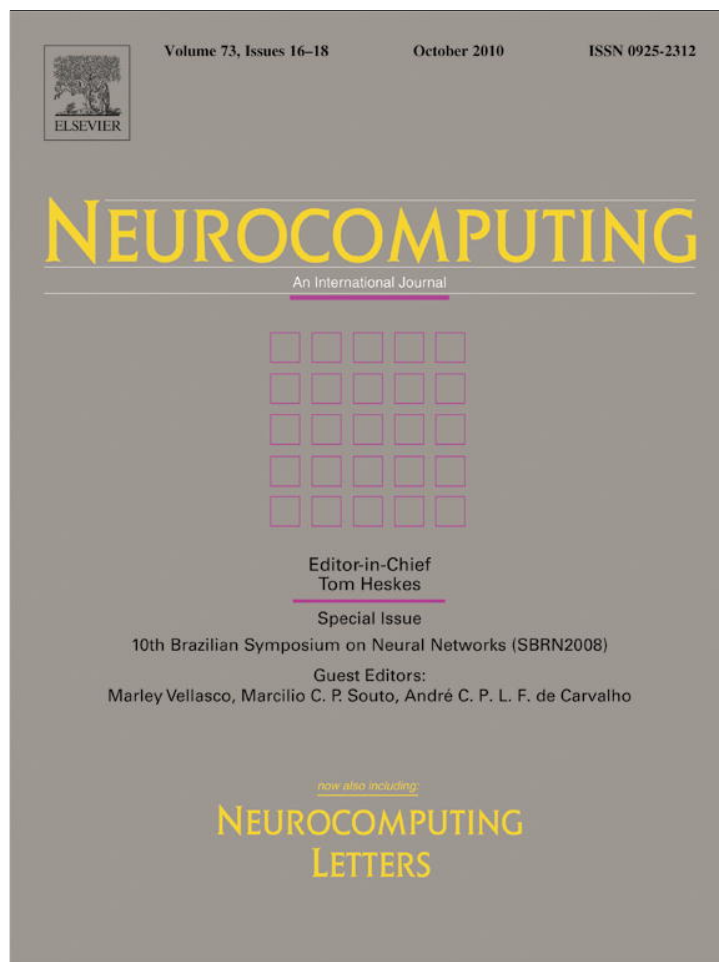


Provided for non-commercial research and education use.
Not for reproduction, distribution or commercial use.



This article appeared in a journal published by Elsevier. The attached copy is furnished to the author for internal non-commercial research and education use, including for instruction at the authors institution and sharing with colleagues.

Other uses, including reproduction and distribution, or selling or licensing copies, or posting to personal, institutional or third party websites are prohibited.

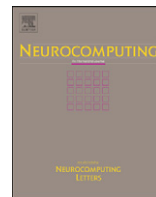
In most cases authors are permitted to post their version of the article (e.g. in Word or Tex form) to their personal website or institutional repository. Authors requiring further information regarding Elsevier's archiving and manuscript policies are encouraged to visit:

<http://www.elsevier.com/copyright>



Contents lists available at ScienceDirect

Neurocomputing

journal homepage: www.elsevier.com/locate/neucom

Large margin classifiers based on affine hulls

Hakan Cevikalp^{a,*}, Bill Triggs^b, Hasan Serhan Yavuz^a, Yalçın Küçük^c, Mahide Küçük^c, Atalay Barkana^d^a Electrical and Electronics Engineering Department of Eskisehir Osmangazi University, Meselik, 26480 Eskisehir, Turkey^b Laboratoire Jean Kuntzmann, Grenoble, France^c Mathematics Department of Anadolu University, Eskisehir, Turkey^d Electrical and Electronics Engineering Department of Anadolu University, Eskisehir, Turkey

ARTICLE INFO

Article history:

Received 29 November 2009

Received in revised form

20 April 2010

Accepted 16 June 2010

Communicated by G.-B. Huang

Available online 29 July 2010

Keywords:

Affine hull

Classification

Convex hull

Kernel methods

Large margin classifier

Quadratic programming

Support vector machines

ABSTRACT

This paper introduces a geometrically inspired large margin classifier that can be a better alternative to the support vector machines (SVMs) for the classification problems with limited number of training samples. In contrast to the SVM classifier, we approximate classes with affine hulls of their class samples rather than convex hulls. For any pair of classes approximated with affine hulls, we introduce two solutions to find the best separating hyperplane between them. In the first proposed formulation, we compute the closest points on the affine hulls of classes and connect these two points with a line segment. The optimal separating hyperplane between the two classes is chosen to be the hyperplane that is orthogonal to the line segment and bisects the line. The second formulation is derived by modifying the ν -SVM formulation. Both formulations are extended to the nonlinear case by using the kernel trick. Based on our findings, we also develop a geometric interpretation of the least squares SVM classifier and show that it is a special case of the proposed method. Multi-class classification problems are dealt with constructing and combining several binary classifiers as in SVM. The experiments on several databases show that the proposed methods work as good as the SVM classifier if not any better.

© 2010 Elsevier B.V. All rights reserved.

1. Introduction

The support vector machine (SVM) classifier is a successful binary classification method that simultaneously minimizes the empirical classification error and maximizes the geometric margin, which is defined as the distance between the separating hyperplane and closest samples from the classes [2,8]. To do so, SVM first approximates each class with a convex hull and finds the closest points in these convex hulls [1]. Then, these two points are connected with a line segment. The hyperplane, orthogonal to the line segment that bisects the line, is chosen to be the separating hyperplane. From the geometrical point of view, in the separable case, the two closest points on the convex hulls determine the separating hyperplane, and the SVM margin is merely equivalent to the minimum distance between the convex hulls that represent classes. However, convex hull approximations tend to be unrealistically tight in high-dimensional spaces since the classes typically extend beyond the convex hulls of their training samples. For example, a convex hull constructed by randomly sampled points from a high-dimensional hypersphere can include only a negligible fraction of the volume of the sphere

even if the chosen samples are well spaced and close to the surface of the sphere [3]. This situation may also be observed when the low-dimensional data samples are mapped to a higher-dimensional feature space through kernel mapping during estimation of the nonlinear decision boundaries between classes.

As opposed to the convex hulls, affine hulls (i.e., spanning linear subspaces that have been shifted to pass through the centroids of the classes) give rather loose approximations to the class regions, because they do not constrain the positions of the training points within the affine subspaces. Therefore, they may be better alternatives to convex hulls for some pattern classification problems especially when the data samples lie in high-dimensional spaces. In the context of classification, affine hulls were first used as global classifiers of isolated word and hand-written digits giving good classification performance [11,14]. In these methods, each class is approximated with an affine hull constructed from its training samples, and the label of a test sample is determined based on the distance to the nearest affine hull. Vincent and Bengio [26] used affine/convex hulls in a local sense by constructing them using k -nearest neighbors of a test sample for classification problems with complex nonlinear decision boundaries. They report that affine hulls usually give higher classification accuracy than convex hulls and that using both models for classification significantly improves the k -nearest neighbor classification performance [26]. We extended local linear affine/convex hull classifiers to the nonlinear case in [4].

* Corresponding author.

