

Multi-view learning with privileged weighted twin support vector machine

Ruxin Xu^a, Huiru Wang^{a,*}

^a*Department of Mathematics, College of Science, Beijing Forestry University, No.35 Qinghua East Road, 100083 Haidian, Beijing, China*

Abstract

Weighted twin support vector machines (WLTSVM) mines as much potential similarity information in samples as possible to improve the common shortcoming of non-parallel plane classifiers. Compared with twin support vector machines (TWSVM), it reduces the time complexity by deleting the superfluous constraints using the inter-class K-Nearest Neighbor (KNN). Multi-view learning (MVL) is a newly developing direction of machine learning, which focuses on learning acquiring information from the data indicated by multiple feature sets. In this paper, we propose multi-view learning with privileged weighted twin support vector machines (MPWTSVM). It not only inherits the advantages of WLTSVM but also has its characteristics. Firstly, it enhances generalization ability by mining intra-class information from the same perspective. Secondly, it reduces the redundancy constraints with the help of inter-class information, thus improving the running speed. Most importantly, it can follow both the consensus and the complementarity principle simultaneously as a multi-view classification model. The consensus principle is realized by minimizing the coupling items of the two views in the original objective function. The complementary principle is achieved by establishing privileged information paradigms and MVL. A standard quadratic programming solver is used to solve the problem. Compared with multi-view classification models such as SVM-2K, MVTSVM, MCPK, and PSVM-2V, our model has better accuracy and classification efficiency. Experimental results on 45 binary data sets prove the effectiveness of our method.

*Corresponding author

Email address: whr2019@bjfu.edu.cn (Huiru Wang)

Keywords: Multi-view learning, Weighted-TWSVM, Privileged information, Consensus principle, Complementarity principle

1. Introduction

Support vector machines (SVMs) [1, 2] are supervised learning models with relevant learning algorithms which is used to analyze data for classification and regression analysis. Till now, many improvements to SVM have been proposed. The multi-surface proximal SVM via generalized eigenvalues (GEPSVM) [3] makes each plane the closest to the samples of its category and the farthest away from the samples of other categories. Compared with SVM, the GEPSVM possesses better XOR performance and lower computational complexity. GEPSVM has been developed into a series of new non-parallel plane classifiers. The twin support vector machines (TWSVM) [4] has gained wide attention due to its good generalization capacity and short calculation time [5] as a kind of GEPSVM. TWSVM obtains two non-parallel planes by working out two quadratic programming problems (QPPs), and every QPP has a smaller size than the standard SVM.

To further improve the solving speed of TWSVM, Ye et al. proposed a non-parallel plane classifier called weighted TWSVM with local information (WLTSVM) [6], which digs as much potential similarity information in the samples as possible. It can discover the geometric and discriminative structures of data manifolds by constructing intra-class and inter-class graphs. By weighting the samples [7], the model discovers information about the intrinsic similarity of samples in the same class and derives as many support vectors that reside in the other class as possible. Based on WLTSVM, a KNN(K-Nearest Neighbor)-based weighted rough ν -TSVM [8], least-squares KNN-based weighted multiclass TSVM [9], enhanced regularized KNN-based TSVM (RKNN-TSVM) [10] were proposed so that redundant samples can be deleted and the running speed can be improved. This idea further reduces the effect of outliers on the model.

Multi-view data are feature data of the same object obtained from different ways or levels, with polymorphism characteristics, multi-source, multi-descriptive, and high-dimensional heterogeneity, etc [11]. For example, SIFT features, color histogram features, texture features, and text descriptions constitute a multi-view of the image in the image recognition problem. Multi-view data use features distributed in different feature spaces from different

perspectives to describe an object. These various features disclose different attributes of objects from distinct perspectives, thus enabling a more comprehensive and accurate description of objects compared to a single perspective. Multi-view learning (MVL), also considered data integration of multiple feature sets, is a rising direction in machine learning. Nowadays, MVL has been widely used in various fields and researches [12, 13, 14, 15]. In addition, a multi-view canonical correlation analysis method based on variational graph neural network is proposed [16]. This method is an advanced model based on multi-view data and adopts multi-view representation learning techniques. MVL has proven effective in different application scenarios, such as improving image classification, annotation, and retrieval performance [17], financial distress prediction [18], predicting the multiple stages of AD progression [19] and mining product adoption intentions from social media [20]. MVL methods can be deduced in following three major classes [21]: co-regularization style algorithms, co-training style algorithms, and margin consistency style algorithms [22, 23].

To better mine the information, MVL generally needs to follow two principles: the principle of consistency and the principle of complementarity [24, 25]. The consistency principle aims to maximize the consistency of multiple views. The principle of complementarity indicates that complementary information from multiple views ought to be used to provide a more comprehensive and accurate description of the object. Under these two principles, MVL algorithms can be sorted into co-regularization algorithms and co-training algorithms. The co-regularization algorithm fuses the regularization term of the discriminant function into the objective function to insure consensus information between different views. SVM-2K [26], multi-view twin SVMs (MVTSSVM) [27], regularized multi-view least squares SVM (RMVLSSVM) [28] and multi-view maximum margin of twin spheres SVM [29] are representative multi-view co-regularization learning algorithms. The co-training algorithm maximizes the mutual consistency of multiple views by iterations and exchanges complementary information to generate a classifier on each view [30]. Algorithms that satisfy these two principles tend to have better generalization capability and performance. Most existing classification models either satisfy the complementarity principle or the consistency principle, but fewer models satisfy both principles like multi-view least squares SVM [31].

Inspired by the above-mentioned theories and conclusions, we propose a novel classification algorithm for MVL called multi-view learning with priv-

ileged weighted twin support vector machines (MPWTSVM). It is achieved by solving two QPPs which makes MPWTSVM work faster. The model we propose mines the potential similarity information between different perspectives and categories by using two graphs (intra-class and inter-class graph) to represent intra-class compactness and inter-class separability. Intuitively, support vectors exist in the closest relationship between samples that sharing different labels. Therefore, by considering possible support vectors, the time complexity can be significantly reduced. We ensure the consistency principle by minimizing the coupling terms of the two views in the objective function. Through the establishment of a privileged information paradigm and MVL, the principle of complementarity is realized. Consequently, the proposed model coordinates all views' information abundantly during the learning procedure and preserves different views' characteristics and inherent similarity information. We use a standard quadratic programming solver in order to seek the solution of MPWTSVM. In addition, numerical experiments are conducted to verify the performance.

In summary, our contributions can be outlined as follows.

- 1) MPWTSVM uses the weighted idea of WLTSVM and incorporates this ideology into multiple views. It is worth noting that in the two-classification process, we weigh the samples separately under two perspectives and obtain the weights of different types and the same type at each view. With the help of KNN, redundant samples are deleted. Therefore, the operation efficiency is greatly improved.
- 2) Our model can satisfy both the two principles of MVL, namely, the consensus principle and the complementarity principle. The consensus principle is realized by minimizing the coupling items of the two views in the objective function. The complementary principle is achieved by establishing privileged information paradigms from different perspectives and MVL. Therefore, the proposed MPWTSVM has a better classification ability.
- 3) We compare our method with five state-of-the-art algorithms and performed numerical experiments on 45 binary multi-view classification data sets. The results demonstrate that MPWTSVM has better accuracy and efficiency than other similar algorithms.

The remainder of this paper is presented below. Section 2 retrospects related work about WLTSVM and PSVM-2V. Section 3 provides a detailed

description of our proposed MPWTSVM, where we derive the method’s dual optimization problem and kernel trick. In Section 4, we compare our algorithm with four state-of-the-art algorithms. Section 5 gives the experimental results, and we provide conclusions and future work in Section 6.

2. Related Works

In this section, we give a brief review on a single view learning algorithm WLTSVM and a multi-view learning algorithm PSVM-2V.

2.1. WLTSVM

In [6], a TSVM-based model WLTSVM with local information was proposed. The WLTSVM uses two graphs, i.e., intra-class and inter-class graph, to describe the tightness within a class and the separability between classes, so as to dig out as much basic similarity information as possible in the samples. WLTSVM finds two nonparallel hyperplanes $f(x)$, one for positive class, and the other for negative class:

$$f_1(x) = w_1^\top x_1 + b_1, f_2(x) = w_2^\top x_2 + b_2,$$

where w_1 and w_2 are two nonparallel hyperplanes’ weights, b_1 and b_2 are the biases. The model classifies the samples according to the hyperplane to which the given sample is close.

Suppose that we have N_1 positive training samples $\{x_i, y_i\}, i = 1, 2, \dots, N_1$, and N_2 negative training samples $\{x_i, y_i\}, i = 1, 2, \dots, N_2$, where $x_i \in \mathbb{R}$ and y_i is the class label. For any pair of points $(x_i, x_j), (i = 1, 2, \dots, N_1, j = 1, 2, \dots, N_1)$ in positive samples and an arbitrary point $x_l (l = 1, 2, \dots, N_2)$ in negative samples, we define the weight matrices of within-class graph $G_s (W_{s,ij}^1)$ and between-class graph $G_d (f_j^1)$ of positive samples as below:

$$W_{s,ij}^1 = \begin{cases} 1, & \text{if } x_i \text{ is the } k\text{-nearest neighbors of } x_j \\ & \text{or } x_j \text{ is the } k\text{-nearest neighbors of } x_i \\ 0, & \text{otherwise,} \end{cases} \quad (1)$$

$$f_j^1 = \begin{cases} 1, & \exists j, W_{d,ij}^1 \neq 0 \\ 0, & \text{otherwise,} \end{cases} \quad (2)$$

where

$$W_{d,ij}^1 = \begin{cases} 1, & \text{if } x_l \text{ is the } k\text{-nearest neighbors of } x_i \\ 0, & \text{otherwise.} \end{cases} \quad (3)$$

Similarly, we can get the weight matrices of negative samples and name them $W_{s,ij}^2$ and f_j^2 .

Let $d_j = \sum_{j=1}^{l_1} W_{s,ij}^1, j = 1, 2, \dots, N_1; q_j = \sum_{j=1}^{l_2} W_{s,ij}^2, j = 1, 2, \dots, N_2$, then the formulation of WLTSVM can be written as follows:

$$\begin{aligned} \min \quad & \frac{1}{2} \sum_{i=1}^{l_1} d_j (\omega_1^\top x_j^+ + b_1)^2 + C \sum_{j=1}^{l_2} \xi_j, \\ \text{s.t.} \quad & -f_j^1 (\omega_1^\top x_j^- + b_1) + \xi_j \geq f_j^1 \cdot 1, \quad \xi_j \geq 0, \end{aligned} \quad (4)$$

and

$$\begin{aligned} \min \quad & \frac{1}{2} \sum_{i=1}^{l_2} q_j (\omega_2^\top x_j^+ + b_2)^2 + C \sum_{j=1}^{l_1} \xi_j, \\ \text{s.t.} \quad & -f_j^2 (\omega_2^\top x_j^+ + b_2) + \xi_j \geq f_j^2 \cdot 1, \quad \xi_j \geq 0, \end{aligned} \quad (5)$$

where $(w_i, b_i) \in (R_n \times R)(i = 1, 2)$, C is the penalty coefficient, and ξ is the nonnegative slack variable.

