}^T$  is equivalent to maximizing  $\text{tr}(\mathbf{B}^T\mathbf{DRD}^T\mathbf{B}) = \sum_{i=1}^n \sum_{j=1}^n r_{ij}\mathbf{d}_i^T\mathbf{B}\mathbf{B}^T\mathbf{d}_j$ . If  $\mathbf{R} = \mathbf{I}$ , it is equivalent to standard PCA. However, if  $\mathbf{R}$  is  $\mathbf{G}(\mathbf{G}^T\mathbf{G})^{-1}\mathbf{G}^T$ , where  $\mathbf{G}$  is the indicator matrix (or an approximation), the weighted covariance only maximizes the covariance within each cluster. This effectively maximizes the correlation between each pair of points in the same class. Figure 1.1 shows a toy problem with two oriented Gaussian classes. The first eigenvector in PCA finds a direction of maximum variance that does not necessarily correspond to maximum discrimination. In fact, by projecting the data into the first principal component, the clusters overlap. If  $\mathbf{R}$  is the initial matrix of neighbors, the first step of DCA finds a more suitable projection that maximizes class separability (see fig. 1.1).

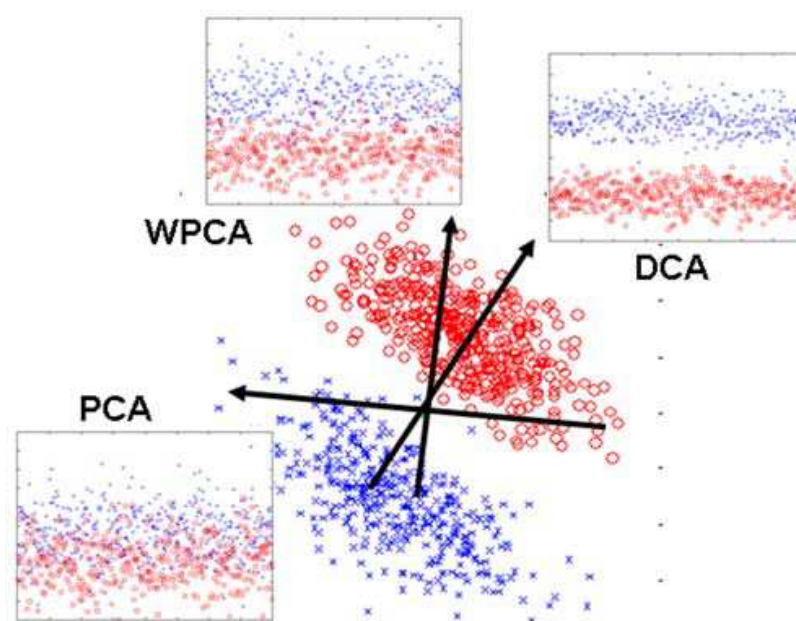


Fig. 1.1 Two class toy problem. PCA, WPCA, and DCA projections in one dimensional space.



## 4. Experiments

This section describes three experiments using synthetic and real data that demonstrate the effectiveness of DCA for clustering.

### 4.1 Clustering with DCA

In the first experiment, we show how the DCA error function is able to correctly cluster oriented clusters.