E-mail addresses: hakan.cevikalp@gmail.com (H. Cevikalp), Bill.Triggs@imag.fr (B. Triggs), hsyavuz@ogu.edu.tr (H.S. Yavuz), ykucuk@anadolu.edu.tr (Y. Küçük), mkucuk@anadolu.edu.tr (M. Küçük), atalaybarkan@anadolu.edu.tr (A. Barkana).

More recently, we compared different convex class models for high-dimensional classification problems, then found that affine hull approximations are typically more accurate than convex hull approximations [3]. These results are not surprising due to the fact that high-dimensional approximations tend to be simple: For a fixed sample size, the amount of geometric details that can be resolved usually decreases rapidly as the dimensionality increases. Therefore, affine hulls tend to be better models for high-dimensional data approximations.

Besides the classification, approximations based on affine hulls have also been used for dimensionality reduction. Mixtures of principal component analyzers [12] use local affine hulls to estimate nonlinear data manifolds. Similarly, locally linear embedding [18] approximates the nonlinear structure of the high-dimensional data by exploiting local affine hull reconstructions. Verbeek [25] combined several locally valid affine hulls to obtain a global nonlinear mapping between the high-dimensional sample space and low-dimensional manifold. Many applications of affine hulls in the context of classification and dimensionality reduction can be attributed in part to their simplicity and computational efficiency. Finding distances from test samples to affine hulls requires only simple linear algebra. On the other hand, computing distances to nonlinear complex models can be problematic. Even if the models are restricted to being convex hulls, distance computations require the solution of a quadratic optimization problem.

The classification methods using affine hulls or other convex sets described above are “nearest convex model” classifiers and they are instance-based in nature. In other words, decision boundaries are not explicitly created during a training phase. Instead, the decision boundaries remain implicit, and new examples are classified online based on the distances to the nearest convex class models. This paper investigates an alternative “margin between convex model” strategy that is based on explicitly building maximum margin separators between pairs of affine hulls. As a first example of the power of this approach, note that the SVM itself is the maximum margin separator between the convex hulls of the training samples of the two classes. One motivation for replacing nearest-convex-model approaches with margin-based ones is that for all of the above cited nearest-convex-model classifiers, the decision boundaries (surfaces equidistant from the two convex models) are generically at least quadratic or piecewise quadratic in complexity. For example, for affine hulls they are generically hyperboloids. Such decision boundaries are more flexible than linear ones, but in high dimensions when the training data are scarce this may lead to overfitting, thus damaging generalization to unseen examples. Linear margin-based approaches have fewer degrees of freedom, so they are typically less sensitive to the precise arrangement of the training samples. For example, for an SVM classifier, motions of the SVM support vectors parallel to the SVM decision surface do not alter the margin and hence do not invalidate the classifier, whereas they do typically change piecewise quadratic decision surface of the equivalent nearest convex hull classifier. Another motivation for studying margin between affine hulls approach is their potential flexibility and compactness. In linear case, affine models allow each class to be fitted individually and represented compactly, following which the linear separator between any two classes can be found quickly by simple linear algebra.

In our preliminary work [5] we showed how to construct maximum margin classifier that separates linear affine hulls. Another study addressing the same problem was independently given in [29]. Here we extend the method such that it can be used when the class samples lie on nonlinear manifolds that cannot be modeled with linear affine hulls. To this end, we map the samples in each class into a much higher-dimensional feature space

through kernel mapping and then construct the linear affine hulls in the mapped space. Since the problem is cast in a much higher-dimensional space, it is more likely that class regions can now be approximated with linear affine hulls. Although the constructed linear affine hulls in the mapped space correspond to nonlinear manifolds in the input space, finding the maximum separating hyperplane between these nonlinear manifolds is still straightforward because of their linear nature in the mapped space. In case of outliers, to allow soft margin solutions, we first reduce affine hulls in order to alleviate the effects of those outliers and then search for the best separating hyperplane between these reduced robust models.

The rest of the paper is organized as follows: In Section 2, we introduce the proposed method. Section 3 describes the experimental results. Concluding remarks are given in Section 4.

2. Method

Consider a binary classification problem with the training data given in the form $\{\mathbf{x}_i, y_i\}$, $i=1, \dots, n$, $y_i \in \{-1, +1\}$, $\mathbf{x}_i \in \mathbb{R}^d$. To separate classes, SVM classifier finds a separating hyperplane that maximizes the margin, which is defined as the distance between the hyperplane and closest samples from the classes. To do so, SVM first approximates each class with a convex hull [1]. A convex hull consists of all points that can be written as a convex combination of the points in the original set, and a convex combination of points is a linear combination of data points where all coefficients are nonnegative and sum up to 1. More formally, the convex hull of samples $\{\mathbf{x}_i\}_{i=1, \dots, n}$ can be written as

$$H^{\text{convex}} = \left\{ \mathbf{x} = \sum_{i=1}^n \alpha_i \mathbf{x}_i \mid \sum_{i=1}^n \alpha_i = 1, \alpha_i \geq 0 \right\}. \quad (1)$$

Convex hulls of two classes are illustrated in Fig. 1. Following this approximation, SVM finds the closest points in these convex hulls. Then, these two points are connected with a line segment. The plane, orthogonal to the line segment that bisects the line, is selected to be the separating hyperplane as shown in Fig. 1.

In contrast to the SVM classifier, the proposed method approximates each class (positive and negative classes) with an affine hull of its training samples. An affine hull of a class is the smallest affine subspace containing them. This is an unbounded, and hence typically rather loose model for each class, thus affine hull modeling can be a better choice than convex hull modeling

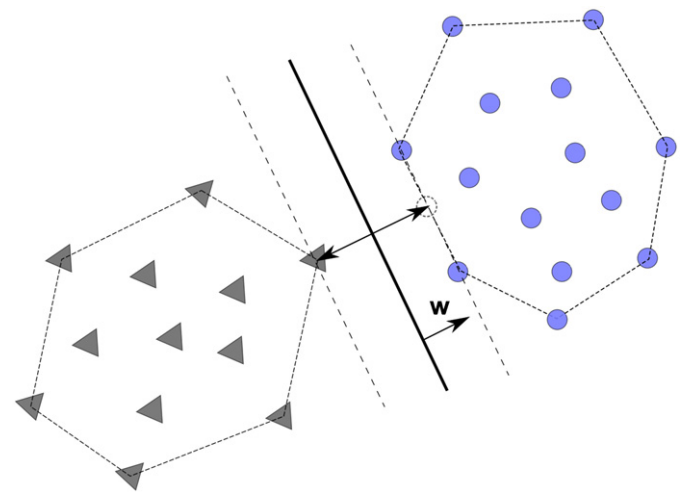


Fig. 1. Two closest points on the convex hulls determine the separating hyperplane.

