# MAXIMUM MARGIN DISCRIMINANT PROJECTIONS FOR FACIAL EXPRESSION RECOGNITION

*Symeon Nikitidis, Anastasios Tefas and Ioannis Pitas*

Department of Informatics, Aristotle University of Thessaloniki, Greece
{nikitidis,tefas,pitas}@aiia.csd.auth.gr

## ABSTRACT

We present a novel dimensionality reduction method which aims to identify a low dimensional projection subspace, where samples form classes that are better discriminated and separated with maximum margin. The proposed method brings certain advantages, both to data embedding and classification. It improves classification performance, reduces the required training time of the SVM classifier, since it is trained over the projected low dimensional samples and also data outliers and the overall data samples distribution inside classes do not affect its performance. The proposed method has been applied for facial expression recognition in Cohn-Kanade database verifying its superiority in this task, against other state-of-the-art dimensionality reduction techniques.

*Index Terms*— Subspace learning, maximum margin projections, support vector machines, facial expression recognition

## 1. INTRODUCTION

One of the most crucial problems that every facial image analysis algorithm encounters is the high dimensionality of the image data, which can range from several hundreds to thousands of extracted image features. Directly dealing with such high dimensional data is not only computational inefficient, but also yields several problems in subsequently performed statistical learning algorithms, due to the so-called *"curse of dimensionality"*. Thus, various techniques have been proposed for efficient data embedding (or dimensionality reduction) that obtain a more manageable problem and alleviate computational complexity. Moreover, reducing the dimensionality of the original data can reveal the actual hidden underlying data structure that can be efficiently described using only a small number of degrees of freedom. Such a popular category of methods is the subspace image representation algorithms which aim to discover the latent image features by projecting linearly or non-linearly the high-dimensional input samples to a low-dimensional subspace, where an appropriately formed criterion is optimized.

Focusing on the underlying optimization criterion, a popular category of subspace learning algorithms are those that attempt to enhance classes discrimination in the reduced dimensional projection space. These algorithms aim to identify a discriminative subspace, in which the data samples from different classes are far apart from each other. Linear Discriminant Analysis (LDA) [1] and its variants, are such representative methods that extract discriminant information by finding projection directions that achieve intra-class compactness and inter-class separability.

Margin maximizing embedding algorithms [2, 3, 4] inspired by the great success of Support Vector Machines (SVMs) [5] also aim to enhance data discrimination in the low dimensional space. In [3] the Maximum Margin Projection (MMP) algorithm has been proposed, which is an unsupervised embedding method that attempts to find different subspace directions that separate data points in different clusters with maximum margin. To do so, MMP seeks for such a data labelling, so that, if an SVM classifier is trained, the resulting separating hyperplanes can separate different data clusters with the maximum margin. He et. al in [2] also exploited the margin maximization concept proposing a semisupervised dimensionality reduction method for image retrieval that aims to discover both geometrical and discriminant structures of the data manifold. This algorithm constructs a within-class and a between-class graph by exploiting both class and neighborhood information and finds a linear transformation matrix that maps image data to a subspace, where, at each local neighborhood, the margin between relevant and irrelevant images is maximized.

In this paper we integrate optimal data embedding and SVM classification in a single framework to be called Maximum Margin Discriminant Projections (MMDP). MMDP algorithm directly operates on the random features extracted using an orthogonal Gaussian random projection matrix and derives an optimal projection matrix such that the separating margin between the projected samples of different classes is maximized, by exploiting the decision hyperplanes obtained from training a SVM classifier. The MMDP approach brings certain advantages, both to data embedding and classification.

Since it is combined with a classification method, MMDP is appropriately tuned towards improving classification performance. Furthermore, the SVM classifier is trained over the projected low dimensional data samples determined by MMDP, thus the required computational effort is significantly reduced. Moreover, since the decision hyperplane identified by SVM training is explicitly determined by the support vectors, data outliers and the overall data samples distribution inside classes do not affect MMDP performance, in contrast to other discrimination enhancing subspace learning algorithms, such as LDA, which assumes a Gaussian data distribution for optimal classes discrimination.

