

# One-Class Support Vector Machine with Relative Comparisons<sup>\*</sup>

GU Hong (顾弘)<sup>\*\*</sup>, ZHAO Guangzhou (赵光宙), QIU Jun (裘君)<sup>†</sup>

College of Electrical Engineering, Zhejiang University, Hangzhou 310027, China;

<sup>†</sup> Ningbo Institute of Technology, Zhejiang University, Ningbo 315000, China

**Abstract:** One-class support vector machines (one-class SVMs) are powerful tools that are widely used in many applications. This paper describes a semi-supervised one-class SVM that uses supervision in terms of relative comparisons. The analysis uses a hypersphere version of one-class SVMs with a penalty term appended to the objective function. The method simultaneously finds the minimum sphere in the feature space that encloses most of the target points and considers the relative comparisons. The result is a standard convex quadratic programming problem, which can be solved by adapting standard methods for SVM training, i.e., sequential minimal optimization. This one-class SVM can be applied to semi-supervised clustering and multi-classification problems. Tests show that this method achieves higher accuracy and better generalization performance than previous SVMs.

**Key words:** one-class support vector machines; semi-supervised learning; relative comparisons; clustering; multiclass classification

## Introduction

One-class support vector machines (one-class SVMs) were first proposed for novelty detection<sup>[1]</sup> and high dimensional density estimation<sup>[2]</sup>. The two kinds of one-class SVMs are the hypersphere version and the hyperplane version. The hypersphere model, also known as the support vector domain description (SVDD)<sup>[1,3]</sup>, tries to find a minimum surrounding sphere in the feature space that encloses most of the given points. The hyperplane model<sup>[2]</sup>, on the other hand, separates the given points from the origin with the maximum margin. Both mechanisms transform the data into the feature space corresponding to a given kernel and use relaxation parameters to control the outliers. Originally, these two methods both belonged to

the unsupervised learning category, with the target class points provided only (or positive samples). The semi-supervised versions employ negative samples to improve the performance. One method closely related to the hyperplane one-class SVMs is the transductive SVMs (TSVMs)<sup>[4,5]</sup>, which focuses on the labeled data. TSVMs try to develop a large margin hyperplane classifier using labeled training data, while simultaneously forcing the hyperplane to be far from the unlabeled data.

Labeled data points are normally used as inequality constraints like  $y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 \geq 0$  in SVMs for classification problems or are transformed to pairwise constraints (must-links and cannot-links) for clustering. This paper presents a new kind of constraint for relative comparisons. A relative comparison constraint is expressed by a triplet  $(\mathbf{x}^1, \mathbf{x}^2, \mathbf{a})$ , which means  $\mathbf{a}$  is closer to  $\mathbf{x}^1$  than to  $\mathbf{x}^2$ . In many applications, relative comparisons are believed to be more easily obtained than other types of constraints<sup>[6,7]</sup>, especially in applications with user feedback. Another advantage of these constraints, which can be used to regularize the

---

Received: 2009-11-16; revised: 2010-01-21

<sup>\*</sup> Supported by the National Natural Science Foundation of China (No. 60872070)

<sup>\*\*</sup> To whom correspondence should be addressed.

E-mail: ghong@zju.edu.cn; Tel: 86-13588705645

local density, is that they indicate which points are closer to the focus point, or the clustering center in one-class data. For example, relevance feedback of content-based image retrieval naturally asks the user which instance is similar to the query even when they do not belong to the same class. Relative comparisons have gained much attention in metric learning, but little research has been done related to instance-based learning theories. This analysis considers these constraints in a one-class SVM. This may be the first study to integrate relative constraints into a one-class SVM. This method appends a new penalty term to the learning objective function, which leads to a convex quadratic programming (QP) problem. The tests demonstrate the performance of this method for clustering and multi-classification problems.

## 1 Related Works

In essence, a one-class SVM can be regarded as a density estimator<sup>[2,3,8]</sup>. However, there are various applications of one-class SVMs. Schölkopf et al.<sup>[9]</sup> applied a one-class SVM to novelty detection. They trained a one-class SVM with normal target data and then used the SVM to detect novel data that deviated from the learned classifier. Manevitz and Yousef<sup>[10]</sup> compared the real performance of one-class SVMs for document classification problems. They showed that the one-class SVMs outperform all other methods except the neural network one.