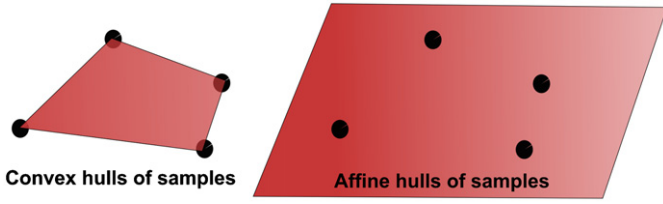


Fig. 2. Comparison of convex and affine hulls of samples.

for high-dimensional data. Affine and convex hulls of four samples are illustrated in Fig. 2. The affine hull of samples $\{\mathbf{x}_i\}_{i=1,\dots,n}$ contains all points of the form $\sum_{i=1}^n \alpha_i \mathbf{x}_i$ with $\sum_{i=1}^n \alpha_i = 1$. More formally affine hull of a class with samples $\{\mathbf{x}_i\}_{i=1,\dots,n}$ can be written as

$$H^{aff} = \left\{ \mathbf{x} = \sum_{i=1}^n \alpha_i \mathbf{x}_i \mid \sum_{i=1}^n \alpha_i = 1 \right\}. \quad (2)$$

Our goal is to find the maximum margin linear separating hyperplane between affine hulls of classes. The points \mathbf{x} which lie on the separating hyperplane satisfy $\langle \mathbf{w}, \mathbf{x} \rangle + b = 0$, where \mathbf{w} is the normal of the separating hyperplane, $|b|/\|\mathbf{w}\|$ is the perpendicular distance from the hyperplane to the origin, and $\|\mathbf{w}\|$ is the Euclidean norm of \mathbf{w} . For any separating hyperplane, all points \mathbf{x}_i in the positive class satisfy $\langle \mathbf{w}, \mathbf{x}_i \rangle + b > 0$ and all points \mathbf{x}_i in the negative class satisfy $\langle \mathbf{w}, \mathbf{x}_i \rangle + b < 0$, so that $y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) > 0$ for all training data points. Finding the best separating hyperplane between affine hulls can be solved by computing the closest points on them. The optimal separating hyperplane will be the one that bisects perpendicularly the line segment connecting the closest points as in SVM classifier. The offset (also called threshold), b , can be chosen as the distance from the origin to the point halfway between the closest points along the normal \mathbf{w} . Once the best separating hyperplane is determined, a new sample \mathbf{x} is classified based on the sign of the decision function, $f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$.

Next, we will first show how to find the best separating hyperplane for linearly separable affine hulls and then extend the idea for inseparable case. After that, we explain kernelization process. This is followed by introducing a second equivalent formulation based on a variation of ν -SVM [20]. Lastly we show the relation between the proposed method and the least squares SVM (LS-SVM) [24,23] and derive a geometric intuition for LS-SVM.

2.1. Linearly separable case

Suppose that affine hulls belonging to the positive and negative classes are linearly separable. The affine hulls of two classes do not intersect, i.e., they are linearly separable, if the affine combinations of their samples satisfy the rule

$$\sum_{i:y_i=+1} \alpha_i \mathbf{x}_i \neq \sum_{j:y_j=-1} \alpha_j \mathbf{x}_j \quad \text{for} \quad \sum_{i:y_i=+1} \alpha_i = \sum_{j:y_j=-1} \alpha_j = 1. \quad (3)$$

It should be noted that linear separability of data points does not necessarily guarantee the separability of corresponding affine hulls of classes. For linearly separable case, it is more convenient to write an affine hull as

$$\{\mathbf{x} = \mathbf{U}\mathbf{v} + \boldsymbol{\mu} \mid \mathbf{v} \in \mathbb{R}^l\}, \quad (4)$$

where $\boldsymbol{\mu} = (1/n)\sum_i \mathbf{x}_i$ is the mean of the samples (or any other reference point in the hull) and \mathbf{U} is an orthonormal basis for the directions spanned by the affine subspace. The vector \mathbf{v} contains the reduced coordinates of the point within the subspace, expressed with respect to the basis \mathbf{U} . Numerically, \mathbf{U} can be

found as the \mathbf{U} -matrix of the “thin” singular value decomposition (SVD) of $[\mathbf{x}_1 - \boldsymbol{\mu}, \dots, \mathbf{x}_n - \boldsymbol{\mu}]$. Here, “thin” indicates that we take only the columns of \mathbf{U} corresponding to “significantly non-zero” singular values λ_k ; l is the number of such non-zero singular values. This subspace estimation process is essentially orthogonal least squares fitting. Discarding near-zero singular values corresponds to discarding directions that appear to be predominantly “noise”. As an alternative, samples can be fitted with some other more robust subspace estimation processes such as L1 norm based subspace fitting procedures described in [10,13]. But, we will consider only the least squares fitting (L2 norm) in this study.

Now suppose that we have two affine hulls with point sets $\{\mathbf{U}_+ \mathbf{v}_+ + \boldsymbol{\mu}_+\}$ and $\{\mathbf{U}_- \mathbf{v}_- + \boldsymbol{\mu}_-\}$. (These can be estimated with either L2 or L1 fitting and they may have different numbers of dimensions l). A closest pair of points between the two hulls can be found by solving

$$\min_{\mathbf{v}_+, \mathbf{v}_-} \|(\mathbf{U}_+ \mathbf{v}_+ + \boldsymbol{\mu}_+) - (\mathbf{U}_- \mathbf{v}_- + \boldsymbol{\mu}_-)\|^2. \quad (5)$$

Defining $\mathbf{U} \equiv (\mathbf{U}_+ \quad -\mathbf{U}_-)$ and $\mathbf{v} \equiv (\mathbf{v}_+^T \quad \mathbf{v}_-^T)^T$, this can be written as the standard least squares problem

$$\min_{\mathbf{v}} \|\mathbf{U}\mathbf{v} - (\boldsymbol{\mu}_- - \boldsymbol{\mu}_+)\|^2. \quad (6)$$

If we take the derivative of the objective function (6) with respect to \mathbf{v} and equate it to zero, then we obtain

$$\mathbf{U}^T \mathbf{U} \mathbf{v} - \mathbf{U}^T (\boldsymbol{\mu}_- - \boldsymbol{\mu}_+) = \mathbf{0}. \quad (7)$$

Subsequently, we get the solution of the problem as $\mathbf{v} = (\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T (\boldsymbol{\mu}_- - \boldsymbol{\mu}_+)$. Taking the decision boundary $f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$,

$$\mathbf{w} = \frac{1}{2} (\mathbf{x}_+ - \mathbf{x}_-) = \frac{1}{2} (\mathbf{I} - \mathbf{P}) (\boldsymbol{\mu}_+ - \boldsymbol{\mu}_-), \quad (8)$$

where $\mathbf{P} = \mathbf{U}(\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T$ is the orthogonal projection onto the joint span of the directions contained in the two subspaces, $\mathbf{I} - \mathbf{P}$ is the corresponding projection onto the orthogonal complement of this span,¹ and \mathbf{x}_+ and \mathbf{x}_- denote the closest points on the positive and negative classes, respectively. Note that \mathbf{w} lies along the line segment joining the two closest points and it is half the line segment’s size. The offset b of the separating hyperplane is given by

$$b = -\mathbf{w}^T (\mathbf{x}_+ + \mathbf{x}_-) / 2. \quad (9)$$

2.2. Inseparable case

A problem arises if the affine hulls of classes intersect, i.e., affine hulls are not linearly separable. If the affine hulls of classes are close to being linearly separable and they overlap because of a few outliers, we can restrict the influence of outlying points by reducing affine hulls. Note that ignoring directions corresponding to the overly small singular values during affine hull constructions reduces the effects of noise and outliers to some degree. But, we will use a different approach here in order to cope with the outliers. To this end, we use the initial affine hull formulation (2) and introduce upper and lower bounds on coefficients α_i to reduce affine hulls inspired by the idea that is introduced to reduce convex hulls in [1]. It should be noted that the reduced affine hulls are not uniformly scaled versions of the initial complete affine hulls. One may go further and choose different lower and upper bounds, or define a different interval for every sample in the