2.2. PSVM-2V

Suppose there is a dataset with two perspectives of l labeled augmented samples $\{(X_A, X_B, Y)\} = \{(x_i^A, x_i^B, y_i)\}_{i=1}^l = \{(x_i^A; 1), (x_i^B; 1), y_i)\}_{i=1}^l$, where X_A and X_B are two views' feature spaces. The i^{th} sample of two views are represented by the superscripts A and B of x_i . PSVM-2V finds two hyperplanes, one for view A, another for view B:

$$w_A^{*\top} \phi_A(x^A) = 0, w_B^{*\top} \phi_B(x^B) = 0,$$

with the optima w_A^* and w_B^* from problem (6).

The optimization problem of PSVM-2V [32] can be written as follows,

$$\begin{aligned}
\min \quad & \frac{1}{2} (\|\omega_A\|^2 + \gamma\|\omega_B\|^2) + C^A \sum_{i=1}^l \xi_i^{A*} + C^B \sum_{i=1}^l \xi_i^{B*} + C \sum_{i=1}^l \eta_i \\
\text{s.t.} \quad & |(w_A \cdot \phi_A(x_i^A)) - (w_B \cdot \phi_B(x_i^B))| \leq \epsilon + \eta_i \\
& y_i(w_A \cdot \phi_A(x_i^A)) \geq 1 - \xi_i^{A*} \\
& y_i(w_B \cdot \phi_B(x_i^B)) \geq 1 - \xi_i^{B*} \\
& \xi_i^{A*} \geq y_i(w_B \cdot \phi_B(x_i^B)), \quad \xi_i^{A*} \geq 0 \\
& \xi_i^{B*} \geq y_i(w_A \cdot \phi_A(x_i^A)), \quad \xi_i^{B*} \geq 0 \\
& \eta_i \geq 0, \quad i = 1, \dots, l,
\end{aligned} \tag{6}$$

where w_A, w_B are weight vectors for views A and B respectively and ϕ_A, ϕ_B are mappings from inputs to high-dimensional feature spaces. The principle of complementarity is realized by limiting the non-negative slack variables $\xi_A = (\xi_1^A, \xi_2^A, \dots, \xi_l^A)^\top$ and $\xi_B = (\xi_1^B, \xi_2^B, \dots, \xi_l^B)^\top$ by the non-negative correction function. C^A, C^B, C are non-negative penalty parameters, γ is a non-negative trade-off parameter and η is the nonnegative slack variable.

3. The multi-view learning with privileged weighted twin support vector machines

In this paper, we propose a novel MVL method called multi-view learning with privileged weighted TSVM (MPWTSVM) which realizes the consensus and complementarity principle at the same time. It not only inherits the weighting idea of WLTSVM, but also combines MVL to produce better performance. The linear and nonlinear cases of MPWTSVM and the dual formulations are presented in the following sections. Major notations used in this paper are summarized in Table 1.

Table 1
List of notations

Notation	Description
(x_i^A, x_i^B, y_i)	i th training point
l	number of training points
$(x_i \cdot x_j)$	inner product between x_i and x_j as $x_i^\top x_j$
ω^A, ω^B	weight vectors for view A and view B

Table 1
List of notations

Notation	Description
$\phi(\cdot)$	mappings from inputs to high-dimensional feature spaces
$\mathcal{K}(x_i, x_j)$	kernel function ($\phi(x_i) \cdot \phi(x_j)$)
$C_A, C_B, C, C_{A2}, C_{B2}, C_2$	non-negative penalty parameter
γ	non-negative trade-off parameter
W^A, W^B	intra-class weight matrix of view A and view B
$f_i^A, f_i^B, f_j^A, f_j^B$	inter-class weight matrix of view A and view B

3.1. Linear MPWTSVM

To make full use of similarity information in data affinity we define the intra-class weight matrix of view A's positive and negative samples respectively.

$$W_{s,ij}^A = \begin{cases} 1, & \text{if } x_i^A \text{ is the } k\text{-nearest neighbors of } x_j^A \\ & \text{or } x_j^A \text{ is the } k\text{-nearest neighbors of } x_i^A \\ 0, & \text{otherwise,} \end{cases} \quad (7)$$

pairs of points (x_i^A, x_j^A) are in view A's positive and negative samples respectively. Similarly, we can acquire intra-class weight matrix of view B's positive samples and negative samples respectively.

The inter-class weight matrix can be defined as below:

$$f_j^A = \begin{cases} 1, & \exists j, W_{d,ij}^A \neq 0 \\ 0, & \text{otherwise,} \end{cases} \quad (8)$$

$$\text{where } W_{d,ij}^A = \begin{cases} 1, & \text{if } x_i^A \text{ is the } k\text{-nearest neighbors of } x_j^A \\ 0, & \text{otherwise,} \end{cases} \quad (9)$$

and pairs of points (x_i^A, x_j^A) are in view A's positive samples, x_i^A is an arbitrary point in view A's negative samples. In the same way, the inter-class weight matrix of view A and view B can be well defined. We can use the weighted thought of WLTSVM and incorporate this idea into multiple views. By this means, we can obtain the intra-class and inter-class weight matrix weight in the corresponding view.

Fig.1 displays the model construction of MPWTSVM. Multi-view data are gathered from different fields or gained from different feature extractors for effective learning. By weighting the data from different perspectives, MPWTSVM makes full use of the similarity information in data affinity. Combining with privilege information and introducing coupling items, MPWTSVM meets the two principles of multi-view classification.

MPWTSVM finds four hyperplanes, two for view A, two for view B:

$$\begin{aligned} f_1^A(x^A) &= w_1^A x_1^A + b_1^A, & f_2^A(x^A) &= w_2^A x_2^A + b_2^A, \\ f_1^B(x^B) &= w_1^B x_1^B + b_1^B, & f_2^B(x^B) &= w_2^B x_2^B + b_2^B, \end{aligned}$$

where $f_1^A(x^A)$ and $f_2^A(x^A)$ are two nonparallel hyperplanes for positive class and negative class of view A separately, $f_1^B(x^B)$ and $f_2^B(x^B)$ are two nonparallel hyperplanes for positive class and negative class of view B. w_1^t and w_2^t are the weights of two nonparallel hyperplanes, b_1^t and b_2^t are the biases of view t (t=A,B). The view t's model of MPWTSVM classifies the samples relying on which hyperplane (from $f_1^t(x)$ and $f_2^t(x)$) the given view t's sample is close to.

Formally, MPWTSVM can be established as follows:

$$\begin{aligned} \min \quad & \frac{1}{2} \sum_{i=1}^{l_1} \sum_{j=1}^{l_1} W_{s,ij}^A (\omega_+^{A\top} x_j^{A,+} + b_+^A)^2 + \frac{1}{2} \gamma \sum_{i=1}^{l_1} \sum_{j=1}^{l_1} W_{s,ij}^B (\omega_+^{B\top} x_j^{B,+} + b_+^B)^2 \\ & + C_A \sum_{j=1}^{l_2} \xi_j^A + C_B \sum_{j=1}^{l_2} \xi_j^B + C \sum_{j=1}^{l_2} \xi_j^A \xi_j^B \\ \text{s.t.} \quad & -f_j^A(\omega_+^{A\top} x_j^{A,-} + b_+^A) + \xi_j^A \geq f_j^A \cdot 1, \\ & -f_j^B(\omega_+^{B\top} x_j^{B,-} + b_+^B) + \xi_j^B \geq f_j^B \cdot 1, \\ & \xi_j^A \geq -f_j^B(\omega_+^{B\top} x_j^{B,-} + b_+^B), \quad \xi_j^A \geq 0, \\ & \xi_j^B \geq -f_j^A(\omega_+^{A\top} x_j^{A,-} + b_+^A), \quad \xi_j^B \geq 0, \quad (j \in I^-) \end{aligned} \quad (10)$$

and

$$\begin{aligned}
\min \quad & \frac{1}{2} \sum_{i=1}^{l_2} \sum_{j=1}^{l_2} W_{s,ij}^A (\omega_-^{A\top} x_j^{A,-} + b_-^A)^2 + \frac{1}{2} \gamma \sum_{i=1}^{l_2} \sum_{j=1}^{l_2} W_{s,ij}^B (\omega_-^{B\top} x_j^{B,-} + b_-^B)^2 \\
& + C_{A2} \sum_{i=1}^{l_1} \xi_i^A + C_{B2} \sum_{i=1}^{l_1} \xi_i^B + C_2 \sum_{i=1}^{l_1} \xi_i^A \xi_i^B \\
s.t. \quad & f_i^A (\omega_-^{A\top} x_i^{A,+} + b_-^A) + \xi_i^A \geq f_i^A \cdot 1, \\
& f_i^B (\omega_-^{B\top} x_i^{B,+} + b_-^B) + \xi_i^B \geq f_i^B \cdot 1, \\
& \xi_i^A \geq f_i^B (\omega_-^{B\top} x_i^{B,+} + b_-^B), \quad \xi_i^A \geq 0, \\
& \xi_i^B \geq f_i^A (\omega_-^{A\top} x_i^{A,+} + b_-^A), \quad \xi_i^B \geq 0, \quad (i \in I^+)
\end{aligned} \tag{11}$$

where $C_A, C_B, C, C_{A2}, C_{B2}, C_2$ are non-negative parameters and γ is a non-negative trade-off parameter.

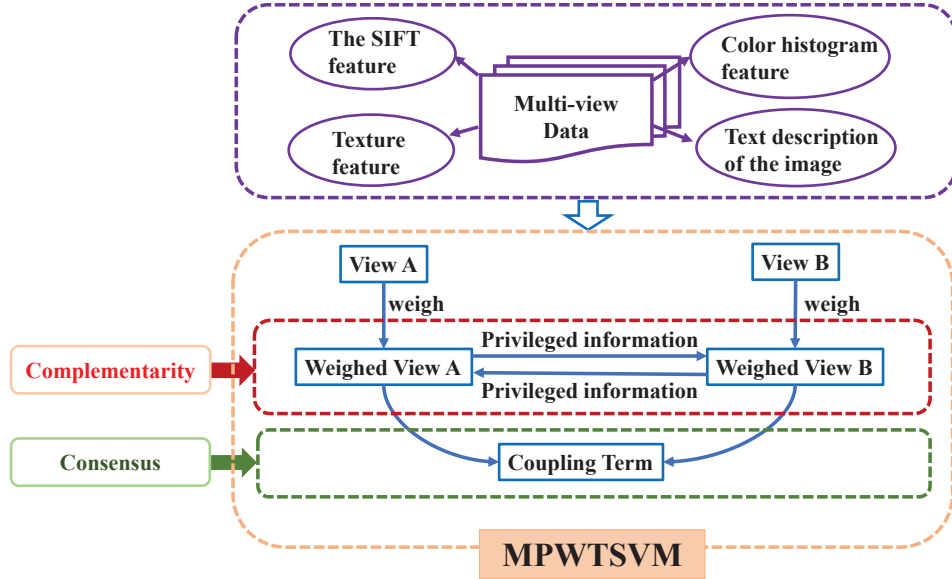


Figure 1: Schematic diagram of MPWTSVM model construction.