Ben-Hur et al.<sup>[11]</sup> presented an unsupervised learning method called support vector clustering (SVC) that yields good results for arbitrary geometrical shapes. SVC is based on the observation that given a pair of data points that belong to different components (clusters), any path that connects them must exit from the sphere in the feature space. The learned hypersphere, when mapped back to the data space, can separate into several components, each enclosing a separate cluster of points. The disadvantage of SVC is that it can hardly get the prototype of each cluster. Recent studies have thus preferred using the multiple sphere representation<sup>[12-14]</sup>. The common basis of these approaches is to train a one-class SVM for each class or cluster. Inspired by metric learning, Wang et al.<sup>[15]</sup> proposed a structured one-class SVM training algorithm with constraints expressed by  $(\mathbf{x} - \mathbf{a})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \mathbf{a}) \leq R$  rather than  $\|(\mathbf{x} - \mathbf{a})\|^2 \leq R$  where  $\boldsymbol{\Sigma}$  is computed like RCA<sup>[16]</sup>

in a cluster.

Relative comparisons were first used to study the Mahalanobis distance<sup>[7]</sup> and for SVaD measures<sup>[6]</sup>. Both results show that their methods take advantage of the relative comparisons and outperform state-of-art methods using pairwise constraints. Frome et al.<sup>[17]</sup> proposed an attractive large-margin distance learning algorithm based on relative comparisons. They trained the local distance functions in an SVM-like optimization framework. Another work related to Frome et al.<sup>[17]</sup> is the CLAD measures proposed by Krishna et al.<sup>[18]</sup> Krishna et al.<sup>[18]</sup> used the linear accumulative error model and gradient-based search. The local image distance and CLAD measures were used as the asymmetric distance functions. The main difference between these two methods is the role of the relative comparisons. In Krishna et al.<sup>[18]</sup>, the comparisons were used to generate the objective function which was the accumulative error sum. Frome et al.<sup>[17]</sup> used the comparisons as constraints in a convex optimization problem. The advantage of this method is the global-consistent of the trained distance functions.

The semi-supervised learning improves the learning results by incorporating some supervision information. Generally speaking, there are two approaches to achieve this goal. The first approach is to learn the metric or kernel functions<sup>[7,16-18]</sup>. The other approach uses constraint-based approaches<sup>[5,19,20]</sup>. One popular semi-supervised clustering algorithm that combines metric learning and constraint-based approach is the HMRF-KMeans proposed by Sugato et al.<sup>[20]</sup>, where pairwise constraints are considered. Semi-supervised learning based on one-class SVMs has not been well studied. Tax and Duin<sup>[3]</sup> gives the original solution when there are negative samples in the training set. However, this semi-supervised form of the SVDD is rarely used in real applications, even in classification problems. As mentioned before, relative comparisons have been investigated in metric learning algorithms but not in other types of algorithms. This paper shows how these constraints can be integrated into one-class SVMs.

## 2 One-Class SVMs with Relative Comparisons

Given the training dataset  $\mathbf{X} = \{\mathbf{x}_i\}, i = 1, \dots, N$ , first introduce the relative comparison triplet  $\tau(\mathbf{x}^1, \mathbf{x}^2, \mathbf{a})$ ,

where  $\mathbf{x}^1 \in X$  and  $\mathbf{x}^2 \in X$  are points in the data space and  $\mathbf{a}$  represents the hypersphere center of the target points. The triplet  $\tau(\mathbf{x}^1, \mathbf{x}^2, \mathbf{a})$  means

$$\tau(\mathbf{x}^1, \mathbf{x}^2, \mathbf{a}) \rightarrow \|\mathbf{x}^1 - \mathbf{a}\|^2 \leq \|\mathbf{x}^2 - \mathbf{a}\|^2 \quad (1)$$

in the linear feature space. Assume a set of triplets  $T = \{\tau(\mathbf{x}^1, \mathbf{x}^2, \mathbf{a})_k\}, k=1, \dots, M$ , then the objective is

$$\min R^2 + C_\xi \sum_i \xi_i + C_\gamma \sum_k \gamma_k$$

$$\text{s.t. } \|\mathbf{x}_i - \mathbf{a}\|^2 \leq R^2 + \xi_i, \quad \forall i, \xi_i \geq 0,$$

$$\|\mathbf{x}_k^1 - \mathbf{a}\|^2 \leq \|\mathbf{x}_k^2 - \mathbf{a}\|^2 + \gamma_k, \quad \forall \tau(\mathbf{x}^1, \mathbf{x}^2, \mathbf{a})_k, \quad \gamma_k \geq 0 \quad (2)$$

where  $\xi_i$  and  $\gamma_k$  are slack variables and  $C_\xi$  and  $C_\gamma$  are penalty regularization parameters. The basic idea of the objective is similar to that of SVDD, to find a minimum sphere enclosing most of the target data points. There are two different kinds of constraints in Eq. (2). The first is the same as the original one in Tax and Duin<sup>[3]</sup>. The second kind,  $\|\mathbf{x}_k^1 - \mathbf{a}\|^2 \leq \|\mathbf{x}_k^2 - \mathbf{a}\|^2 + \gamma_k$ , corresponding to the penalty term  $C_\gamma \sum_k \gamma_k$ , is the

additional relative comparison. The objective then tries to find the most feasible model that fulfills these conditions with  $\mathbf{x}_k^1$  closer to the sphere center than  $\mathbf{x}_k^2$ . Note that the metric regularization form described in Wang et al.<sup>[15]</sup> can be easily extended to this framework.