¹ If the two subspaces share common directions, $\mathbf{U}^T \mathbf{U}$ is not invertible and the solution for $(\mathbf{v}_+, \mathbf{v}_-)$ and $(\mathbf{x}_+, \mathbf{x}_-)$ is non-unique, but the orthogonal complement remains well defined, giving a unique minimum norm separator \mathbf{w} . Numerically all cases can be handled by finding $\tilde{\mathbf{U}}$, the \mathbf{U} matrix of the thin SVD of \mathbf{U} , and taking $\mathbf{P} = \tilde{\mathbf{U}} \tilde{\mathbf{U}}^T$.

training set if a-priori information is available. For instance, if the lower bound is set to zero, then the method will be equivalent to the SVM classifier. Finding the closest points on the reduced affine hulls can be written as a quadratic optimization problem

$$\begin{aligned} \min_{\alpha} & \left\| \sum_{i:y_i=+1} \alpha_i \mathbf{x}_i - \sum_{i:y_i=-1} \alpha_i \mathbf{x}_i \right\|^2 \\ \text{s.t.} & \sum_{i:y_i=+1} \alpha_i = 1, \quad \sum_{i:y_i=-1} \alpha_i = 1, \quad -\tau \leq \alpha_i \leq \tau, \end{aligned} \quad (10)$$

where τ is the user-chosen bound. This optimization problem (10) can be written in a more compact form as

$$\begin{aligned} \min_{\alpha} & \sum_{ij} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\ \text{s.t.} & \sum_i \alpha_i y_i = 0, \quad \sum_i \alpha_i = 2, \quad -\tau \leq \alpha_i \leq \tau. \end{aligned} \quad (11)$$

This is a quadratic programming problem that can be solved using standard optimization techniques. Note that the Hessian matrix, $\mathbf{G} = [G_{ij}] = y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$, is a positive semi-definite matrix, thus the objective function is convex and a global minimum exists as in SVM classifier. Moreover, if the Hessian matrix is strictly positive definite, the solution is unique and it is guaranteed to be the global minimum.

Since the coefficients are bounded between $-\tau$ and $+\tau$, the solution is determined by more points and no extreme point or noisy point can excessively influence the solution for well-chosen τ . Once we compute the optimal values of coefficients α_i , the normal and the offset of the separating hyperplane can be computed as in the linearly separable case

$$\mathbf{w} = \frac{1}{2} \left(\sum_{i:y_i=+1} \alpha_i \mathbf{x}_i - \sum_{i:y_i=-1} \alpha_i \mathbf{x}_i \right), \quad (12)$$

$$b = -\frac{1}{2} \mathbf{w}^T \left(\sum_{i:y_i=+1} \alpha_i \mathbf{x}_i + \sum_{i:y_i=-1} \alpha_i \mathbf{x}_i \right). \quad (13)$$

We call this method large margin classifier of affine hulls (LMCAH) since it uses affine hulls to approximate class regions and finds the optimal separating hyperplane yielding the largest margin between the affine hulls.

If the underlying geometry of the classes is highly complex and nonlinear, and approximating classes with linear affine hulls is not appropriate, we can map the data into a higher-dimensional space, where the classes can be approximated with linear affine hulls. Note that the objective function of (11) is written in terms of the dot products of samples, which allows the use of the kernel trick. Thus, by using kernel trick - i.e., replacing $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$ with the kernel function $k(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$ where $\phi: \mathbb{R}^d \rightarrow \mathfrak{F}$ is the mapping function from the input space to the mapped space \mathfrak{F} - we can find the best separating hyperplane parameters in the mapped space. As a result, more complex nonlinear decision boundaries between classes can be approximated by using this trick.

2.3. An equivalent formulation based on variation of ν -SVM classifier

The ν -SVM formulation has been proposed as an alternative to the classical SVM formulation [20]. A new parameter ρ' is introduced in this formulation, and error penalty term C that appears in classical SVM formulation is removed. Here, we introduce an alternative formulation to find the best separating hyperplane between affine hulls based on a variation of ν -SVM

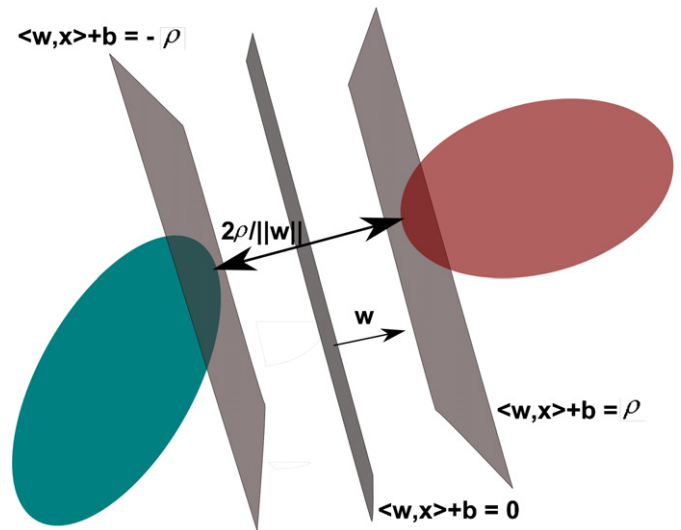


Fig. 3. Illustration of the supporting and the best separating hyperplanes in linearly separable case for ν -SVM classifier.

formulation. A major advantage of new formulation is that one can relate the parameter τ in (11) with the expected error bounds and this may help us to find a more sophisticated procedure for choosing unknown parameters that appear in both formulations.

The ν -SVM optimization is formulated as

$$\begin{aligned} \min_{\mathbf{w}, b, \xi, \rho'} & \frac{1}{2} \|\mathbf{w}'\|^2 - \nu \rho' + \frac{1}{n} \sum_i \xi'_i \\ \text{s.t.} & y_i (\langle \mathbf{w}', \mathbf{x}_i \rangle + b') \geq \rho' - \xi'_i, \quad \xi'_i \geq 0, \quad \rho' \geq 0, \end{aligned} \quad (14)$$

where \mathbf{w}' represents the normal of the separating hyperplane, b' is the offset, ν is a user-chosen parameter between 0 and 1, and $\xi'_i, i=1, \dots, n$, are the positive slack variables. In this formulation, for linearly separable case, there exist two parallel supporting hyperplanes positioned such that all points in the positive class satisfy $\langle \mathbf{w}', \mathbf{x} \rangle + b' \geq \rho'$ and all points in the negative class satisfy $\langle \mathbf{w}', \mathbf{x} \rangle + b' \leq -\rho'$ as shown in Fig. 3. Therefore, classes are separated by the margin $2\rho'/\|\mathbf{w}'\|$ and it is shown that ν acts as an upper bound on the fraction of margin errors and a lower bound on the fraction of support vectors [20]. Moreover, the decision function produced by ν -SVM can also be produced by classical SVM for appropriate choice of error penalty term C .