The rest of the paper is organized as follows. Section 2 presents the proposed linear dimensionality reduction algorithm, while Section 3 discusses its initialization. Section 4 describes the conducted experiments for facial expression recognition and concluding remarks are drawn in Section 5.

## 2. MAXIMUM MARGIN DISCRIMINANT PROJECTIONS

Given a set $\mathcal{X} = \{(\mathbf{x}_1, y_1), ..., (\mathbf{x}_N, y_N)\}$ of $N$ training data pairs, where $\mathbf{x}_i \in \mathcal{R}^m, i = 1, ..., N$ are the $m$-dimensional input feature vectors and $y_i \in \{-1, 1\}$ is the class label associated with each sample $\mathbf{x}_i$, a binary SVM classifier attempts to find the separating hyperplane that separates training data points of the two classes with maximum margin, while minimizes the classification error defined according to which side of the decision hyperplane training samples of each class fall in. Considering that each training sample of $\mathcal{X}$ is firstly projected to a low-dimensional subspace using a projection matrix $\mathbf{R} \in \mathcal{R}^{r \times m}$, where $r \ll m$ and performing the linear transformation $\acute{\mathbf{x}}_i = \mathbf{R}\mathbf{x}_i$, the binary SVM optimization problem is formulated as follows:

$$\min_{\mathbf{w}, \xi_i} \frac{1}{2}\mathbf{w}^T\mathbf{w} + C\sum_{i=1}^{N}\xi_i \tag{1}$$

subject to the constraints:

$$y_i\left(\mathbf{w}^T\mathbf{R}\mathbf{x}_i + b\right) \geq 1 - \xi_i \tag{2}$$
$$\xi_i \geq 0, \quad i = 1, \ldots, N, \tag{3}$$

where $\mathbf{w} \in \mathcal{R}^r$ is the $r$-dimensional normal vector of the separating hyperplane, $b \in \mathcal{R}$ is its bias term, $\boldsymbol{\xi} = [\xi_1, \ldots, \xi_N]^T$ are the slack variables, each one associated with a training sample and $C$ is the term that penalizes the training error.

The MMDP algorithm attempts to learn a projection matrix $\mathbf{R}$, such that the low-dimensional data sample projection is performed efficiently, thus enhancing the discrimination between the two classes. To quantify the discrimination power of the projection matrix $\mathbf{R}$, we formulate our MMDP algorithm based on geometrical arguments. To do so, we employ a combined iterative optimization framework, involving the

simultaneous optimization of the separating hyperplane normal vector $\mathbf{w}$ and the projection matrix $\mathbf{R}$, performed by successively updating the one variable, while keeping the other fixed. Next we first discuss the derivation of the optimal separating hyperplane normal vector $\mathbf{w}_o$, in the projection subspace determined by $\mathbf{R}$ and subsequently, we demonstrate the projection matrix update with respect to the fixed $\mathbf{w}_o$.

### 2.0.1. Finding the optimal $\mathbf{w}_o$ in the projection subspace determined by $\mathbf{R}$

The optimization with respect to $\mathbf{w}$, is essentially the conventional binary SVM training problem performed in the projection subspace determined by $\mathbf{R}$, rather than in the input space. To solve the constrained optimization problem in (1) with respect to $\mathbf{w}$, we introduce positive Lagrange multipliers $\alpha_i$ and $\beta_i$ each associated with one of the constraints in (2) and (3), respectively and formulate the Lagrangian function:

$$
\begin{aligned}
\mathcal{L}(\mathbf{w}, \boldsymbol{\xi}, \mathbf{R}, \boldsymbol{\alpha}, \boldsymbol{\beta}) &= \frac{1}{2}\mathbf{w}^T\mathbf{w} + C\sum_{i=1}^{N}\xi_i \\
&- \sum_{i=1}^{N}\alpha_i\left[y_i\left(\mathbf{w}^T\mathbf{R}\mathbf{x}_i + b\right) - 1 + \xi_i\right] \\
&- \sum_{i=1}^{N}\beta_i\xi_i.
\end{aligned}
\tag{4}
$$

