



Consistency and diversity neural network multi-view multi-label learning

Dawei Zhao^{a,b}, Qingwei Gao^{a,b,*}, Yixiang Lu^b, Dong Sun^b, Yusheng Cheng^c

^a School of Computer Science and Technology, Anhui University, Hefei 230601, PR China

^b School of Electrical Engineering and Automation, Anhui University, Hefei 230601, PR China

^c School of Computer and Information, Anqing Normal University, Anqing 246011, PR China

ARTICLE INFO

Article history:

Received 4 September 2020

Received in revised form 20 December 2020

Accepted 2 February 2021

Available online 10 February 2021

Keywords:

Multi-view learning

Multi-label learning

Consistency

Diversity

Neural network

ABSTRACT

In multi-view multi-label learning, each object is represented by multiple heterogeneous data and is simultaneously associated with multiple class labels. Previous studies usually use shared subspaces to fuse multi-view representations. However, as the number of views increases, it is more challenging to capture the high-order relationships among multiple views. Therefore, a novel neural network multi-view multi-label learning framework is proposed, which is intended to solve the problem of consistency and diversity among views through a simple and effective method (CDMM). First, we build a separate classifier for each view based on the neural network method of the nonlinear kernel mapping function and require each view to learn a consistent label result. Then, we consider the diversity of individual views while learning a consistent representation among views. For this reason, we combine the Hilbert–Schmidt Independence Criterion with exploring the diversity among different views. Finally, the label correlation factor is in addition to the classification model, and the view contribution factor is added to the prediction model. A large number of comparative experiments with existing state-of-the-art solutions on benchmark multi-view multi-label learning data sets show the effectiveness of this method.

© 2021 Elsevier B.V. All rights reserved.

1. Introduction

Multi-label (ML) learning solves the problem of label ambiguity by associating a single instance with a set of labels. For example, a picture can be labeled as the “lake”, “mountain”, and “forest” at the same time, and there may be a strong correlation among these labels. The main task of ML learning is to build a learning model, which can effectively predict the possible label set of unknown objects. In the past few decades, scholars have proposed many effective ML learning algorithms in various domains, such as image annotation [1,2], video annotation [3], and bioinformatics [4,5] etc.

Traditional ML learning is to learn knowledge from a single data structure. However, in the real world, due to the increasing diversity of data collection and feature extraction methods, a single example has multiple views. This type of data is usually associated with numerous heterogeneous feature representations simultaneously, and each feature representation provides a different view of the data [6–9]. For example, in image classification, natural scene images can usually be reproduced by visual features

or described by a specific text. The main challenge of this type of task is how to effectively learn the heterogeneity among multiple views while accurately classifying data. Therefore, in the face of more complex data classification problems in real scenes, a multi-view multi-label (MVML) learning framework has emerged.

Due to the widespread existence of MVML datasets, MVML learning has become an active research area in many practical applications [10,11]. In MVML learning, each instance is represented by multiple heterogeneous feature data, which is also associated with multiple class labels. Both solutions based on multi-view (MV) learning or ML learning have their fundamental problems to be resolved. The main issues to be solved urgently in the method that focuses on MV learning are:

1. There should be consistent information representation among different views, so how to effectively solve the consistency problem among mining views and fusing the correlation between high-dimensional heterogeneous data is of paramount importance.
2. Information observed among views is different, so there is individual diversity in the information obtained. Individual diversified information mining contributes to enhancing communication among views, thereby improving the performance of the algorithm.

* Corresponding author at: School of Computer Science and Technology, Anhui University, Hefei 230601, PR China.

E-mail address: qingweigao@ahu.edu.cn (Q. Gao).

3. The structural differences of the data among the various views lead to different importance of each view, so the contribution degree of each view is also different.

Based on these problems, we divide these algorithms into two types. The first type of method is a two-step learning strategy. The first step directly uses the MV learning method to solve the MV learning problem. The second step utilizes the existing ML learning method to solve the ML learning problem. For example, Liu et al. [12] proposed a multi-view framework IrMMC based on matrix decomposition. The framework first seeks the shared representation of multiple views and then completes the classification based on the matrix of the shared feature space. Furthermore, [11] maps each view to a shared space to eliminate noise and redundancy while maintaining the sparse and manifold structure of the image data, respectively. This two-step learning strategy learning often results in sub-optimal results.

The second type of method is the joint learning strategy. They established a unified, MVML learning model to solve the problem. For example, Zhao et al. [13] introduced a predictive reliability measurement method to select samples for sharing labels to share information with other opinions in a co-training manner; Luo et al. [14] jointly extract the consistency and specificity of heterogeneous features for subspace learning; Zhang et al. [15] proposed an MVML method based on matrix factorization, which uses the complementarity among different views to obtain a common semantic representation. The complexity of this type of method is high, but the performance of the model has been significantly improved.

The main problems that need to be solved urgently in the method that focuses on ML learning are:

1. In the face of high-dimensional heterogeneous data, how to effectively predict the label set of unknown instances;
2. Effectively fusion of information among each view and mining label correlation to improve the performance of the classifier.

Based on this, some scholars have proposed solutions. One type of method is tantamount to connect all the data views in series to one data view, and then use the ML learning method to solve it. However, this concatenation strategy has the following problems: it ignores the different physical interpretations of the view data features; the concatenation strategy causes the feature dimension of the data to be too large, and the model training will overfit when the training; the other type establishes an ML classification model for each MV heterogeneous data, and unify the results of these models to construct the final prediction model. For example, Ren et al. [16] fuse multiple views into a mixed feature matrix and use low-rank structure and manifold regularization to utilize global label correlation and local smoothness. Nevertheless, parallel strategy forces all views to output consistent results, ignoring the diversity among different views. Another common problem in MVML learning is that there is an individual difference in the degree of contribution among other views.

Based on the above analysis, it can be seen that the current MVML learning mainly faces the following challenges: (1) the problems of consistency, diversity, and contribution degree among views. (2) the problem of label correlation in ML learning. Some existing learning methods that adopt step-by-step strategies ignore the information exchange between MV and ML when solving MVML problems, so they often get sub-optimal results. However, the learning method that adopts a unified strategy will have excessive model complexity or insufficient consideration of critical issues. To effectively solve the current problems facing MVML learning, a consistency and diversity neural network multi-view multi-label learning method was proposed and

named CDMM. First, we design a random single-hidden layer feedforward neural network(SLFN) to perform ML learning for each view to ensure the consistency of all views. Then, we use the Hilbert–Schmidt Independence Criterion(HSIC) [17–19] to induce diversity among different views to learn the diversity information represented by different views. Finally, in the classification, we combine the advantages of the proposed classifier to enrich the original label space with the idea of label dependence propagation; in the prediction, we make the final prediction according to the different effects of each view on the contribution of the MVML learning task.

The main contributions of the paper are as follows:

1. An MVML feedforward neural network model is constructed, different from the traditional neural network model. The CDMM does not require iteration, it is efficient and straightforward, and its parts are closely integrated.
2. CDMM has a unified framework to jointly study the view consistency and diversity issues in MVML learning, and at the same time, integrate the correlation of labels and different contributing factors of views into the classification and prediction models. In addition, a similar ensemble learning method is used to predict the final model.
3. A large number of empirical results of CDMM on the benchmark data set proves that it has certain advantages compared with some related and competitive methods.

The remainder of this paper is organized as follows. In Section 2, the related work of MVML learning is briefly introduced. Section 3 introduces the technical details of CDMM. The results of comparative experiments and specific analyses are illustrated in Section 4. Finally, Section 5 concludes this paper.

2. Related work

2.1. Multi-label learning

The difference from traditional single-label learning tasks is that the goal of ML learning is to assign multiple class labels for a single instance, which has attracted the attention and research of a large number of scholars in different machine learning tasks. According to different types of label correlations used, the existing ML methods can usually be divided into three categories. First-order: Consider that each label has its unique attributes and ignores the correlation among labels, such as Binary Relevance(BR) [1], ML-kNN [20], and LIFT [21]. FRS-LIFT [22] is an extended algorithm on LIFT, which implements label-specific feature reduction through fuzzy rough sets. However, first-order strategies often get sub-optimal results; Second-order: Consider the correlation between paired labels, such as Calibrated Label Ranking(CLR) [23], MLRL [24], but the relationship among labels in the real world is usually more complicated; high-level: Mining the correlation among all category labels. For example, Xu et al. [25] proposed an ensemble framework for training multi-label models while using a low-rank structure to capture the complex correlation among labels. MLMF [26] learns the label correlation and assumes that the learning of each label depends on the approximate output of the feature vector and other labels; LSML [27] improves the performance of multi-label classification by jointly learning label correlation and label-specific features; WML-LSC [28] divides the noisy feature matrix into a feature information-rich matrix and outliers through matrix factorization and uses a linear self-recovery model to reconstruct the original label matrix. It can not only learn robust multi-label models but also refine the noise data of features and labels to a certain extent. In addition, most of the existing ML methods on the hypothesis of label correlation are considered global label correlation. Huang

et al. [29] first proposed the concept of local label correlation when solving ML learning problems.

