



# On minimum class locality preserving variance support vector machine

Xiaoming Wang<sup>a</sup>, Fu-lai Chung<sup>b</sup>, Shitong Wang<sup>a,b,\*</sup>

<sup>a</sup> School of Information Technology, Jiangnan University, WuXi, JiangSu, China

<sup>b</sup> Department of Computing, Hong Kong Polytechnic University, Hong Kong, China

## ARTICLE INFO

### Article history:

Received 7 July 2009

Received in revised form

14 January 2010

Accepted 13 February 2010

### Keywords:

Supervised learning

Support vector machine

Minimum class variance support machine

Locality preserving projections

## ABSTRACT

In this paper, a so-called minimum class locality preserving variance support machine (MCLPV\_SVM) algorithm is presented by introducing the basic idea of the locality preserving projections (LPP), which can be seen as a modified class of support machine (SVM) and/or minimum class variance support machine (MCVSVM). MCLPV\_SVM, in contrast to SVM and MCVSVM, takes the intrinsic manifold structure of the data space into full consideration and inherits the characteristics of SVM and MCVSVM. We discuss in the paper the linear case, the small sample size case and the nonlinear case of the MCLPV\_SVM. Similar to MCVSVM, the MCLPV\_SVM optimization problem in the small sample size case is solved by using dimensionality reduction through principal component analysis (PCA) and one in the nonlinear case is transformed into an equivalent linear MCLPV\_SVM problem under kernel PCA (KPCA). Experimental results on real datasets indicate the effectiveness of the MCLPV\_SVM by comparing it with SVM and MCVSVM.

© 2010 Elsevier Ltd. All rights reserved.

## 1. Introduction

In the past decade, kernel methods [1] are widely studied and applied [2,3]. Support vector machine (SVM), as a kernel method, is a powerful machine learning method based on Vapnik's Statistical Learning Theory [4]. Different from other pattern recognition methods which usually attempt to minimize the misclassification errors on the training set (empirical risk minimization), SVM minimizes the structural risk, which is the probability of misclassifying a previously unseen sample [4–6]. The essential point of SVM is to find a linear separating hyperplane which achieves the maximal margin among two different classes of data in the linear case. Furthermore, SVM can be extended to build nonlinear separating decision hyperplane by exploiting kernelization techniques [1]. The optimization problem of SVM can be formulated as a quadratic programming problem which can be solved very efficiently by its dual optimization problem.

However, SVM solution does not take into consideration the class distribution and may result in a non-robust solution [7]. In order to overcome the drawback of SVM, a modified class of SVM called minimum class variance support vector machine (MCVSVM) is presented in [7] which is inspired from the optimization of Fisher's discriminant ratio [8,9]. When the

training set contains fewer samples than the dimensionality of the training vectors, it has been proved that the solution of MCVSVM problems in such cases can be found through principal component analysis (PCA) dimensionality reduction [9,10]. In the nonlinear case, it has also been shown that, under kernel PCA (KPCA) [11], the nonlinear optimization problem can be transformed into an equivalent linear MCVSVM problem. Unlike SVM, the solution of MCVSVM takes into consideration both the samples in the boundaries and the distribution of the classes and gives a robust solution. However, the intrinsic manifold structure of the data space has not been taken into full consideration in MCVSVM.

Recently, a novel linear dimensionality reduction algorithm called locality preserving projections (LPP) [12–14] is proposed. LPP aims to preserve the local manifold structure of the samples space. To be specific, the manifold structure is modeled by a nearest-neighbor graph which preserves the local structure of the data space. The LPP method is the linear approximation to the eigenfunctions of the Laplace Beltrami operator on the samples manifold. In [15], the author presented the Laplacian support vector machine (LSVM) algorithm by combining SVM and LPP. However, LSVM aimed at the semi-supervised learning method.

Aiming at the drawback of SVM and MCVSVM that the intrinsic manifold structure of the data space is ignored, in this paper, we propose a novel learning algorithm called minimum class locality preserving variance support machine (MCLPV\_SVM) in which the manifold structure within each class is explicitly considered. First, an adjacency graph in the same class is built, which can best reflect the geometry structure of the data manifold and the class

\* Corresponding author at: School of Information Technology, Jiangnan University, WuXi, JiangSu, China.

E-mail address: [wxwangst@yahoo.com.cn](mailto:wxwangst@yahoo.com.cn) (S. Wang).

relationship between the sample points. Then, by using the basic idea of the LPP, we define the locality preserving within-class scatter matrix. Finally, the optimization problem of MCLPV\_SVM is formulated by using the locality preserving within-class scatter matrix. MCLPV\_SVM are very closely related to MCVSVM and share some properties with the SVM and the MCVSVM. Both MCVSVM and MCLPV\_SVM can also be seen as the large margin classifier but they employ the Mahalanobis distance metric rather than the Euclidean distance metric in SVM. However, a significant difference between MCVSVM and MCLPV\_SVM is that the latter explicitly considers the data manifold structure. Moreover, the experimental results show that the intrinsic manifold structure is helpful to improve the classification performance. Besides, as is pointed out in [13], recently, a number of research efforts have shown that the face images possibly reside on a nonlinear submanifold. Therefore, in this case MCLPV\_SVM can usually achieve better performance in contrast to SVM and MCVSVM.

The rest of this paper is organized as follows. The related works will be reviewed in Section 2. In Section 3, the linear case of MCLPV\_SVM is presented and the small sample size case is discussed where the number of the training vectors is smaller than the dimensions of samples. In Section 4, the nonlinear decision hyperplanes will be defined and solved. In Section 5, a discussion is carried out about the relationship of the proposed method with SVM and MCVSVM and constructing the weight matrix. The experimental results are reported in Section 6. Finally, conclusions are drawn in Section 7.

## 2. Related work

In this section, SVM, MCVSVM and LPP will be briefly introduced. For simplicity, only binary classification tasks are considered here. Multi-class case can be solved by one-against-one [16]. Let a training dataset contain two classes of  $N$  samples, represented by  $\{(\mathbf{x}_1, z_1), \dots, (\mathbf{x}_N, z_N)\}$  with training samples  $\mathbf{x}_i = [\mathbf{x}_{i1}, \dots, \mathbf{x}_{iM}]^T \in \mathbf{X}$  and labels  $z_i \in \{1, -1\}$ , where  $i = 1, \dots, N$ ,  $T$  denotes transpose and the dimensionality of the sample space is denoted by  $M$  (i.e.  $\mathbf{X} \in \mathbf{R}^M$ ).

### 2.1. Minimum class variance support vector machine (MCVSVM)

Rather than simply minimizing the training error, support vector machine (SVM) [4] minimizes the structural risk which expresses an upper bound on generalization error [17]. However, actually, SVM is a local method in the sense that solution is exclusively determined by support vectors, whereas all other data points are irrelevant to the decision hyperplane [17]. Thus, a modified class of SVM called minimum class variance support machine (MCVSVM) which takes into consideration both the samples in the boundaries and the distribution of the classes and gives a robust solution is proposed in [7]. In the case where the training vectors are linearly separable MCVSVM optimization problem is defined as

$$\min_{\mathbf{w}, b} \mathbf{w}^T \mathbf{S}_W \mathbf{w} \quad (1)$$

$$\text{s.t. } z_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \quad i = 1, \dots, N \quad (2)$$

where the matrix  $\mathbf{S}_W$  is the within-class scatter matrix defined as

$$\mathbf{S}_W = \sum_{K=1}^2 \sum_{\mathbf{x} \in \mathbf{X}_K} (\mathbf{x} - \mathbf{u}_K)(\mathbf{x} - \mathbf{u}_K)^T \quad (3)$$

Here,  $\mathbf{u}_K = (1/N_K) \sum_{\mathbf{x} \in \mathbf{X}_K} \mathbf{x}$  is the mean sample vector for the  $K$ th class ( $K=1, 2$ ). MCVSVM can be seen as a compromise between SVM and FLDA [9]. Similar to SVM, the optimization problem

(1) subject to the separability constraints (2) of MCVSVM can be efficiently solved by switching to its Wolfe dual problem using a Lagrangian formulation of the problem when the within-class scatter matrix  $\mathbf{S}_W$  is nonsingular. However, MCVSVM encounter the same small size sample problem [18] as FLDA where the training set contains fewer samples than the dimensionality of the training vectors. It has been proven that the solution of the MCVSVM optimization problems in such cases can be found through PCA dimensionality reduction. It has also showed that the nonlinear MCVSVM problem is equivalent to a linear one, subject to an initial KPCA embedding of the training data.