The solution can be found from the saddle point of the Lagrangian function, which has to be maximized with respect to the dual variables $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ and minimized with respect to the primal ones $\mathbf{w}, \boldsymbol{\xi}$ and $b$. According to the Karush-Kuhn-Tucker (KKT) conditions the partial derivatives of $\mathcal{L}(\mathbf{w}, \boldsymbol{\xi}, \mathbf{R}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ with respect to the primal variables $\mathbf{w}, \boldsymbol{\xi}$ and $b$ vanish deriving the following equalities:

$$\frac{\partial\mathcal{L}(\mathbf{w}, \boldsymbol{\xi}, \mathbf{R}, \boldsymbol{\alpha}, \boldsymbol{\beta})}{\partial\mathbf{w}} = 0 \quad \Rightarrow \quad \mathbf{w} = \sum_{i=1}^{N}\alpha_i y_i \mathbf{R}\mathbf{x}_i, \tag{5}$$

$$\frac{\partial\mathcal{L}(\mathbf{w}, \boldsymbol{\xi}, \mathbf{R}, \boldsymbol{\alpha}, \boldsymbol{\beta})}{\partial b} = 0 \quad \Rightarrow \quad \sum_{i=1}^{N}\alpha_i y_i = 0, \tag{6}$$

$$\frac{\partial\mathcal{L}(\mathbf{w}, \boldsymbol{\xi}, \mathbf{R}, \boldsymbol{\alpha}, \boldsymbol{\beta})}{\partial\xi_i} = 0 \quad \Rightarrow \quad \beta_i = C - \alpha_i y_i. \tag{7}$$

By substituting the terms from the above equalities into (4), we switch to the dual formulation, where the optimization problem in (1) is reformulated to the maximization of the following Wolfe dual problem:

$$\max_{\boldsymbol{\alpha}} \sum_{i=1}^{N}\alpha_i - \frac{1}{2}\sum_{i,j}^{N}\alpha_i\alpha_j y_i y_j \mathbf{x}_i^T\mathbf{R}^T\mathbf{R}\mathbf{x}_j. \tag{8}$$

subject to the constraints:

$$\sum_{i=1}^{N}\alpha_i y_i = 0, \quad \alpha_i \geq 0, \quad \forall \ i = 1, \ldots, N. \tag{9}$$

Consequently, solving (8) for $\boldsymbol{\alpha}$ the optimal separating hyperplane normal vector $\mathbf{w}_o$ in the reduced dimensional space determined by $\mathbf{R}$, is subsequently derived from (5).

### 2.0.2. Maximum margin projection matrix update for fixed $\mathbf{w}_o$

At each optimization round $t$ we seek to update the projection matrix $\mathbf{R}^{(t-1)}$, so that its new estimate $\mathbf{R}^{(t)}$ maximizes the separating margin between the two classes. To do so, we first project the high dimensional training samples $\mathbf{x}_i$ from the input space to a low dimensional subspace, using the projection matrix $\mathbf{R}^{(t-1)}$ derived during the previous step, and subsequently, train the binary SVM classifier in order to obtain the optimal Lagrange multipliers $\boldsymbol{\alpha}_o$ specifying the normal vector of the separating hyperplane $\mathbf{w}_o^{(t)}$.

To formulate the optimization problem for the projection matrix $\mathbf{R}$ update, we exploit the dual form of the binary SVM cost function defined in (8). However, since term $\sum_{i=1}^{N} \alpha_i$ is constant with respect to $\mathbf{R}$, we can remove it from the cost function. Moreover, in order to retain the geometrical correlation between samples in the projection subspace, we constrain the derived updated projection matrix $\mathbf{R}^{(t)}$ to be orthogonal. Consequently, the constrained optimization problem for the projection matrix $\mathbf{R}$ update can be summarized as follows:

$$\max_{\mathbf{R}} \mathcal{O}(\mathbf{R}) = \frac{1}{2} \sum_{i,j}^{N} \alpha_{i,o} \alpha_{j,o} y_i y_j \mathbf{x}_i^T \mathbf{R}^T \mathbf{R} \mathbf{x}_j, \quad (10)$$

subject to the constraint:

$$\mathbf{R}\mathbf{R}^T = \mathbf{I}, \quad (11)$$

where $\mathbf{I}$ is a $r \times r$ identity matrix.