Consider the DCA optimization expression, eq. (1.10), when  $\mathbf{B} = \mathbf{I}_d$  (i.e. no projection); in this case, eq. (1.10) becomes  $tr((\mathbf{G}^T\mathbf{G})^{-1}\mathbf{G}^T\mathbf{D}^T(\mathbf{D}\mathbf{D}^T)^{-1}\mathbf{D}\mathbf{G})$ . This error function, due to the term  $(\mathbf{D}\mathbf{D}^T)^{-1}$ , provides affine invariance to clustering. To illustrate this property, we have generated three examples of three two-dimensional random Gaussian clusters. Figure 1.2.a shows three clusters of 300 samples each, generated from three two-dimensional Gaussians:  $N_2(\mathbf{x}; [-4;3], 0.25\mathbf{I}_2)$ ,  $N_2(\mathbf{x}; [-4;2], 0.25\mathbf{I}_2)$  and  $N_2(\mathbf{x}; [7;3], 0.25\mathbf{I}_2)$ . Similarly, fig. 1.2.b illustrates 300 samples generated from three two-dimensional Gaussians  $N_2(\mathbf{x}; [-10;-10], 0.25\mathbf{I}_2)$ ,  $N_2(\mathbf{x}; [-10;-5], 0.25\mathbf{I}_2)$  and  $N_2(\mathbf{x}; [30;15], 0.25\mathbf{I}_2)$ . Analogously, fig. 1.2.c shows  $N_2(\mathbf{x}; [-4;3], 2[1 \ 0.8; 0.1 \ 1])$ ,  $N_2(\mathbf{x}; [-4;2], 0.25[1 \ 0.8; 0.1 \ 1])$  and  $N_2(\mathbf{x}; [3;3], 0.25[1 \ 0.8; 0.1 \ 1])$ .

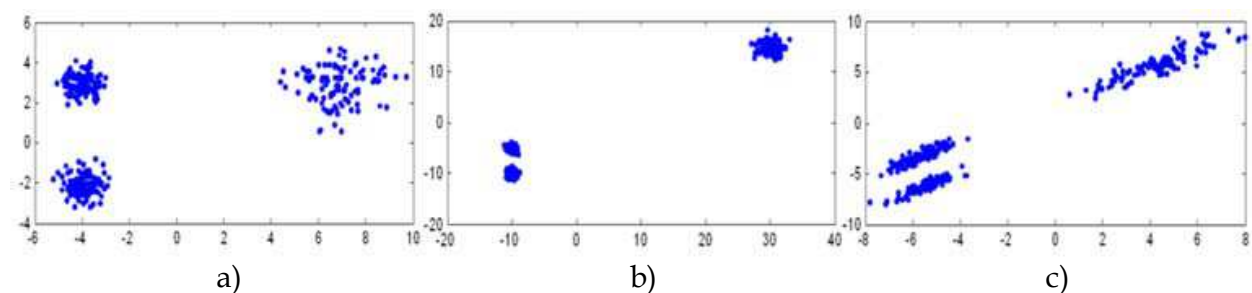


Fig. 1.2 Three examples of three two-dimensional Gaussian clusters.

We run DCA and  $k$ -means with the same random initialization and let both algorithms converge. To compute the accuracy of the results for a  $c$  cluster case, we compute a  $c$ -by- $c$  confusion matrix  $\mathbf{C}$ , where each entry  $c_{ij}$  is the number of data points in cluster  $i$  that belong to class  $j$ . It is difficult to compute the accuracy by strictly using the confusion matrix  $\mathbf{C}$ , because it is unknown which cluster matches with which class. An optimal way to solve this problem is to compute the following maximization problem (Zha et al., 2001; Knuth, 1993):

$$\max tr(\mathbf{C}\mathbf{P}) \mid \mathbf{P} \text{ is a permutation matrix} \quad (1.15)$$

To solve eq. (1.15), we use the classical Hungarian algorithm (Knuth, 1993). Table (1.2) shows the clustering accuracy for the three examples described above. We run the algorithms 1000 times from different random initializations (same for  $k$ -means and DCA).

	$k$ -means	DCA
Fig. 1.2.a	0.713±0.23%	0.990±0.05%
Fig. 1.2.b	0.526±0.07%	0.959±0.09%
Fig. 1.2.c	0.594±0.13%	0.974±0.06%

Table 1.1 Comparison of clustering accuracy for DCA and  $k$ -means.

As we can see from the results in table 1.1, DCA is able to achieve better clustering results starting from the same initial condition as  $k$ -means. Moreover, DCA results in a more stable

(less variance) clustering.  $k$ -means clustering accuracy largely degrades when two clusters are closer together or the clusters are not spherical. DCA is able to keep the accuracy even with oriented clusters (fig. 1.2.c).