Under the supposition that both perspectives are equally important, the MPWTSVM targets on seeking four hyperplanes, which are explained in detail below:

- (1) The variables to be worked out in problem (10) are $w_+^A, w_+^B, b_+^A, b_+^B, \xi^A$ and ξ^B . Problem (11) is similar.
- (2) A larger $W_{s,ij}^A$ indicates a larger weight within the positive class of view A, which can make full use of the structural information. The introduction of this term can make a more compact structure and thus have a better generalization performance. View B is the same. These matrices fully exploit the intra-class information of view A and B.
- (3) Minimizing ξ_i^A, ξ_i^B demonstrates the product of error variables of A and B should be as small as possible. Besides, minimizing the coupling term $C \sum_{i=1}^l \xi_i^A \xi_i^B$ indicates that the high amount of error in one perspective can be recompensed by the low amount of error in the other perspective. It means that a bigger error variable can be allowed in one perspective. In this way, the classification results of the models constructed in different views can converge, and the principle of consensus is realized.
- (4) If the inter-class weight f^A or f^B is 0, it means that the sample constraint is redundant and can be deleted, as a consequence the algorithm efficiency is greatly improved.
- (5) Our model uses each view separately as privileged information to reformulate the slack variables. By limiting the non-negative slack variables ξ^A and ξ^B through the unknown non-negative correcting functions determined by view A and B, therefore, MPWTSVM realizes the complementary principle.

3.2. The dual problem

For simplicity, we let $d_i^{A/B} = \sum_{i=1}^{l_1} W_{s,ij}^{A/B}$, $\mathbf{w}_\pm^{A/B} = \begin{pmatrix} \omega_\pm^{A/B} \\ b_\pm \end{pmatrix}$, $\mathbf{x}^{A/B} = (x^{A/B}, 1)$. Aiming at getting the solution of (10), the corresponding La-

grangian function can be constructed as

$$\begin{aligned}
L(\mathbf{w}_+^A, \mathbf{w}_+^B, \xi_j^A, \xi_j^B) = & \frac{1}{2} \sum_{i=1}^{l_1} d_i^A (\mathbf{w}_+^{AT} \mathbf{x}_i^{A,+})^2 + \frac{1}{2} \gamma \sum_{i=1}^{l_1} d_i^B (\mathbf{w}_+^{BT} \mathbf{x}_j^{B,+})^2 \\
& + C_A \sum_{j=1}^{l_2} \xi_j^A + C_B \sum_{j=1}^{l_2} \xi_j^B + C \sum_{j=1}^{l_2} \xi_j^A \xi_j^B \\
& - \sum_{j=1}^{l_2} \alpha_j^A \left[-f_j^A(\mathbf{w}_+^{AT} \mathbf{x}_j^{A,-}) + \xi_j^A - f_j^A \right] \\
& - \sum_{j=1}^{l_2} \alpha_j^B \left[-f_j^B(\mathbf{w}_+^{BT} \mathbf{x}_j^{B,-}) + \xi_j^B - f_j^B \right] \\
& - \sum_{j=1}^{l_2} \lambda_j^A \left[\xi_j^A + f_j^B(\mathbf{w}_+^{BT} \mathbf{x}_j^{B,-}) \right] \\
& - \sum_{j=1}^{l_2} \lambda_j^B \left[\xi_j^B + f_j^A(\mathbf{w}_+^{AT} \mathbf{x}_j^{A,-}) \right] \\
& - \sum_{j=1}^{l_2} \beta_j^A \xi_j^A - \sum_{j=1}^{l_2} \beta_j^B \xi_j^B, \tag{12}
\end{aligned}$$

where $\alpha^A = (\alpha_1^A, \dots, \alpha_{l_2}^A)^\top$, $\alpha^B = (\alpha_1^B, \dots, \alpha_{l_2}^B)^\top$, $\lambda^A = (\lambda_1^A, \dots, \lambda_{l_2}^A)^\top$, $\lambda^B = (\lambda_1^B, \dots, \lambda_{l_2}^B)^\top$, $\beta^A = (\beta_1^A, \dots, \beta_{l_2}^A)^\top$, $\beta^B = (\beta_1^B, \dots, \beta_{l_2}^B)^\top$ are the non-negative Lagrange multipliers vectors.

Differentiating the Lagrangian function L with respect to variables $\mathbf{w}_+^A, \mathbf{w}_+^B$,

ξ_j^A, ξ_j^B yields the following Karush-Kuhn-Tucker (KKT) conditions:

$$\frac{\partial L}{\partial \mathbf{w}_+^A} = \sum_{i=1}^{l_1} d_i^A \mathbf{x}_i^{A,+} \mathbf{x}_i^{A,+ \top} \mathbf{w}_+^A + \sum_{j=1}^{l_2} \alpha_j^A f_j^A \mathbf{x}_j^{A,-} - \sum_{j=1}^{l_2} \alpha_j^B f_j^A \mathbf{x}_j^{A,-} = 0, \quad (13)$$

$$\frac{\partial L}{\partial \mathbf{w}_+^B} = \gamma \sum_{i=1}^{l_1} d_i^B \mathbf{x}_i^{B,+} \mathbf{x}_i^{B,+ \top} \mathbf{w}_+^B + \sum_{j=1}^{l_2} \alpha_j^B f_j^B \mathbf{x}_j^{B,-} - \sum_{j=1}^{l_2} \alpha_j^A f_j^B \mathbf{x}_j^{B,-} = 0, \quad (14)$$

$$\frac{\partial L}{\partial \xi_j^A} = C_A + C \cdot \xi_j^B - \alpha_j^A - \lambda_j^A - \beta_j^A = 0, \quad (15)$$

$$\frac{\partial L}{\partial \xi_j^B} = C_B + C \cdot \xi_j^A - \alpha_j^B - \lambda_j^B - \beta_j^B = 0, \quad (16)$$

$$\alpha_j^A \left(-f_j^A \mathbf{w}_+^{A \top} \mathbf{x}_j^{A,-} - f_j^A + \xi_j^A \right) = 0, \quad (17)$$

$$\alpha_j^B \left(-f_j^B \mathbf{w}_+^{B \top} \mathbf{x}_j^{B,-} - f_j^B + \xi_j^B \right) = 0, \quad (18)$$

$$\lambda_j^A \left(\xi_j^A + f_j^B \mathbf{w}_+^{B \top} \mathbf{x}_j^{B,-} \right) = 0, \quad (19)$$

$$\lambda_j^B \left(\xi_j^B + f_j^A \mathbf{w}_+^{A \top} \mathbf{x}_j^{A,-} \right) = 0, \quad (20)$$

$$\beta_j^A \xi_j^A = 0, \quad (21)$$

$$\beta_j^B \xi_j^B = 0. \quad (22)$$

Expressing (13) and (14) in matrix form, we get:

$$X_+^{A \top} D_+^A X_+^A w_+^A + X_-^{A \top} F_-^A \alpha_-^A - X_-^{A \top} F_-^A \lambda_-^B = 0, \quad (23)$$

$$\gamma X_+^{B \top} D_+^B X_+^B w_+^B + X_-^{B \top} F_-^B \alpha_-^B - X_-^{B \top} F_-^B \lambda_-^A = 0, \quad (24)$$

where $D_+^A = \text{diag}(d_1^A, d_2^A, \dots, d_{l_1}^A)$, $D_+^B = \text{diag}(d_1^B, d_2^B, \dots, d_{l_1}^B)$, $F_-^A = \text{diag}(f_1^A, f_2^A, \dots, f_{l_2}^A)$ and $F_-^B = \text{diag}(f_1^B, f_2^B, \dots, f_{l_2}^B)$.

Thereupon we can obtain the optimal solutions w_+^A and w_+^B of (10):

$$w_+^A = - \left(X_+^{A \top} D_+^A X_+^A \right)^{-1} \left(X_-^{A \top} F_-^A (\alpha_-^A - \lambda_-^B) \right), \quad (25)$$

$$w_+^B = - \left(\gamma X_+^{B \top} D_+^B X_+^B \right)^{-1} \left(X_-^{B \top} F_-^B (\alpha_-^B - \lambda_-^A) \right). \quad (26)$$

If the matrix $X_+^{A\top} D_+^A X_+^A$ or $\gamma X_+^{B\top} D_+^B X_+^B$ is irreversible, we can introduce a regularization term ϵI , $\epsilon > 0$, then (25) and (26) can be written as

$$w_+^A = -\left(X_+^{A\top} D_+^A X_+^A + \epsilon I\right)^{-1} \left(X_-^{A\top} F_-^A (\alpha_-^A - \lambda_-^B)\right), \quad (27)$$

$$w_+^B = -\left(\gamma X_+^{B\top} D_+^B X_+^B + \epsilon I\right)^{-1} \left(X_-^{B\top} F_-^B (\alpha_-^B - \lambda_-^A)\right). \quad (28)$$

The same method is used if we encounter the problem of matrix integrability in the later part.

By substituting the above equation into (12), we can derive the dual formulation as follows,

$$\begin{aligned} \max \quad & -\frac{1}{2}(\alpha_-^A - \lambda_-^B)^\top F_-^A X_-^A \left(X_+^{A\top} D_+^A X_+^A\right)^{-1} X_-^{A\top} F_-^A (\alpha_-^A - \lambda_-^B) \\ & -\frac{1}{2\gamma}(\alpha_-^B - \lambda_-^A)^\top F_-^B X_-^B \left(X_+^{B\top} D_+^B X_+^B\right)^{-1} X_-^{B\top} F_-^B (\alpha_-^B - \lambda_-^A) \\ & + \alpha_-^{A\top} F_-^A \mathbf{e}_- + \alpha_-^{B\top} F_-^B \mathbf{e}_- - C \xi_-^{A\top} \xi_-^B \\ \text{s.t.} \quad & \alpha_-^A, \alpha_-^B, \lambda_-^A, \lambda_-^B, \beta_-^A, \beta_-^B \geq 0, \end{aligned} \quad (29)$$

where $\xi_-^A = \frac{1}{C}(\alpha_-^B + \lambda_-^B + \beta_-^B - C_B \cdot \mathbf{e}_-)$ and $\xi_-^B = \frac{1}{C}(\alpha_-^A + \lambda_-^A + \beta_-^A - C_A \cdot \mathbf{e}_-)$.