Essentially, however, the optimization problem (1) subject to the separability constraints (2) of MCVSVM can be also written as

$$\max_{\mathbf{w}, b} \rho \quad (4)$$

$$\text{s.t. } \frac{z_i(\mathbf{w}^T \mathbf{x}_i + b)}{\sqrt{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}} \geq \rho, \quad i = 1, \dots, N \quad (5)$$

Note that the optimization problem (4) subject to the constraints (5) can be changed equivalently to (1) subject to the constraints (2). However, in (4) we can find easy that MCVSVM is also a large margin classifier but it employs the Mahalanobis distance metric [19] when calculating the distance from the separating hyperplane  $\mathbf{w}^T \mathbf{x} + b = 0$  to the data point. Therefore, MCVSVM reflect the class distribution through introducing the matrix  $\mathbf{S}_W$  in the Mahalanobis distance metric.

### 2.2. Locality preserving projections (LPP)

Locality preserving projections (LPP) is a linear dimensionality reduction algorithm by feature extraction or projection. It builds an adjacency graph incorporating neighborhood information of the data set. Using the Laplacian graph, LPP then computes a transformation matrix which maps the data points into a subspace. This linear transformation optimally preserves local neighborhood information in a certain sense. The representation map generated by this method may be viewed as a linear discrete approximation to a continuous map that naturally arises from the geometry of the manifold [12].

Let  $N_k(\mathbf{x}_i)$  denotes  $k$  nearest neighbors of node  $i$  and  $G$  denote the adjacency graph of dataset  $\mathbf{X}$  with  $N$  nodes. Here the  $i$ th node corresponds to the data point  $\mathbf{x}_i$ . Nodes  $i$  and  $j$  are connected by an edge if  $i$  is among  $k$  nearest neighbors of  $j$  or  $j$  is among  $k$  nearest neighbors of  $i$ , i.e.  $\mathbf{x}_j \in N_k(\mathbf{x}_i)$  or  $\mathbf{x}_i \in N_k(\mathbf{x}_j)$ . In order to weigh the edges of the adjacency graph  $G$ , we generally need to calculate the weight matrix  $\mathbf{W}$  in different ways. One common choice is weight of the Gaussian kernel as follows:

$$W_{ij} = \begin{cases} \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{t}\right) & \text{if } \mathbf{x}_i \in N_k(\mathbf{x}_j) \text{ or } \mathbf{x}_j \in N_k(\mathbf{x}_i) \\ 0 & \text{other} \end{cases} \quad (6)$$

where  $\|\mathbf{x}\| = (\sum_{i=1}^M \mathbf{x}_i^2)^{1/2}$  is the usual Euclidean ( $L_2$ ) norm in  $\mathbf{R}^M$ ,  $t > 0$  is the Gaussian kernel parameter and can be empirically determined. Note, the weight matrix  $\mathbf{W}$  of the graph  $G$  models the local structure of the data manifold. LPP finds the transformation vector  $\mathbf{w} \in \mathbf{R}^M$  by minimizing the following objective function:

$$\min_{\mathbf{w}} \mathbf{w}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{w} \quad (7)$$

$$\text{s.t. } \mathbf{w}^T \mathbf{D} \mathbf{X} \mathbf{D}^T \mathbf{w} = 1 \quad (8)$$

where  $\mathbf{D}$  is a diagonal matrix and its entries are column (or row, since  $\mathbf{S}$  is symmetric) sum of  $\mathbf{W}$ , i.e.  $D_{ii} = \sum_j W_{ij}$ ,  $\mathbf{L} = \mathbf{D} - \mathbf{W}$  is the Laplacian matrix. The transformation vector  $\mathbf{w}$  that minimizes the objective function is given by the minimum eigenvalue solution to

the generalized eigenvalue problem. Note that the two matrices  $\mathbf{X}\mathbf{L}\mathbf{X}^T$  and  $\mathbf{X}\mathbf{D}\mathbf{X}^T$  are both symmetric and positive semidefinite since the matrix  $\mathbf{L}$  and the diagonal matrix  $\mathbf{D}$  are both symmetric and positive semidefinite [13].

LPP seeks to preserve the intrinsic geometry of the data and local structure by minimizing the above objective function (7). LPP preserves well the intrinsic geometry of the data and local structure and has been successfully applied to face recognition. More detail can be found in [13,14].

### 3. Minimum class locality preserving variance support vector machine

In this section, MCLPV\_SVM will be presented. Firstly, the locality preserving scatter matrix and the locality preserving within-class scatter matrix are defined. Secondly, the optimization problem formulation of the linear MCLPV\_SVM is given by using the locality preserving within-class scatter matrix. Finally, we discuss MCLPV\_SVM in the small sample size case.

#### 3.1. Locality preserving within-class scatter matrix

Before continuing our discussion, here we would like to give the following definitions by using the basic idea of LPP.

**Definition 2.1.** Let  $\mathbf{L}$  be the Laplacian matrix of the dataset  $\mathbf{X}$ , the matrix  $\mathbf{Z} = \mathbf{X}\mathbf{L}\mathbf{X}^T = \mathbf{X}(\mathbf{D} - \mathbf{W})\mathbf{X}^T$  is called the locality preserving scatter matrix.

**Definition 2.2.** Suppose the dataset  $\mathbf{X}$  be separated into two different classes, the matrix

$$\mathbf{Z}_W = \sum_{K=1}^2 \mathbf{Z}_K \quad (9)$$

is called the locality preserving within-class scatter matrix. Here,  $\mathbf{Z}_K$  ( $K=1,2$ ) is the locality preserving scatter matrix of the  $K$ th class  $\mathbf{X}_K$ , i.e.  $\mathbf{Z}_K = \mathbf{X}_K(\mathbf{D}^K - \mathbf{W}^K)\mathbf{X}_K^T$ ,  $\mathbf{D}^K$  is a diagonal matrix and  $\mathbf{D}_{ij}^K = \sum_j W_{ij}^K$ ,  $\mathbf{W}^K$  is the weight matrix of the  $K$ th class  $\mathbf{X}_K$  which can be defined as

$$W_{ij}^K = \begin{cases} \exp\left(-\frac{\|\mathbf{x}_{Ki} - \mathbf{x}_{Kj}\|^2}{t}\right) & \text{if } \mathbf{x}_{Ki} \in N_k(\mathbf{x}_{Kj}) \text{ or } \mathbf{x}_{Kj} \in N_k(\mathbf{x}_{Ki}) \\ 0 & \text{other} \end{cases} \quad (10)$$

where  $\mathbf{x}_{Kj}$  refers to the  $j$ th sample point in the  $K$ th class ( $K=1,2$ ).

It is worthwhile to note that the locality preserving within-class scatter matrix  $\mathbf{Z}_W$  is symmetric and positive semidefinite and formally similar to the within-class scatter matrix  $\mathbf{S}_W$ . However, the locality preserving within-class scatter matrix  $\mathbf{Z}_W$  reflects the intrinsic geometry and local structure of the data. Furthermore,  $\mathbf{Z}_W$  is different from the objective function of LPP when defining the weight matrix. In  $\mathbf{Z}_W$  the weight matrix uses the class labels and carries not only the intrinsic manifold structure information of the data space but also the discriminating information, but in LPP the class labels are not used and the discriminating information is not carried.