The ν -SVM formulation can be interpreted as a maximal separation between the reduced convex hulls of classes [9]. Since we use affine hulls to model classes, we need to revise the optimization problem to accommodate this change. To this end, we first divide the objective function by $\nu^2/2$, the constraints by ν , and make the following substitutions as in [9]

$$\tau = \frac{2}{\nu n}, \quad \mathbf{w} = \frac{\mathbf{w}'}{\nu}, \quad b = \frac{b'}{\nu}, \quad \rho = \frac{\rho'}{\nu}, \quad \xi_i = \frac{\xi'_i}{\nu}. \quad (15)$$

These modifications yield the equivalent formulation²

$$\begin{aligned} \min_{\mathbf{w}, b, \xi, \rho} & \|\mathbf{w}\|^2 - 2\rho + \tau \sum_i \xi_i \\ \text{s.t.} & y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq \rho - \xi_i, \quad \xi_i \geq 0, \end{aligned} \quad (16)$$

with the new decision function $f(\mathbf{x}) \equiv f'(\mathbf{x})/\nu$.

² Crisp and Burges [9] showed that the constraint $\rho' \geq 0$ in (14) is redundant and hence it can be removed.

Note that two affine hulls are linearly separable if they lie parallel to each other in the given input space since affine hulls extend to infinity in all directions. In this case, the supporting hyperplanes yielding the largest margin between affine hulls will entirely include them, so that all affine combinations of samples belonging to the positive class satisfy $\langle \mathbf{w}, \mathbf{x}_+^{aff} \rangle + b = \rho$ and all affine combinations of samples belonging to the negative class satisfy $\langle \mathbf{w}, \mathbf{x}_-^{aff} \rangle + b = -\rho$ as illustrated in Fig. 4. Fig. 4 illustrates affine hulls of two classes where the affine hull of the first class is a line and the affine hull of the second class is a plane. Note that affine hulls are linearly separable if they lie parallel to each other. Therefore, all samples of classes and their affine combinations lie on the supporting hyperplanes, which yield the largest margin between the affine hulls. In case of outliers, we must construct reduced compact affine hulls that will fit the data robustly. Therefore, we should allow errors for outlier samples from all over the input space, not just the ones near the decision boundary as

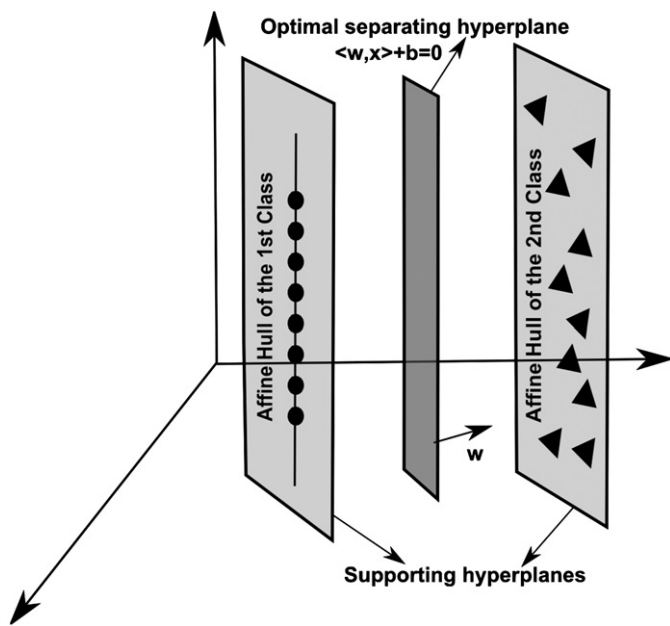


Fig. 4. Optimal separating hyperplane between affine hulls of two classes. Note that affine hulls lie on the supporting hyperplanes.

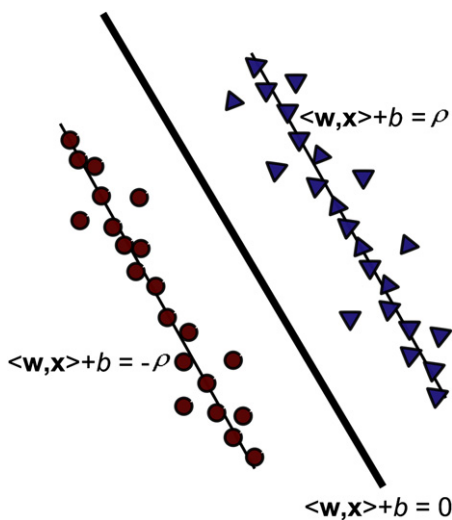


Fig. 5. To obtain better separating hyperplanes between affine hulls, we should allow errors for outlier samples from all over the input space.

illustrated in Fig. 5. To do so, the inequality constraints in (16) is replaced with equality constraints $y_i(\langle \mathbf{w}, \mathbf{x} \rangle + b) = \rho - \delta_i \xi_i$, where δ_i is a term which takes values +1 or -1 based on the location of outliers with respect to the supporting hyperplanes. This leads to the new optimization problem

$$\begin{aligned} \min_{\mathbf{w}, b, \xi, \rho} \quad & \|\mathbf{w}\|^2 - 2\rho + \tau \sum_i \xi_i \\ \text{s.t.} \quad & y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) = \rho - \delta_i \xi_i, \xi_i \geq 0. \end{aligned} \quad (17)$$

To derive the dual, we consider the Lagrangian

$$\begin{aligned} L(\mathbf{w}, b, \xi, \rho, \alpha, \beta) = \quad & \|\mathbf{w}\|^2 - 2\rho + \tau \sum_i \xi_i - \sum_i \alpha_i [y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \\ & - \rho + \delta_i \xi_i] - \sum_i \beta_i \xi_i, \end{aligned} \quad (18)$$

where $\beta_i \geq 0$. The Lagrangian L has to be maximized with respect to α_i , β_i and minimized with respect to \mathbf{w}, b, ξ , and ρ . The optimality conditions yield

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{w}} = 0 \rightarrow \mathbf{w} &= \frac{1}{2} \sum_i \alpha_i y_i \mathbf{x}_i, \\ \frac{\partial L}{\partial b} = 0 \rightarrow \sum_i \alpha_i y_i &= 0, \\ \frac{\partial L}{\partial \rho} = 0 \rightarrow \sum_i \alpha_i &= 2, \\ \frac{\partial L}{\partial \xi_i} = 0 \rightarrow \alpha_i &= \frac{\tau - \beta_i}{\delta_i} \rightarrow -\tau \leq \alpha_i \leq \tau. \end{aligned} \quad (19)$$

Thus, the dual of the optimization problem becomes

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{4} \sum_{ij} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\ \text{s.t.} \quad & \sum_i \alpha_i y_i = 0, \sum_i \alpha_i = 2, -\tau \leq \alpha_i \leq \tau. \end{aligned} \quad (20)$$

This optimization is equivalent to the one given in (11)—here 1/4 appears in the objective function, but rescaling objective function with a positive constant does not change the solution. Therefore, the new formulation based on modified v -SVM is in fact equivalent to finding the best separating hyperplane between the reduced affine hulls that represent classes. We call this method v -LMC-AH. Due to the Karush-Kuhn-Tucker (KKT) conditions, slack variables can occur only when $\alpha_i = \pm \tau$. To compute offset b , we use the primal constraints and take equal number of samples with coefficients $\alpha_i \neq \pm \tau$ from positive and negative classes. Assume that there are l selected samples. By using KKT conditions, we know that $\xi_i = 0$ for the samples with $\alpha_i \neq \pm \tau$. Thus, the offset will be

$$b = -\frac{1}{2l} \sum_{i=1}^l \sum_{j=1}^n \alpha_j y_j \langle \mathbf{x}_j, \mathbf{x}_i \rangle. \quad (21)$$

This offset is not necessarily equivalent to the one given in (13). Therefore, using geometrically inspired formulation and v -LMC-AH formulation create separating hyperplanes with the same normal, but the positions (perpendicular distances from the origin) of these hyperplanes may be different. It is not a priori evident that which offset is the best and one can use other principled methods to determine the best b for a given problem, e.g., given \mathbf{w} , b can be computed as value yielding the smallest classification error on a validation set. As in the previous case, extension to the nonlinear case can be done by using the kernel trick.