#### 4.2 Removing undesirable dimensions

The second experiment demonstrates the ability of DCA to deal with undesired dimensions not relevant for clustering. A synthetic problem is created as follows: 200 samples from a two-dimensional Gaussian distribution with mean  $[-5,-5]$  and another 200 samples from another Gaussian distribution with mean  $[5,5]$  are generated ( $x$  and  $y$  dimensions). We add a third dimension generated with uniform noise between  $[0,35]$  ( $z$  dimension). Figure 1.3 shows 200 samples of each class in the original space (fig. 1.3.a), as well as the projection (fig. 1.3.b) onto  $x$  and  $y$ . The  $k$ -means algorithm is biased by the noise (fig. 1.4.a). Similarly, projecting the data into the first two principal components also produces the wrong clustering because PCA preserves the energy of the uniform noise, which is not relevant for clustering. However, DCA is able to remove the noise and achieve the correct clustering as evidenced in fig. 1.4.b. In this particular example 15 neighbors were initially selected and  $\mathbf{B} \in \mathfrak{R}^{3 \times 2}$ .

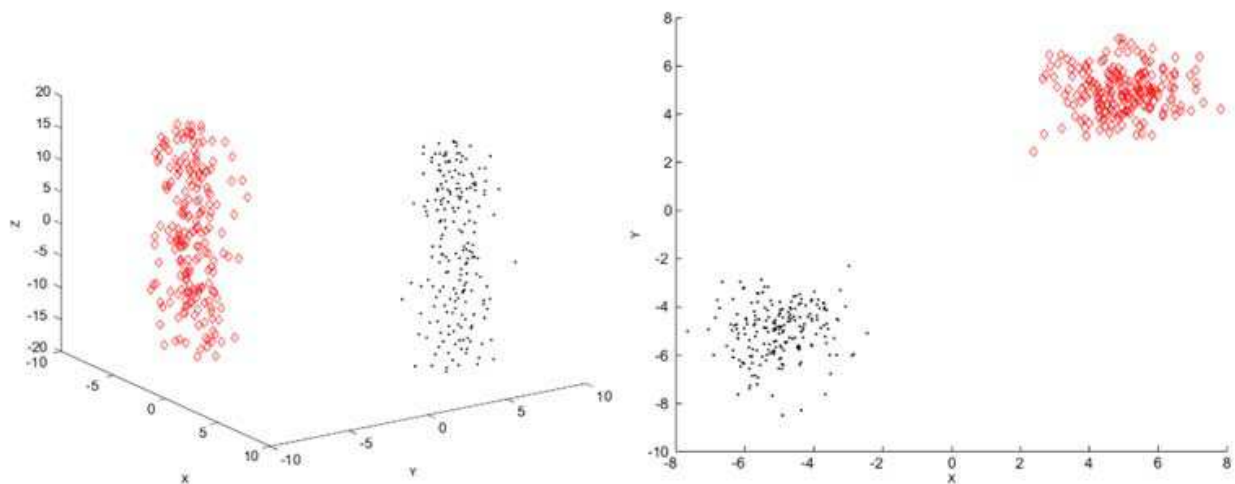


Fig. 1.3 a) 2 classes of 3 dimensional data. b) Projection onto XY space.