#### 3.2. Minimum class locality preserving variance support vector machine

In this subsection, we will present MCLPV\_SVM formulation. Assuming that the data is linearly separable, similar to the optimization problem (4) of MCVSVM, we define the optimization

problem of MCLPV\_SVM as follows:

$$\max_{\mathbf{w}, b} \rho \quad (11)$$

$$\text{s.t. } \frac{z_i(\mathbf{w}^T \mathbf{x}_i + b)}{\sqrt{\mathbf{w}^T \mathbf{Z}_W \mathbf{w}}} \geq \rho, \quad i = 1, \dots, N. \quad (12)$$

Compared with MCVSVM, MCLPV\_SVM incorporate the intrinsic geometry and local structure of the data in the similar way. Here, the essential difference between SVM, MCVSVM and MCLPV\_SVM can be clearly seen: when we define the distance for the sample point  $\mathbf{x}_i$  to the decision hyperplane  $\mathbf{w}^T \mathbf{x} + b = 0$ , SVM is  $z_i(\mathbf{w}^T \mathbf{x}_i + b) / \sqrt{\mathbf{w}^T \mathbf{w}}$ , MCVSVM is  $z_i(\mathbf{w}^T \mathbf{x}_i + b) / \sqrt{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$ , and MCLPV\_SVM is  $z_i(\mathbf{w}^T \mathbf{x}_i + b) / \sqrt{\mathbf{w}^T \mathbf{Z}_W \mathbf{w}}$ . However, SVM, MCVSVM and MCLPV\_SVM are all the large margin classifier. SVM employ directly the Euclidean distance metric, but both MCVSVM and MCLPV\_SVM employ the Mahalanobis distance metric. MCVSVM incorporate the distribution of the classes by using the within-class covariance matrix  $\mathbf{S}_W$  in the distance metric, while MCLPV\_SVM do the intrinsic geometry of the data and local structure by using the locality preserving within-class covariance matrix  $\mathbf{Z}_W$ . Note, the optimization problem (11) subject to the constraints (12) can be changed equivalently to

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \mathbf{w}^T \mathbf{Z}_W \mathbf{w} \quad (13)$$

$$\text{s.t. } z_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \quad i = 1, \dots, N \quad (14)$$

Fig. 1 describes the decision hyperplanes of SVM, MCVSVM, and MCLPV\_SVM on an artificial dataset. As can be seen from the case illustrated in Fig. 1, the MCLPV\_SVM decision hyperplane reflects the intrinsic manifold structure of the data and shows it is more reasonable. The MCSVM decision hyperplane reflects the average information of the class distribution and the SVM decision hyperplane does neither the intrinsic manifold structure of the data, nor the class distribution.

In the case where the training vectors are not linearly separable, the optimum decision hyperplane is found by using the soft margin [4] formulation and solving the following optimization problem:

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \mathbf{w}^T \mathbf{Z}_W \mathbf{w} + C \sum_{i=1}^N \xi_i \quad (15)$$

$$\text{s.t. } z_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i = 1, \dots, N. \quad (16)$$

where  $\xi_i$  denotes the non-negative slack variables for data point  $\mathbf{x}_i$ ,  $C$  is a given constant that defines the cost of the errors after the classification.  $C$  is also called regularization parameter. Larger values of  $C$  correspond to higher penalty assigned to errors. The linearly separable case can be achieved when choosing  $C = \infty$ . Obviously, (15) is a quadratic programming problem. As in SVM and MCVSVM, we can transform this problem into its corresponding dual problem as follows. The primal Lagrangian is

$$L(\mathbf{w}, b, \alpha, \beta, \xi) = \frac{1}{2} \mathbf{w}^T \mathbf{Z}_W \mathbf{w} + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i [z_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 + \xi_i] - \sum_{i=1}^N \beta_i \xi_i \quad (17)$$

where the vectors  $\alpha = [\alpha_1, \dots, \alpha_N]^T$  and  $\beta = [\beta_1, \dots, \beta_N]^T$  ( $\alpha, \beta \in R^N$ ) are the Lagrangian multipliers for the constraints (16). By differentiating with respect to  $\mathbf{w}$ ,  $b$  and  $\beta$  and using the Karush–Kuhn–Tucker (KKT) conditions [20], the following holds:

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{Z}_W \mathbf{w} - \sum_{i=1}^N \alpha_i z_i \mathbf{x}_i = 0$$

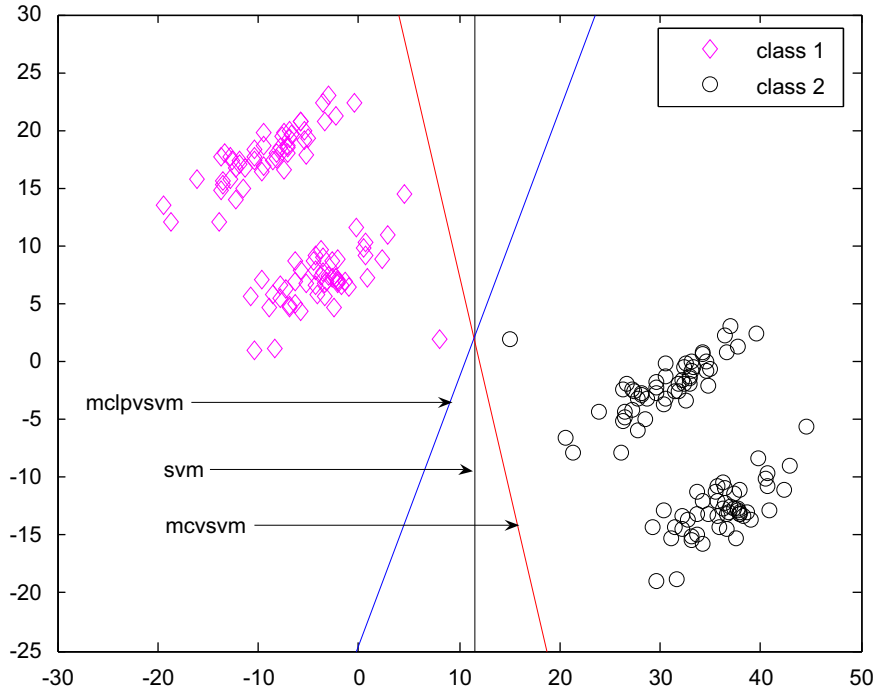


Fig. 1. Illustration of the decision hyperplanes generated by SVM, MCVSVM, and MCLPV\_SVM.

$$\frac{\partial L}{\partial b} = -\sum_{i=1}^N \alpha_i z_i = 0$$

$$\frac{\partial L}{\partial \beta} = C\mathbf{e} - \alpha - \beta = 0$$

$$\alpha_i(z_i(\mathbf{w}^T \mathbf{x}_i + b) - 1) = 0, \quad i = 1, \dots, N \quad (18)$$

where  $\mathbf{e}$  is a  $N$ -dimensional vector of ones, i.e.,  $\mathbf{e} = [1, \dots, 1]^T$ ,  $\mathbf{0} = [0, \dots, 0]^T$ . If the matrix  $\mathbf{Z}_W$  is nonsingular or invertible, we have

$$\mathbf{w} = \mathbf{Z}_W^{-1} \sum_{i=1}^N \alpha_i z_i \mathbf{x}_i \quad (19)$$

By replacing (19) into (17) and using the KKT conditions, the constraint optimization problem (15) is reformulated to the Wolf dual problem

$$\max_{\alpha} -\frac{1}{2} \alpha^T \mathbf{H} \alpha + \mathbf{e}^T \alpha \quad (20)$$

$$\text{s.t.} \quad \sum_{i=1}^N \alpha_i z_i = 0, \quad 0 \leq \alpha_i \leq C, \quad i = 1, \dots, N. \quad (21)$$

where  $\mathbf{H}_{ij} = z_i z_j \mathbf{x}_i^T \mathbf{Z}_W^{-1} \mathbf{x}_j$ . Suppose  $\alpha^* = [\alpha_1^*, \dots, \alpha_N^*]^T$  can be used to solve the above optimization problem, then the optimal weight vector

$$\mathbf{w}^* = \mathbf{Z}_W^{-1} \sum_{i=1}^N z_i \alpha_i^* \mathbf{x}_i \quad (22)$$

If  $0 < \alpha_i^* < C$ , the corresponding data point  $\mathbf{x}_i$  can be called a support vector. The optimal threshold  $b^*$  can be found by exploiting the fact that for all support vectors  $\mathbf{x}_i$  their corresponding slack variables are zero, according to the KKT condition (18). However, averaging over all support vectors yields usually a numerically stable solution. We can calculate the

optimal threshold  $b^*$  as follows:

$$b^* = \frac{1}{N_{SV}} \sum_{i=1}^{N_{SV}} (z_i - z_i \sum_{j=1}^N \alpha_j^* \mathbf{H}_{ij}), \quad \mathbf{x}_i \in D_{SV} \quad (23)$$

where  $D_{SV}$  consist of all support vectors and  $N_{SV}$  is the number of the support vectors. So, the corresponding decision function of MCLPV\_SVM will be

$$f(\mathbf{x}) = \text{sgn} \left( \sum_{i=1}^N z_i \alpha_i^* (\mathbf{x}_i^T \mathbf{Z}^{-1} \mathbf{x}) + b^* \right) \quad (24)$$