Equation (2) is a standard convex quadratic programming problem which can be solved using the Lagrangian:

$$\begin{aligned} L(R, \mathbf{a}, \xi_i, \gamma_j, \alpha_i, \beta_i, \chi_k, \eta_k) = & R^2 + C_\xi \sum_i \xi_i + C_\gamma \sum_k \gamma_k - \\ & \sum_i \alpha_i \{R^2 + \xi_i - (\mathbf{x}_i \cdot \mathbf{x}_i - 2\mathbf{a} \cdot \mathbf{x}_i + \mathbf{a} \cdot \mathbf{a})\} - \sum_i \beta_i \xi_i - \\ & \sum_k \chi_k \{(\mathbf{x}_k^2 \cdot \mathbf{x}_k^2 - 2\mathbf{a} \cdot \mathbf{x}_k^2) + \gamma_k - (\mathbf{x}_k^1 \cdot \mathbf{x}_k^1 - 2\mathbf{a} \cdot \mathbf{x}_k^1)\} - \\ & \sum_k \eta_k \gamma_k \end{aligned} \quad (3)$$

For simplicity, the symbol  $\sum_i$  is used here to traverse all points in  $X$  and  $\sum_k$  is used to traverse all triplets in  $T$ .  $\alpha_i, \beta_i, \chi_k$ , and  $\eta_k$  are Lagrangian multipliers associated with the variables in Eq. (2) which satisfies  $\alpha_i \geq 0, \beta_i \geq 0, \chi_k \geq 0, \eta_k \geq 0$ .

From the condition

$$\frac{\partial L}{\partial R} = \frac{\partial L}{\partial \mathbf{a}} = \frac{\partial L}{\partial \xi_i} = \frac{\partial L}{\partial \gamma_k} = 0 \quad (4)$$

we get

$$\sum_i \alpha_i = 1 \quad (5)$$

$$\mathbf{a} = \sum_i \alpha_i \mathbf{x}_i + \sum_k \chi_k (\mathbf{x}_k^1 - \mathbf{x}_k^2) \quad (6)$$

$$\alpha_i = C_\xi - \beta_i, \quad \forall i \quad (7)$$

$$\chi_k = C_\gamma - \gamma_k, \quad \forall k \quad (8)$$

Using Eqs. (5)-(8), the solution of Eq. (2) can be obtained by maximizing the Wolfe dual form of the Lagrangian in Eq. (3):

$$W = \sum_i \alpha_i (\mathbf{x}_i \cdot \mathbf{x}_i) - \sum_i \sum_j \alpha_i \alpha_j (\mathbf{x}_i \cdot \mathbf{x}_j) +$$

$$\sum_k \chi_k \{(\mathbf{x}_k^1 \cdot \mathbf{x}_k^1) - (\mathbf{x}_k^2 \cdot \mathbf{x}_k^2)\} -$$

$$\sum_k \sum_s \chi_k \chi_s \{(\mathbf{x}_k^1 \cdot \mathbf{x}_s^1) - (\mathbf{x}_k^1 \cdot \mathbf{x}_s^2) - (\mathbf{x}_k^2 \cdot \mathbf{x}_s^1) + (\mathbf{x}_k^2 \cdot \mathbf{x}_s^2)\} - \\ 2 \sum_i \sum_k \alpha_i \chi_k \{(\mathbf{x}_k^1 \cdot \mathbf{x}_i) - (\mathbf{x}_k^2 \cdot \mathbf{x}_i)\} \quad (9)$$

with the following constraints:

$$\sum_i \alpha_i = 1,$$

$$0 \leq \alpha_i \leq C_\xi, \quad \forall i,$$

$$0 \leq \chi_k \leq C_\gamma, \quad \forall k \quad (10)$$

Finally, the distance from a point to the learned data center is given by

$$\begin{aligned} D^2(\mathbf{x}) = & \|\mathbf{x} - \mathbf{a}\|^2 = \\ & (\mathbf{x} \cdot \mathbf{x}) - 2 \sum_i \alpha_i (\mathbf{x} \cdot \mathbf{x}_i) + \sum_i \sum_j \alpha_i \alpha_j (\mathbf{x}_i \cdot \mathbf{x}_j) - \\ & 2 \sum_k \chi_k \{(\mathbf{x} \cdot \mathbf{x}_k^1) - (\mathbf{x} \cdot \mathbf{x}_k^2)\} + \sum_k \sum_s \chi_k \chi_s \{(\mathbf{x}_k^1 \cdot \mathbf{x}_s^1) - \\ & (\mathbf{x}_k^1 \cdot \mathbf{x}_s^2) - (\mathbf{x}_k^2 \cdot \mathbf{x}_s^1) + (\mathbf{x}_k^2 \cdot \mathbf{x}_s^2)\} \end{aligned} \quad (11)$$

where  $i$  and  $j$  are used to traverse all points in  $X$  and  $k$  and  $s$  are used to traverse all triplets in  $T$ . The radius of the hypersphere is given by  $R = D(\mathbf{x}^{\text{sv}})$  for any point belonging to the spherical support vectors.