2.4. Geometric interpretation of least squares SVM classifier

Least squares support vector machines (LS-SVM) was initially proposed by Suykens and Vandewalle [24] for classification and nonlinear function estimation and then new variants of this method have been introduced [7,21,22]. The basic motivation was to simplify the classical SVM formulation without losing much generalization performance. To this end, inequality constraints in the SVM classification formulation are heuristically replaced with equality constraints,³ and the square of the slack variables are used in the objective function. More formally the optimization problem is defined as

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{2} \sum_i \xi_i^2 \\ \text{s.t.} \quad & y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) = 1 - \xi_i, \end{aligned} \quad (22)$$

where C is a user-chosen error penalty term as in classical SVM classifier. It is shown that the solution is obtained by solving a set of linear equations rather than solving a quadratic programming problem [24]. Note that, in this formulation, for the linearly separable case, there exist two parallel supporting hyperplanes positioned such that all points in the positive class satisfy $\langle \mathbf{w}, \mathbf{x} \rangle + b = 1$ and all points in the negative class satisfy $\langle \mathbf{w}, \mathbf{x} \rangle + b = -1$ where the margin between these hyperplanes is given by $2/\|\mathbf{w}\|$. As in v -LMC – AH formulation, this corresponds to approximating each class with an affine hull instead of a convex hull since all samples and their affine combinations are forced to lie on the supporting hyperplanes. In LS-SVM, L2 norm (squares) of the slack variables are used in the optimization, but using L1 norm of the slack variables in the objective function may be more appropriate for some applications since it allows more robust fitting of data samples. If we use L1 norm of the slack variables, new optimization problem becomes

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i \\ \text{s.t.} \quad & y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) = 1 - \delta_i \xi_i, \quad \xi_i \geq 0, \end{aligned} \quad (23)$$

where δ_i is a term which takes values $+1$ or -1 . The Lagrangian will be

$$\begin{aligned} L(\mathbf{w}, b, \xi, \alpha, \beta) = \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i - \sum_i \alpha_i [y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \\ & - 1 + \delta_i \xi_i] - \sum_i \beta_i \xi_i, \end{aligned} \quad (24)$$

under the constraint $\beta_i \geq 0$. The optimality conditions yield

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{w}} = \mathbf{0} \rightarrow \mathbf{w} &= \frac{1}{2} \sum_i \alpha_i y_i \mathbf{x}_i, \\ \frac{\partial L}{\partial b} = 0 \rightarrow \sum_i \alpha_i y_i &= 0, \\ \frac{\partial L}{\partial \xi_i} = 0 \rightarrow \alpha_i &= \frac{C - \beta_i}{\delta_i} \rightarrow -C \leq \alpha_i \leq C. \end{aligned} \quad (25)$$

Thus, the dual of the optimization problem becomes

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle - \sum_i \alpha_i \\ \text{s.t.} \quad & \sum_i \alpha_i y_i = 0, \quad -C \leq \alpha_i \leq C. \end{aligned} \quad (26)$$

³ In fact, using equality constraints for nonlinear function estimation was introduced earlier in [19].

Similar to the previous cases, this is a convex quadratic optimization problem with a global minimum. Due to the Karush–Kuhn–Tucker (KKT) conditions, slack variables can occur only when $\alpha_i = \pm C$. To compute offset b , we use the primal constraints and take equal number of samples with coefficients $\alpha_i \neq \pm C$ from positive and negative classes as in v -LMC – AH. Assume that there are l selected samples. By using KKT conditions, we know that $\xi_i = 0$ for the samples with $\alpha_i \neq \pm C$. Thus, the offset will be

$$b = -\frac{1}{2l} \sum_{i=1}^l \sum_{j=1}^n \alpha_j y_j \langle \mathbf{x}_j, \mathbf{x}_i \rangle. \quad (27)$$

A new sample is classified based on the sign of the decision function $f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$. Nonlinearization can be done by replacing the dot products $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$ with the kernel function $k(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$. We call this method as C -LMC-AH since it uses error penalty term C . It follows from Proposition 6 of [20] that for appropriate choices of C , the C -LMC-AH algorithm will yield identical results to LMC-AH and v -LMC – AH classifiers. More precisely, if v -LMC – AH classification leads to $\rho \geq 0$, then C -LMC-AH method with C set a-priori to $1/\rho^n$ (or $1/\rho v n$), leads to the same decision function as v -LMC – AH [6,20].

2.5. Extension to the multi-class classification problems

To use the proposed methods in multi-class classification problems, we can use most of the strategies adopted for extending binary SVM classifiers to the multi-class cases. Here we will discuss the most popular two strategies: one-against-one (OAO) and one-against-rest (OAR). For a c -class classification problem, the OAR strategy trains c binary classifiers, in which each classifier separates one class from the remaining $c - 1$ classes. All classifiers are needed to be trained on the entire training set, and the class label of a test sample is determined according to the highest output of the classifiers in the ensemble. On the other hand, the OAO strategy constructs all possible $c(c-1)/2$ binary classifiers out of c classes. The decision of the ensemble is decided by max wins algorithm: Each OAO classifier casts one vote for its preferred class, and the final decision is the class with the most votes. In addition to these we can also use directed acyclic graphs [17] or binary decision trees [28] for multi-class classification.

3. Experiments

We tested⁴ the linear and kernelized versions of the proposed methods LMC-AH, v -LMC – AH and C -LMC-AH (L1 norm based LS-SVM) on a number of datasets and compared them to the SVM classifier. For the linearly separable case, linear separator is determined by using affine subspace estimation formulation, and subspace dimensions are set by retaining enough leading eigenvectors to account for 95–98% of the total energy in the eigen-decomposition. For the inseparable and nonlinear cases, we used quadratic programming formulations. Both one-against-rest (OAR) and one-against-one (OAO) approaches are used for multi-class classification problems and we report the results of whichever yields the best.

We first tested the linear LMC-AH method on multiple and single shot face recognition problems to demonstrate that affine hull approximations are more appropriate than convex hull approximations when the dimensionality of the input space is high. To assess the generalization performances of kernelized

⁴ For software see <http://www2.ogu.edu.tr/~mlcv/software.html>.



Fig. 6. Some detected face images from videos belonging to two subjects.

Table 1
Classification rates (%) on the Honda/UCSD database.

Methods	Clean	Corrup. training	Corrupted test	Corrup. training+test
LMC-AH	97.44	97.44	92.31	87.18
SVM	94.87	92.31	92.31	82.05

versions of the methods, we tested them on seven low-dimensional databases chosen from the UCI repository.