### 3.3. MCLPV\_SVM for small sample size problem

Similar to MCVSVM, the major drawback of MCLPV\_SVM is that it may encounter the so-called small sample size problem [18]. The small sample size problem occurs whenever the number of samples is smaller than the dimensionality of the samples. In this case, the locality preserving within-class scatter matrix  $\mathbf{Z}_W$  and the within-class scatter matrix  $\mathbf{S}_W$  which are both  $M \times M$  matrix will be singular. To deal with the problem, through dimensionality reduction using PCA, the optimization problem of MCVSVM is reformulated into an equivalent one in a lower dimensional space, where the within-class scatter matrix  $\mathbf{S}_W$  is nonsingular and MCVSVM optimization problem can be efficiently solved. For MCLPV\_SVM, we employ the same way, i.e. the data in the sample space is first transformed to a low-dimension space where the matrix  $\mathbf{Z}_W$  is nonsingular through dimensionality reduction using PCA and then MCLPV\_SVM is applied in the transformed space. In PCA, the total scatter matrix be defined as

$$\mathbf{S}_t = \sum_{\mathbf{x} \in \mathbf{X}} (\mathbf{x} - \mathbf{u})(\mathbf{x} - \mathbf{u})^T \quad (25)$$

where  $\mathbf{u} = (1/N) \sum_{\mathbf{x} \in \mathbf{X}} \mathbf{x}$  is the total sample mean vector. Let  $\Psi$  and  $\Pi$  be the complementary dimensional spaces spanned by the orthonormal eigenvectors of  $\mathbf{S}_t$  that correspond to nonzero eigenvalues and to zero eigenvalues, respectively. For MCLPV\_SVM, similar to MCVSVM, we have the following theorem.

**Theorem.** Let  $\mathbf{w} = \mathbf{v} + \gamma$  where  $\mathbf{w} \in \mathbf{R}^M$ ,  $\mathbf{v} \in \Psi$ ,  $\gamma \in \Pi$ , thus the optimization problem (15) subject to the constraints (16) of MCLPV\_SVM is equivalent to

$$\min_{\mathbf{v}, \mathbf{b}, \xi} \frac{1}{2} \mathbf{v}^T \mathbf{Z}_W \mathbf{v} + C \sum_{i=1}^N \xi_i \quad (26)$$

$$\text{s.t. } z_i(\mathbf{v}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i = 1, \dots, N \quad (27)$$

A proof of the above Theorem can be found in Appendix.

According to the above theorem, we can conclude that the optimization problem of MCLPV\_SVM can be derived from  $\Psi$  without any loss of the optimal discriminatory information. Suppose  $\mathbf{S}_i$  has  $m$  nonzero eigenvectors and let the column vectors of  $\mathbf{P}$  are eigenvectors corresponding to nonzero eigenvectors of  $\mathbf{S}_i$ , by linear algebra theory,  $\Psi$  is isomorphic to  $m$ -dimensional Euclidean space  $\mathbf{R}^m$  [21,22]. And the corresponding isomorphic mapping is

$$\mathbf{v} = \mathbf{P}\eta, \quad \mathbf{v} \in \Psi, \quad \eta \in \mathbf{R}^m \quad (28)$$

Thus, in  $\mathbf{R}^m$  the optimization problem of MCLPV\_SVM can be written as

$$\min_{\eta, \mathbf{b}, \xi} \frac{1}{2} \eta^T \tilde{\mathbf{Z}}_W \eta + C \sum_{i=1}^N \xi_i \quad (29)$$

$$\text{s.t. } z_i(\eta^T \mathbf{y}_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i = 1, \dots, N \quad (30)$$

where  $\eta \in \mathbf{R}^m$ ,  $\mathbf{y}_i = \mathbf{P}^T \mathbf{x}_i$  and  $\tilde{\mathbf{Z}}_W = \mathbf{P}^T \mathbf{Z}_W \mathbf{P}$ . However, in  $\mathbf{R}^m$  the locality preserving within-class scatter matrix  $\tilde{\mathbf{Z}}_W$  may be still singular. In order to find the MCLPV\_SVM hyperplane, we can transform the data into a lower dimension space through using PCA. So, the column vectors of  $\mathbf{P}$  can be eigenvectors corresponding to the largest  $\tilde{m}$  nonzero eigenvectors of  $\mathbf{S}_i$ . Here  $\tilde{m}$  must be small enough to make  $\tilde{\mathbf{Z}}_W = \mathbf{P}^T \mathbf{Z}_W \mathbf{P}$  nonsingular. Suppose  $\{\eta^*, \mathbf{b}^*, \xi^*\}$  solves the above optimization problem, the decision function is

$$f(\mathbf{x}) = \text{sgn}(\eta^{*T} \mathbf{P}^T \mathbf{x} + b^*) \quad (31)$$

#### 4. MCLPV\_SVM in the nonlinear case

In the previous discussion, the derived decision function (or hyperplane) is derived in a linear form. In order to handle with nonlinear classification problems, we can seek to use the kernelization trick [1] to map the  $M$ -dimensional data points into a high-dimensional feature space, where a linear hyperplane corresponds to a nonlinear hyperplane in the original space. However, MCVSVM or MCLPV\_SVM could not directly get the hyperplane because the within-class scatter matrix  $\mathbf{S}_W^\phi$  or the locality preserving within-class scatter matrix  $\mathbf{Z}_W^\phi$  is generally singular in the feature space. In order to overcome the difficulty, MCVSVM employ PCA to transform the data in the feature space into a low-dimension space where the matrix  $\mathbf{S}_W^\phi$  is nonsingular, and then the linear MCVSVM algorithm is applied in the low-dimension space. In [7] the author pointed out that this process is in nature to transform the data in the sample space using KPCA into a new space where the matrix  $\mathbf{S}_W^\phi$  is nonsingular, and then the linear MCVSVM algorithm is applied in the space.

Similarly, for the nonlinear MCLPV\_SVM, we can first transform the data in the original space into a new space where the matrix  $\mathbf{Z}_W^\phi$  is nonsingular using KPCA, and the linear MCLPV\_SVM is used in the transformed space. In the transformed space, the

optimization problem is reformulated as

$$\min_{\mathbf{w}, \mathbf{b}, \xi} \frac{1}{2} \mathbf{w}^T \tilde{\mathbf{Z}}_W^\phi \mathbf{w} + C \sum_{i=1}^N \xi_i \quad (32)$$

$$\text{s.t. } z_i(\mathbf{w}^T \mathbf{y}_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i = 1, \dots, N \quad (33)$$

where  $\mathbf{y}_i$  is the corresponding projected sample point in the transformed space using KPCA for the sample point  $\mathbf{x}_i$  in the original space,  $\tilde{\mathbf{Z}}_W^\phi$  is the locality preserving within-class scatter matrix in the transformed space. Assume the optimum case of (32) is  $\{\mathbf{w}^*, \mathbf{b}^*, \xi^*\}$ , the decision function can be written as

$$f(\mathbf{x}) = \text{sgn}(\mathbf{w}^{*T} \mathbf{y} + b^*) \quad (34)$$

where  $\mathbf{y}$  is the projected vector in the transformed space using the KPCA transform for the data point  $\mathbf{x}$  which is unlabeled.

## 5. Discussion

### 5.1. Connection to SVM

When the locality preserving within-class scatter matrix  $\mathbf{Z}_W$  is nonsingular, let  $\mathbf{P} = (\mathbf{Z}_W)^{-1/2}$ , we have  $\mathbf{P}^T = ((\mathbf{Z}_W)^{-1/2})^T = (\mathbf{Z}_W)^{-1/2} = \mathbf{P}$  since  $\mathbf{Z}_W$  is invertible and symmetric. Therefore, (15) subject to the constraints (16) can be written as

$$\min_{\mathbf{v}, \mathbf{b}, \xi} \mathbf{v}^T \mathbf{v} + C \sum_{i=1}^N \xi_i \quad (35)$$

$$\text{s.t. } z_i(\mathbf{v}^T \mathbf{y}_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i = 1, \dots, N \quad (36)$$

where  $\mathbf{y}_i = \mathbf{P}^T \mathbf{x}_i$ ,  $\mathbf{v} = \mathbf{P} \mathbf{w}$ . This is the primal optimization problem of SVM. It can be shown that the MCLPV\_SVM formulation can be solved using the existing SVM software, thus making the solution easy to be computed. Assume the optimum case of (35) is  $\{\mathbf{v}^*, \mathbf{b}^*, \xi^*\}$ , thus, the decision function can be written as

$$f(\mathbf{x}) = \text{sgn}(\mathbf{v}^{*T} \mathbf{P}^T \mathbf{x} + b^*) \quad (37)$$

### 5.2. Connection to MCVSVM

In MCLPV\_SVM, if the weight matrix  $\mathbf{W}^K$  in the  $K$ th ( $K=1,2$ ) class  $\mathbf{X}_K$  is defined as follows:

$$\mathbf{W}_{ij}^K = \frac{1}{N_K} \quad (38)$$

according to [13], the locality preserving scatter matrix  $\mathbf{Z}_K = \mathbf{X}_K \mathbf{L}_K (\mathbf{X}_K)^T$  is equivalent to the scatter matrix  $\mathbf{S}_K$ . So, we have

$$\mathbf{Z}_W = \sum_{K=1}^2 \mathbf{Z}_K = \sum_{K=1}^2 \mathbf{S}_K = \mathbf{S}_W \quad (39)$$

This suggests that in this case the locality preserving within-class scatter matrix  $\mathbf{Z}_W$  is equivalent to the within-class scatter matrix  $\mathbf{S}_W$ . So, it can be easily concluded that MCVSVM can be obtained from the special affinity (or adjacency) matrix in MCLPV\_SVM. MCLPV\_SVM can be also seen as the generalized version of MCVSVM since its affinity (or adjacency) matrix is more general.