The one-class SVM works in the feature space by replacing the dot product  $(\mathbf{x}_i \cdot \mathbf{x}_j)$  with a kernel function  $g(\cdot, \cdot)$ :

$$g(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j) \quad (12)$$

Typically, when the Gaussian kernel is used, the one-class SVM can model arbitrary cluster shapes of the data.

### 3 Training the One-Class SVMs

Sequential minimal optimization (SMO)<sup>[21]</sup> was used to solve the quadratic optimization problem (QP). SMO breaks up the constrained minimization of Eq. (9) into a set of local optimization steps. In this setting,

both  $\alpha$  and  $\chi$  need to be optimized. The basic optimizing step is as follows:

$$\Delta = \alpha_1^{\text{new}} + \alpha_2^{\text{new}}, \alpha_2^{\text{new}} = \alpha_2 + \frac{D^2(\mathbf{x}_2) - D^2(\mathbf{x}_1)}{2\Gamma(\mathbf{x}_1, \mathbf{x}_2)},$$

$$\alpha_1^{\text{new}} = \Delta - \alpha_2^{\text{new}} \quad (13)$$

$$\chi_k^{\text{new}} = \chi_k + \frac{D^2(\mathbf{x}_k^1) - D^2(\mathbf{x}_k^2)}{2\Gamma(\mathbf{x}_k^1, \mathbf{x}_k^2)} \quad (14)$$

where  $\Gamma(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1 \cdot \mathbf{x}_1 + \mathbf{x}_2 \cdot \mathbf{x}_2 - 2\mathbf{x}_1 \cdot \mathbf{x}_2$ . In the training process, any point that deviates from the Karush-Kuhn-Tucker (KKT) conditions of QP in Eq. (2) is selected and optimized. The KKT for Eq. (2) is

$$\begin{cases} \alpha_i \{R^2 + \xi_i - (\mathbf{x}_i \cdot \mathbf{x}_i - 2\mathbf{a} \cdot \mathbf{x}_i + \mathbf{a} \cdot \mathbf{a})\} = 0 \quad \forall i, \\ \beta_i \xi_i = 0 \quad \forall i, \\ \chi_k \{(\mathbf{x}_k^2 \cdot \mathbf{x}_k^2 - 2\mathbf{a} \cdot \mathbf{x}_k^2) + \gamma_k - (\mathbf{x}_k^1 \cdot \mathbf{x}_k^1 - 2\mathbf{a} \cdot \mathbf{x}_k^1)\} = 0 \quad \forall k, \\ \eta_k \gamma_k = 0 \quad \forall k \end{cases} \quad (15)$$

Points corresponding to  $0 < \alpha_i < C_\xi$  are called support vectors (SVs) during the optimization. They lie on the spherical surface of the hypersphere, so can be referred to as spherical SVs as well. Points corresponding to  $\alpha_i = 0$  and  $\alpha_i = C_\xi$  are called inner points and boundary support vectors (BSVs). Analogously, the relative comparison triplets with  $\chi_k = 0$  are called legal triplets while the triplets with  $0 < \chi_k < C_\gamma$  are called critical support triplets (CSTs) and triplets with  $\chi_k = C_\gamma$  are deviation support triplets (DSTs). After the optimization, the inner points and legal triplets are not needed for the classifier since only the terms with  $\alpha_i \neq 0$  and  $\chi_k \neq 0$  need to be restored in Eq. (11). During the iterations, the value of  $\alpha$  is reduced if it violates Eq. (10) or Eq. (13). The general strategy is to first ensure the optimization of  $\alpha$  first. Once all the  $\alpha$  satisfy the KKT,  $\chi$  is optimized.

## 4 Application: Semi-Supervised Clustering

Semi-supervised clustering aims to improve the clustering results using limited supervision. This has become a topic of significant interest in the recent machine learning literature. Existing methods take advantage of metric learning and background or feedback constraints. Sugato et al.<sup>[20]</sup> proposed a probabilistic framework based on hidden Markov random fields (HMRFs) for semi-supervised clustering that combines

the constraint-based and distance-based approaches. An extended kernel version of HMRF-KMeans, known as SS-KERNEL-KMeans was presented by Kulis et al.<sup>[19]</sup> HMRF-KMeans and SS-KERNEL-KMeans use pairwise constraints. SSSVaD<sup>[6]</sup> is another recently proposed semi-supervised clustering algorithm. SSSVaD learns the Mahalanobis distance using relative comparisons which gives better performance than HMRF-KMeans. Xing et al.<sup>[22]</sup> described a convex programming approach to learn the Mahalanobis distance for K-Means with side information. Other methods have also been used in Refs. [23-25].