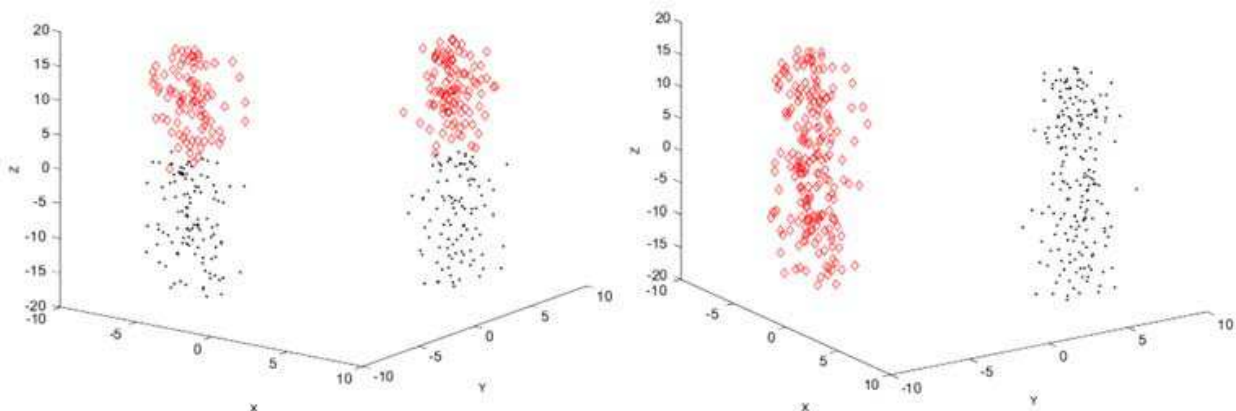


Fig. 1.4 a)  $k$ -means clustering. b) DCA clustering.

### 4.3 Clustering faces

The final experiment shows results on clustering faces from the ORL face database (Samaria & Harter, 1994). The ORL face database is composed of 40 subjects and 10 images per subject. We randomly selected  $c$  subjects from the database and add the 10 images for each subject to  $\mathbf{D} \in \mathfrak{R}^{d \times 10c}$  (e.g. fig. 1.5.a). Afterwards, we compute PCA, weighted PCA (WPCA), PCA+LDA (preserving 95% of the energy in PCA), and DCA. After computing PCA, WPCA (with the initial matrix  $\mathbf{R}$ ), and PCA+LDA, we run the  $k$ -means algorithm 10 times and the solution with smallest error is chosen. This procedure is repeated 40 times for different number of classes (between 4 and 40 subjects). To perform a fair comparison, we project the data into the number of classes minus ones ( $c-1$ ) dimensions for all methods.

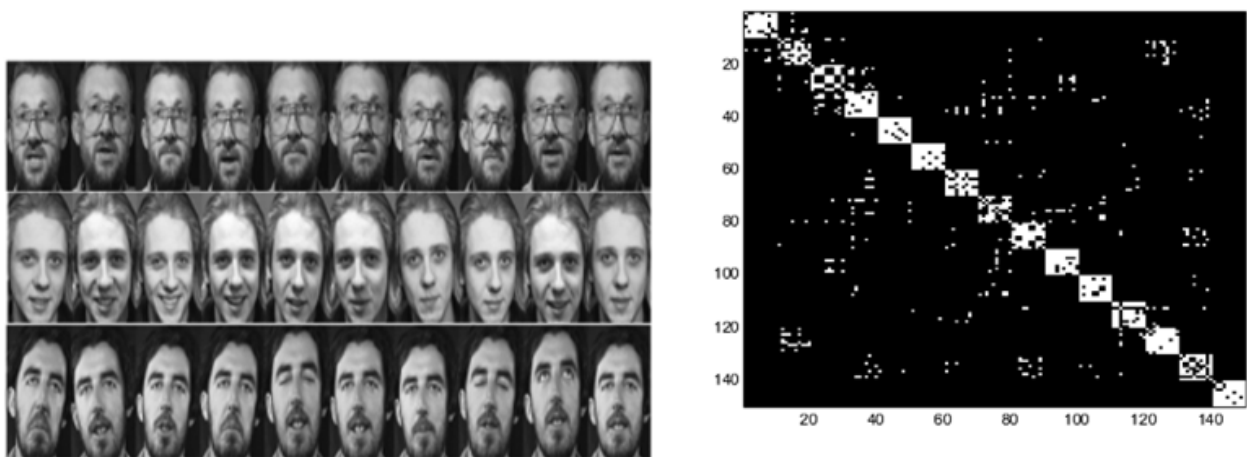


Fig. 1.5 a) Some faces of the ORL data base. b) Estimate of  $\mathbf{R}$  for 15 clusters (people), each cluster has 10 samples. The samples are ordered by clusters.