In order to apply the constraint $\mathbf{R}\mathbf{R}^T = \mathbf{I}$ we first solve (10), without the orthogonality constraints on its rows and obtain $\acute{\mathbf{R}}$. Thus, we solve (10) for $\mathbf{R}$ keeping $\mathbf{w}_o^{(t)}$ fixed, by applying a steepest ascent optimization algorithm, which, at a given iteration $t$, invokes the following update rule:

$$\acute{\mathbf{R}}^{(t)} = \mathbf{R}^{(t-1)} + \lambda_t \nabla \mathcal{O}(\mathbf{R}^{(t-1)}), \quad (12)$$

where $\lambda_t$ is the learning step parameter for the $t$-th iteration evaluated using the methodology presented in [6] and $\nabla \mathcal{O}(\mathbf{R}^{(t-1)})$ is the partial derivative of the objective function in (10) with respect to $\mathbf{R}^{(t-1)}$, evaluated as:

$$\begin{aligned} \nabla \mathcal{O}(\mathbf{R}^{(t-1)}) &= \sum_{i,j}^{N} \alpha_{i,o} \alpha_{j,o} y_i y_j \mathbf{R}^{(t-1)} \mathbf{x}_i \mathbf{x}_j^T \\ &= \sum_{i=1}^{N} \alpha_{i,o} y_i \mathbf{w}_o^{(t)} \mathbf{x}_i^T. \end{aligned} \quad (13)$$

Thus, $\acute{\mathbf{R}}^{(t)}$ is derived as:

$$\acute{\mathbf{R}}^{(t)} = \mathbf{R}^{(t-1)} + \lambda_t \left( \sum_{i=1}^{N} \alpha_{i,o} y_i \mathbf{w}_o^{(t)} \mathbf{x}_i^T \right). \quad (14)$$

Obtaining the projection matrix $\acute{\mathbf{R}}^{(t)}$ that increases the separating margin between the two classes in the projection subspace, we subsequently orthonormalize its rows by performing a Gram-Schmidt procedure, to derive $\mathbf{R}^{(t)}$.

After deriving the new projection matrix $\mathbf{R}^{(t)}$, the previously identified separating hyperplane is no longer optimal, since it has been evaluated in the projection subspace determined by $\mathbf{R}^{(t-1)}$. Consequently, it is required to re-project the training samples using $\mathbf{R}^{(t)}$ and retrain the SVM classifier to obtain the current optimal separating hyperplane and its normal vector. Thus, MMDP algorithm iteratively updates the projection matrix and evaluates the normal vector of the optimal separating hyperplane $\mathbf{w}_o$ in the projection subspace determined by $\mathbf{R}$, until the algorithm converges. In order to determine algorithms convergence we track the partial derivative value in (13) to identify stationarity. The following stationarity check step is performed, which examines whether the following termination condition is satisfied:

$$||\nabla \mathcal{O}(\mathbf{R}^{(t)})||_F \leq e_{\mathbf{R}} ||\nabla \mathcal{O}(\mathbf{R}^{(0)})||_F, \quad (15)$$

where $e_{\mathbf{R}}$ is a predefined stopping tolerance. In our conducted experiments, we considered that $e_{\mathbf{R}} = 10^{-3}$. The combined iterative optimization process of the MMDP algorithm is summarized in Algorithm 1.