Kernel Grower<sup>[13,26]</sup> is a new unsupervised clustering algorithm based on one-class SVMs and KMeans. The present method adopts a similar mechanism but uses a generative model and the semi-supervised one-class SVM. At the beginning of the clustering, some of the data points are selected to construct the initial clusters. For each cluster, a one-class SVM is trained with some relative comparisons generated from prior knowledge. Formally, let  $\mathbf{CL} = \{\mathbf{CL}_1, \mathbf{CL}_2, \dots, \mathbf{CL}_z\}$  be the clusters and the data points  $\mathbf{X} = \{\mathbf{x}_i\}$ . The points re-assignment rule is assumed to be  $\arg \max_z p(\mathbf{CL}_z | \mathbf{x}_i)$ , where

$$p(\mathbf{CL}_z | \mathbf{x}_i) = \frac{p(\mathbf{x}_i | \mathbf{CL}_z) p(\mathbf{CL}_z)}{p(\mathbf{x}_i)} \quad (16)$$

Using the Bayesian rule, Eq. (16) can be rewritten as

$$p(\mathbf{CL}_z | \mathbf{x}_i) = \frac{p(\mathbf{x}_i | \mathbf{CL}_z) p(\mathbf{CL}_z)}{\sum_{j=1, \dots, Z} p(\mathbf{x}_i | \mathbf{CL}_j) p(\mathbf{CL}_j)} \quad (17)$$

For a given cluster, calculate  $p(\mathbf{x}_i | \mathbf{CL}_z)$  by

$$p(\mathbf{x}_i | \mathbf{CL}_z) = \max(R_{\mathbf{CL}_z}^2 - D_{\mathbf{CL}_z}^2(\mathbf{x}_i), 0) \quad (18)$$

where  $R_{\mathbf{CL}_z}$  denotes the radius of the learned  $z$ -th hypersphere and  $D_{\mathbf{CL}_z}(\mathbf{x})$  denotes the distance from the sphere center to  $\mathbf{x}$ , Eq. (11). One issue in the clustering algorithm is the generation of the triplets set  $\mathbf{T} = \{\tau(\mathbf{x}^1, \mathbf{x}^2, \mathbf{a})_k\}, k=1, \dots, M$ . The format used here differs from ones described in Refs. [6,7]. Given a cluster  $\mathbf{CL}_z$ , dynamically transform  $(\mathbf{x}^1, \mathbf{x}^2, \mathbf{x}^3)$  to the present format. If  $\mathbf{x}^1$  and  $\mathbf{x}^2$  belong to  $\mathbf{CL}_z$ , construct two triplets  $\tau(\mathbf{x}^1, \mathbf{x}^3, \mathbf{a}_z)$  and  $\tau(\mathbf{x}^2, \mathbf{x}^3, \mathbf{a}_z)$ . A prior label does not need to be assigned to start the clustering, with this simple heuristic working well in the tests.

Two datasets were used to validate the system performance. Letters-IJL is a subset of the Letter Recognition Dataset (Letters) in the UCI data repository<sup>[27]</sup>. Letters contains 20 000 characters, each a 16-dimensional vector. A subset of 10% of the data points was chosen randomly from three classes  $\{I, J, L\}$  of the Letters. The second dataset was Protein<sup>[22]</sup>, which contains 20-dimensional feature vectors representing 116 proteins from 6 classes. The clustering results were evaluated using the standard normalized mutual information (NMI) defined as

$$NMI = \frac{I(U, V)}{(H(U) + H(V)) / 2} \quad (19)$$

where  $U$  is a random variable denoting the cluster assignments of the points and  $V$  is a random variable denoting the underlying class labels of the points.  $I(U, V) = H(U) - H(U|V)$  is the mutual information between the random variables  $U$  and  $V$ .  $H(U)$  is the Shannon entropy of  $U$  and  $H(U|V)$  is the conditional entropy of  $U$  given  $V$ . More information about the NMI measures can be found in Ref. [28].