$c$	PCA	WPCA	DCA	PCA+LDA
4	73±0%	1±0%	87±2%	<b>1±0%</b>
10	88±6%	95±6%	<b>97±4%</b>	88±8%
15	86±5%	88±4%	<b>96±1%</b>	82±6%
20	80±4%	84±4%	<b>87±2%</b>	83±4%
25	77±3%	80±4%	<b>87±2%</b>	80±4%
30	75±3%	79±3%	<b>81±3%</b>	81±4%
35	73±4%	77±3%	78±4%	<b>81±3%</b>
40	71±2%	74±3%	73±3%	<b>80±4%</b>

Table 1.2 Comparison of the clustering accuracy for several projection methods (same number of bases).

Fig. 1.6 shows the accuracy in clustering for PCA+ $k$ -means versus DCA. For a given number of clusters, we show the mean and variance over 40 realizations. DCA always outperforms PCA+ $k$ -means. Table 1.2 shows some numerical values for the clustering accuracy. DCA outperforms most of the methods when there are between 5 and 30 classes. For more classes, PCA+LDA performs marginally better. In addition, the accuracy of the PCA+ $k$ -means method drops as the number of classes increases (as expected).

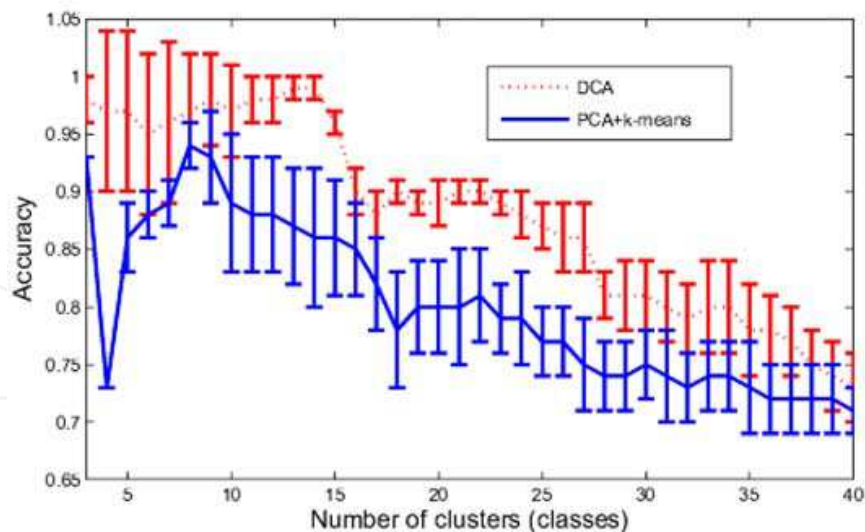


Fig. 1.6 Accuracy of clustering versus the number of classes. Blue PCA and red DCA (dotted line).

## 5. Discussion and future work

In this paper, we have proposed DCA, a technique that jointly performs dimensionality reduction and clustering. In synthetic and real examples, DCA outperforms standard  $k$ -means and PCA+ $k$ -means, for clustering high dimensional data. DCA provides a discriminative embedding that minimizes cluster separation and is less prone to local minima. Additionally, we have proposed an unbiased least-squares formulation for LDA. Although DCA has shown promising preliminary results, several issues still need to be addressed. It remains unclear how to select the optimal number of clusters. Several model order selection (e.g. Minimum Description Length or Akaike information criterion) could be applied towards this end. On the other hand, DCA assumes that all the clusters have the same orientation (not necessarily spherical). This limitation could be easily address by using kernel extensions of eq. (1.10) to deal with non-Gaussian clusters.