---

**Algorithm 1** Maximum Margin Discriminant Projections Algorithm Considering a Binary Classification Problem.

---

1: **Input:** The set $\mathcal{X} = \{(\mathbf{x}_i, y_i), \quad i = 1, \ldots, N\}$ of $N$ $m$-dimensional two class train data samples.
2: **Output:** The optimal maximum margin projection matrix $\mathbf{R}_o$ and the optimal separating hyperplane normal vector $\mathbf{w}_o$.
3: **Initialize:** $t = 1$ and $\mathbf{R}^{(0)} \in \mathcal{R}^{r \times m}$ as an orthogonal Gaussian random projection matrix.
4: **repeat**
5:     **Project** $\mathbf{x}_i$ to a low dimensional subspace performing the linear transformation:
    $\acute{\mathbf{x}}_i = \mathbf{R}^{(t-1)} \mathbf{x}_i \quad \forall i = 1, \ldots, N$.
6:     **Train** the binary SVM classifier in the projection subspace by solving the optimization problem in (8) subject to the constraints in (9) to obtain the optimal Lagrange multipliers $\boldsymbol{\alpha}_o$.
7:     **Obtain** the normal vector of the optimal separating hyperplane as:
    $\mathbf{w}_o^{(t)} = \sum_{i=1}^{N} \alpha_{i,o} y_i \mathbf{R}^{(t-1)} \mathbf{x}_i$.
8:     **Determine** learning rate $\lambda_t$.
9:     **Evaluate** $\nabla \mathcal{O}(\mathbf{R}^{(t-1)}) = -\sum_{i=1}^{N} \alpha_{i,o} y_i \mathbf{w}_o^{(t)} \mathbf{x}_i^T$.
10:     **Update** projection matrix $\mathbf{R}^{(t-1)}$ given $\mathbf{w}_o^{(t)}$ as:
    $\mathbf{R}^{(t)} = \text{Orthogonalize}(\mathbf{R}^{(t-1)} - \lambda_t \sum_{i=1}^{N} \alpha_{i,o} y_i \mathbf{w}_o^{(t)} \mathbf{x}_i^T)$.
11:     $t = t + 1$, $\mathbf{R}_o = \mathbf{R}^{(t)}$ and $\mathbf{w}_o = \mathbf{w}^{(t)}$.
12: **until** $||\nabla \mathcal{O}(\mathbf{R}^{(t)})||_F \leq 10^{-3} ||\nabla \mathcal{O}(\mathbf{R}^{(0)})||_F$

---

## 3. MMDP ALGORITHM INITIALIZATION

To initialize MMDP, it is first required to train the binary SVM classifier and obtain the optimal $\mathbf{w}_o$ in a low dimensional subspace determined by an initial projection matrix $\mathbf{R}^{(0)}$, used in order to perform dimensionality reduction and form the basis of the projection subspace. To do so, we construct $\mathbf{R}^{(0)}$ as an orthogonal Gaussian random projection matrix. To derive $\mathbf{R}^{(0)}$, the following procedure is applied. We create a $m \times m$ matrix $\mathbf{G}$ of i.i.d., zero-mean, unit variance Gaussian random variables and partition it into the $m \times r$ matrix $\mathbf{Q}$ and the $m \times (m-r)$ matrix $\mathbf{P}$, thus $\mathbf{G} = [\mathbf{Q} \ \mathbf{P}]$. Consequently, we orthonormalize the columns of $\mathbf{G}$ and create an orthonormal matrix $\mathbf{G}_\perp = [\mathbf{Q}_\perp \ \mathbf{P}_\perp]$. To do so, we normalize the first column of $\mathbf{G}$ and orthogonalize the remaining columns with respect to the first, via a Gram-Schmidt procedure. This procedure results in the Gaussian random projection matrix $\mathbf{R}^{(0)} = \mathbf{Q}_\perp^T$ having orthonormal rows that can be used for the initialization of the iterative optimization framework.

## 4. EXPERIMENTAL RESULTS

To visualize the ability of MMDP algorithm to estimate useful subspaces that enhance data discrimination, we applied the proposed algorithm in a two class toy classification problem using artificial data, aiming to learn a 2D projection space. To generate our toy dataset we collected 500 300-dimensional samples for each class, with the first class features drawn randomly from a standard normal distribution $\mathcal{N}(0, 1)$ and the second class drawn from a $\mathcal{N}(0.2, 1)$ normal distribution and used 100 of them for training, while the rest were used to compose the toy test set. Figure 1 shows the 2D projection of the two classes data samples after different iterations of the MMDP algorithm. As can be observed, the proposed algorithm was able, after a few iterations, to perfectly separate linearly the two classes, by continuously maximizing the separating margin.