The experiments used 2-fold cross validation. The data points were separated into two groups for training and for testing. Figures 1 and 2 show the average clustering results for Letters-IJL and Protein. The present method is referred as to the semi-supervised K multi-sphere support vector clustering (SS-KMSVC) method in the figures. The clustering results of some other algorithms are also given. All the KMeans-based algorithms found in the literature suffer from a dependence on the initialization algorithm. Therefore, to get a fair comparison, all the methods used the neighborhood inference algorithm provided by Sugato et al.<sup>[20]</sup> to initialize the clusters. The SS-KMSVC used the Gaussian kernel  $g(\mathbf{x}_1, \mathbf{x}_2) = \exp(-q \cdot \|\mathbf{x}_1 - \mathbf{x}_2\|^2)$ . The width parameter  $q$  was adjusted to get the best performance, typically  $1.2E-3$  for Letters-IJL and  $2.5E-3$  for Protein. The penalty regularization parameters  $C_\xi$  and  $C_\gamma$  were both set to 1 in the present method. These two variables are robust to the clustering results. In the testing with very few relative comparison triplets are the DSTs where most  $\chi_k$  were less than 0.1.

The semi-supervised clustering results obtained by the SSSVaD<sup>[6]</sup> algorithm consistently outperform the other methods when there is a small number of constraints as shown in Figs. 1 and 2. However, other

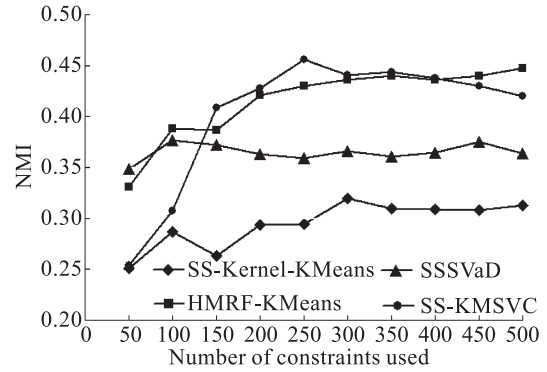


Fig. 1 Clustering results for Letters-IJL

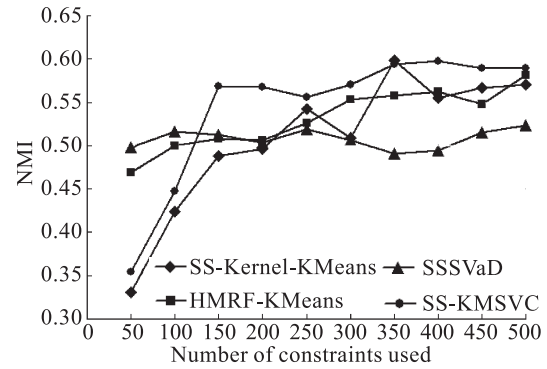


Fig. 2 Clustering results for Protein

methods had better performance than SSSVaD with a large number of constraints. The kernel methods, SS-KERNEL-KMeans did poorly on the Letters-IJL dataset, and yielded unstable results on the Protein dataset. The clustering results obtained by HMRF-KMeans are very good on both datasets, with better performance as the number of constraints increases. The present method achieves the best overall performance for a modest number of constraints. With the Protein dataset, SS-KMSVC did significantly better with more than 200 constraints. On the Letters-IJL dataset, the best clustering result occurred with about 250 constraints. The performance decreases slightly as the number of constraints increases due to the triplets conversion collisions and overfitting.

## 5 Application: Multi-Classification

This section applies the one-class SVMs for multi-classification problems. Generally speaking, the two basic strategies to solve multi-classification problem are to classify all the data at once based on a posterior probability and to use binary classification methods and then combine the binary classification results. A well known algorithm is the one-versus-one

method (also referred to as the pairwise method) based on SVMs. If there are  $z$  classes, the pairwise method takes  $z(z-1)/2$  two-class SVMs to construct the final classification tree when the whole problem is decomposed to a set of binary classification problems. Another popular algorithm is the one-versus-all (OVA) method. A one-class classifier is trained by classifying points outside of this class as negative samples, thus  $z$  SVMs need be trained for the  $z$  class problem. Lee and Lee<sup>[14]</sup> proposed an OVA method based on the Bayesian optimal decision theory and SVDD. For multi-classification problems, they used the basic SVDD to train a one-class classifier for each class using only the target (positive) data points during the training. The final multiclass classifier was constructed by maximizing the pseudo posterior probability function:

$$\arg \max_{z=1, \dots, Z} p(\text{CL}_z) \cdot p(\mathbf{x} | \text{CL}_z) \quad (20)$$

where  $p(\text{CL}_z)$  is the priori and  $p(\mathbf{x} | \text{CL}_z)$  is the

density function obtained in Eq. (18).