Due to $\beta_j^A, \beta_j^B \geq 0$ in (15) and (16), we have $\alpha_j^A + \lambda_j^A - C \cdot \xi_j^B \leq C_A$, $\alpha_j^B + \lambda_j^B - C \cdot \xi_j^A \leq C_B$. Owing to the complexity of objective function (29), we can work out the following unanimous dual problem as a substitute for simplification,

$$\begin{aligned} \min \quad & \frac{1}{2}(\alpha_-^A - \lambda_-^B)^\top F_-^A X_-^A \left(X_+^{A\top} D_+^A X_+^A\right)^{-1} X_-^{A\top} F_-^A (\alpha_-^A - \lambda_-^B) \\ & + \frac{1}{2\gamma}(\alpha_-^B - \lambda_-^A)^\top F_-^B X_-^B \left(X_+^{B\top} D_+^B X_+^B\right)^{-1} X_-^{B\top} F_-^B (\alpha_-^B - \lambda_-^A) \\ & - \alpha_-^{A\top} F_-^A \mathbf{e}_- - \alpha_-^{B\top} F_-^B \mathbf{e}_- + C \xi_-^{A\top} \xi_-^B \\ \text{s.t.} \quad & \alpha_-^A + \lambda_-^A - C \cdot \xi_-^B \leq C_A \cdot \mathbf{e}_- \\ & \alpha_-^B + \lambda_-^B - C \cdot \xi_-^A \leq C_B \cdot \mathbf{e}_- \\ & \alpha_-^A, \alpha_-^B, \lambda_-^A, \lambda_-^B, \beta_-^A, \beta_-^B \geq 0 \cdot \mathbf{e}_-. \end{aligned} \quad (30)$$

Next, we define $\pi_+ = (\alpha_-^{A\top}, \alpha_-^{B\top}, \lambda_-^{A\top}, \lambda_-^{B\top}, \xi_-^{A\top}, \xi_-^{B\top})^\top$. Concisely, (30)

can be further reformulated as

$$\begin{aligned}
\min_{\pi_+} \quad & \frac{1}{2} \pi_+^\top H_+ \pi_+ + p_+^\top \pi_+ \\
\text{s.t.} \quad & A_+ \pi_+ \leq b_+, \quad \pi_+ \geq 0,
\end{aligned} \tag{31}$$

where

$$H_+ = \begin{pmatrix} H_1^+ & 0_{l_2} & 0_{l_2} & -H_1^+ & 0_{l_2} & 0_{l_2} \\ 0_{l_2} & H_2^+ & -H_2^+ & 0_{l_2} & 0_{l_2} & 0_{l_2} \\ 0_{l_2} & -H_2^+ & H_2^+ & 0_{l_2} & 0_{l_2} & 0_{l_2} \\ -H_1^+ & 0_{l_2} & 0_{l_2} & H_1^+ & 0_{l_2} & 0_{l_2} \\ 0_{l_2} & 0_{l_2} & 0_{l_2} & 0_{l_2} & 0_{l_2} & C \cdot E_{l_2} \\ 0_{l_2} & 0_{l_2} & 0_{l_2} & 0_{l_2} & C \cdot E_{l_2} & 0_{l_2} \end{pmatrix}_{6l_2 \times 6l_2},$$

$$H_1^+ = \left(F_-^A X_-^A \left(X_+^{AT} D_+^A X_+^A \right)^{-1} X_-^{AT} F_-^A \right),$$

$$H_2^+ = \left(\frac{1}{\gamma} F_-^B X_-^B \left(X_+^{BT} D_+^B X_+^B \right)^{-1} X_-^{BT} F_-^B \right),$$

$$p_+^T = \left(-e_-^T F_-^A \quad -e_-^T F_-^B \quad 0_{1 \times l_2} \quad 0_{1 \times l_2} \quad 0_{1 \times l_2} \right)_{1 \times 6l_2},$$

$$A_+ = \begin{pmatrix} E_{l_2} & 0_{l_2} & E_{l_2} & 0_{l_2} & 0_{l_2} & -C \cdot E_{l_2} \\ 0_{l_2} & E_{l_2} & 0_{l_2} & E_{l_2} & -C \cdot E_{l_2} & 0_{l_2} \end{pmatrix}_{2l_2 \times 6l_2},$$

$$b_+^T = (C_A \cdot e_-^T \quad C_B \cdot e_-^T)_{1 \times 2l_2},$$

E_{l_2} is the $l_2 \times l_2$ identity matrix, 0_{l_2} is the $l_2 \times l_2$ matrix with all entries be 0 and e_- is a column vector with the proper dimension with element 1.

Using a similar process, the second optimization problem (11) can be

written as:

$$\begin{aligned} \min_{\pi_-} \quad & \frac{1}{2} \pi_-^\top H_- \pi_- + p_-^\top \pi_- \\ \text{s.t.} \quad & A_- \pi_- \leq b_-, \quad \pi_- \geq 0, \end{aligned} \tag{32}$$

where $\pi_- = (\alpha_+^{A^\top}, \alpha_+^{B^\top}, \lambda_+^{A^\top}, \lambda_+^{B^\top}, \xi_+^{A^\top}, \xi_+^{B^\top})^\top$,

$$H_- = \begin{pmatrix} H_1^- & 0_{l_1} & 0_{l_1} & -H_1^- & 0_{l_1} & 0_{l_1} \\ 0_{l_1} & H_2^- & -H_2^- & 0_{l_1} & 0_{l_1} & 0_{l_1} \\ 0_{l_1} & -H_2^- & H_2^- & 0_{l_1} & 0_{l_1} & 0_{l_1} \\ -H_1^- & 0_{l_1} & 0_{l_1} & H_1^- & 0_{l_1} & 0_{l_1} \\ 0_{l_1} & 0_{l_1} & 0_{l_1} & 0_{l_1} & 0_{l_1} & C_2 \cdot E_{l_1} \\ 0_{l_1} & 0_{l_1} & 0_{l_1} & 0_{l_1} & C_2 \cdot E_{l_1} & 0_{l_1} \end{pmatrix}_{6l_1 \times 6l_1},$$

$$H_1^- = \left(F_+^A X_+^A \left(X_-^{A^\top} D_-^A X_-^A \right)^{-1} X_+^{A^\top} F_+^A \right),$$

$$H_2^- = \left(\frac{1}{\gamma_2} F_+^B X_+^B \left(X_-^{B^\top} D_-^B X_-^B \right)^{-1} X_+^{B^\top} F_+^B \right),$$

$$p_-^\top = \left(-e_+^\top F_+^A \quad -e_+^\top F_+^B \quad 0_{1 \times l_1} \quad 0_{1 \times l_1} \quad 0_{1 \times l_1} \right)_{1 \times 6l_1},$$

$$A_- = \begin{pmatrix} E_{l_1} & 0_{l_1} & E_{l_1} & 0_{l_1} & 0_{l_1} & -C_2 \cdot E_{l_1} \\ 0_{l_1} & E_{l_1} & 0_{l_1} & E_{l_1} & -C_2 \cdot E_{l_1} & 0_{l_1} \end{pmatrix}_{2l_1 \times 6l_1},$$

$$b_-^\top = \left(C_{A2} \cdot e_+^\top \quad C_{B2} \cdot e_+^\top \right)_{1 \times 2l_1},$$

E_{l_1} is the $l_1 \times l_1$ identity matrix, 0_{l_1} is the $l_1 \times l_1$ matrix with all entries be 0 and e_+ is a column vector with the proper dimension with element 1.

A new data point $x \in \mathbb{R}^n$ is assigned to class r ($r=1,2$), depending on which plane it is nearer to. We first have the decision function of view A and view B respectively:

$$\text{class}^A(x_A) = \arg \min_{r=1,2} \left(\frac{|x_A^\top w_r^A + b_r^A|}{\|w_r^A\|} (x_A) \right) \tag{33}$$

and

$$\text{class}^B(x_B) = \arg \min_{r=1,2} \left(\frac{|x_B^\top w_r^B + b_r^B|}{\|w_r^B\|} (x_B) \right). \quad (34)$$

Then the decision function combining two views can be given below:

$$\text{class}(x) = \arg \min_{r=1,2} (d_r(x)), \quad (35)$$

where

$$d_r(x) = \frac{1}{2} \left(\frac{|x_A^\top w_r^A + b_r^A|}{\|w_r^A\|} + \frac{|x_B^\top w_r^B + b_r^B|}{\|w_r^B\|} \right).$$

For the sake of perspicuity, we explicitly express the MPWTSVM algorithm in Algorithm 1.

Algorithm 1 QP Algorithm for MPWTSVM.

Input: Training datasets $S = \{(\mathbf{x}_i^A, \mathbf{x}_i^B, y_i)\}_{i=1}^l = \{((x_i^A; 1), (x_i^B; 1), y_i)\}_{i=1}^l$, where label $y_i \in \{-1, 1\}$ and the testing sample, \mathbf{x} ;

Initial parameters $\gamma, C_A, C_B, C, D \geq 0$.

Output: Decision function as in (33),(34),and(35).

1. Choose two appropriate kernels $\mathcal{K}_A(x_i^A, x_j^A)$, $\mathcal{K}_B(x_i^B, x_j^B)$ and initialize the kernel parameters.
 2. Establish and solve QPPs of (31) and (32) by using 5-fold cross-validation and choose the best parameters.
 3. Construct the separating hyperplanes $\mathbf{w}_+^{A\top} \phi_A(x_A) = 0$, $\mathbf{w}_-^{A\top} \phi_A(x_A) = 0$, $\mathbf{w}_+^{B\top} \phi_B(x_B) = 0$, $\mathbf{w}_-^{B\top} \phi_B(x_B) = 0$, and $0.5(\mathbf{w}_+^{A\top} \phi_A(x_A) + \mathbf{w}_-^{A\top} \phi_A(x_A)) + 0.5(\mathbf{w}_+^{B\top} \phi_B(x_B) + \mathbf{w}_-^{B\top} \phi_B(x_B)) = 0$.
 4. For a new testing point \mathbf{x} , predict its label according to the decision functions (33) or (34) respectively and (35) collectively.
-

3.3. Nonlinear MPWTSVM

In this section, we extend linear MPWTSVM to the nonlinear case. The kernel-generated hyperplanes are:

$$\mathcal{K}(x_+^A, C^A)w_+^A + b_+^A = 0; \quad \mathcal{K}(x_-^A, C^A)w_-^A + b_-^A = 0; \quad (36)$$

$$\mathcal{K}(x_+^B, C^B)w_+^B + b_+^B = 0; \quad \mathcal{K}(x_-^B, C^B)w_-^B + b_-^B = 0; \quad (37)$$

where \mathcal{K} is a chosen kernel function defined by $\mathcal{K}(x_i, x_j) = (\phi(x_i) \cdot \phi(x_j))$. $\phi(\cdot)$ is a nonlinear mapping that maps the low-dimensional feature space to the high-dimensional feature space in a non-linear manner. C denotes training examples from view A and view B respectively. $C^A = [X_1^A; X_2^A]$ and $C^B = [X_1^B; X_2^B]$, so that positive examples from view t ($t=A, B$) are denoted as Ω_+^t and negative examples from view t are denoted as Ω_-^t .