3.1. Experiments on the honda/UCSD database

Honda/UCSD database [15] has been collected for video-based face recognition and it consists of 59 video sequences belonging to 20 individuals. Each video consists of approximately 300–500 frames. It is a fixed database, so that 20 of the videos are allocated for training and the remaining 39 for testing. Here, we consider face recognition based on multiple images. In this scenario, face recognition problem is defined as taking a set of face images from an unknown person and finding the most similar set among the database of labeled image sets. We used the cascade face detector of Viola and Jones [27] to detect faces in each video sequence and resized the detected face images to the gray images of size 40×40 followed by histogram equalization. Then, these images are used to construct image sets of individuals. Some of the detected face images are shown in Fig. 6. We used affine hulls and convex hulls to model image sets, and used the distances between these models as a similarity measure. In other words, computed margins between linear affine hulls and convex hulls are used to determine the label of the tested image sets.⁵ We performed several experiments in order to test the robustness against outliers. In the first experiment, we computed classification rates based on the clean image sets. Then, we systematically corrupted training and test sets by adding images from other classes to each set. These images can be seen as outliers and the changes in classification rates reflect the robustness of the methods against these outliers. The results are given in Table 1. As can be seen from the table, affine hull approximations yield better results than convex hull approximations in all cases except for the corrupted test where both models give same results. Thus, affine hulls seem better and more robust models for approximating image sets.

⁵ For the affine hull case we simply used the minimum distance between the estimated affine subspaces using Eq. (6), whereas soft margin linear SVM algorithm is used to determine the distances between convex hulls.

3.2. Experiments on the AR face database

The AR face dataset [16] contains 26 frontal images with different facial expressions, illumination conditions and occlusions for each of 126 subjects, recorded in two 13-image sessions spaced by 14 days. For this experiment, we randomly selected 20 male and 20 female subjects. The images were down-scaled (from 768×576), aligned, so that centers of the two eyes fell at fixed coordinates, then cropped to size 105×78 . Some pre-processed images are shown in Fig. 7. Raw pixel values were used as features. For training we randomly selected $n=7,13,20$ samples for each individual, keeping the remaining $26-n$ for testing. This process was repeated 10 times, with the final classification rate being obtained by averaging the 10 results. The results are presented in Table 2. Best results are obtained by OAR strategy for both tested methods. The proposed method gives better classification rates than soft margin linear SVM classifier in all cases. The performance difference is more apparent for $n=7$. These results support our claims, suggesting that affine hulls can be better models for representing classes in high-dimensional spaces when the number of samples is limited.

3.3. Experiments on the UCI databases

In this group of experiments, we tested the kernelized versions of the methods (quadratic programming formulations) on seven lower-dimensional datasets from the UCI repository: Ionosphere, Iris, Letter Recognition (LR), Multiple Features (MF)-pixel averages, Pima Indian Diabetes (PID), Wine, and Wisconsin Diagnostic Breast Cancer (WDBC). The key parameters of these datasets are summarized in Table 3. We used the Gaussian kernels, and all design parameters are set based on random partitions of datasets into training and test sets. OAO strategy was used for multi-class problems. Reported classification rates given in Table 4 are computed by fivefold cross-validation. Although being quite mixed, results indicate that generalization performances of LMC-AH and v -LMC-AH methods compare favorably with SVM classifier, whereas C-LMC-AH generally yields the worst classification accuracy.



Fig. 7. Aligned images of one subject from the AR face database.

Table 2
Classification rates (%) on the AR face database.

Methods	$n=7$	$n=13$	$n=20$
LMC-AH	95.19 \pm 0.6	98.95 \pm 0.3	99.62 \pm 0.3
SVM	94.54 \pm 0.6	98.66 \pm 0.2	99.58 \pm 0.3

Table 3

Low-dimensional databases selected from UCI repository.

Databases	Number of classes	Dataset size	Dimensionality
Ionosphere	2	351	34
Iris	3	150	4
LR	26	20000	16
MF	10	2000	256
PID	2	768	8
Wine	3	178	13
WDBC	2	569	30

Table 4

Classification rates (%) on the UCI datasets.

UCI	LMC-AH	ν -LMC – AH	C-LMC-AH	SVM
Ionosphere	93.7 \pm 2.9	93.7 \pm 2.9	93.4 \pm 2.3	92.9 \pm 3.2
Iris	94.7 \pm 2.9	94.7 \pm 2.9	95.3 \pm 3.8	95.3 \pm 3.8
LR	99.98 \pm 0.02	99.98 \pm 0.02	99.89 \pm 0.13	99.64 \pm 0.12
MF	98.4 \pm 0.4	98.4 \pm 0.4	97.8 \pm 0.3	98.0 \pm 0.4
PID	99.9 \pm 0.3	99.9 \pm 0.3	99.9 \pm 0.3	99.9 \pm 0.3
Wine	98.8 \pm 1.6	98.8 \pm 1.6	94.8 \pm 2.6	98.2 \pm 1.6
WDBC	96.0 \pm 2.5	96.0 \pm 2.5	94.9 \pm 3.0	97.6 \pm 0.7

4. Summary and conclusion

We investigated the idea of basing large margin classifiers on affine hulls of classes as an alternative to the SVM (convex hull large margin classifier). Given two affine hull models, their corresponding large margin classifier is easily determined by finding a closest pair of points on these two models and bisecting the displacement between them. We also investigated another formulation obtained by revising the ν -SVM classifier. This formulation yields a separating hyperplane with the same normal as in our first formulation, but the offset is not necessarily the same. This suggests that for a fixed hyperplane normal \mathbf{w} in a specific problem, there may be principled procedures to determine the best offset b . To allow soft margin solutions, we first reduce affine hulls to alleviate the effects of outliers and then find the best separating hyperplanes between these reduced models. Such classifiers can also be kernelized, and extension to the multi-class classification is straightforward using any of the standard approaches such as OAO or OAR.

The experimental results provided useful insights on the potential application areas of the proposed method. The proposed method is much more efficient than SVM classifier in terms of classification accuracy and real-time performance (testing time) when the dimensionality of the sample space is high and affine hulls are linearly separable (in this case solution is easily determined based on subspace estimation which requires simple linear algebra, whereas SVM formulation requires solving a quadratic programming). For the low-dimensional databases generalization performances of the proposed methods compare favorably with SVM classifier but SVM is more efficient in terms of testing time. This is because of the fact that all training data points contribute to the affine hull models (almost all computed α_i coefficients are non-zero), thus the proposed quadratic optimization solutions lack sparseness, and we need more computations to evaluate decision functions. Nevertheless, some pruning techniques can be employed to overcome this problem.

Acknowledgment

This work was supported by the Young Scientists Award Programme (TÜBA-GEBİP/2010) of the Turkish Academy of Sciences.