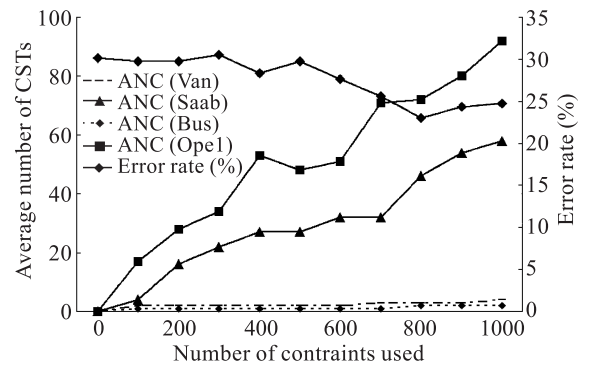
The current tests follow the tests in Lee and Lee<sup>[14]</sup>. The five datasets used in the tests were Iris, Vehicle, Vowel, Segment, and Letters. All these datasets can be freely downloaded from the UCI machine learning repository. The generation of relative comparisons is quite straightforward. Suppose that a one-class SVM classifier is trained for class A with  $m$  positive points in A and other  $n$  negative points (belong to other classes) to generate  $m \times n$  triplets. This can be a very large number for some datasets; however, not all the triplets are needed. In general, we randomly select  $(m+n)$  triplets as described in the following. The parameter  $q$  is set to the same value as in Lee and Lee<sup>[14]</sup> and  $C_\xi$  and  $C_\gamma$  are set to 1. The test results are shown in Table 1 which indicates that the current method outperforms that of Lee and Lee<sup>[14]</sup> on most datasets except for Segment, with a slightly higher error rate.

**Table 1 Multi-classification test results for several datasets based on the error rate (%). The Vehicle dataset was normalized while the others were not. The parameter  $q$  for the five datasets was set to 3, 600, 0.9, 0.0037, and 0.0012. The boldface terms are the best results (the lowest error rates) for each dataset and the term in the parentheses is the differences of error rate between the present method and Lee and Lee<sup>[14]</sup>. LDA and QDA denote linear discriminant analysis and quadratic discriminant analysis<sup>[29]</sup>.**

Data	Error rate (%)					
	LDA	QDA	1-1 SVMs	1-all SVMs	Lee's <sup>[14]</sup>	Our method
Iris	<b>2</b>	<b>2</b>	6	10	<b>2</b>	<b>2</b> (↓0)
Vehicle	22.7	<b>20.21</b>	39.36	33.33	30.14	23.05 (↓7.09)
Vowel	55.63	52.81	38.74	34.85	37.23	<b>30.7</b> (↓6.53)
Segment	8.83	22.21	13.25	6.88	<b>4.94</b>	5.03 (↑0.09)
Letters	30.14	11.28	6.2	N/A	4.92	<b>4.76</b> (↓0.16)

As mentioned, many relative comparisons can be generated using all the original data points. The results in Fig. 3 show this effect for the Vehicle dataset. The error rate decreased for less than 800 relative comparisons and then increased slightly. The average number of CSTs (ANC) for each Vehicle class show that the class Opel had the most CSTs, about 10% of the number of generalized relative comparisons, with Saab having the next highest. Classes Van and Bus have small numbers of CSTs which shows good discrimination of the features. For the hard classification classes like Opel and Saab, the number of CSTs increases with the number of relative comparison triplets. However, with too many triplets the training process is time-consuming and does not bring the desired performance boosting, as shown in Table 2 which gives a

detailed comparison of the training times needed for increasing numbers of relative comparisons. A compromise is to use about  $(m+n)$  triplets in the training which gives better results without excessively long training times.



**Fig. 3 Multi-classification results for the Vehicle dataset for various numbers of relative comparisons**

**Table 2 Training times for various numbers of relative comparisons on dataset Vehicle**

Number of constraints used	Training time (ms)
0	1010
100	1964
500	10 511
800	22 798
1000	28 819

## 6 Conclusions

A one-class SVM was developed for relative comparisons as perhaps the first to incorporate this type of constraint into one-class SVMs. The SVM was then used for semi-supervised clustering and multiclass classification problems. The method yields better results than other state-of-the-art algorithms. The results shown in Section 4 suggest that the metric learning plays an important role for learning global attributes of feature vectors while the constraint-based methods are more suitable for regularizing the local density. Further work is needed to integrate metric learning algorithms into one-class SVMs.