### 5.3. Constructing the weight matrix

Generally, we construct the weight matrix  $\mathbf{W}_K$  in the  $K$ th ( $K=1,2$ ) class according to (10). However, the neighbor parameter  $k$  must firstly be determined when constructing the adjacency graph of the  $K$ th ( $K=1,2$ ) class. Furthermore, even though two data points in the same class are close to each other but not connected in the adjacency graph, the weight of their edge is 0. This means



that they do not have similarity and are seen as data points which have different class labels. So, we suggest that here the weight matrix can be constructed as follows:

$$W_{ij}^k = \exp\left(-\frac{\|\mathbf{x}_i^k - \mathbf{x}_j^k\|^2}{t}\right) \quad (40)$$

This means that two data points in the same class have always similarity. If the above way does not severely influence the classification accuracy, it will obviously bring two aspects of benefit: reduction in the parameter and time. This is in that we do not require to construct the adjacency graph.

## 6. Experiments

In this section, experimental results will be reported. In the first experiment we investigate the influence of parameters on MCLPV\_SVM performance when using (10) to construct the weight matrix. In the second experiment we report the evaluation results of the linear MCLPV\_SVM in comparison with SVM and MCVSVM on several datasets of the well-known University of California at Irvine (UCI) Repository of machine learning databases [23]. At last, we report the experimental results on a face dataset in the small sample size case and the nonlinear case. For the SVM algorithm, we employ LIBSVM [24].

At present, choosing the algorithms parameters for the kernel methods such as SVM is an open problem. In general, the algorithms parameters are manually set. In order to evaluate the performance, a strategy, as is pointed out and adopted in [16], is that a set of the parameters is first given and then the best cross-validation mean rate among the set is used to estimate the generalized accuracy. In this work we adopted this strategy.

### 6.1. Parameter influence on performance of MCLPV\_SVM

Compared with SVM and MCVSVM on the single parameter, however, MCLPV\_SVM introduce two additional parameters—the heat parameter  $t$  and the neighborhood parameter  $k$  as shown in (10). In order to investigate the influence of the parameters on classification accuracy in MCLPV\_SVM, we test MCLPV\_SVM with different parameters on the heart dataset [23]. In this experiment, we use 40% data points of the heart dataset for training and the rest for testing. The data have been normalized (that is, mean zero and standard deviation one). For SVM, MCVSVM and MCLPV\_SVM, the regularization parameter  $C$  is selected from the set {0.001, 0.01, 0.1, 1, 10, 100}. In addition, in MCLPV\_SVM, the heat kernel parameter  $t$  and the neighborhood parameter  $k$  are, respectively, selected from the set {0.5, 1.0, 1.5, 2.0, 2.5} and {3, 6, 9, 12, 15, 30, 60, 90}.

The classification accuracy of SVM and MCVSVM on different regularization parameter  $C$  is reported in Table 1. Table 2 presents the classification accuracy of MCLPV\_SVM on different regularization parameter  $C$  and different neighborhood parameter  $k$  when constructing the weight matrix according to (10). Note, here the heat kernel parameter  $t$  vary in its set for each pair  $(C, k)$ , and we give the best classification accuracy on different  $t$ . When constructing the weight matrix according to (40) the

**Table 2**

Classification accuracy (%) on different  $C$  and  $k$  in MCLPV\_SVM.

	$C=0.001$	$C=0.01$	$C=0.1$	$C=1$	$C=10$	$C=100$
$k=3$	87.037 $t=0.5$	90.741 $t=1$	90.741 $t=2.5$	85.185 $t=0.5$	85.185 $t=1$	85.185 $t=1$
$k=6$	87.037 $t=0.5$	90.741 $t=1$	90.741 $t=2.5$	85.185 $t=0.5$	87.037 $t=0.5$	85.185 $t=1$
$k=9$	87.037 $t=0.5$	88.889 $t=1.5$	92.593 $t=2.5$	87.037 $t=0.5$	85.185 $t=1$	85.185 $t=1$
$k=12$	87.037 $t=0.5$	88.889 $t=1.5$	92.593 $t=2.5$	87.037 $t=2.5$	85.185 $t=1$	85.185 $t=1$
$k=15$	87.037 $t=0.5$	88.889 $t=1.5$	90.741 $t=2.5$	87.037 $t=2.5$	87.037 $t=0.5$	85.185 $t=1$
$k=30$	87.037 $t=0.5$	88.889 $t=1.5$	90.741 $t=2.5$	87.037 $t=2.5$	85.185 $t=1$	85.185 $t=1$
$k=60$	87.037 $t=0.5$	88.889 $t=1.5$	90.741 $t=2.5$	87.037 $t=2.5$	87.037 $t=0.5$	85.185 $t=1$
$k=90$	87.037 $t=0.5$	88.889 $t=1.5$	90.741 $t=2.5$	87.037 $t=2.5$	85.185 $t=1$	85.185 $t=1$

**Table 3**

Classification accuracy (%) on different  $C$  and  $t$  in MCLPV\_SVM.

	$C=0.001$	$C=0.01$	$C=0.1$	$C=1$	$C=10$	$C=100$
$t=0.5$	87.037	87.037	87.037	81.481	87.037	85.185
$t=1$	77.778	88.889	87.037	85.185	85.185	85.185
$t=1.5$	55.556	88.889	88.889	85.185	85.185	85.185
$t=2$	55.556	55.556	90.741	85.185	85.185	85.185
$t=2.5$	55.556	55.556	88.889	87.037	85.185	85.185

classification accuracy on different regularization parameter  $C$  and different heat kernel  $t$  are reported in Table 3.

From the experimental results, it can be found that the choice of  $C$  plays an important role in terms of the accuracy. Furthermore, if the best is selected as the experimental result, MCLPV\_SVM shows better classification accuracy no matter what the weight matrix is constructed according to (10) or (40). In addition, we can find that it is feasible to construct the weight matrix according to (40) and this way does not influence severely the performance in MCLPV\_SVM. However, it reduces the parameters and the running time when employing (40) to construct the weight matrix since we do not need to construct the adjacency graph. Consequently, in the following experiments, we give the classification performance based on (40).

### 6.2. Performance comparison for linear cases

In this subsection, experimental results are presented for several benchmarking datasets from the well-known University of California at Irvine (UCI) Repository of machine learning databases [23]. A summary of the characteristics of the selected datasets are presented in Table 4. In these datasets, two-class cases and multi-class cases are both included. In multi-class cases, one-against-one strategy [16] is used. All data have been normalized before experiment.

The criteria used to estimate the generalized performance is the 5-fold cross validation accuracy [16] on the whole dataset, according to the size of the dataset, in order to ensure good statistical behavior. In 5-fold cross validation test, a dataset is divided into five subsets. Each time, one of the five subsets is used as the test set and all other four subsets are put together to form a training set. This procedure is repeated five times and then the average accuracy or error across all trials is computed. We give here the mean accuracy and standard deviation of the 5-fold cross

**Table 1**

Classification accuracy (%) on different  $C$  in SVM and MCVSVM.

	$C=0.001$	$C=0.01$	$C=0.1$	$C=1$	$C=10$	$C=100$
SVM	59.259	87.037	85.185	85.185	85.185	85.185
MCVSVM	55.556	55.556	55.556	87.037	85.185	85.185

validation. The regularization parameter  $C$  is selected from the set  $\{0.001, 0.01, 0.1, 1, 10, 100\}$ .