Although the above methods have reached the most advanced level in processing ML data, they mainly focus on single-view data. There have been some ML learning methods to solve the MV problem. However, they did not explore the heterogeneous relationship among each view, which is of great importance to the successful construction of multi-view learning models [16,30]. Besides, label correlations as an essential influencing factor are also widely used in MVML learning [31,32].

2.2. Multi-view multi-label learning

Since a single-view can no longer effectively address the issue of data on diversification, MVML learning has become an active research field in many practical applications [10,11]. In this section, we classify the MVML algorithm from the different view fusion periods: early fusion and late fusion [30].

Early fusion: The typical way is to learn a shared space representation for all heterogeneous feature data, and then construct an ML classification model based on the shared information. For example, Zhu et al. [33] proposed an MVML learning method for incomplete views. First, learn label-specific features through global and local label correlation, then use a low-rank hypothesis matrix to recover incomplete views, and finally use consistent multi-view representations for different complementary information of the view is encoded. LSA-MML [15] solves the MVML learning problem based on the premise that there is a common representation among different views and obtains undiscovered latent semantics through alignment among different views in the kernel space. The goal of MVMLE [8] and LSA-MML is to use the HSIC during the mapping process to maintain consensus on the multi-view potential space. MVD-ELM [34] is a multi-view deep neural network based on extreme learning machine local receptive fields(ELM-LRF) to realize a fast and high-quality projection feature learning method for three-dimensional shapes. It is used to solve the problem of deep learning applications in ML 3D shape segmentation. Tan et al. [35] proposed an individual and commonality-based MVML learning(ICM2L) method to explicitly explore the individuality and commonality information about MVML data in a unified model.

The above method has the following two different problems. (1) This type of method maps each view to a shared subspace to find shared information among all views. However, in this process, there is usually no communication among the various views, and it is challenging to ensure that the shared semantic information among the views is fully utilized. (2) The potential fact that each view has a different specific contribution to ML prediction is ignored.

Therefore, some scholars attempt to solve the above problems. SIMM [36] optimizes a confusing adversarial loss and an ML loss to extract the shared information among views. Secondly, orthogonal constraints are in addition to use the view-private discrimination information, and finally, pass the synergy of shared and private information for semantic learning. TMV-LE [37] explores the high-order relationship among multiple views by constructing the mutual constraints between the tensor factorization and the mapping matrix. Then, the multi-view is utilized for more comprehensively mine the topological structure in the feature space and migrate it to the label space to obtain the label distribution. The key to the effectiveness of the early fusion strategy is how to learn an effective common representation, but it is not easy to learn an accurate, common representation as to the number and types of views increases.

Late fusion: this type of method usually constructs an ML classifier for all views that combine the classification results

Table 1

Notations and their corresponding definitions.

Notation	Description
$\mathbf{X} = \{\mathbf{X}^v\}_{v=1}^V$	Multi-view matrices
\mathbf{Y}	The label matrix
\mathbf{A}	Label correlation matrix
\mathbf{T}	Label completion matrix
\mathbf{C}	The Laplacian matrix of \mathbf{A}
\mathbf{H}	The hidden layer output matrix
β	The neural network weight coefficient
\mathbf{K}	The Gram matrices
θ^v	Contribution weight coefficient of different views
N	Number of samples
d	Number of features
m	Number of labels
V	Number of views

of each view to make the final label prediction. For example, VLSF [30] learns the label specific representation of each data view, and uses label correlation and view consensus to build an ML classification model. MLSO [38] builds support vector machines(SVM) classifier for each heterogeneous data and jointly learns multi-source ML learning tasks under a unified optimization framework. The problem with the later fusion strategy is that it constructs a separate classification model for each view, which improves the accuracy. However, the complexity of the algorithm also increases accordingly. Besides, the redundancy and the influence of noisy data among the views are not considered.

Based on the above research, we found that the early fusion method has the problem that it is difficult to obtain an accurate, common representation of all views. Moreover, there is no algorithm to simultaneously explore the consistency and complementary knowledge between the views in the later fusion. Based on this, we consider the excellent performance of the neural network method on MVML and propose a novel single hidden layer feedforward neural network MVML algorithm named CDMM. A large number of empirical results on the benchmark data set show that the CDMM method is superior to these related competitive methods.

3. Proposed approach

3.1. Problem statement and notations

Let $\mathbf{X} = \{\mathbf{X}^v\}_{v=1}^V$ represents a feature space dataset with v views, where $\mathbf{X}^v = [\mathbf{x}_1, \dots, \mathbf{x}_N]^T \in \mathbb{R}^{N \times d}$ represents the entire feature space of the v th view with N samples. $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_m] \in \mathbb{R}^{N \times m}$ represents the corresponding label space, where $\mathbf{y}_i \in \{-1, 1\}^{N \times m}$ is the label vector of \mathbf{x}_i , and m represents the number of labels. Table 1 summarizes the definitions of some notations used in this paper.

3.2. Label matrix reconstruction

Label correlation is a crucial factor to improve the performance of ML learning, so in this section, we embed the label correlation into the original label matrix. After reconstructing the original label matrix, the member labels in the enhanced label matrix have strong interdependence and weak dependence on other non-member labels. Specifically, inspired by the literature [39], we use two different sources when constructing the label correlation matrix \mathbf{A} : label space \mathbf{Y} and instance space \mathbf{X} . In the label space, we use Jaccard similarity to measure the dependency between two labels:

$$A_{jk}^{(L)} = \frac{\sum_{i=1}^N \mathbf{y}_{ij} \cdot \mathbf{y}_{ik}}{\sum_{i=1}^N (\mathbf{y}_{ij} + \mathbf{y}_{ik} - \mathbf{y}_{ij} \cdot \mathbf{y}_{ik})} \quad (1)$$

In the instance space, to model the affinity matrix $A_{jk}^{v(I)}$ for each view, we define:

$$A_{jk}^{v(I)} = \begin{cases} \exp\left(-\frac{\|\mu_j^v - \mu_k^v\|_2^2}{\sigma^2}\right) & \text{if } |A_{jk}^{v(I)}| \geq \varepsilon \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

where $\mu_j^v = \frac{\sum_{i=1}^N x_i^v y_{ij}}{\sum_{i=1}^N y_{ij}}$. Furthermore, the parameter σ in Eq. (2) is fixed to be 1. We combine the above two measurement methods into a weight matrix \mathbf{A} in the following way:

$$A_{jk} = \begin{cases} \eta A_{jk}^{(L)} + (1 - \eta) \frac{1}{V} \sum_{v=1}^V A_{jk}^{v(I)} & (A_{jk} > \varepsilon) \\ 0 & (A_{jk} \leq \varepsilon) \end{cases} \quad (3)$$

Finally, inspired by the idea of label propagation dependence [40], we have:

$$\mathbf{T} = \mathbf{Y}\mathbf{C} \quad (4)$$

where $\eta \in [0, 1]$ is a balance factor, $\mathbf{C} = \mathbf{D} - \mathbf{A}$ is the Laplacian matrix of \mathbf{A} , \mathbf{D} is the diagonal matrix $\mathbf{D} = (\mathbf{D}_{ij}) = (\sum_k A_{jk})$, $\mathbf{T} = [\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_m] \in \mathbb{R}^{N \times m}$.