In addition we compared the performance of the proposed method for facial expression recognition, on the Cohn-Kanade database [7], with that of several state-of-the-art dimensionality reduction techniques, such as Eigenfaces (PCA) [8], Fisherfaces (LDA), Laplacianfaces (LPP) [9] and Randomfaces (RP) resulting by projecting facial images using random projections. For baseline comparison we also directly feed the initial high dimensional samples to a linear SVM classifier. In our implementation we have combined our optimization algorithm with LIBSVM [10], which provides an efficient implementation for solving several binary linear SVMs for multiclass classification problems. Similarly, the discriminant low-dimensional facial representations derived from the other examined algorithms were also fed to LIBSVM for classification.

Each subject in each video sequence of the Cohn-Kanade database poses a facial expression, starting from the neutral emotional state and finishing at the expression apex. To form our data collection we considered only the last video frame depicting each formed facial expression at its highest intensity. Face detection was performed on these images and the resulting facial regions of interest were manually aligned with respect to the eyes position, anisotropically scaled to a fixed size of $150 \times 200$ pixels and converted to grayscale. Thus, we used in total 407 images depicting 100 subjects, posing 7 different expressions. To measure the facial expression recognition accuracy, we randomly partitioned the available samples into 5-folds and a cross-validation has been performed. Figure 2 shows example images from the Cohn-Kanade dataset, depicting the 7 recognized facial expressions arranged in the following order: anger, fear, disgust, happiness, sadness, surprise and the neutral emotional state.



**Fig. 2**. Sample images depicting facial expressions in the Cohn-Kanade database.

Table 1 summarizes the best average facial expression recognition rates achieved by each examined embedding method, across different subspace dimensionalities varying from 3 to 500. The best recognition rate attained by MMDP is 80.4% using 150-dimensional discriminant representations of the initial 30,000-dimensional input samples. MMDP outperforms all other competing embedding algorithms by more than 3% compared against the second best performing method, which is PCA. The best average expression recognition rate attained by PCA, LDA, LPP and RP were 77.3%, 74.2%, 76.6% and 75.2%, respectively.

**Table 1**. Best average expression recognition accuracy rates (%) in Cohn-Kanade database. In parentheses is shown the dimension that results in the best performance for each method.

| SVM | PCA | LDA | LPP | RP | MMDP |
|---|---|---|---|---|---|
| 73.4 | 77.3 | 74.2 | 76.6 | 75.2 | **80.4** |
| $(30,000)$ | $(325)$ | $(6)$ | $(6)$ | $(500)$ | $(150)$ |

## 5. CONCLUSION

We proposed a discrimination enhancing subspace learning method called MMDP that aims to identify a low dimensional projection subspace where samples form classes that are separated with maximum margin. MMDP directly works with
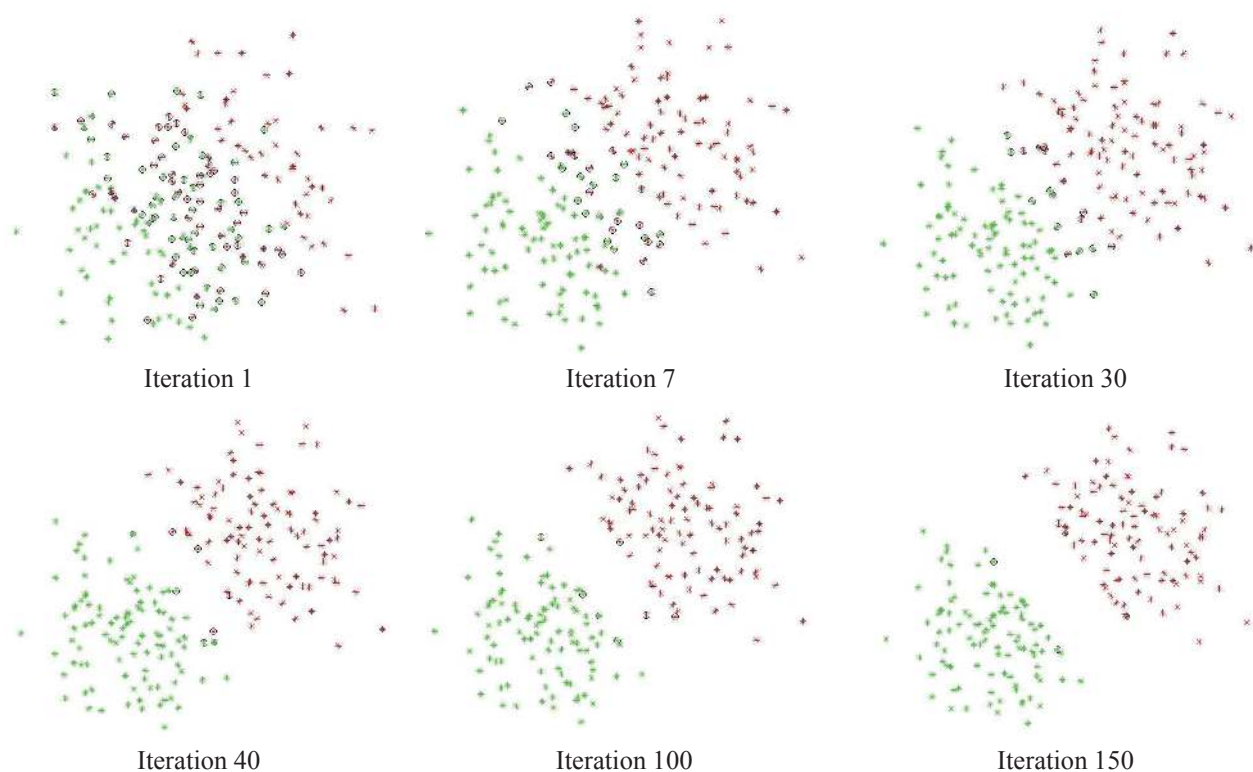
| Iteration 1 | Iteration 7 | Iteration 30 |