Table 5 reports the experimental results of SVM, MCVSVM and MCLPV\_SVM on the selected datasets. Please note, here we give the best results on different parameters. From Table 5, on the whole, it can be found that there is an improvement in the generalized performance of MCLPV\_SVM over SVM and MCVSVM. MCVSVM also has the comparable performance to SVM. These experimental results indicate that the generalized performance in the large margin classifier can be improved when the characteristics of data, especially the intrinsic manifold structure of the data space, are considered.

For a rigorous comparison of SVM, MCVSVM and MCLPV\_SVM, we further performed the paired two-tailed  $t$ -tests [27] on these algorithms. The  $p$ -value of a  $t$ -test represents the probability that two sets of compared samples come from distributions with equal means. The smaller the  $p$ -value, the more significant the difference of the two average values is, and a  $p$ -value of 0.05 is a typical threshold which is considered statistically significant. Table 6 reports the experimental results of the  $t$ -tests. For example, the  $p$ -value of the  $t$ -test when comparing MCVSVM and SVM on the Newthyroid dataset is **0.039301** ( $< 0.05$ ), meaning that MCVSVM performs significantly better than SVM on this dataset at the 0.05 significant level. From Table 6, although MCLPV\_SVM have on the whole better generalized performance, compared with MCVSVM, its improvement in the generalized performance is not significant. However, MCLPV\_SVM

**Table 6**

$P$ -value of  $t$ -test on the selected datasets.

Dataset	MCVSVM/SVM	MCLPV_SVM /SVM	MCLPV_SVM/ MCVSVM
Heart	0.62131	0.47666	0.3739
Breast	–	–	–
Pima	0.86447	0.99724	0.70453
Wdbc	0.18003	0.7163	0.18228
Ionosphere	0.12102	<b>0.04509</b>	<b>0.0010826</b>
Newthyroid	<b>0.039301</b>	<b>0.039301</b>	1
Wine	0.76549	<b>0.03172</b>	<b>0.046973</b>
Iris	<b>0.020484</b>	<b>0.020484</b>	1
Vehicle	0.53438	<b>0.019329</b>	0.67022
Glass	0.99695	<b>0.011586</b>	<b>0.029963</b>

significantly outperforms SVM on six of 10 datasets. Please note, for the Breast dataset, since SVM, MCVSVM, and MCLPV\_SVM always keep the same accuracy in each classification test, so the  $t$ -test cannot be performed for these algorithms.

### 6.3. Performance comparison for small size sample case and nonlinear cases

In the subsection, we report the experimental results in the small sample size case and the nonlinear case. In this study, the Yale face databases [25] were tested. The Yale face database was constructed at the Yale Center for Computational Vision and Control. It contains 165 gray scale images of 15 individuals. The images demonstrate variations in lighting condition, facial expression (normal, happy, sad, sleepy, surprised, and wink). In the experiments, preprocessing to locate the faces was applied. Original images were manually aligned (two eyes were aligned at the same position), cropped, and then re-sized to  $32 \times 32$  pixels, with 256 gray levels per pixel. Each image is represented by a 1,024-dimensional vector in image space. More details can be found in [26]. The databases in Matlab format after being preprocessed is available at: <http://www.cs.uiuc.edu/homes/dengcai2/Data/data.html>. Fig. 2 depicts some sample images after being preprocessed.

The typical kernel used in our experiments is the Gaussian kernel, i.e.  $\exp(-(\mathbf{u}-\mathbf{v})^T(\mathbf{u}-\mathbf{v})/2\sigma^2)$ , where  $\sigma$  is the spread of the Gaussian kernel. Here, we let  $\sigma=10$ . As we described previously, in this case both the locality preserving within-class scatter matrix  $\mathbf{Z}_W$  and the within-class scatter matrix  $\mathbf{S}_W$  are singular. So, PCA or KPCA is used to project the data in the original space into a subspace where the locality preserving within-class scatter matrix  $\mathbf{Z}_W$  in MCLPV\_SVM or the within-class scatter matrix  $\mathbf{S}_W$  in MCVSVM is nonsingular. In order to make our experimental results fair, we took the same dimension reduction before all three algorithms run.

The experimental results in the small sample size case are reported in Table 7 and those in the nonlinear case do in Table 8. Note, here we give the experimental results after reducing to different dimensions using PCA or KPCA. Since the training samples are linearly independent, we can project the data onto the  $N_{tr}-2$  dimension space. Here  $N_{tr}$  is the number of the training samples. The experimental results obtained by directly applying SVM are reported in the bottom row of the table. From Tables 7 and 8, it can be found that MCLPV\_SVM outperforms SVM and MCVSVM on the whole. This suggests that the performance can indeed be improved when the intrinsic manifold structure of the data space is taken into full consideration. This characteristic is embodied in MCLPV\_SVM.

Tables 9 and 10 report respectively the experimental results of  $t$ -tests in the linear and nonlinear cases. It can be found that in the linear case MCLPV\_SVM performs significantly better in

**Table 4**

Characteristics of the selected datasets.

Dataset	No. of patterns	No. of features	No. of classes
Heart	270	13	2
Breast	699	9	2
Pima	768	8	2
Wdbc	569	30	2
Ionosphere	351	34	2
Newthyroid	215	5	3
Wine	178	13	3
Iris	150	4	3
Vehicle	846	8	4
Glass	214	9	6

**Table 5**

Mean accuracy (%) and standard deviation of the cross validation on the selected datasets.

Dataset	SVM	MCVSVM	MCLPV_SVM
Heart	$84.815 \pm 0.0428$ $C=0.01$	$85.185 \pm 0.0370$ $C=10$	<b><math>85.556 \pm 0.0318</math></b> $C=1, t=2$
Breast	<b><math>96.855 \pm 0.0209</math></b> $C=1$	<b><math>96.855 \pm 0.0209</math></b> $C=100$	<b><math>96.855 \pm 0.0209</math></b> $C=100, t=0.1$
Pima	<b><math>77.473 \pm 0.0107</math></b> $C=0.01$	$77.341 \pm 0.0174$ $C=10$	$77.471 \pm 0.0155$ $C=10, t=2$
Wdbc	$97.534 \pm 0.00675$ $C=0.1$	$96.484 \pm 0.0097$ $C=100$	<b><math>97.699 \pm 0.0120</math></b> $C=10, t=3$
Ionosphere	$87.167 \pm 0.0316$ $C=10$	$84.889 \pm 0.0459$ $C=100$	<b><math>90.020 \pm 0.0434</math></b> $C=100, t=0.8$
Newthyroid	$96.279 \pm 0.0237$ $C=10$	<b><math>98.140 \pm 0.0271</math></b> $C=100$	<b><math>98.140 \pm 0.0271</math></b> $C=100, t=0.8$
Wine	$93.235 \pm 0.0455$ $C=0.1$	$93.268 \pm 0.0375$ $C=100$	<b><math>95.523 \pm 0.0375</math></b> $C=100, t=5$
Iris	$96.667 \pm 0.0421$ $C=0.1$	<b><math>98.667 \pm 0.0266</math></b> $C=1$	<b><math>98.667 \pm 0.0266</math></b> $C=1, t=5$
Vehicle	$79.787 \pm 0.0142$ $C=10$	$80.498 \pm 0.0184$ $C=1$	$81.086 \pm 0.0155$ $C=100, t=6$
Glass	$60.299 \pm 0.0695$ $C=10$	$60.310 \pm 0.0706$ $C=10$	<b><math>63.577 \pm 0.0655</math></b> $C=100, t=2$



Fig. 2. The sample cropped face image in the Yale face image dataset after being preprocessed.

Table 7

Mean accuracy (%) and standard deviation of the 5-fold cross validation in the linear case on the Yale face image dataset.

Dim	SVM	MCVSVM	MCLPV_SVM
3	60.889 ± 0.0585 C=1	60.000 ± 0.0421 C=100	<b>64.000 ± 0.0388</b> C=100, t=10
6	72.889 ± 0.1226 C=0.1	72.444 ± 0.1075 C=1	<b>76.000 ± 0.1103</b> C=1, t=10
9	74.222 ± 0.0997 C=0.1	73.778 ± 0.1273 C=1	<b>76.667 ± 0.1192</b> C=1, t=15
12	77.556 ± 0.0989 C=0.1	76.889 ± 0.1149 C=1	<b>80.889 ± 0.0935</b> C=1, t=25
15	76.889 ± 0.1026 C=0.1	76.222 ± 0.1105 C=1	<b>78.444 ± 0.0939</b> C=0.1, t=50
$N_{tr}-2$	77.556 ± 0.0989 C=0.1	<b>79.333 ± 0.1103</b> C=1	<b>79.333 ± 0.1103</b> C=1, t=50
All	76.889 ± 0.1026 C=1	–	–

Table 8

Mean accuracy (%) and standard deviation of the 5-fold cross validation in the nonlinear case on the Yale face image dataset.