3.3. Problem formulation

Different from the widely adopted multi-layer network embedding structure [41], we use a random SLFN method to construct the initial model, similar to the extreme learning machines (ELM) [42], Schmidt et al. [43], and random vector connection function (Random Vector Functional-Link, RVFL) [44] input the weight parameters of the neural network randomly. The details are as following:

$$\sum_{i=1}^L g_i(\mathbf{x}_j)^T \beta_i = \sum_{i=1}^L g(\mathbf{w}_i \mathbf{x}_j + \mathbf{b}_i)^T \beta_i = \mathbf{o}_j \quad (5)$$

$j = 1, \dots, N$

where $\mathbf{x}_j = [\mathbf{x}_{j1}, \mathbf{x}_{j2}, \dots, \mathbf{x}_{jd}]^T \in \mathbb{R}^{d \times N}$ is the input instance feature vector, $\mathbf{w}_i = [\mathbf{w}_{1i}, \mathbf{w}_{2i}, \dots, \mathbf{w}_{di}]^T$ is the weight vector that connects the node of the input feature layer with the node of the i th hidden layer, $\beta_i = [\beta_{1i}, \beta_{2i}, \dots, \beta_{mi}]^T$ is the weight vector connecting the i th hidden node and the output node, and \mathbf{b}_i is the bias of the i th hidden node. The standard SLFN of L hidden nodes with activation function g can approximate these N samples with zero error, which means that $\sum_{j=1}^N \|\mathbf{o}_j - \mathbf{t}_j\| = 0$, that is, there are β_i , \mathbf{w}_i , and \mathbf{b}_i :

$$\sum_{i=1}^L g_i(\mathbf{x}_j)^T \beta_i = \sum_{i=1}^L g(\mathbf{w}_i \mathbf{x}_j + \mathbf{b}_i)^T \beta_i = \mathbf{t}_j \quad (6)$$

$j = 1, \dots, N$

Eq. (6) can be written succinctly as:

$$\mathbf{H}\boldsymbol{\beta} = \mathbf{T} \quad (7)$$

$$\mathbf{H} = \begin{bmatrix} g(\mathbf{w}_1 \mathbf{x}_1 + \mathbf{b}_1)^T & \cdots & g(\mathbf{w}_L \mathbf{x}_1 + \mathbf{b}_L)^T \\ \vdots & \ddots & \vdots \\ g(\mathbf{w}_1 \mathbf{x}_N + \mathbf{b}_1)^T & \cdots & g(\mathbf{w}_L \mathbf{x}_N + \mathbf{b}_L)^T \end{bmatrix}_{N \times L}$$

where $\boldsymbol{\beta} = \begin{bmatrix} \beta_1^T \\ \vdots \\ \beta_L^T \end{bmatrix}_{L \times m}$, and $\mathbf{T} = \begin{bmatrix} \mathbf{t}_1^T \\ \vdots \\ \mathbf{t}_N^T \end{bmatrix}_{N \times m}$. \mathbf{H} is called the hidden layer output matrix of the neural network:

$$\min_{\boldsymbol{\beta}} \|\mathbf{H}\boldsymbol{\beta} - \mathbf{T}\|_F^2 \quad (8)$$

Further extending the network to MVML data includes:

$$\min_{\boldsymbol{\beta}} \sum_{v=1}^V \left(\|\mathbf{H}^v \boldsymbol{\beta}^v - \mathbf{T}\|_F^2 \right); \quad (9)$$

$$\mathbf{H}^v = \begin{bmatrix} g(\mathbf{w}_1^v \mathbf{x}_1^v + \mathbf{b}_1^v)^T & \cdots & g(\mathbf{w}_L^v \mathbf{x}_1^v + \mathbf{b}_L^v)^T \\ \vdots & \ddots & \vdots \\ g(\mathbf{w}_1^v \mathbf{x}_N^v + \mathbf{b}_1^v)^T & \cdots & g(\mathbf{w}_L^v \mathbf{x}_N^v + \mathbf{b}_L^v)^T \end{bmatrix}_{N \times L}$$

where $v = 1, \dots, V$, V represents the number of views. Different from the traditional unsupervised multi-view view-consistency learning, a supervised learning model is adopted, which uses consistent information among different views by requiring different views to predict the same label result.

3.4. View diversity learning

In this section, we will enhance the diversity of all views to explore the view-diversity information among different views. We consider that higher independence means the higher diversity between the two variables, the lower correlation between the variables at this time to approximate the quantification of diversity based on the dependence among these variables. Here, the HSIC a simple and solid theoretical foundation and the ability to measure linear and nonlinear dependence among variables so that we use it. HSIC estimates the correlation of calculating the square norm of the cross-covariance operator on $\mathbf{H}^v \boldsymbol{\beta}^v$ and $\mathbf{H}^s \boldsymbol{\beta}^s$ in Hilbert space, and it is given by:

$$HSIC(\mathbf{H}^v \boldsymbol{\beta}^v, \mathbf{H}^s \boldsymbol{\beta}^s) = (n-1)^{-2} \text{tr}(\mathbf{K}^v \mathbf{P} \mathbf{K}^s \mathbf{P}) \quad (10)$$

In this paper, we use the inner product kernel to specify $\mathbf{K}^s = \mathbf{H}^s \mathbf{H}^{sT}$. Then minimize the overall HSIC on the v th individuality matrix to reduce the redundancy between them. The formula is as following:

$$\Phi(\{\mathbf{H}^v \boldsymbol{\beta}^v\}_{v=1}^V) = \sum_{v=1, v \neq s}^V HSIC(\mathbf{H}^v \boldsymbol{\beta}^v, \mathbf{H}^s \boldsymbol{\beta}^s) \quad (11)$$

$$= \sum_{v=1, v \neq s}^V (n-1)^{-2} \text{tr}((\mathbf{H}^v \boldsymbol{\beta}^v)^T \mathbf{P} \mathbf{K}^s \mathbf{P} \mathbf{H}^v \boldsymbol{\beta}^v)$$

$$= \sum_{v=1}^V \text{tr}((\mathbf{H}^v \boldsymbol{\beta}^v)^T \tilde{\mathbf{K}}^v \mathbf{H}^v \boldsymbol{\beta}^v)$$

where $\tilde{\mathbf{K}}^v = (n-1)^{-2} \sum_{s=1, s \neq v}^V \mathbf{P} \mathbf{K}^s \mathbf{P}$, \mathbf{K}^v , \mathbf{K}^s , $\mathbf{P} \in \mathbb{R}^{N \times N}$, \mathbf{K}^v and \mathbf{K}^s are used to measure the kernel-induced similarity between $\mathbf{H}^v \boldsymbol{\beta}^v$ and $\mathbf{H}^s \boldsymbol{\beta}^s$ vectors, respectively. $\mathbf{P} = \delta_{ij} - 1/N$, $\delta_{ij} = 1$, if $i = j$, $\delta_{ij} = 0$ otherwise. Adding Eq. (9) to Eq. (11), the final optimization model is as following:

$$\min_{\boldsymbol{\beta}} \sum_{v=1}^V \left(\|\mathbf{H}^v \boldsymbol{\beta}^v - \mathbf{T}\|_F^2 + \alpha \text{tr}((\mathbf{H}^v \boldsymbol{\beta}^v)^T \tilde{\mathbf{K}}^v \mathbf{H}^v \boldsymbol{\beta}^v) \right) \quad (12)$$

where α represents a non-negative trade-off parameter, which controls the diversity among views.

3.5. Model optimization and prediction

Different from the traditional SLN method which uses iterative backpropagation to solve the weight, we adopt the same scheme as the standard ELM [42,45], replacing backpropagation with a one-time matrix inverse operation and canceling the iterative operation of the neural network. Specifically, the derivation of Eq. (12) w.r.t $\boldsymbol{\beta}$ is as following:

$$\boldsymbol{\beta}^v = (\mathbf{H}^v)^{\dagger} \mathbf{T} \quad (13)$$

$$\begin{cases} (\mathbf{H}^v)^\dagger = \left((\mathbf{H}^v)^T \mathbf{H}^v + (\mathbf{H}^v)^T \tilde{\mathbf{K}}^v \mathbf{H}^v \right)^{-1} (\mathbf{H}^v)^T \\ (\mathbf{H}^v)^\dagger = (\mathbf{H}^v)^T \left(\mathbf{H}^v (\mathbf{H}^v)^T + \tilde{\mathbf{K}}^v \mathbf{H}^v (\mathbf{H}^v)^T \right)^{-1} \end{cases} \quad (14)$$