We can define:

$$\Omega_+^t = \mathcal{K}(x_+^t, C^t), \quad \Omega_-^t = \mathcal{K}(x_-^t, C^t),$$

$$\mathbf{w}_+^t = \begin{pmatrix} \omega_+^t \\ b_+^t \end{pmatrix}, \quad \mathbf{w}_-^t = \begin{pmatrix} \omega_-^t \\ b_-^t \end{pmatrix}.$$

In order to simplify the calculation, we update the matrices above:

$$\Omega_+^t = (\mathcal{K}(x_+^t, C^t), e_+); \quad \Omega_-^t = (\mathcal{K}(x_-^t, C^t), e_-).$$

Then the optimization problems for non-linear MPWTSVM can be formulated as

$$\begin{aligned} \min \quad & \frac{1}{2} \sum_{i=1}^M \sum_{j=1}^M W_{s,ij}^{A_1} (\mathbf{w}_+^{A\top} \Omega_i^{A,+})^2 + \frac{1}{2} \gamma \sum_{i=1}^M \sum_{j=1}^M W_{s,ij}^{B_1} (\mathbf{w}_+^{B\top} \Omega_i^{B,+})^2 \\ & + C_A \sum_{j=1}^N \xi_j^A + C_B \sum_{j=1}^N \xi_j^B + C \sum_{j=1}^N \xi_j^A \xi_j^B \\ \text{s.t.} \quad & -f_j^A (\mathbf{w}_+^{A\top} \Omega_j^{A,-}) + \xi_j^A \geq f_j^A \cdot 1, \\ & -f_j^B (\mathbf{w}_+^{B\top} \Omega_j^{B,-}) + \xi_j^B \geq f_j^B \cdot 1, \\ & \xi_j^A \geq -f_j^B (\mathbf{w}_+^{B\top} \Omega_j^{B,-}), \quad \xi_j^A \geq 0, \\ & \xi_j^B \geq -f_j^A (\mathbf{w}_+^{A\top} \Omega_j^{A,-}), \quad \xi_j^B \geq 0, \quad (j \in I^-); \end{aligned} \tag{38}$$

and

$$\begin{aligned}
\min \quad & \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N W_{s,ij}^{A_2} (\mathbf{w}_-^{A\top} \Omega_j^{A,-})^2 + \frac{1}{2} \gamma \sum_{i=1}^N \sum_{j=1}^N W_{s,ij}^{B_2} (\mathbf{w}_-^{B\top} \Omega_j^{B,-})^2 \\
& + C_{A_2} \sum_{i=1}^M \xi_i^A + C_{B_2} \sum_{i=1}^M \xi_i^B + C_2 \sum_{i=1}^M \xi_i^A \xi_i^B \\
s.t. \quad & f_i^A(\mathbf{w}_-^{A\top} \Omega_i^{A,+}) + \xi_i^A \geq f_i^A \cdot 1, \\
& f_i^B(\mathbf{w}_-^{B\top} \Omega_i^{B,+}) + \xi_i^B \geq f_i^B \cdot 1, \\
& \xi_i^A \geq f_i^B(\mathbf{w}_-^{B\top} \Omega_i^{B,+}), \quad \xi_i^A \geq 0, \\
& \xi_i^B \geq f_i^A(\mathbf{w}_-^{A\top} \Omega_i^{A,+}), \quad \xi_i^B \geq 0, \quad (i \in I^+). \tag{39}
\end{aligned}$$

4. Comparison with other algorithms

In this section, we compare our MPWTSVM with SVM-2K [26], MVTSVM [27], MCPK [33] and PSVM-2V [32]. The complexity analysis is also included in the following comparisons. For simplicity, we suppose the number of samples of each class are equal, namely $l_1 = l_2 = l/2$, where l represents the number of training samples. Problem (31) and (32) involve two convex QPPs. Both of them can be worked out by the classical QP solver with a time complexity less than $2\mathcal{O}((3l)^3)$ for the reason that the inter-class weight matrix can remove redundant samples and greatly reduce the time complexity.

4.1. MPWTSVM vs. SVM-2K

SVM-2K solves a QPP, which combines two-stage learning and SVM into a single optimization. Moreover, it only satisfies the consistency principle of multi-view training by using Kernel Canonical Correlation Analysis (KCCA) theory [34] and does not satisfy the principle of complementarity. The time complexity of SVM-2K is $\mathcal{O}((4l)^3)$. Our model has better efficiency compared to SVM-2K. Our MPWTSVM solves two QPPs which makes it work faster than SVM-2K. Besides, the introduction of KNN makes it more efficient to identify the potential support vector. Based on satisfying the principle of consistency, we met the principle of complementarity with the help of privileged information.

4.2. MPWTSVM vs. MVT SVM

Similar to MVT SVM, our model solves two QPPs. The time complexity of MVT SVM is about $2 \times \mathcal{O}((2l)^3)$. MVT SVM combines two views by bringing in the similarity constraint between the two-dimensional projections of two different TSVMs from the two feature spaces. It is only applied to the classification of the two views, which cannot solve the general multi-view problem. The supplementary information between different views cannot be effectively used so that the model does not satisfy the principle of complementarity. Our model makes full use of each perspective as privileged information to redefine slack variables, thereby satisfying the principle of complementarity.

4.3. MPWTSVM vs. MCPK

Compared with MCPK, our model shares with it that they both satisfy both consistency and complementarity principles. The difference is that MCPK solves one QPP, while MPWTSVM solves two problems. The time complexity of MCPK is $\mathcal{O}((6l)^3)$. Our model has better efficiency compared to MCPK. In addition, our model extends the weighting idea of WLTSVM to different perspectives to measure as much similarity information between samples as possible and obtain higher accuracy.

4.4. MPWTSVM vs. PSVM-2V

Both PSVM-2V and our model satisfy the complementarity and the consistency principle simultaneously. PSVM-2V solves a QPP. The time complexity of PSVM-2V is $\mathcal{O}((6l)^3)$. Our model has better efficiency compared to PSVM-2V. PSVM-2V uses regularization terms to limit the differences of prediction results from different perspectives to achieve the consistency principle, which is achieved by the coupling terms in the objective function of our model. PSVM-2V and our model realize the principle of complementarity both with the help of privileged information. However, compared to PSVM-2V, our model also draws on the idea of weighting, which enables better preservation of the inter and intra connections and differences of different views in the data.

5. Experiments

In this section, we make comparisons between our MPWTSVM and five benchmark methods, SVM+, SVM-2K, MVT SVM, PSVM-2V, and MCPK.

We verify the performance of MPWTSVM for binary classification on 45 datasets obtained from *Animals with Attributes (AwA)*¹[35]. In order to eliminate the influence of size and simplify the numerical calculation, we scale all the features to the range of [0, 1] in the data preprocessing. The experiments are conducted in Matlab R2015a on Windows 7 running on a PC with system configuration Inter(R) Core(TM) i7-6700 CPU (3.40GHz) with 8.00 GB of RAM.

5.1. Experimental setup

Dataset. The AwA dataset consists of 30,475 images in 50 animal categories, each image has six pre-extracted features re-represented. The detailed characteristics of the data set are demonstrated in Table 2. Similar to the study[36], we use the ten classes in AwA dataset, i.e. *antelope, grizzly bear, killer whale, beaver, dalmatian, Persian cat, horse, German shepherd, blue whale and Siamese cat*. Through the one-to-one strategy, we randomly select 200 samples in each class for training and train 45 binary classifiers for each class pair combination.

Table 2

Comprehensive information of datasets we use in the experiments.

Data set	#Data	#Classes	#Features (View A)	#Features (View B)	#Binary data sets
AwA	6249	10	252	2000	45 (one-versus-one)

Kernels. In these experiments, we choose the Gaussian radial basis function (RBF) $\mathcal{K}(x_i, x_j) = \exp(-\|x_i - x_j\|^2/\sigma^2)$ for all the algorithms since the RBF is most widely used in the classification problem [31, 37, 33, 38, 39].

Measures. We assess the performance of the classifier by the test accuracy. We perform the grid search strategy and 5-fold cross-validation for all datasets to select the optimal parameters. For methods other than SVM+, we ponder on the mixed prediction function $\text{sign}(0.5(f_A(x_A) + f_B(x_B)))$ in addition to two views' prediction functions $\text{sign}(f_A(x_A))$ and $\text{sign}(f_B(x_B))$, and choose the one who has the highest precision.

Benchmark method. We make comparisons between the proposed method and five of the most recent methods:

¹Available at <https://cvml.ist.ac.at/AwA2/>.

- 1) SVM₊ : The SVM+ algorithm [40] replaces the standard SVM’s slack variable by using the non-negative correction function determined by the privilege information. During the training process, we separately use the two views as privileged information for each other.
- 2) SVM-2K: SVM-2K combines two-stage learning—KCCA followed by SVM, into a single optimization [26].
- 3) MVTSVM: MVTSVM merges two perspectives by introducing similarity constraints between two one-dimensional projections [27]. The model learns two hyperplanes and solves a pair of QPPs rather than one.
- 4) PSVM-2V: PSVM-2V extends LUPI (learning using privileged information) to MVL [32].
- 5) MCPK: MCPK is a simple and effective approach to MVL coupling privilege kernels and satisfies two principles for MVL [33].

Parameters. For all algorithms, the optimal parameters are decided by five-fold cross validation. The parameter C in SVM₊ is selected from the set $\{10^{-3}, 10^{-2}, 10^{-1}, 1, 10^1, 10^2, 10^3\}$. Penalty parameters C_A, C_B and C of SVM-2K and PSVM-2V are selected from $\{10^{-3}, 10^{-2}, 10^{-1}, 1, 10^1, 10^2, 10^3\}$. For MVTSVM, we let $C_1 = C_2 = C_3 = C_4$ and $D = H$, and both of them are turned over the set $\{10^{-3}, 10^{-2}, 10^{-1}, 1, 10^1, 10^2, 10^3\}$. For MCPK, we set $C_A = C_B = D$ varying in the set $\{10^{-3}, 10^{-2}, 10^{-1}, 1, 10^1, 10^2, 10^3\}$. For MPWTSVM, we set $C_A = C_B = C = D$ from $\{10^{-3}, 10^{-2}, 10^{-1}, 1, 10^1, 10^2, 10^3\}$. The neighborhood size k is searched within $\{3, 5, 7, 9, 11\}$. Moreover, the trade-off parameter γ in SVM₊, PSVM-2V, MCPK and MPWTSVM is chosen from the set $\{10^{-3}, 10^{-2}, 10^{-1}, 1, 10^1, 10^2, 10^3\}$. The kernel parameter σ for RBF kernel function is chosen from $\{10^{-3}, 10^{-2}, 10^{-1}, 1, 10^1, 10^2, 10^3\}$. For the sake of simplicity, in the multi-view model, the kernel parameters of the two views are set to be the same.