Dim	SVM	MCVSVM	MCLPV_SVM
3	62.000 ± 0.0884 C=100	62.444 ± 0.0656 C=100	<b>63.778 ± 0.0413</b> C=100, t=5
6	70.667 ± 0.0646 C=100	<b>75.111 ± 0.1296</b> C=1	74.444 ± 0.1169 C=10, t=20
9	76.222 ± 0.1119 C=1	75.778 ± 0.1048 C=100	<b>76.889 ± 0.1259</b> C=1, t=5
12	75.556 ± 0.1307 C=10	77.556 ± 0.1211 C=100	<b>79.556 ± 0.1170</b> C=10, t=5
15	<b>77.333 ± 0.1388</b> C=10	76.000 ± 0.1254 C=100	77.111 ± 0.1125 C=10, t=5
$N_{tr}-2$	76.000 ± 0.1388 C=10	77.333 ± 0.1388 C=100	<b>78.667 ± 0.1240</b> C=10, t=5
All	76.667 ± 0.1333 C=10	–	–

Table 9

P-value of t-test in the linear case on the Yale face image dataset.

Dim	MCVSVM/SVM	MCLPV_SVM/SVM	MCLPV_SVM/MCVSVM
3	0.55426	<b>0.03598</b>	<b>0.012685</b>
6	0.71743	<b>0.03166</b>	<b>0.029925</b>
9	0.82464	<b>0.04937</b>	<b>0.039657</b>
12	0.70405	<b>0.04628</b>	<b>0.01781</b>
15	0.50228	<b>0.04016</b>	<b>0.032046</b>
$N_{tr}-2$	0.077737	0.077737	–

comparison with SVM and MCVSVM, while in the nonlinear case it does not so. This suggests that although in the nonlinear case MCLPV\_SVM has insignificant performance difference in contrast to SVM and MCVSVM, it performs on the whole better in the sense of the average accuracy.

Table 10

P-value of t-test in the linear case on the Yale face image dataset.

Dim	MCVSVM /SVM	MCLPV_SVM /SVM	MCLPV_SVM/MCVSVM
3	0.82464	0.54475	0.3739
6	0.33432	0.33102	0.62116
9	0.71759	0.62124	0.60047
12	0.20799	<b>0.02799</b>	<b>0.02080</b>
15	0.39201	0.22049	0.095771
$N_{tr}-2$	0.3739	<b>0.049301</b>	0.17781

## 7. Conclusion

In this paper, we propose a novel minimum class locality preserving variance support machines (MCLPV\_SVM). Different from SVM and MCVSVM, MCLPV\_SVM take the intrinsic manifold structure of the data space into full consideration and inherit the characteristics of SVM and MCVSVM. In small sample size cases and nonlinear cases, similar to MCVSVM, PCA or KPCA is used to transform the data in original space into a low dimension space where the optimization problem of linear MCLPV\_SVM can be efficiently solved. Experimental results indicate the effectiveness of MCLPV\_SVM by comparing it with SVM and MCVSVM. Although the proposed method here demonstrates our initial attempt to solve small sample size problems and nonlinear classification tasks from a new perspective, we still hope to study its more effective version to deal well with these problems and tasks.

## Acknowledgments

This work is supported by the Hong Kong Polytechnic University Grants (Grant no.Z-08R and G-U296), National Science Foundation of China (Grant nos.60773206, 60704047 and 90820002), 2007 Cultivation Fund of the key Scientific and Technical Innovation Project of Ministry of Education of China, National Key Lab. Of CAD & CG at Zhejiang University, Key Lab. of Computer Information Technologies at JiangSu Province, 2008 Postgraduate student's Research Fund at JiangSu Province.

## Appendix

Proof of Theorem in Section 3.3.

**Proof.** Since  $\gamma \in \Pi$ ,  $\gamma^T S_t \gamma = 0$ , then under the projection  $\gamma$ , for all training vectors  $\mathbf{x}_i$ ,  $\mathbf{x}_j$  with  $\mathbf{x}_i \neq \mathbf{x}_j$  then  $\gamma^T \mathbf{x}_i = \gamma^T \mathbf{x}_j$ . In other words, all the training vectors  $\mathbf{x}_i$  fall in the same point under the projection  $\gamma$ . Thus,  $\gamma^T \mathbf{x}_i = r$  is a constant,  $\forall \mathbf{x}_i$ . So, we have

$$\mathbf{w}^T \mathbf{x}_i = (\mathbf{v} + \gamma)^T \mathbf{x}_i = \mathbf{v}^T \mathbf{x}_i + r \quad (41)$$

For the  $K$ th ( $K=1,2$ ) class  $\mathbf{X}_K$ , according to Definition 2.2, we have

$$\mathbf{Z}_K = \mathbf{X}_K \mathbf{L}_K \mathbf{X}_K^T = \mathbf{X}_K (\mathbf{D}_K - \mathbf{W}_K) \mathbf{X}_K^T \quad (42)$$



where  $\mathbf{W}_K$  is the weight matrix and  $\mathbf{L}_K$  is the Laplacian matrix in the  $K$ th ( $K=1,2$ ) class,  $\mathbf{D}^K$  is a diagonal matrix and  $D_{ij}^K = \sum_j \mathbf{W}_{ij}^K$ . Since  $\gamma^T \mathbf{x}_i = r$ , it follows that

$$\begin{aligned} \mathbf{Z}_K \gamma &= \mathbf{X}_K (\mathbf{D}_K - \mathbf{S}_K) \mathbf{X}_K^T \gamma \\ &= \left( \sum_{i=1}^{N_K} D_{ii}^K \mathbf{x}_{Ki} \mathbf{x}_{Ki}^T - \sum_{i=1}^{N_K} \sum_{j=1}^{N_K} \mathbf{W}_{ij}^K \mathbf{x}_{Ki} \mathbf{x}_{Kj}^T \right) \gamma \\ &= \sum_{i=1}^{N_K} \sum_{j=1}^{N_K} \mathbf{W}_{ij}^K \mathbf{x}_{Ki} \mathbf{x}_{Kj}^T \gamma - \sum_{i=1}^{N_K} \sum_{j=1}^{N_K} \mathbf{W}_{ij}^K \mathbf{x}_{Ki} \mathbf{x}_{Kj}^T \gamma \\ &= r \sum_{i=1}^{N_K} \sum_{j=1}^{N_K} \mathbf{W}_{ij}^K \mathbf{x}_{Ki} - r \sum_{i=1}^{N_K} \sum_{j=1}^{N_K} \mathbf{W}_{ij}^K \mathbf{x}_{Ki} = 0 \end{aligned} \quad (43)$$

where  $\mathbf{x}_{Ki}$  is the  $i$ th sample point in the  $K$ th class,  $N_K$  is the number of samples in the  $K$ th class,  $\mathbf{0}$  is a  $M$ -dimensional vector of zeros. Similarly, we have

$$\gamma^T \mathbf{Z}_K = \mathbf{0}^T \quad (44)$$

$$\gamma^T \mathbf{Z}_K \gamma = 0 \quad (45)$$

According to (44) and (45), we can conclude that

$$\mathbf{v}^T \mathbf{Z}_K \gamma = r \mathbf{v}^T \mathbf{0} = 0 \quad (46)$$

$$\gamma^T \mathbf{Z}_K \mathbf{v} = r \mathbf{0}^T \mathbf{v} = 0 \quad (47)$$

Thus, it follows that

$$\begin{aligned} \mathbf{w}^T \mathbf{Z}_W \mathbf{w} &= (\mathbf{v} + \gamma)^T \left( \sum_{K=1}^2 \mathbf{Z}_K \right) (\mathbf{v} + \gamma) \\ &= \sum_{K=1}^2 (\mathbf{v}^T \mathbf{Z}_K \mathbf{v} + \mathbf{v}^T \mathbf{Z}_K \gamma + \gamma^T \mathbf{Z}_K \mathbf{v} + \gamma^T \mathbf{Z}_K \gamma) \\ &= \sum_{K=1}^2 \mathbf{v}^T \mathbf{Z}_K \mathbf{v} = \mathbf{v}^T \mathbf{Z}_W \mathbf{v} \end{aligned} \quad (48)$$