When $N \ll L$, then $(\mathbf{H}^v)^\dagger$ is the Moore–Penrose generalized inverse of matrix \mathbf{H} , and \mathbf{H} is the hidden layer output matrix. Besides, according to the ridge regression theorem, it is recommended to add a positive value $1/E$ to the diagonal of $\mathbf{H}^T \mathbf{H}$ or $\mathbf{H} \mathbf{H}^T$ to calculate the output weight [46], which can get a more stable solution and has better generalization performance. In other words, to improve the stability of the Eq. (14), we can have:

$$\boldsymbol{\beta}^v = (\mathbf{H}^v)^T \left(\mathbf{H}^v (\mathbf{H}^v)^T + \tilde{\mathbf{K}}^v \mathbf{H}^v (\mathbf{H}^v)^T + \frac{\mathbf{I}}{E} \right)^{-1} \mathbf{T} \quad (15)$$

According to Eq. (15), our final output function can be expressed as:

$$\begin{aligned} f(\mathbf{x}^v) &= h(\mathbf{x}^v) \boldsymbol{\beta} \\ &= h(\mathbf{x}^v) (\mathbf{H}^v)^T \left((\mathbf{I} + \tilde{\mathbf{K}}^v) \mathbf{H}^v (\mathbf{H}^v)^T + \frac{\mathbf{I}}{E} \right)^{-1} \mathbf{T} \end{aligned} \quad (16)$$

Besides, when the hidden layer feature map $h(\mathbf{x})$ is unknown, a kernel matrix can be defined by Eq. (16) (in this paper, we select the Radial Basis Function(RBF) kernel):

$$\boldsymbol{\Omega}_{\text{ELM}} = \mathbf{H} \mathbf{H}^T : \Omega_{ij} = h(\mathbf{x}_j) \cdot h(\mathbf{x}_i) = \text{Kernel}(\mathbf{x}_i, \mathbf{x}_j) \quad (17)$$

At this time, the number of hidden nodes L (the dimension of the hidden layer feature space) does not need to be specified. Eq. (16) can be abbreviated as:

$$\begin{aligned} f(\mathbf{x}^v) &= h(\mathbf{x}^v) \boldsymbol{\beta} \\ &= h(\mathbf{x}^v) (\mathbf{H}^v)^T \left((\mathbf{I} + \tilde{\mathbf{K}}^v) \mathbf{H}^v (\mathbf{H}^v)^T + \frac{\mathbf{I}}{E} \right)^{-1} \mathbf{T} \\ &= \begin{bmatrix} \text{Kernel}(\mathbf{x}^v, \mathbf{x}_1^v) \\ \vdots \\ \text{Kernel}(\mathbf{x}^v, \mathbf{x}_N^v) \end{bmatrix} \left((\mathbf{I} + \tilde{\mathbf{K}}^v) \boldsymbol{\Omega}_{\text{ELM}} + \frac{\mathbf{I}}{E} \right)^{-1} \mathbf{T} \end{aligned} \quad (18)$$

Considered the different contributions of each view in MVML learning, when making predictions, the output of other views should be weighed, as showed below:

$$\begin{aligned} \min_{\boldsymbol{\theta}^v} \quad & \frac{1}{2} \sum_{v=1}^V \boldsymbol{\theta}^v \text{PreLoss}^v + \frac{\lambda}{2} \|\boldsymbol{\theta}^v\|_2^2 \\ \text{s.t.} \quad & \boldsymbol{\theta}^v > 0, \sum_{v=1}^V \boldsymbol{\theta}^v = 1 \end{aligned} \quad (19)$$

Solve the objective function Eq. (19) by Lagrangian multiplier method:

$$\boldsymbol{\theta}^v = \begin{cases} \frac{\sum_{v=1}^V \text{PreLoss}^v + \lambda - V \text{PreLoss}^v}{V \lambda} & \text{if } \sum_{v=1}^V \text{PreLoss}^v + \lambda > V \text{PreLoss}^v \\ \varepsilon & \text{if } \sum_{v=1}^V \text{PreLoss}^v + \lambda \leq V \text{PreLoss}^v \end{cases} \quad (20)$$

where $\text{PreLoss}^v = \|\mathbf{H}_{\text{train}}^v \boldsymbol{\beta}^v - \mathbf{T}_{\text{train}}\|_F^2$, and ε is a very small non-negative value.

Finally, our prediction function is expressed as:

$$\mathbf{T}_{\text{pre}} = \sum_{v=1}^V \boldsymbol{\theta}^v \mathbf{H}_{\text{test}}^v \boldsymbol{\beta}^v \quad (21)$$

Based on the above optimization and prediction of the model, the main process of CDMM is summarized in Algorithm 1.

Algorithm 1 Consistency and Difference Neural Network Multi-view Multi-label Learning(CDMM).

Input:

The training data set: $\mathbf{X}^v \in \mathbb{R}^{N \times d}$;

The label data set: $\mathbf{Y} \in \mathbb{R}^{N \times m}$;

The non-negative trade-off parameters: $\alpha, \lambda, C, \sigma$, and η ;

Output:

Final prediction objective function: \mathbf{T}_{pre} ;

1: Calculate the label correlation matrix by Eq. (3)

2: Reconstruction of the original label matrix by Eq. (4)

3: **for** Train Data $v = 1$ to V **do**

4: Calculate the kernel matrix $\boldsymbol{\Omega}_{\text{ELM}}$ by Eq. (17);

5: Calculate neural network output weights by Eq. (15);

6: Calculate the prediction function of the training set by Eq. (18);

7: Calculate view contribution weight $\boldsymbol{\theta}^v$ by Eq. (20);

8: **end for**

9: Calculate the prediction function of the test set by Eq. (21)

10: **return** \mathbf{T}_{pre} .

3.6. Complexity analysis

Specifically, the complexity of the HSIC function is $\mathcal{O}(N^3)$, but because the matrix \mathbf{P} is a sparse matrix, the overall complexity will be lower than $\mathcal{O}(N^3)$. The complexity of solving weight $\boldsymbol{\beta}$ is $\mathcal{O}(N^2 m + N^2 d_{\text{max}})$. Since $N \gg m$ and $N \gg d_{\text{max}}$, the overall time complexity of CDMM is $\mathcal{O}(N^3)$. In practice, the time cost mainly comes from the size of the sample size.

4. Experiments

4.1. Datasets

To verify the effectiveness of the CDMM algorithm, we use a total of 7 benchmark MVML data sets for performance evaluation, which can be download from MULAN¹ and LEAR². The details of the data sets are summarized in Table 2.

4.2. Comparing algorithms

To verify the performance of CDMM, we selected three state-of-the-art MVML learning algorithms and one incomplete view MVML weak-label learning algorithm as the comparison algorithm.

1. iMvWL³ [32]: Learning incomplete views and weak labels by sharing subspace information. In the experiment, the complete view information can be obtained. The parameters α and β are both searched in $\{10^{-5}, 10^{-4}, \dots, 10^0\}$.

2. ICM2L⁴: Explicitly explore the individuality and common information of MVML data in the unified model. Set α, β and k within the range suggested in the literature [35].

3. VLSF⁵: Multi-View Multi-Label Learning With View-Label-Specific Features. All parameters are searched within the range given in the literature [30].

¹ data sets: <http://mulan.sourceforge.net/datasets-mlc.html>.

² data sets: <http://lear.inrialpes.fr/people/guillaumin/data.php>.

³ code: <http://mlda.swu.edu.cn/codes.php?name=iMvWL>.

⁴ code: <http://mlda.swu.edu.cn/codes.php?name=ICM2L>.

⁵ code: <http://www.escience.cn/people/huangjun/index.html>.

Table 2
MVML data sets.

Views	Emotions	Yeast	Pascal07	Corel5k	Espgame	laprtc12	mirflickr
1	rhythmic attributes (8)	Genetic Expression (79)	DenseSift (1000)	DenseHue (100)	DenseHue (100)	DenseHue (100)	DenseHue (100)
2	timbre attributes (64)	phylogenetic profile (24)	HarrisSift (1000)	DenseSift (1000)	DenseSift (1000)	DenseSift (1000)	DenseSift (1000)
3	–	–	Gist(512)	Gist(512)	Gist(512)	Gist(512)	Gist(512)
4	–	–	HSV(4096)	HSV(4096)	HSV(4096)	HSV(4096)	HSV(4096)
5	–	–	RGB(4096)	Lab(4096)	Lab(4096)	Lab(4096)	Lab(4096)
6	–	–	Tags(804)	RGB(4096)	RGB(4096)	RGB(4096)	RGB(4096)
Domain	music	biology	image	image	image	image	image
m	6	14	20	260	268	291	457
N	593	2417	9963	4999	20770	19627	25000

4. SIMM⁶: Multi-View Multi-Label Learning with View-Specific Information Extraction. The number of hidden layers is fixed to 64. The parameter α is set to 1, and β is searched in $\{0.1, 0.01, 0.001, 0.0001\}$.