References

- [1] K.P. Bennett, E.J. Bredensteiner, Duality and geometry in SVM classifiers, in: International Conference on Machine Learning, 2000.
- [2] C.J.C. Burges, Tutorial on support vector machines for pattern recognition, Data Min. Knowl. Discovery 2 (1998) 121–167.
- [3] H. Cevikalp, B. Triggs, R. Polikar, Nearest hyperdisk methods for high-dimensional classification, in: International Conference on Machine Learning, 2008.
- [4] H. Cevikalp, D. Larlus, M. Neamtu, B. Triggs, F. Jurie, Manifold based local classifiers: linear and nonlinear approaches, J. Signal Process. Syst. 61 (1) (2010) 61–73.
- [5] H. Cevikalp, B. Triggs, Large margin classifiers based on convex class models, in: International Conference on Computer Vision Workshops, 2009.
- [6] C.C. Chang, C.J. Lin, Training ν -support vector classifiers: theory and algorithms, Neural Comput. 13 (9) (2001) 2119–2147.
- [7] W. Chu, C.J. Ong, S. Keerthi, An improved conjugate gradient scheme to the solution of least squares SVM, IEEE Trans. Neural Networks 16 (2005) 498–501.
- [8] C. Cortes, V. Vapnik, Support vector networks, Mach. Learn. 20 (1995) 273–297.
- [9] D.J. Crisp, C.J. Burges, A geometric interpretation of ν -SVM classifiers, in: Neural Information Processing Systems, 1999.
- [10] C. Ding, D. Zhou, X. He, H. Zha, R1-pca: rotational invariant l1-norm principal component analysis for robust subspace factorization, in: International Conference on Machine Learning, 2006.
- [11] M.B. Gulmezoglu, V. Dzhafarov, A. Barkana, The common vector approach and its relation to principal component analysis, IEEE Trans. Speech Audio Process. 9 (2001) 655–662.
- [12] G.E. Hinton, P. Dayan, M. Revow, Modeling the manifolds of images of handwritten digits, IEEE Trans. Neural Networks 18 (1997) 65–74.
- [13] Q. Ke, T. Kanade, Robust L1 norm factorization in the presence of outliers and missing data by alternative convex programming, in: IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2005.
- [14] J. Laaksonen, Subspace classifiers in recognition of handwritten digits, Technical Report, 1997.
- [15] K.C. Lee, J. Mo, M.H. Yang, D. Kriegman, Video-based face recognition using probabilistic appearance manifolds, in: IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2003.
- [16] A.M. Martinez, R. Benavente, The AR face database, Technical Report, Computer Vision Center, Barcelona, Spain, 1998.
- [17] J.C. Platt, N. Cristianini, J. Shawe-taylor, Large margin dags for multiclass classification, in: Advances in Neural Information Processing Systems, 2000.
- [18] S.T. Roweis, L.K. Saul, Nonlinear dimensionality reduction by locally linear embedding, Science 290 (2000) 2323–2326.
- [19] C. Saunders, A. Gammerman, V. Vovk, Ridge regression learning algorithm in dual variables, in: International Conference on Machine Learning, 1998.
- [20] B. Schölkopf, A.J. Smola, R.C. Williamson, P.L. Bartlett, New support vector algorithms, Neural Comput. 12 (2000) 1207–1245.
- [21] J.A.K. Suykens, L. Lukas, J. Vandewalle, Sparse least squares support vector machine classifiers, in: Proceedings of the European Symposium on Artificial Neural Networks (ESANN'2000), Bruges, Belgium, 2000, pp. 37–42.
- [22] J.A.K. Suykens, J. De Brabanter, L. Lukas, J. Vandewalle, Weighted least squares support vector machines: robustness and sparse approximation, Neurocomputing 48 (1–4) (2002) 85–105.
- [23] J.A.K. Suykens, T. Van Gestel, J. De Brabanter, B. De Moor, J. Vandewalle, Least Squares Support Vector Machines, World Scientific Publishing Co. Pte. Ltd., 2002.
- [24] J.A.K. Suykens, J. Vandewalle, Least squares support vector machine classifiers, Neural Process. Lett. 9 (1999) 293–300.
- [25] J. Verbeek, Learning non-linear image manifolds by global alignment of local linear models, IEEE Trans. Pattern Anal. Mach. Intell. 28 (2006) 1236–1250.
- [26] P. Vincent, Y. Bengio, K-local hyperplane and convex distance nearest neighbor algorithms, in: Advances in Neural Information Processing Systems, 2001.
- [27] P. Viola, M. Jones, Robust real-time face detection, Int. J. Comput. Vision 57 (2004) 137–154.
- [28] V. Vural, J.G. Dy, A hierarchical method for multi-class support vector machines, in: International Conference on Machine Learning, 2004.
- [29] Z. Xiaofei, S. Yong, Affine subspace nearest points classification algorithm for wavelet face recognition, in: 2009 WRI World Congress on Computer Science and Information Engineering, 2009.



pattern recognition, neural networks, image and signal processing, optimization, and computer vision. He is a member of the IEEE.

Hakan Cevikalp received the M.S. degree from the Department of Electrical and Electronics Engineering, Eskisehir Osmangazi University, Eskisehir, Turkey, in 2001 and the Ph.D. degree from the Department of Electrical Engineering and Computer Science, Vanderbilt University, Nashville, TN, in 2005. He worked as a post-doctoral researcher at LEAR team of INRIA Rhone-Alpes in France in 2007 and Rowan University in USA in 2008. He is currently working as an assistant professor in the Department of Electrical and Electronics Engineering, Eskisehir Osmangazi University, Eskisehir, Turkey. His research interests include pat-



Yalçın Küçük graduated from the Department of Mathematics of Hacettepe University, Ankara, Turkey in 1980, and started to work as a research assistant in the same department. He received M.Sc. degree in 1983 and Ph.D. degree in 1987 from the Department of Mathematics of Hacettepe University. He is currently working as a professor in the Department of Mathematics of Anadolu University, Eskisehir, Turkey. His research interests include continuity and differentiability of set valued functions, convex and nonconvex set valued optimization and vector optimization.



Bill Triggs is a CNRS researcher who works mainly on machine learning based approaches to understand images and other sensed data. He leads the AI (Apprentissage et Interfaces) team in the Laboratoire Jean Kuntzmann (LJK) in Grenoble, France, and he is also the deputy director of LJK, coordinator of the EU research project CLASS on unsupervised image and text understanding, and coordinator of the CNRS partner of the EU network of excellence PASCAL 2.



Mahide Küçük graduated from the Department of Mathematics of Hacettepe University, Ankara, Turkey in 1979, and started to work as a research assistant in the same department. She received M.Sc. degree in 1982 and Ph.D. degree in 1987 from the Department of Mathematics of Hacettepe University. She is currently working as a professor in the Department of Mathematics of Anadolu University, Eskisehir, Turkey. Her research interests include bitopological spaces and differentiability of set valued functions, convex and nonconvex set valued optimization and vector optimization.



Hasan Serhan Yavuz received the B.S., M.S. and Ph.D. degrees from the Electrical and Electronics Engineering Department of Eskisehir Osmangazi University, Eskisehir, Turkey, in 1999, 2002 and 2008, respectively. He is in the academic staff of the Electrical and Electronics Engineering Department of Eskisehir Osmangazi University. His research interests include pattern recognition, image and signal processing, fuzzy logic and computer vision.



Atalay Barkana received B.S. degree from Robert College of İstanbul in 1969, M.S. and Ph.D. degrees from University of Virginia in 1971 and 1974, respectively, all in Electrical Engineering. He is currently working in the Electrical and Electronics Engineering Department of Anadolu University, Eskisehir, Turkey. His research interests include pattern recognition, computer vision, and optimal control.