Now, using the KKT condition  $\alpha^T \mathbf{z} = \sum_{i=1}^N \alpha_i z_i = 0$ , the following holds:

$$\sum_{i=1}^N \alpha_i z_i r = r \sum_{i=1}^N \alpha_i z_i = 0 \quad (49)$$

Using the previous facts the Lagrangian in (17) can be written as

$$\begin{aligned} L(\mathbf{w}, \mathbf{b}, \alpha, \beta, \xi) &= \frac{1}{2} \mathbf{w}^T \mathbf{Z}_W \mathbf{w} + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i [z_i (\mathbf{w}^T \mathbf{x}_i + \mathbf{b}) - 1 + \xi_i] - \sum_{i=1}^N \beta_i \xi_i \\ &= \frac{1}{2} \mathbf{v}^T \mathbf{Z}_W \mathbf{v} + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i [z_i (\mathbf{v}^T \mathbf{x}_i + r + \mathbf{b}) - 1 + \xi_i] - \sum_{i=1}^N \beta_i \xi_i \\ &= \frac{1}{2} \mathbf{v}^T \mathbf{Z}_W \mathbf{v} + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i [z_i (\mathbf{v}^T \mathbf{x}_i + \mathbf{b}) - 1 + \xi_i] - \sum_{i=1}^N \alpha_i z_i r - \sum_{i=1}^N \beta_i \xi_i \\ &= \frac{1}{2} \mathbf{v}^T \mathbf{Z}_W \mathbf{v} + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i [z_i (\mathbf{v}^T \mathbf{x}_i + \mathbf{b}) - 1 + \xi_i] - \sum_{i=1}^N \beta_i \xi_i \end{aligned} \quad (50)$$

Note, this is the primal Lagrangian of (26). Then using the chain rule we can easily prove that

$$\left. \frac{\partial L}{\partial \mathbf{w}} \right|_{\mathbf{w}=\mathbf{w}^*} = \left. \frac{\partial L}{\partial \mathbf{v}} \right|_{\mathbf{v}=\mathbf{v}^*} = \mathbf{0} \Leftrightarrow \mathbf{Z} \mathbf{v}^* - \sum_{i=1}^N \alpha_i z_i \mathbf{x}_i = \mathbf{0} \quad (51)$$

Thus, the decision hyperplane depends only on  $\mathbf{v} \in \Psi$  (an arbitrary vector  $\gamma \in \Pi$  can be chosen). So, the theorem has been proven.  $\square$

## References

- [1] B. Scholkopf, A. Smola, Learning With Kernels, MIT Press, Cambridge, MA, 2002.
- [2] I.W. Tsang, J.T. Kwok, P.M. Cheung, Core machines: fast svm training on very large data sets, Journal of Machine Learning Research 6 (2005) 363–392.
- [3] Z.H. Deng, F.L. Chung, S.T. Wang, FRSD: fast reduced set density estimator using minimal enclosing ball approximation, Pattern Recognition 41 (2008) 1363–1372.
- [4] V. Vapnik, The Nature of Statistical Learning Theory, Springer-Verlag, New York, 1995.
- [5] C.J.C. Burges, A tutorial on support machines for pattern recognition, Data Mining and Knowledge Discovery 2 (2) (1998) 121–167.
- [6] B. Scholkopf, S. Mika, C.J.C. Burges, P. Knirsch, K.R. Muller, G. Ratsch, A.J. Smola, Input space vs. feature space in Kernel-based methods, IEEE Transactions on Neural Network 10 (5) (1999) 1000–1017.
- [7] S. Zafeiriou, A. Tefas, I. Pitas, Minimum class variance support vector machines, IEEE Transactions on Image Processing 16 (10) (2007) 2551–2564.
- [8] K. Fukunaga, Statistical Pattern Recognition, Academic, San Diego, CA, 1990.
- [9] R.O. Duda, P.E. Hart, D.G. Stork, Pattern Classification, second ed., Wiley, New York, 2001.
- [10] K.I. Diamantaras, S.Y. Kung, Principal Component Neural Networks, Wiley, New York, 1996.
- [11] A. Scholkopf, B. Smola, K.R. Muller, Nonlinear component analysis as a Kernel eigenvalue problem, Neural Computation 10 (1998) 1299–1319.
- [12] X. He, P. Niyogi, Locality preserving projections, in: Proceedings of the Conference on Advances in Neural Information Processing Systems, 2003, pp. 585–591.
- [13] X.F. He, S.C. Yan, Y.X. Hu, P. Niyogi, H.J. Zhang, Face recognition using Laplacianfaces, IEEE Transactions on Pattern Analysis and Machine Intelligence 27 (3) (2005) 328–340.
- [14] E. Kokiopoulou, Y. Saad, Orthogonal neighborhood preserving projections: a projection-based dimensionality reduction technique, IEEE Transactions on Pattern Analysis and Machine Intelligence 29 (12) (2007) 2143–2156.
- [15] M. Belkin, P. Niyogi, V. Sindhwani, Manifold regularization: a geometric framework for learning from labeled and unlabeled examples, The Journal of Machine Learning Research 7 (2006) 2399–2434.
- [16] L.G. Abril, C. Angulo, F. Velasco, J.A. Ortega, A note on the bias in SVMs for multiclassification, IEEE Transactions on Neural Networks 19 (4) (2008) 723–725.
- [17] T. Xiong, V. Cherkassky, A combined SVM and LDA approach for classification, in: Proceedings of International Joint Conference on Neural Networks, Montreal, Canada, July 31–August 4, 2005, pp. 1455–1459.
- [18] P. Howland, J. Wang, H. Park, Solving the small sample size problem in face recognition using generalized discriminant analysis, Pattern Recognition 39 (2) (2006) 277–287.
- [19] S.M. Xiang, F.P. Nie, C.S. Zhang, Learning a Mahalanobis distance metric for data clustering and classification, Pattern Recognition 41 (2008) 3600–3612.
- [20] R. Fletcher, Practical Methods of Optimization, second ed., Wiley, New York, 1987.
- [21] J. Yang, J.Y. Yang, Why can LDA be performed in PCA transformed space? Pattern Recognition 36 (2) (2003) 563–566.
- [22] J. Yang, A.F. Frangi, J. Yang, D. Zhang, Z. Jin, KPDA plus LDA: a complete kernel Fisher discriminant framework for feature extraction and recognition, IEEE Transactions on Pattern Analysis and Machine Intelligence 27 (2) (2005) 230–244.
- [23] C. Blake, C. Merz, UCI Repository of machine learning databases. Available: <http://www.ics.uci.edu/~mllearn/MLRepository.html>.
- [24] C.C. Chang, C.J. Lin, LIBSVM: a library for support machines, 2001. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- [25] Yale Univ, Face Database, <http://cvc.yale.edu/projects/yalefaces/yalefaces.html>, 2002.
- [26] D. Cai, X. He, J. Han, H.J. Zhang, Orthogonal Laplacianfaces for face recognition, IEEE Transactions on Image Processing 15 (11) (2006) 3608–3614.
- [27] E. Alpaydin, Introduction to Machine Learning, The MIT Press, Cambridge, MA, USA, 2004.

**About the author**—XIAOMING WANG received B.S. and M.S. degrees in the school of information from JiangNan University in 2001 and 2007, respectively. He is currently a Ph.D. candidate in the same school of the same university. His current research interests include artificial intelligence, pattern recognition, data mining and fuzzy system.

**About the author**—FU-LAI CHUNG received the B.Sc. degree from the University of Manitoba, Canada, in 1987, and the M.Phil. and Ph.D. degrees from the Chinese University of Hong Kong, in 1991 and 1995, respectively. He joined the Department of Computing, Hong Kong Polytechnic University, in 1994, where he is currently an associate professor. He has published widely in the areas of fuzzy systems, neural networks, and pattern recognition. His current research interests include fuzzy data mining, fuzzy neural network modeling, and fuzzy techniques for multimedia applications.

**About the author**—SHITONG WANG received the M.S. degree in computer science from Nanjing University of Aeronautics and Astronautics, China, in 1987. He visited London University and Bristol University in U.K., Hiroshima International University in Japan, Hong Kong University of Science and Technology, Hong Kong Polytechnic University, as a research scientist, for over 5 years. Currently, he is a full professor in The School of Information of Southern Yangtze University, China. His research interests include AI, neuron-fuzzy systems, pattern recognition, and image processing. He has published about 80 papers in international/national journals and has authored seven books. (Tel.: +86 510 8373140; e-mail: wxwangst@yahoo.com.cn).