5.2. Parameter analysis

To study the influence of parameters, we examine the parameter sensitivity of MPWTSVM on nine datasets of AwA. The ranges of variation parameters C, γ, σ are the same as that given in Section 5.1, and we calculate the test data’s accuracy.

The consequences are shown in Fig.2. It depicts the values of the main parameter k for each combination of parameters C and the kernel parameter corresponding to the highest accuracy rate, in which we control the remaining parameters to be consistent with C to simplify the calculation and graphing. These graphs show that the accuracy of MPWTSVM has something to do with the parameters k , C , $\sigma(ker)$ on all datasets and is sensitive to them all. Therefore, these parameters should be carefully adjusted.

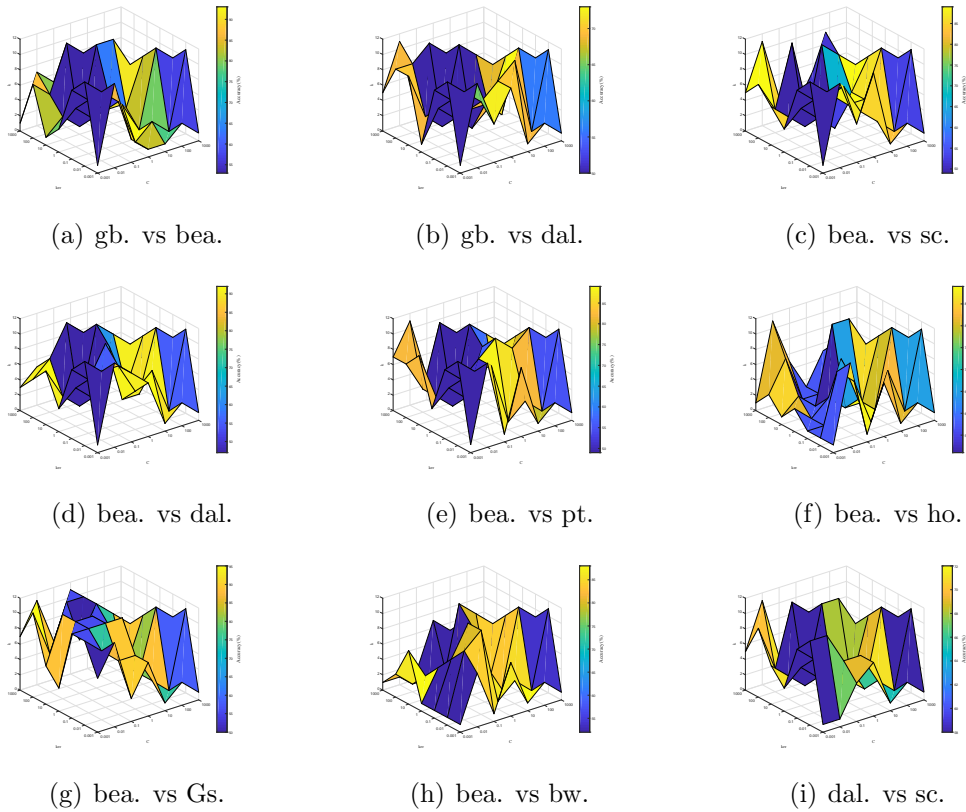


Figure 2: The figure indicates the highest accuracy k value corresponding to each combination of C and ker on the datasets from AWA. The colors indicate the degree of accuracy.

5.3. Experimental results

Below, we make the capability of MPWTSVM and all benchmark algorithms under the case of nonlinearity in comparison. Table 3 gives a summary

of the detailed information of the overall experimental results on 45 multi-view datasets, including the accuracy, time, and average rank. ‘Accuracy’ represents the average of five test results, plus or minus the standard deviation. ‘Time’ means the average training time of five experiments. The bold values in Table 3 demonstrate the best accuracy.

Table 3

Performance on *AwA* dataset (average accuracy \pm standard deviation of accuracy(%)).

		SVM+A	SVM+B	SVM-2K	MVTSVM	MCPK	PSVM-2V	MPWTSVM
1	ante. vs. sc.	58.50 \pm 3.791	84.50 \pm 3.710	68.00 \pm 11.374	63.50 \pm 4.183	82.50 \pm 6.374	84.50 \pm 4.108	86.00\pm6.021
2	ante. vs. gb.	67.00 \pm 3.260	82.50 \pm 7.071	75.00 \pm 3.062	72.50 \pm 6.847	82.00 \pm 4.472	85.00 \pm 2.500	85.50\pm4.108
3	ante. vs. kw.	79.00 \pm 8.023	88.50 \pm 5.477	83.50 \pm 5.184	77.00 \pm 12.550	89.00 \pm 3.791	88.00 \pm 1.118	91.50\pm3.791
4	ante. vs. bea.	88.50 \pm 5.477	98.00 \pm 2.092	92.00 \pm 5.420	91.00 \pm 2.850	98.50\pm2.236	96.50 \pm 4.183	98.00 \pm 1.118
5	ante. vs. dal.	73.50 \pm 6.982	80.50 \pm 6.708	73.00 \pm 4.809	72.50 \pm 8.839	84.00 \pm 7.416	84.00\pm3.354	81.00 \pm 7.624
6	ante. vs. pt.	72.50 \pm 4.677	88.50 \pm 8.023	75.50 \pm 5.701	76.00 \pm 4.873	88.00 \pm 3.708	87.00 \pm 2.092	90.00\pm5.590
7	ante. vs. ho.	68.00 \pm 8.551	84.00\pm7.416	76.00 \pm 8.023	69.50 \pm 7.786	83.50 \pm 4.183	82.00 \pm 5.123	83.50 \pm 3.791
8	ante. vs. Gs.	73.00 \pm 9.747	92.50\pm3.953	78.50 \pm 17.375	74.50 \pm 2.739	89.00 \pm 1.369	91.50 \pm 4.183	92.00 \pm 2.739
9	ante. vs. bw.	75.50 \pm 4.809	84.50 \pm 9.585	75.50 \pm 5.969	73.00 \pm 5.420	83.50 \pm 5.184	82.00 \pm 8.178	84.50\pm6.471
10	gb. vs. sc.	71.00 \pm 10.548	87.00 \pm 6.225	79.50 \pm 7.159	75.00 \pm 6.124	86.50 \pm 2.850	84.50 \pm 4.472	88.50\pm2.850
11	gb. vs. kw.	74.00 \pm 9.117	82.00 \pm 9.906	77.00 \pm 4.809	80.50 \pm 7.583	84.50 \pm 4.809	83.00 \pm 5.420	85.50\pm2.092
12	gb. vs. bea.	86.50 \pm 4.183	91.50 \pm 5.755	84.00 \pm 5.755	87.50 \pm 6.374	92.50\pm2.500	91.50 \pm 4.541	92.50 \pm 3.536
13	gb. vs. dal.	68.00 \pm 8.551	76.00 \pm 6.755	73.00 \pm 5.420	74.00 \pm 7.202	76.00 \pm 5.184	77.00 \pm 6.471	79.50\pm7.159
14	gb. vs. pt.	73.50 \pm 8.944	85.00 \pm 7.071	81.00 \pm 8.944	83.50 \pm 4.873	86.50 \pm 8.768	87.50 \pm 4.677	90.50\pm5.969
15	gb. vs. ho.	59.00 \pm 8.944	70.00 \pm 8.292	62.50 \pm 15.910	62.50 \pm 7.500	65.00 \pm 10.155	73.50 \pm 6.519	76.50\pm6.755
16	gb. vs. Gs.	76.00 \pm 2.236	83.00 \pm 4.108	67.50 \pm 8.478	76.00 \pm 5.755	83.50 \pm 6.021	84.50 \pm 6.937	87.50\pm5.000
17	gb. vs. bw.	68.00 \pm 7.374	81.00 \pm 3.791	77.50 \pm 3.062	77.00 \pm 4.809	81.50 \pm 6.021	81.50 \pm 6.755	88.50\pm3.354
18	kw. vs. sc.	80.50 \pm 4.809	87.00 \pm 2.739	83.00 \pm 6.471	81.00 \pm 7.202	86.00 \pm 1.369	89.00 \pm 3.354	89.50\pm4.108
19	kw. vs. bea.	81.50 \pm 6.519	87.00 \pm 4.809	86.50 \pm 6.519	85.00 \pm 5.590	89.00 \pm 4.183	90.50 \pm 2.092	92.00\pm2.092
20	kw. vs. dal.	68.00 \pm 6.708	70.50 \pm 6.708	70.50 \pm 4.108	72.00 \pm 2.092	72.00 \pm 13.393	74.00 \pm 3.354	75.50\pm2.739
21	kw. vs. pt.	74.00 \pm 6.275	83.50 \pm 2.850	77.00 \pm 4.472	71.50 \pm 8.944	86.00 \pm 2.850	86.50\pm4.183	86.00 \pm 3.354
22	kw. vs. ho.	70.50 \pm 7.159	83.50 \pm 3.791	82.50 \pm 8.101	76.50 \pm 6.021	86.00 \pm 6.275	84.50 \pm 7.159	88.00\pm2.739
23	kw. vs. Gs.	69.00 \pm 5.184	87.50\pm5.863	74.50 \pm 5.701	73.00 \pm 3.26	84.50 \pm 5.420	85.50 \pm 3.260	87.00 \pm 5.701
24	kw. vs. bw.	75.50 \pm 3.708	84.00 \pm 6.755	78.00 \pm 6.225	77.00 \pm 5.701	83.00 \pm 8.178	82.00 \pm 4.809	86.50\pm2.236
25	bea. vs. sc.	89.50 \pm 4.108	95.00 \pm 4.330	88.50 \pm 2.236	93.00 \pm 2.092	96.50 \pm 1.369	95.00 \pm 3.536	97.00\pm1.118
26	bea. vs. dal.	91.50 \pm 5.184	96.00 \pm 2.850	95.50 \pm 3.708	93.00 \pm 3.260	96.50 \pm 2.236	95.50 \pm 2.092	97.50\pm3.062
27	bea. vs. pt.	88.50 \pm 3.354	93.50 \pm 2.850	92.50 \pm 3.953	90.00 \pm 3.062	95.50\pm3.260	93.50 \pm 3.791	95.50\pm3.260
28	bea. vs. ho.	88.50 \pm 5.184	98.00 \pm 2.092	91.00 \pm 5.755	91.00 \pm 2.236	98.00 \pm 2.092	97.50 \pm 2.500	98.50\pm1.369