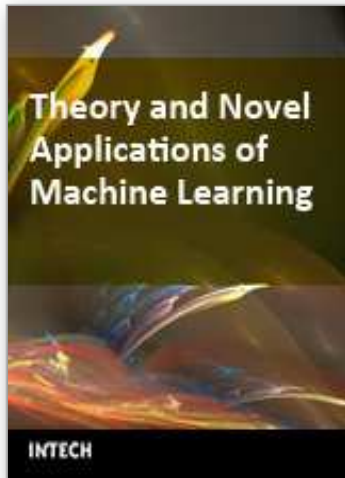
## 6. Acknowledgements

This work has been partially supported by NIH R01 51435 from the National Institute of Mental Health, N000140010915 from the Naval Research Laboratory, the Department of the Interior National Business Center contract no. NBCHD030010, and SRI International subcontract no. 03-000211.

## 7. References

- Baldi, P., & Hornik, K. (1989). Neural networks and principal component analysis: Learning from examples without local minima. *Neural Networks*, 2, 53–58.
- de la Torre, F., Gross, R., Baker, S., & Kumar, V. (2005). Representational oriented component analysis for face recognition with one sample image per training class. *Computer Vision and Pattern Recognition*.
- de la Torre, F., & Kanade, T. (2005). Multimodal oriented discriminant analysis. *International Conference on Machine Learning* (pp. 177–184).

- Dhillon, I. S., Guan, Y., & Kulis, B. (2004). A unified view of kernel k-means, spectral clustering and graph partitioning. *UTCS Tech. Report TR-04-25*.
- Ding, C., & He, X. (2004). K-means clustering via principal component analysis. *International Conference on Machine Learning* (pp. 225–232).
- Ding, C., He, X., & Simon, H. (2005). On the equivalence of nonnegative matrix factorization and spectral clustering. *Siam International Conference on Data Mining (SDM)*.
- Ding, C., & Li, T. (2007). Adaptive dimension reduction using discriminant analysis and k-means clustering. *International Conference on Machine Learning*.
- Fletcher, R. (1987). *Practical methods of optimization*. John Wiley and Sons.
- Fukunaga, K. (1990). *Introduction to statistical pattern recognition, second edition*. Academic Press, Boston, MA.
- Gallinari, P., Thiria, S., Badran, F., & Fogelman-Soulie, F. (1991). On the relations between discriminant analysis and multilayer perceptrons. *Neural Networks*, 4, 349–360.
- Golub, G., & Loan, C. F. V. (1989). *Matrix computations*. 2nd ed. The Johns Hopkins University Press.
- He, X., & Niyogi, P. (2003). Locality preserving projections. *Neural Information Processing Systems*.
- Jain, A., Murty, M., & Flynn, P. (1999). Data clustering: A review. *ACM Computing Surveys*.
- Jain, A. K. (1988). *Algorithms for clustering data*. Prentice Hall.
- Jolliffe, I. T. (1986). *Principal component analysis*. New York: Springer-Verlag.
- Knuth, D. E. (1993). *The standford graphbase*. Addison-Wesley Publishing Company.
- Lee, D., & Seung, H. (2000). Algorithms for non-negative matrix factorization. *Neural Information Processing Systems* (pp. 556–562).
- Liu, W., & Yi, J. (2003). Existing and new algorithms for nonnegative matrix factorization. *University of Texas at Austin*.
- Lowe, D. G., & Webb, A. (1991). Optimized feature extraction and the bayes decision in feed-forward classifier networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 355–364.
- MacQueen, J. B. (1967). Some methods for classification and analysis of multivariate observations. *5-th Berkeley Symposium on Mathematical Statistics and Probability*. Berkeley, University of California Press. (pp. 1:281–297).
- Moghaddam, B., & Pentland, A. (1997). Probabilistic visual learning for object representation. *Pattern Analysis and Machine Intelligence*, 19, 137–143.
- Roweis, S., & Ghahramani, Z. (1999). A unifying review of linear gaussian models. *Neural Computation*, 11, 305–345.
- Samaria, F., & Harter, A. (1994). Parameterization of a stochastic model for human face identification. *Proceedings of the 2nd IEEE Workshop on Applications of Computer Vision*.
- Shi, J., & Malik, J. (2000). Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22, 888–905.
- Tipping, M., & Bishop, C. M. (1999). Probabilistic principal component analysis. *Journal of the Royal Statistical Society B*, 61, 611–622.
- Ye, J. (2005). Generalized low rank approximation of matrices. *Machine Learning*.
- Ye, J., Zhao, Z., & Wu, M. (2007). Discriminative k-means for clustering. *Advances in Neural Information Processing Systems*.
- Zass, R., & Shashua, A. (2005). A unifying approach to hard and probabilistic clustering. *International Conference on Computer Vision*. Beijing.
- Zha, H., Ding, C., Gu, M., He, X., & Simon, H. (2001). Spectral relaxation for k-means clustering. *Neural Information Processing Systems* (pp. 1057–1064).
- Zhang, S., & Sim, T. (2007). Discriminant subspace analysis: A fukunaga-koontz approach. *PAMI*, 29, 1732–1745.