## References

- [1] Tax D M J, Duin R P W. Support vector domain description. *Pattern Recognition Letters*, 1999, **20**(11-13): 1191-1199.
- [2] Scholkopf B, Platt J C, Shawe-Taylor J, et al. Estimating the support of a high-dimensional distribution. *Neural Computation*, 2001, **13**(7): 1443-1471.
- [3] Tax D M J, Duin R P W. Support vector data description. *Machine Learning*, 2004, **54**(1): 45-66.
- [4] Vapnik V. *The Nature of Statistical Learning Theory*. Springer, 1995.
- [5] Collobert R, Sinz F, Weston J, et al. Large scale transductive SVMs. *The Journal of Machine Learning Research*, 2006, **7**(1): 1712.
- [6] Kumar N, Kummamuru K, Paranjpe D. Semisupervised clustering with metric learning using relative comparisons. *IEEE Transactions on Knowledge and Data Engineering*, 2008, **20**(4): 496-503.
- [7] Schultz M, Joachims T. Learning a distance metric from relative comparisons. In: *Proceeding of the Advances in Neural Information Processing Systems*. Vancouver, Canada: The MIT Press, 2004: 41.
- [8] Vert R, Vert J P. Consistency and convergence rates of one-class SVMs and related algorithms. *The Journal of Machine Learning Research*, 2006, **7**(1): 854.
- [9] Scholkopf B, Williamson R C, Smola A J, et al. Support vector method for novelty detection. *Advances in Neural Information Processing Systems*, 2000, **12**(4): 582-588.
- [10] Manevitz L M, Yousef M. One-class SVMs for document classification. *The Journal of Machine Learning Research*, 2002, **2**(1): 154.
- [11] Ben-Hur A, Horn D, Siegelmann H T, et al. Support vector clustering. *The Journal of Machine Learning Research*, 2002, **2**: 125-137.
- [12] Chiang J H, Hao P Y. A new kernel-based fuzzy clustering approach: Support vector clustering with cell growing. *IEEE Transactions on Fuzzy Systems*, 2003, **11**(4): 518-527.
- [13] Camastra F, Verri A. A novel kernel method for clustering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2005, **3**(1): 801-805.
- [14] Lee D, Lee J. Domain described support vector classifier for multi-classification problems. *Pattern Recognition*, 2007, **40**(1): 41-51.
- [15] Wang D, Yeung D S, Tsang E C C. Structured one-class classification. *IEEE Transactions on Systems, Man, and Cybernetics, Part B*, 2006, **36**(6): 1283-1295.
- [16] Bar-Hillel A, Hertz T, Shental N, et al. Learning a mahalanobis metric from equivalence constraints. *Journal of Machine Learning Research*, 2006, **6**(1): 937.
- [17] Frome A, Singer Y, Sha F, et al. Learning globally-consistent local distance functions for shape-based image retrieval and classification. In: *Proceeding of the Eleventh IEEE International Conference on Computer Vision (ICCV' 07)*. Rio de Janeiro, Brazil: IEEE, 2007: 1255-1263.
- [18] Krishna K, Raghu K, Rakesh A. On learning asymmetric dissimilarity measures. In: *Proceedings of the Fifth IEEE International Conference on Data Mining*. Washington, DC, USA, 2005: 697-700.
- [19] Kulis B, Basu S, Dhillon I, et al. Semi-supervised graph clustering: A kernel approach. *Machine Learning*, 2009, **74**(1): 1-22.
- [20] Sugato B, Mikhail B, Raymond J M. A probabilistic framework for semi-supervised clustering. In: *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. Seattle, WA, USA, 2004: 59-68.
- [21] Platt J. Sequential minimal optimization: A fast algorithm for training support vector machines. *Advances in Kernel*

*Methods-Support Vector Learning*, 1999, **208**: 41-65.

- [22] Xing E P, Ng A Y, Jordan M I, et al. Distance metric learning with application to clustering with side-information. *Advances in Neural Information Processing Systems*, 2003: 521-528.
- [23] Cohn D, Caruana R, McCallum A. Semi-supervised clustering with user feedback. *Constrained Clustering: Advances in Algorithms, Theory, and Applications*, 2003, **4**(1): 17.
- [24] Chapelle O, Scholkopf B, Zien A. Semi-Supervised Learning. Cambridge, MA, USA: The MIT Press, 2006.
- [25] Grira N, Crucianu M, Boujemaa N. Unsupervised and semi-supervised clustering: A brief survey. A Review of Machine Learning Techniques for Processing Multimedia Content. Tech. Rep. The MUSCLE European Network of Excellence (FP6), 2004.
- [26] Chang L, Deng X M, Zheng S W, et al. Scaling up kernel grower clustering method for large data sets via core-sets. *Acta Automatica Sinica*, 2008, **34**(3): 376-382.
- [27] Asuncion A, Newman D J. {UCI} Machine Learning Repository. <http://www.ics.uci.edu/~mllearn/MLRepository.html>. Irvine, CA: University of California, School of Information and Computer Science, 2007.
- [28] Strehl A, Ghosh J, Mooney R. Impact of similarity measures on web-page clustering. In: Proceedings of Workshop on Artificial Intelligence for Web Search (AAAI 2000). AAAI, 2000: 58-64.
- [29] Hastie T, Tibshirani R, Friedman J, et al. The Elements of Statistical Learning (2nd edition). The Mathematical Intelligencer, Springer-Verlag, 2008.