5. CDMM: The parameter α is searched in $\{10^{-10}, 10^{-9}, \dots, 10^{-5}\}$, λ is searched in $\{10^1, 10^2, \dots, 10^7\}$, the regularization coefficient C is searched in $\{10^{-5}, 10^{-1}, \dots, 10^5\}$, the kernel parameter σ is searched in $\{10^{-2}, 10^{-1}, \dots, 10^2\}$, and the label correlation balance parameter η is explored within the range $\{0.5 - 0.8\}$.

4.3. The experimental environment and evaluation metrics

In this paper, the experiment was performed on Windows 10, Intel(R) Core(TM) i7-7700K, and 32GB RAM, and the method was implemented in MATLAB 2016b. We used (1) Hamming Loss(HL), (2) Average Precision(AP), (3) One Error(OE), (4) Ranking Loss(RL), (5) Coverage(CV), (6) Macro-F1, (7) Micro-F1, and (8) Subset Accuracy(SA), eight widely used ML metrics for performance evaluation. They evaluate the performance of the algorithm from the perspective of ranking and classification, respectively. The specific measurement definition can be found in [47]. For Average Precision, Macro-F1, Micro-F1, and Subset Accuracy, the larger the value, the better the performance. For the other four metrics, the smaller the value, the better the performance.

4.4. Experimental results and analysis

For each MVML data set, we randomly select 80% of the total data as the training set, and the remaining data as the test set, and repeat the experiment for all comparison algorithms five times. Tables 3, to 4 respectively list the average results (mean \pm std) of the 5 comparison algorithms on the 7 MVML benchmark data sets, and the best results are shown in black. \downarrow means that the smaller the evaluation metrics value is, the better, and \uparrow means that the larger the evaluation metrics value is, the better.

Furthermore, a statistical hypothesis test was utilized to verify and compare the relative performance of various algorithms. The Friedman test [48] was utilized for performance analysis. Table 5 summarizes the Friedman statistics F_F and the corresponding critical values of the various evaluation metrics. As showed in Table 5, at the significance level $\rho = 0.05$, each evaluation metric is rejected when the null hypothesis that all comparison algorithms are equivalently executed. Therefore, the Nemenyi test [48] is utilized as a post-hoc test to compare the performance of each

algorithm and observe whether the CDMM algorithm is competitive. There is a significant difference in performance between the two classifiers if the corresponding average ranking reaches at least a critical difference (CD):

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

Analysis of experimental results: according to the experimental results reported in Tables 3 and 4, the following observations can be drawn:

- Among 56 configurations (7 data sets and 8 evaluation metrics), CDMM ranks first and second at 64.3% and 26.9%, respectively.
- It is worth noting that CDMM outperforms all comparison algorithms on all metrics on *Emotions* and *Yeast*.
- Except for Pascal07 and Mirflickr, CDMM achieved the best performance on AP, OE, Macro-F1, Micro-F1 and Subset Accuracy.
- All data sets: the HL value of CDMM is almost the same as the optimal result, and the CV and RL are also not much different from the optimal result.

For Nemenyi test, $q = 2.728$ at significance level $\rho = 0.05$, and thus $CD = 2.31(k = 5, N = 7)$. Fig. 1 indicates the CD diagrams of each algorithm under different evaluation metrics, respectively. In each subfigure, two or more algorithms are connected by solid red lines indicate that there is no significant difference in performance among them. For each approach, there are 40 comparative results (5 parallel approaches and 8 evaluation metrics).

Observing Fig. 1, we can see that under the confidence interval of significance level $\rho = 0.05$, CDMM has a significant advantage over other comparison algorithms in 30% of the cases, ranking first. And SIMM ranked second with 27.5%.

Through the above experimental results, it can be analyzed:

- SIMM, iMvWL, and ICM2L are all designed to use subspace learning to solve the problem of MVML learning, but the performance of the SIMM is always better than the other two. The reason is that SIMM builds a model through a neural network, and while mining the shared subspace, it also extracts the private information of the view.
- The VLSF algorithm performs better than iMvWL and ICM2L, because it pays more attention to the consensus learning among views, and at the same time mines the view-label-Specific to improve the effectiveness of the algorithm. The performance of VLSF is inferior to CDMM mainly because VLSF ignores the diversity of views, and there are limitations in extracting view information.

⁶ code: <http://palm.seu.edu.cn/zhangml/>

Table 3The results of each comparison algorithm on HL, OE, RL, and CV evaluation metrics (\downarrow) on 7 data sets.

Dataset	metric	SIMM	ICML	iMvWL	VLSF	CDMM
Emotions	HL	0.246 \pm 0.008	0.375 \pm 0.015	0.395 \pm 0.011	0.293 \pm 0.007	0.207 \pm 0.014
	OE	0.310 \pm 0.056	0.530 \pm 0.030	0.521 \pm 0.021	0.539 \pm 0.073	0.304 \pm 0.034
	RL	0.178 \pm 0.024	0.443 \pm 0.013	0.414 \pm 0.012	0.330 \pm 0.022	0.174 \pm 0.013
	CV	0.307 \pm 0.014	0.530 \pm 0.036	0.506 \pm 0.014	0.441 \pm 0.022	0.304 \pm 0.018
Yeast	HL	0.207 \pm 0.005	0.278 \pm 0.008	0.269 \pm 0.005	0.258 \pm 0.004	0.189 \pm 0.006
	OE	0.225 \pm 0.028	0.235 \pm 0.024	0.292 \pm 0.020	0.343 \pm 0.015	0.211 \pm 0.013
	RL	0.165 \pm 0.008	0.215 \pm 0.011	0.214 \pm 0.008	0.320 \pm 0.006	0.151 \pm 0.008
	CV	0.450 \pm 0.004	0.503 \pm 0.006	0.494 \pm 0.009	0.601 \pm 0.008	0.426 \pm 0.008
Pascal07	HL	0.046 \pm 0.001	0.115 \pm 0.003	0.086 \pm 0.002	0.050 \pm 0.000	0.049 \pm 0.001
	OE	0.255 \pm 0.008	0.589 \pm 0.002	0.397 \pm 0.021	0.290 \pm 0.014	0.308 \pm 0.011
	RL	0.066 \pm 0.002	0.241 \pm 0.040	0.138 \pm 0.011	0.070 \pm 0.002	0.070 \pm 0.002
	CV	0.106 \pm 0.002	0.308 \pm 0.048	0.189 \pm 0.015	0.110 \pm 0.001	0.111 \pm 0.003
Corel5k	HL	0.011 \pm 0.000	0.022 \pm 0.000	0.022 \pm 0.000	0.013 \pm 0.000	0.011 \pm 0.000
	OE	0.363 \pm 0.011	0.697 \pm 0.007	0.687 \pm 0.003	0.438 \pm 0.017	0.362 \pm 0.006
	RL	0.059 \pm 0.002	0.149 \pm 0.002	0.130 \pm 0.003	0.076 \pm 0.004	0.069 \pm 0.005
	CV	0.148 \pm 0.006	0.334 \pm 0.000	0.286 \pm 0.004	0.187 \pm 0.009	0.179 \pm 0.011
ESPgame	HL	0.017 \pm 0.000	0.029 \pm 0.000	0.028 \pm 0.000	0.017 \pm 0.000	0.018 \pm 0.000
	OE	0.476 \pm 0.022	0.713 \pm 0.030	0.674 \pm 0.000	0.533 \pm 0.011	0.465 \pm 0.009
	RL	0.120 \pm 0.006	0.203 \pm 0.002	0.190 \pm 0.002	0.142 \pm 0.003	0.124 \pm 0.002
	CV	0.308 \pm 0.012	0.479 \pm 0.001	0.447 \pm 0.004	0.365 \pm 0.005	0.337 \pm 0.005
laprtc12	HL	0.019 \pm 0.000	0.032 \pm 0.000	0.031 \pm 0.000	0.019 \pm 0.000	0.019 \pm 0.000
	OE	0.447 \pm 0.021	0.720 \pm 0.005	0.624 \pm 0.002	0.465 \pm 0.006	0.439 \pm 0.007
	RL	0.089 \pm 0.007	0.189 \pm 0.001	0.165 \pm 0.002	0.106 \pm 0.001	0.089 \pm 0.002
	CV	0.270 \pm 0.019	0.497 \pm 0.001	0.435 \pm 0.003	0.327 \pm 0.002	0.284 \pm 0.006
Mirflickr	HL	0.006 \pm 0.000	0.013 \pm 0.000	0.013 \pm 0.000	0.006 \pm 0.000	0.006 \pm 0.000
	OE	0.865 \pm 0.006	0.908 \pm 0.004	0.887 \pm 0.003	0.891 \pm 0.003	0.869 \pm 0.004
	RL	0.2222 \pm 0.013	0.288 \pm 0.001	0.285 \pm 0.009	0.197 \pm 0.003	0.184 \pm 0.003
	CV	0.3006 \pm 0.016	0.499 \pm 0.002	0.492 \pm 0.007	0.344 \pm 0.005	0.332 \pm 0.002

Table 4The results of AP, Macro-F1, Micro-F1, and SA evaluation metrics (\uparrow) of each comparison algorithm on 7 data sets.