continued table 3

29	bea. vs. Gs.	92.00±5.701	94.00±6.755	94.00±3.354	93.00±3.708	95.50±3.260	96.50±1.369	97.00±2.092
30	bea. vs. bw.	90.50±2.092	95.00±4.677	89.00±2.236	92.00±2.739	97.00±1.118	96.00±2.85	98.50±2.236
31	dal. vs. sc.	76.00±6.275	85.50±7.583	74.50±5.969	73.50±5.477	87.50±6.374	85.50±4.809	89.00±2.236
32	dal. vs. pt.	74.00±4.183	86.50±3.354	79.50±9.906	77.00±4.108	88.00±2.092	89.00±6.275	89.50±4.809
33	dal. vs. ho.	67.00±11.096	76.00±5.184	76.00±6.755	68.00±7.984	82.00±4.472	78.50±6.519	78.50±6.519
34	dal. vs. Gs.	66.50±5.477	84.00±5.755	69.50±8.178	71.00±6.982	81.50±3.791	81.50±5.755	85.00±3.953
35	dal. vs. bw.	65.00±5.303	78.50±6.021	73.50±16.919	69.50±11.646	84.50±7.159	84.50±3.26	82.00±3.260
36	pt. vs. sc.	60.00±4.677	78.50±3.354	65.00±1.768	65.00±5.303	77.00±7.583	76.50±4.541	78.00±5.701
37	pt. vs. ho.	70.50±4.809	92.50±3.062	87.50±0.000	76.50±5.755	89.50±4.472	92.00±2.739	91.50±5.755
38	pt. vs. Gs.	68.00±2.739	85.00±6.374	81.00±10.093	72.00±7.583	84.00±4.183	82.50±6.124	87.50±5.303
39	pt. vs. bw.	71.50±3.791	84.00±7.416	75.50±3.708	72.00±7.374	88.00±3.708	87.00±8.551	89.00±4.183
40	ho. vs. sc.	67.00±3.260	83.50±4.541	76.00±7.202	67.50±5.590	83.00±3.260	80.50±4.472	86.00±4.873
41	ho. vs. Gs.	64.00±6.275	85.00±3.953	70.50±8.551	71.00±6.755	84.50±9.906	85.00±3.953	84.50±5.969
42	ho. vs. bw.	66.50±3.791	73.50±7.202	68.00±4.108	69.00±8.944	75.50±8.178	75.00±5.000	76.00±4.183
43	Gs. vs. sc.	75.00±8.292	90.50±3.708	71.00±5.477	74.00±6.755	90.50±3.260	90.50±3.260	93.00±2.092
44	Gs. vs. bw.	65.00±5.303	85.50±9.253	68.00±7.786	70.50±7.159	83.50±4.873	83.00±4.472	87.50±4.677
45	bw. vs. sc.	64.50±8.178	73.00±6.225	70.50±7.374	72.50±6.124	77.00±9.083	73.50±5.477	73.50±4.873
	Avg.Acc.	73.59	85.22	77.94	76.50	85.72	85.64	87.57
	Avg.Time.	0.111	0.138	0.347	0.067	0.088	0.858	0.042
	Avg.Rank.	6.678	3.011	5.500	5.756	2.678	2.911	1.467
	*W/D/L	45/0/0	37/2/6	45/0/0	45/0/0	35/5/5	38/2/5	0/45/0

*W/D/L is short for Win/Draw/Loss.

The average accuracy of 45 datasets of MPWTSVM is 87.57%, which is the highest among the seven algorithms followed by MCPK (85.72%), PSVM-2V (85.64%), SVM+A (85.22%), SVM-2K (77.94%), MVT SVM (76.50%) and SVM+B (73.59%). MPWTSVM, MCPK and PSVM-2V can follow both the consensus principle and the complementarity principle. Therefore, the three algorithms perform better than the other four. It is worth noting that our model has better generalization performance compared to MCPK and PSVM-2V. The reason is that it utilizes the weighting idea to better exploit the intrinsic connections and differences within and between different perspectives. SVM-2K and MVT SVM only satisfy the principle of consensus, so the generalization ability is slightly insufficient compared with MPWTSVM, PSVM-2V and MCPK. Among the several algorithms, the one with the lowest accuracy is SVM+, since the model considers only one perspective and

uses the other perspective as its privileged information.

The average time of MPWTSVM is 0.042s which is the shortest one compared to MVTSSVM (0.067s), MCPK (0.088s), SVM+A (0.111s), SVM+B (0.138s), SVM-2K (0.347s), and PSVM-2V (0.858s). The reason is that the proposed MPWTSVM not only solves two small-scale QPPs, but also uses inter-class weights to remove redundant samples. Therefore, the time complexity can be greatly reduced.

For the same dataset, the accuracy of seven algorithms is ranked and the optimal one is assigned as 1, the sub optimal algorithm is designated as 2, and so on. From this, the average rank of each algorithm in 45 datasets is obtained. Our model has an lowest average rank of 1.467 followed by MCPK(2.678) and PSVM-2V(2.911) among the seven algorithms, which shows its good classification performance.

In addition, Table 3 shows the wins and losses of our algorithm compared to other algorithms. If the accuracy of our MPWTSVM is higher than the other algorithm, it is marked as ‘Win’; if it is equal to the other one, it is marked as ‘Draw’; and if it is lower than the other one, it is marked as ‘Loss’. After comparing our algorithm with other algorithms in 45 data sets, we can count the win-loss situation of each algorithm. It can be seen that compared with SVM+A, SVM-2K, and MVTSSVM, our algorithm wins 45 times in 45 data sets; compared with SVM+B, MPWTSVM wins 37 times, loses six times, and draws twice; compared with MCPK, MPWTSVM wins 35 times, loses five times and draws five times; compared with PSVM-2V, MPWTSVM wins 38 times, loses five times and draws twice. It denotes that our algorithm has more prominent advantages compared with other parallel algorithms.

Fig.3 describes the differences separating the winning algorithm’s accuracy and the remaining algorithms’ average accuracy. The length of each bar reflects the generalization performance of the algorithm, and different colors correspond to different algorithms. We can see that MPWTSVM has a better accuracy among these models.

Fig.4 depicts the comparison between the training time of the winning algorithm and the remaining algorithms. The value of the ordinate represents each method’s training time. The numbers on the abscissa indicate the number of data sets. It is due to the deletion of redundant samples with the help of the KNN idea, which reduces the running time and improves the computational efficiency. It shows the superior efficiency of MPWTSVM. These results strengthen the fact that MPWTSVM itself can make full use of intra-

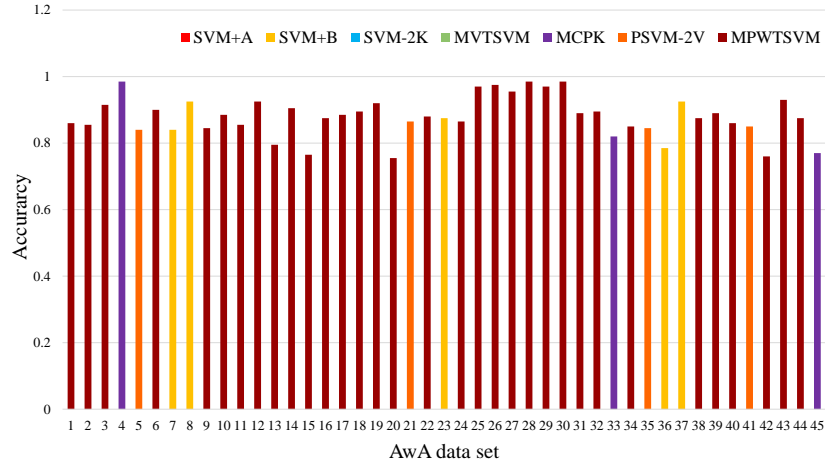


Figure 3: The plot denotes the differences separating the winning method MPWTSVM’s (dark red) classification accuracy and the remaining algorithms’ average accuracy.

class and inter-class information to obtain better classification performance, which means that MPWTSVM has good generalization ability.

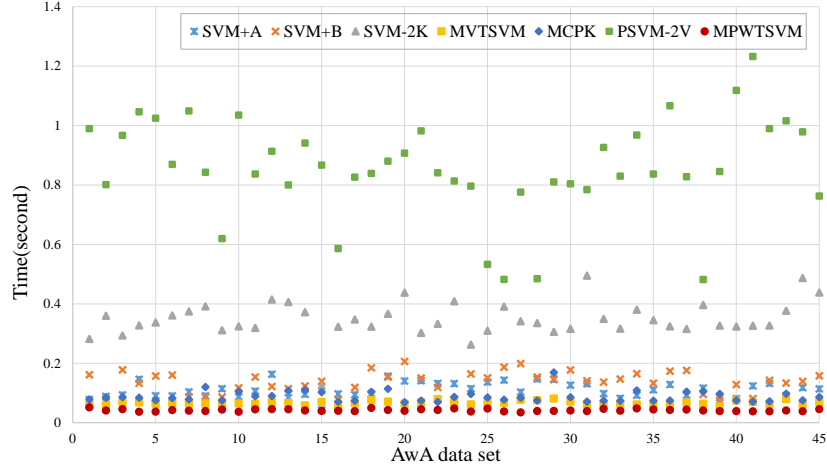


Figure 4: The plot denotes the training time of the following algorithms: SVM+A, SVM+B, SVM-2K, MVTSVM, MCPK, PSVM-2V, and MPWTSVM.

5.4. Friedman test

We use Friedman test [41] for further analysis since the averaged accuracy of our MPWTSVM in Table 3 is not always optimal. For each dataset, the

highest accuracy ranking is 1, followed by 2, and so on. Table 3 also shows the average ranking of the six comparison algorithms. As can be seen from the Avg.Rank line in Table 3, the average ranking of MPWTSVM is 1.467, which ranks the bottom of the six methods. The outcome shows that MPWTSVM proposed by us has the best performance among the six comparison methods.

The following null hypothesis is made: the 7 methods are identical. Friedman statistics can be computed using the following formula:

$$\chi_F^2 = \frac{12N}{k(k+1)} \left[\sum_j R_j^2 - \frac{k(k+1)^2}{4} \right] \quad (40)$$

where $R_j = \frac{1}{N} \sum_i r_i^j$ and r_i^j is the j th of k algorithms on the i th of N datasets.