## **Theory and Novel Applications of Machine Learning**

Edited by Meng Joo Er and Yi Zhou

ISBN 978-953-7619-55-4

Hard cover, 376 pages

**Publisher** InTech

**Published online** 01, January, 2009

**Published in print edition** January, 2009

Even since computers were invented, many researchers have been trying to understand how human beings learn and many interesting paradigms and approaches towards emulating human learning abilities have been proposed. The ability of learning is one of the central features of human intelligence, which makes it an important ingredient in both traditional Artificial Intelligence (AI) and emerging Cognitive Science. Machine Learning (ML) draws upon ideas from a diverse set of disciplines, including AI, Probability and Statistics, Computational Complexity, Information Theory, Psychology and Neurobiology, Control Theory and Philosophy. ML involves broad topics including Fuzzy Logic, Neural Networks (NNs), Evolutionary Algorithms (EAs), Probability and Statistics, Decision Trees, etc. Real-world applications of ML are widespread such as Pattern Recognition, Data Mining, Gaming, Bio-science, Telecommunications, Control and Robotics applications. This book reports the latest developments and futuristic trends in ML.

### **How to reference**

In order to correctly reference this scholarly work, feel free to copy and paste the following:

Fernando De la Torre and Takeo Kanade (2009). Discriminative Cluster Analysis, Theory and Novel Applications of Machine Learning, Meng Joo Er and Yi Zhou (Ed.), ISBN: 978-953-7619-55-4, InTech, Available from:

[http://www.intechopen.com/books/theory\\_and\\_novel\\_applications\\_of\\_machine\\_learning/discriminative\\_cluster\\_analysis](http://www.intechopen.com/books/theory_and_novel_applications_of_machine_learning/discriminative_cluster_analysis)

**INTECH**  
open science | open minds

### **InTech Europe**

University Campus STeP Ri  
Slavka Krautzeka 83/A  
51000 Rijeka, Croatia  
Phone: +385 (51) 770 447  
Fax: +385 (51) 686 166  
[www.intechopen.com](http://www.intechopen.com)

### **InTech China**

Unit 405, Office Block, Hotel Equatorial Shanghai  
No.65, Yan An Road (West), Shanghai, 200040, China  
中国上海市延安西路65号上海国际贵都大饭店办公楼405单元  
Phone: +86-21-62489820  
Fax: +86-21-62489821

© 2009 The Author(s). Licensee IntechOpen. This chapter is distributed under the terms of the [Creative Commons Attribution-NonCommercial-ShareAlike-3.0 License](#), which permits use, distribution and reproduction for non-commercial purposes, provided the original is properly cited and derivative works building on this content are distributed under the same license.

IntechOpen

IntechOpen