Dataset	metric	SIMM	ICML	iMvWL	VLSF	CDMM
Emotions	AP	0.780 \pm 0.027	0.578 \pm 0.022	0.584 \pm 0.015	0.621 \pm 0.029	0.790 \pm 0.019
	Macro-F1	0.386 \pm 0.038	0.182 \pm 0.005	0.183 \pm 0.004	0.105 \pm 0.009	0.643 \pm 0.028
	Micro-F1	0.430 \pm 0.040	0.404 \pm 0.020	0.394 \pm 0.014	0.246 \pm 0.034	0.665 \pm 0.022
	SA	0.233 \pm 0.032	0.106 \pm 0.097	0.009 \pm 0.000	0.037 \pm 0.017	0.294 \pm 0.043
Yeast	AP	0.765 \pm 0.016	0.708 \pm 0.014	0.704 \pm 0.011	0.631 \pm 0.006	0.781 \pm 0.009
	Macro-F1	0.282 \pm 0.006	0.280 \pm 0.005	0.340 \pm 0.013	0.324 \pm 0.014	0.450 \pm 0.015
	Micro-F1	0.575 \pm 0.010	0.580 \pm 0.014	0.594 \pm 0.008	0.442 \pm 0.018	0.693 \pm 0.011
	SA	0.105 \pm 0.019	0.004 \pm 0.002	0.008 \pm 0.004	0.074 \pm 0.009	0.207 \pm 0.011
Pascal07	AP	0.786 \pm 0.005	0.460 \pm 0.025	0.660 \pm 0.013	0.767 \pm 0.007	0.759 \pm 0.006
	Macro-F1	0.538 \pm 0.009	0.073 \pm 0.032	0.413 \pm 0.024	0.523 \pm 0.008	0.474 \pm 0.018
	Micro-F1	0.616 \pm 0.009	0.340 \pm 0.020	0.502 \pm 0.010	0.623 \pm 0.007	0.603 \pm 0.014
	SA	0.394 \pm 0.010	0.038 \pm 0.010	0.099 \pm 0.012	0.356 \pm 0.003	0.359 \pm 0.013
Corel5k	AP	0.534 \pm 0.006	0.258 \pm 0.004	0.274 \pm 0.003	0.476 \pm 0.011	0.545 \pm 0.007
	Macro-F1	0.098 \pm 0.006	0.008 \pm 0.001	0.017 \pm 0.001	0.116 \pm 0.003	0.173 \pm 0.014
	Micro-F1	0.320 \pm 0.009	0.214 \pm 0.007	0.235 \pm 0.004	0.387 \pm 0.018	0.485 \pm 0.017
	SA	0.032 \pm 0.007	0.000 \pm 0.000	0.001 \pm 0.001	0.034 \pm 0.008	0.069 \pm 0.010
ESPgame	AP	0.378 \pm 0.013	0.219 \pm 0.013	0.237 \pm 0.002	0.336 \pm 0.003	0.400 \pm 0.003
	Macro-F1	0.072 \pm 0.013	0.013 \pm 0.004	0.014 \pm 0.000	0.029 \pm 0.004	0.169 \pm 0.008
	Micro-F1	0.194 \pm 0.017	0.206 \pm 0.012	0.216 \pm 0.002	0.182 \pm 0.013	0.347 \pm 0.007
	SA	0.009 \pm 0.001	0.000 \pm 0.000	0.000 \pm 0.000	0.007 \pm 0.002	0.013 \pm 0.001
laprtc12	AP	0.401 \pm 0.017	0.204 \pm 0.000	0.242 \pm 0.001	0.382 \pm 0.002	0.432 \pm 0.004
	Macro-F1	0.065 \pm 0.010	0.010 \pm 0.001	0.021 \pm 0.000	0.060 \pm 0.002	0.183 \pm 0.012
	Micro-F1	0.165 \pm 0.014	0.197 \pm 0.002	0.223 \pm 0.002	0.254 \pm 0.006	0.384 \pm 0.011
	SA	0.002 \pm 0.001	0.000 \pm 0.000	0.000 \pm 0.000	0.005 \pm 0.001	0.018 \pm 0.002
Mirflickr	AP	0.142 \pm 0.009	0.093 \pm 0.001	0.094 \pm 0.005	0.091 \pm 0.003	0.102 \pm 0.002
	Macro-F1	0.001 \pm 0.000	0.007 \pm 0.000	0.004 \pm 0.001	0.011 \pm 0.001	0.009 \pm 0.001
	Micro-F1	0.005 \pm 0.002	0.082 \pm 0.002	0.077 \pm 0.004	0.046 \pm 0.005	0.066 \pm 0.006
	SA	0.250 \pm 0.007	0.000 \pm 0.000	0.000 \pm 0.000	0.207 \pm 0.012	0.198 \pm 0.009

- Compared with SIMM, CDMM has 6 metrics advantages in 8 metrics. The main reason is that CDMM pays more attention to the consistency and diversity of learning among various views when constructing a unified neural network model, and strengthens the information communication among views.

- Although the performance of CDMM is not optimal in the HL and CV metrics, combined with the SA metric, we can find that although CDMM is not as good as SIMM in the prediction results on all labels, it has achieved better results on SA. These show that CDMM has a better ability to recognize rare labels than SIMM.

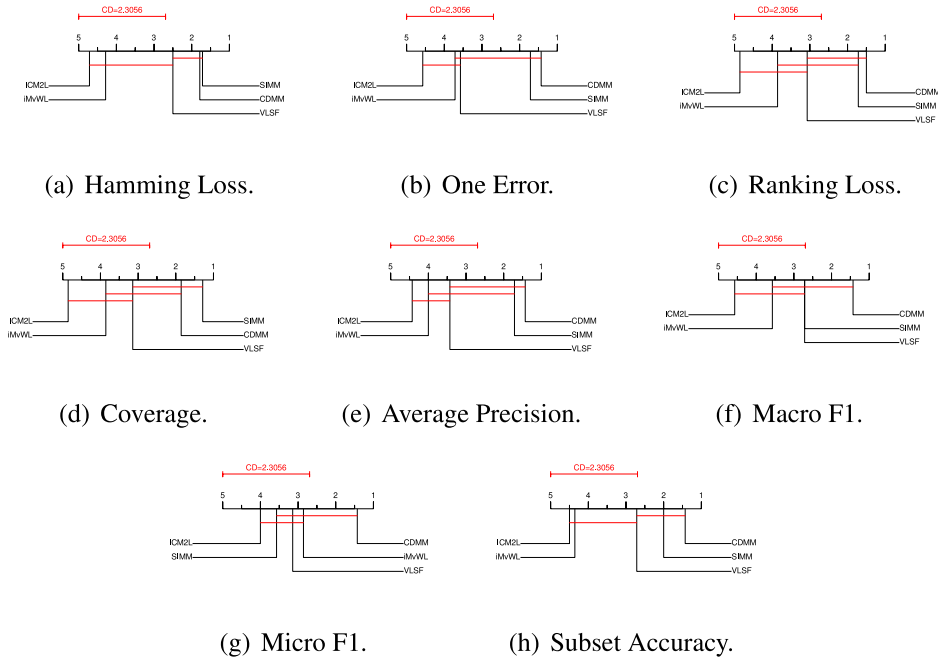


Fig. 1. Nemenyi test results for different evaluation metrics.