From (40) and with the help of R language, we can get χ_F^2 is 222.11 and the P-value of the hypothesis test is 2.2×10^{-16} . If $\alpha = 0.05$, P-value is much less than α . This means that it denies the null hypothesis above. In other words, the differences between these seven methods are obvious.

Since the null hypothesis is rejected, Nemenyi test [42] can be further conducted. The hollow circle represents the average ranking of the seven algorithms, and the critical difference CD is shown in the straight lines centered on "o" in Fig.5, where

$$CD = q_\alpha \sqrt{\frac{k(k+1)}{6N}}. \quad (41)$$

We can get $q_\alpha = 2.949$ and $CD = 1.343$ under the significance level $\alpha = 0.05$ by calculation. The consequence in Fig.5 shows that the proposed MPWTSVM has remarkably better performance than the other models at the confidence level of 95% .

6. Conclusions

In this paper, we propose a weighted MVL kernel method based on privileged information termed MPWTSVM, which satisfies the principle of consensus and complementarity at the same time. The consensus principle is ensured by minimizing the coupling items of the two views in the original goal. We implement the principle of complementarity by establishing a privileged information paradigm and by drawing on multi-view learning. In the two-classification process, we weigh the samples from two perspectives to obtain

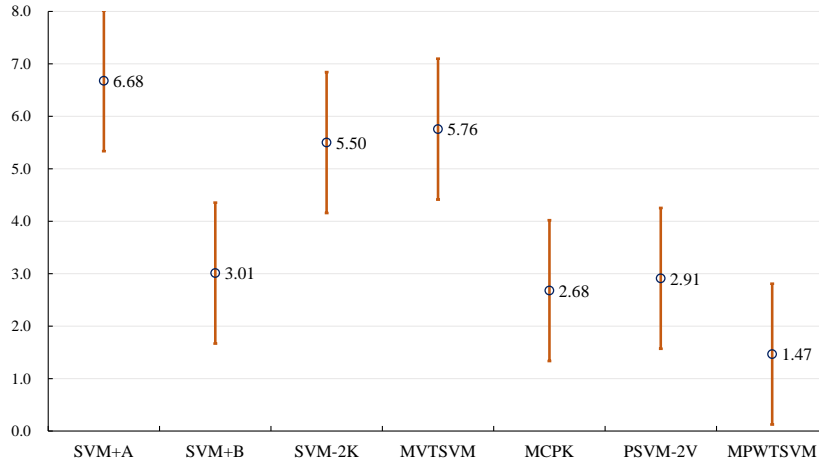


Figure 5: Friedman test.

the weights of different types and the same type in each view. The internal similarity information of samples of the same category is captured. During the learning process, our model entirely integrates the information of diverse views and maintains the characteristics of diverse views to a certain extent. Therefore, the proposed MPWTSVM has better classification accuracy and faster solving speed. We use the standard QP solver to solve MPWTSVM and compare it with the other five algorithms. Through numerical experiments on 45 sets of multi-view binary datasets demonstrate the validity and efficiency of our MPWTSVM. In the future, we will extend MPWTSVM to the case of multiple (more than two) views under the guideline of two principles. Other effective optimization algorithms like the alternating direction methods of multipliers(ADMM) and other solution ways can also be put into consideration before long.

Acknowledgments

This work was supported in part by the Fundamental Research Funds for the Central Universities (No. BLX201928).

References

- [1] V. N. Vapnik, The nature of statistical learning theory, Springer, Berlin, 1995.

- [2] V. Cherkassky, The nature of statistical learning theory , IEEE Transactions on Neural Networks 8 (6) (1997) 1564–1564.
- [3] O. Mangasarian, E. Wild, Multisurface proximal support vector machine classification via generalized eigenvalues, IEEE transactions on pattern analysis and machine intelligence 28 (2006) 69–74.
- [4] Jayadeva, R. Khemchandani, S. Chandra, Twin support vector machines for pattern classification, IEEE Trans. Pattern Anal. Mach. Intell. 29 (5) (2007) 905–910.
- [5] S. Ghorai, A. Mukherjee, P. K. Dutta, Nonparallel plane proximal classifier, Signal Processing 89 (4) (2009) 510–522.
- [6] Q. Ye, C. Zhao, S. Gao, H. Zheng, Weighted twin support vector machines with local information and its application, Neural Networks 35 (2012) 31–39.
- [7] Q. Ye, C. Zhao, N. Ye, X. Chen, Localized twin svm via convex minimization, Neurocomputing 74 (4) (2011) 580–587.
- [8] Y. Xu, J. Yu, Y. Zhang, Knn-based weighted rough ν -twin support vector machine, Knowledge-Based Systems 71 (2014) 303–313.
- [9] M. Tanveer, A. Sharma, P. N. Suganthan, Least squares knn-based weighted multiclass twin svm, Neurocomputing 459 (2021) 454–464.
- [10] J. A. Nasiri, A. M. Mir, An enhanced knn-based twin support vector machine with stable learning rules, Neural Computing and Applications 32 (16) (2020) 12949–12969.
- [11] X. Zhang, L. Zhao, L. Zong, X. Liu, H. Yu, Multi-view clustering via multi-manifold regularized nonnegative matrix factorization, in: 2014 IEEE International Conference on Data Mining, 2014, pp. 1103–1108.
- [12] X. Liu, L. Wang, J. Zhang, J. Yin, Sample-adaptive multiple kernel learning, Proceedings of the National Conference on Artificial Intelligence 3 (2014) 1975–1981.
- [13] J. Chen, G.-B. Huang, Dual distance adaptive multiview clustering, Neurocomputing 441 (2021) 311–322.

- [14] J. Ma, Y. Zhang, L. Zhang, Discriminative subspace matrix factorization for multiview data clustering, *Pattern Recognition* 111 (2021) 107676.
- [15] J. Xu, J. Han, F. Nie, X. Li, Multi-view scaling support vector machines for classification and feature selection, *IEEE Transactions on Knowledge and Data Engineering* 32 (7) (2020) 1419–1430.
- [16] Y. Kaloga, P. Borgnat, S. P. Chepuri, P. Abry, A. Habrard, Variational graph autoencoders for multiview canonical correlation analysis, *Signal Processing* 188 (2021) 108182.
- [17] N. Chen, J. Zhu, E. Xing, Predictive subspace learning for multi-view data: a large margin approach, 2010, pp. 361–369.
- [18] J. Sun, H. Fujita, Y. Zheng, W. Ai, Multi-class financial distress prediction based on support vector machines integrated with the decomposition and fusion methods, *Information Sciences* 559 (2021) 153–170.
- [19] X. Zhang, Y. Yang, T. Li, Y. Zhang, H. Wang, H. Fujita, Cmc: A consensus multi-view clustering model for predicting alzheimer’s disease progression, *Computer Methods and Programs in Biomedicine* 199 (2021) 105895.
- [20] Z. Zhang, X. Wei, X. Zheng, D. D. Zeng, Predicting product adoption intentions: An integrated behavioral model-inspired multiview learning approach, *Information & Management* 58 (7) (2021) 103484.
- [21] J. Zhao, X. Xie, X. Xu, S. Sun, Multi-view learning overview: Recent progress and new challenges, *Information Fusion* 38 (2017) 43–54.
- [22] S. Sun, G. Chao, Multi-view maximum entropy discrimination, 2013, pp. 1706–1712.
- [23] G. Chao, S. Sun, Consensus and complementarity based maximum entropy discrimination for multi-view classification, *Information Sciences* 367-368 (2016) 296–310.
- [24] S. Sun, A survey of multi-view machine learning, *Neural Computing and Applications* 23 (2013) 2031–2038.

- [25] C. Xu, D. Tao, Large-margin multi-view information bottleneck, *Pattern Analysis and Machine Intelligence*, IEEE Transactions on 36 (2014) 1559–1572.
- [26] J. D. R. Farquhar, H. Meng, S. Szedmak, D. R. Hardoon, J. Shawe-taylor, Two view learning: Svm-2k, theory and practice, in: *Advances in Neural Information Processing Systems*, MIT Press, 2006.
- [27] X. Xie, S. Sun, Multi-view twin support vector machines, *Intelligent Data Analysis* 19 (4) (2015) 701–712.
- [28] X. Xie, Regularized multi-view least squares twin support vector machines, *Applied Intelligence* 48 (9) (2018) 3108–3115.
- [29] H. Wang, Z. Zhou, Multi-view learning based on maximum margin of twin spheres support vector machine, *Journal of Intelligent and Fuzzy Systems* 40 (6) (2021) 11273–11286.
- [30] A. Blum, T. Mitchell, Combining labeled and unlabeled data with co-training, in: *Proceedings of the Eleventh Annual Conference on Computational Learning Theory, COLT' 98*, Association for Computing Machinery, New York, NY, USA, 1998, p. 92–100.
- [31] L. Houthuys, R. Langone, J. A. Suykens, Multi-view least squares support vector machines classification, *Neurocomputing* 282 (2018) 78–88.
- [32] J. Tang, Y. Tian, P. Zhang, X. Liu, Multiview privileged support vector machines, *IEEE Transactions on Neural Networks and Learning Systems* 29 (8) (2018) 3463–3477.
- [33] J. Tang, Y. Tian, D. Liu, G. Kou, Coupling privileged kernel method for multi-view learning, *Information Sciences* 481 (2019) 110–127.
- [34] D. R. Hardoon, S. Szedmak, J. Shawe-Taylor, Canonical correlation analysis: An overview with application to learning methods, *Neural Computation* 16 (12) (2004) 2639–2664.
- [35] Y. Xian, C. H. Lampert, B. Schiele, Z. Akata, Zero-shot learning—a comprehensive evaluation of the good, the bad and the ugly, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41 (9) (2019) 2251–2265.

- [36] S. Motiian, M. Piccirilli, D. A. Adjeroh, G. Doretto, Information bottleneck learning using privileged information for visual recognition, in: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016, pp. 1496–1505.
- [37] L. Yang, H. Dong, Support vector machine with truncated pinball loss and its application in pattern recognition, *Chemometrics and Intelligent Laboratory Systems* 177 (2018) 89–99.
- [38] P. Zhu, W. Zhu, Q. Hu, C. Zhang, W. Zuo, Subspace clustering guided unsupervised feature selection, *Pattern Recognition* 66 (2017) 364–374.
- [39] L. Houthuys, R. Langone, J. A. Suykens, Multi-view kernel spectral clustering, *Information Fusion* 44 (2018) 46–56.
- [40] V. Vapnik, A. Vashist, A new learning paradigm: Learning using privileged information, *Neural Networks* 22 (5) (2009) 544–557, advances in Neural Networks Research: IJCNN2009.
- [41] J. Demšar, Statistical comparisons of classifiers over multiple data sets, *J. Mach. Learn. Res.* 7 (2006) 1–30.
- [42] P. B. Nemenyi, *Distribution-free multiple comparisons.*, Princeton University, 1963.