| Iteration 40 | Iteration 100 | Iteration 150 |

**Fig. 1**. 2D projection of the initial data at different iterations of the MMDP algorithm. Circled data samples denote the identified support vectors which reduce during MMDP algorithms convergence. As a result, the SVM training process converges faster and into a sparser solution, since the number of identified support vector decreases as classes discrimination is enhanced.

random features obtained using an orthogonal random Gaussian projection matrix and exploits the separating hyperplane obtained from training a SVM classifier in the identified low dimensional space. Experimental results showed that the proposed method outperforms current state-of-the-art embedding methods for facial expression recognition on the Cohn-Kanade database.

## 6. REFERENCES

[1] K. Fukunaga, *Introduction to statistical pattern recognition*, Academic Press, second edition, 1990.

[2] Xiaofei He, Deng Cai, and Jiawei Han, "Learning a maximum margin subspace for image retrieval," *IEEE Transactions on Knowledge and Data Engineering*, vol. 20, no. 2, pp. 189–201, February 2008.

[3] F. Wang, B. Zhao, and C. Zhang, "Unsupervised large margin discriminative projection," *IEEE Transactions on Neural Networks*, vol. 22, no. 9, pp. 1446–1456, September 2011.

[4] A. Zien and J.Q. Candela, "Large margin non-linear embedding," in *Proceedings of the 22nd International Conference on Machine Learning (ICML)*, 2005, pp. 1060–1067.

[5] Vladimir Vapnik, *The Nature of Statistical Learning Theory*, Springer-Verlag, New York, 1995.

[6] Chih-Jen Lin, "Projected gradient methods for nonnegative matrix factorization," *Neural Computation*, vol. 19, no. 10, pp. 2756–2779, 2007.

[7] T. Kanade, J.F. Cohn, and Y. Tian, "Comprehensive database for facial expression analysis," in *IEEE International Conference on Automatic Face and Gesture Recognition*, March 2000, pp. 46–53.

[8] M. Turk and A. Pentland, "Eigenfaces for recognition," *Journal of cognitive neuroscience*, vol. 3, no. 1, pp. 71–86, 1991.

[9] X He and P. Niyogi, "Locality preserving projections," in *Advances in Neural Information Processing Systems*, Vancouver, British Columbia, Canada, 2003, vol. 16.

[10] Chih-Chung Chang and Chih-Jen Lin, "LIBSVM: A library for support vector machines," *ACM Transactions on Intelligent Systems and Technology*, vol. 2, pp. 27:1–27:27, 2011.