4.5. Component analysis

To verify the effectiveness of each part of CDMM, we additionally conducted component analysis experiments on the *Emotions*, *Yeast*, and *Corel5k* data sets (the parameter values involved have been given in advance), and reported the values under 8 metrics in Fig. 2. Definition of various variant algorithms corresponding to CDMM: CDMM-WD means that label correlation is not considered; CDMM-WL means ignoring view diversity information (maintain the basic structure of the model $\alpha = 0$); CDMM-DL means ignoring the influence of view contribution weight.

The results of the various sub-graphs in Fig. 2, it can be observed that CDMM is superior to its different variant algorithms on multiple metrics. Specifically, CDMM-WD and CDMM-DL achieve better results in most cases than CDMM-WL, which indicates the effectiveness of view diversity learning. CDMM achieves better results in most cases than CDMM-DL, which proves the importance of considering the difference in view contribution in the later fusion model. CDMM-WD has limited performance improvement compared with CDMM-WL and CDMM-DL, which shows that we have room for further improvement in the modeling of the correlation among labels. These results confirm the rationality and effectiveness of the modeling of each part of our model, and also clarify our motivation to use the consistency and diversity of the data from multiple views.

4.6. Parameter sensitivity analysis

CDMM has two essential parameters α and λ , which respectively control the regularization term of view-diversity and the regularization term of view-contribution. We tested the sensitivity of the parameters α and λ in the range of $\{10^{-10}, 10^{-9}, \dots, 10^{-2}\}$ and $\{10^3, 10^4, \dots, 10^{11}\}$ respectively and reported the HL and SA results under the *Yeast* data in Fig. 3.

From Fig. 3, we can observe that when α is about 10^{-6} and λ is about 10^4 , CDMM obtains relatively good performance. Besides, when α tends to a larger value, the performance of CDMM will decrease, mainly because an excessively large α value digs out more view-diversity information, but also loses some vital

Table 5

Summary of the Friedman Statistics F_F ($k = 5$, $N = 7$) and the critical value in terms of each evaluation metric (k : Comparing Algorithms; N : Data sets).

Metric	F_F	Critical Value($\rho = 0.05$)
Hamming Loss	23.5477	2.776
Average Precision	16.6154	
One Error	17.3333	
Ranking Loss	25.4439	
Coverage	32.6842	
Macro-F1	7.1250	
Micro-F1	3.7351	
Subset Accuracy	19.4545	

view consistency information. The λ value tends to achieve the intermediate value so that each view provides enough sufficient information as much as possible. The results of other data sets and evaluation metrics are similar, and similar conclusions can be obtained.

In addition, we also experimentally studied the sensitivity of η to CDMM. Fig. 4 reports the HL and SA values of CDMM on the *Yeast* dataset, with η values varying from 0.1 to 0.8. It can be seen that the performance of CDMM improves with the increase of η and then decreases when $\eta > 0.6$. Without losing generality, we set the value of η in the range of $\{0.5, 0.6, 0.7, 0.8\}$ on all datasets.

5. Conclusion

In this paper, we have studied how to mine the information of view-consistency and view-diversity in multi-view data to achieve effective multi-view multi-label classification. For this reason, a multi-view multi-label learning framework called CDMM is proposed. It uses a unified feedforward neural network model to find out the consistency and diversity among heterogeneous views and additionally considers label correlation factors and view contribution factors. The difference from previous studies is the case that our classifier introduces a nonlinear kernel mapping function into the model. Experiments on multiple benchmark data sets have verified that the CDMM model is concise and efficient, and overall it is superior to relate to competitive solutions.

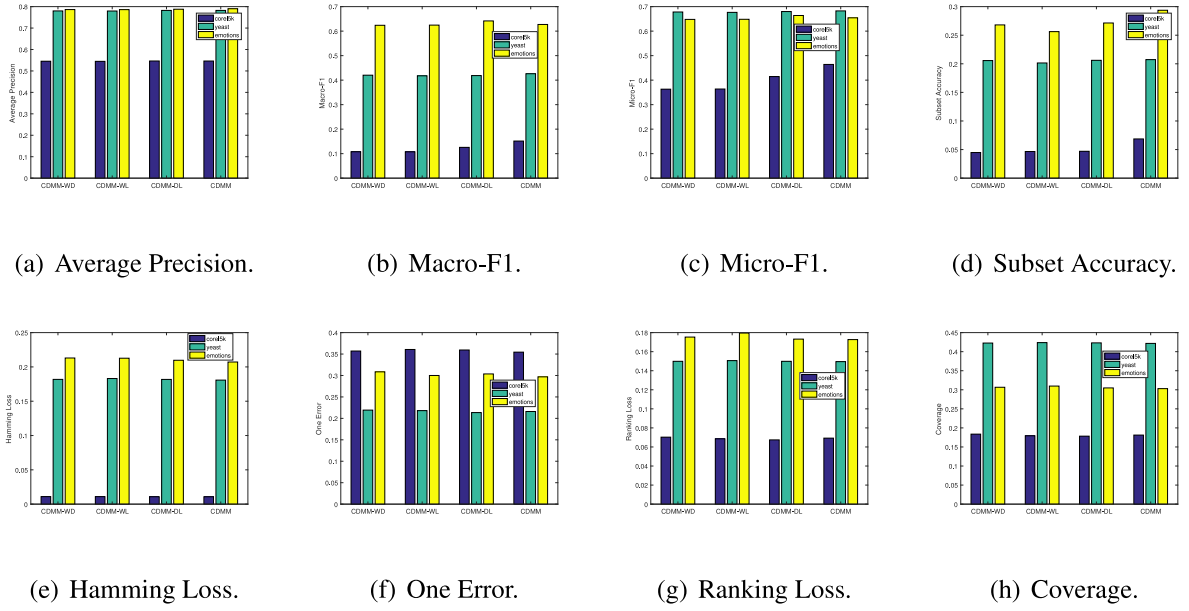


Fig. 2. Eight evaluation metrics results of CDMM and its variants in the *Emotions*, *Yeast* and *Core15k* datasets.

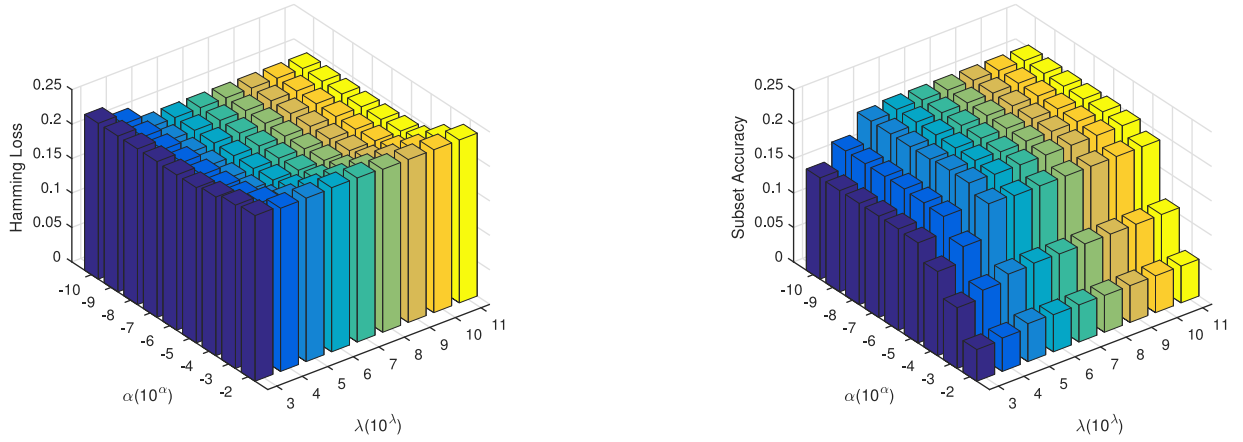


Fig. 3. Parameter sensitivity analysis w.r.t. α and λ on *Yeast*.

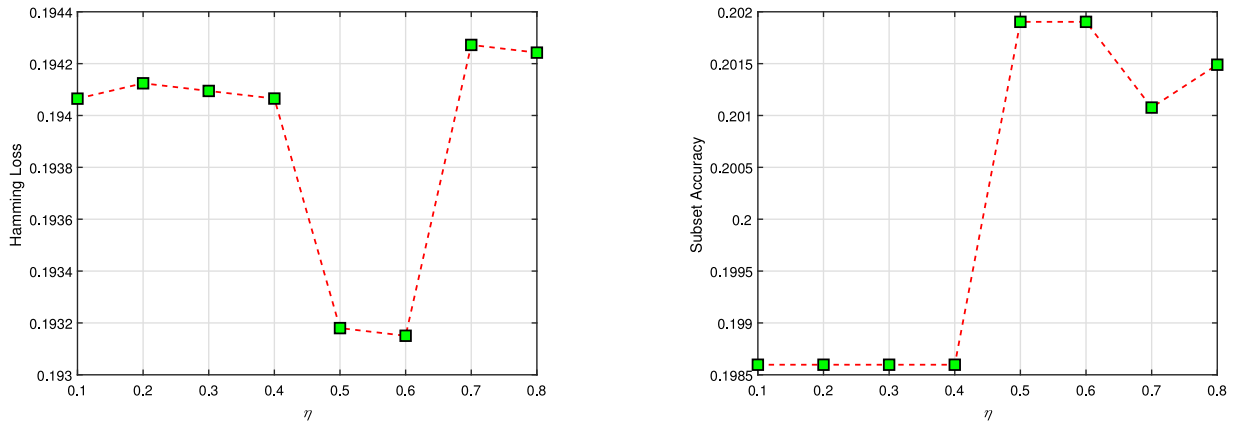


Fig. 4. Parameter sensitivity analysis w.r.t. η on *Yeast*.

The disadvantage of CDMM is apparent that it is not suitable for learning in the context of incomplete data. In the future, we

will expand CDMM to the application of missing feature data and label data.

CRediT authorship contribution statement

Dawei Zhao: Conceptualization, Methodology, Software, Investigation, Data curation, Writing - original draft, Writing - review & editing. **Qingwei Gao:** Validation, Supervision, Project administration, Funding acquisition. **Yixiang Lu:** Visualization, Investigation. **Dong Sun:** Formal analysis, Validation. **Yusheng Cheng:** Writing - review & editing.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgments

This research is supported by Key Laboratory of Intelligent Computing and Signal Processing, Ministry of Education (Anhui University), China (2020A003), and the Nature Science Foundation of Anhui (2008085MF183).

References

- [1] Matthew R. Boutell, Jiebo Luo, Xipeng Shen, Christopher M. Brown, Learning multi-label scene classification, *Pattern Recognit.* 37 (9) (2004) 1757–1771.
- [2] S. Wen, W. Liu, Y. Yang, P. Zhou, Z. Guo, Z. Yan, Y. Chen, T. Huang, Multilabel image classification via feature/label co-projection, *IEEE Transactions on Systems, Man, and Cybernetics: Systems* (2020) 1–10, <http://dx.doi.org/10.1109/TSMC.2020.2967071>.
- [3] Feng Kang, Rong Jin, Rahul Sukthankar, Correlated label propagation with application to multi-label learning, in: 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition vol. 2, CVPR'06, IEEE, 2006, pp. 1719–1726.
- [4] André Elisseeff, Jason Weston, A kernel method for multi-labelled classification, in: *Advances in Neural Information Processing Systems*, 2002, pp. 681–687.
- [5] Minling Zhang, Zhihua Zhou, Multilabel neural networks with applications to functional genomics and text categorization, *IEEE Trans. Knowl. Data Eng.* 18 (10) (2006) 1338–1351.
- [6] Hao Wang, Yan Yang, Bing Liu, Hamido Fujita, A study of graph-based system for multi-view clustering, *Knowl.-Based Syst.* 163 (2019) 1009–1019.
- [7] Zhe Xue, Guorong Li, Qingming Huang, Joint multi-view representation and image annotation via optimal predictive subspace learning, *Inform. Sci.* 451 (2018) 180–194.
- [8] Pengfei Zhu, Qi Hu, Qinghua Hu, Changqing Zhang, Zhizhao Feng, Multi-view label embedding, *Pattern Recognit.* 84 (2018) 126–135.
- [9] Zesen Chen, Xuan Wu, Qingguo Chen, Yao Hu, Minling Zhang, Multi-view partial multi-label learning with graph-based disambiguation, in: *AAAI*, 2020, pp. 3553–3560.
- [10] Y. Zhang, J. Wu, Z. Cai, P.S. Yu, Multi-view multi-label learning with sparse feature selection for image annotation, *IEEE Transactions on Multimedia* 22 (11) (2020) 2844–2857, <http://dx.doi.org/10.1109/TMM.2020.2966887>.
- [11] Xiaofeng Zhu, Xuelong Li, Shichao Zhang, Block-row sparse multiview multilabel learning for image classification, *IEEE Trans. Cybern.* 46 (2) (2015) 450–461.
- [12] Meng Liu, Yong Luo, Dacheng Tao, Chao Xu, Yonggang Wen, Low-rank multi-view learning in matrix completion for multi-label image classification, in: *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015, pp. 2778–2784.
- [13] Xuran Zhao, Nicholas Evans, Jean-Luc Dugelay, A subspace co-training framework for multi-view clustering, *Pattern Recognit. Lett.* 41 (2014) 73–82.
- [14] Shirui Luo, Changqing Zhang, Wei Zhang, Xiaochun Cao, Consistent and specific multi-view subspace clustering, in: *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [15] Changqing Zhang, Ziwei Yu, Qinghua Hu, Pengfei Zhu, Xinwang Liu, Xiaobo Wang, Latent semantic aware multi-view multi-label classification, in: *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [16] Weijieying Ren, Lei Zhang, Bo Jiang, Zhefeng Wang, Guangming Guo, Guquan Liu, Robust mapping learning for multi-view multi-label classification with missing labels, in: *International Conference on Knowledge Science, Engineering and Management*, Springer, 2017, pp. 543–551.
- [17] Xiaochun Cao, Changqing Zhang, Huazhu Fu, Si Liu, Hua Zhang, Diversity-induced multi-view subspace clustering, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 586–594.
- [18] Shixing Yao, Guoxian Yu, Jun Wang, Carlotta Domeniconi, Xiangliang Zhang, Multi-view multiple clustering, 2019, arXiv preprint [arXiv:1905.05053](https://arxiv.org/abs/1905.05053).
- [19] Arthur Gretton, Olivier Bousquet, Alex Smola, Bernhard Schölkopf, Measuring statistical dependence with Hilbert-Schmidt norms, in: *International Conference on Algorithmic Learning Theory*, Springer, 2005, pp. 63–77.
- [20] Minling Zhang, Zhihua Zhou, ML-KNN: A lazy learning approach to multi-label learning, *Pattern Recognit.* 40 (7) (2007) 2038–2048.
- [21] Minling Zhang, Lei Wu, Lift: Multi-label learning with label-specific features, *IEEE Trans. Pattern Anal. Mach. Intell.* 37 (1) (2014) 107–120.
- [22] Suping Xu, Xibei Yang, Hualong Yu, Dong-Jun Yu, Jingyu Yang, Eric C.C. Tsang, Multi-label learning with label-specific feature reduction, *Knowl.-Based Syst.* 104 (2016) 52–61.
- [23] Johannes Fürnkranz, Eyke Hüllermeier, Eneldo Loza Mencía, Klaus Brinker, Multilabel classification via calibrated label ranking, *Mach. Learn.* 73 (2) (2008) 133–153.
- [24] Yu Zhang, Dit-Yan Yeung, Multilabel relationship learning, *ACM Trans. Knowl. Discov. Data (TKDD)* 7 (2) (2013) 1–30.
- [25] Linli Xu, Zhen Wang, Zefan Shen, Yubo Wang, Enhong Chen, Learning low-rank label correlations for multi-label classification with missing labels, in: *2014 IEEE International Conference on Data Mining, IEEE*, 2014, pp. 1067–1072.
- [26] Zhifen He, Ming Yang, Yang Gao, Huidong Liu, Yilong Yin, Joint multi-label classification and label correlations with missing labels and feature selection, *Knowl.-Based Syst.* 163 (2019) 145–158.
- [27] Jun Huang, Feng Qin, Xiao Zheng, Zekai Cheng, Zhixiang Yuan, Weigang Zhang, Qingming Huang, Improving multi-label classification with missing labels by learning label-specific features, *Inform. Sci.* 492 (2019) 124–146.
- [28] Lijuan Sun, Ping Ye, Gengyu Lyu, Songhe Feng, Guojun Dai, Hua Zhang, Weakly-supervised multi-label learning with noisy features and incomplete labels, *Neurocomputing* 413 (2020) 61–71.
- [29] Shengjun Huang, Zhihua Zhou, Multi-label learning by exploiting label correlations locally, in: *Twenty-Sixth AAAI Conference on Artificial Intelligence*, 2012.
- [30] Jun Huang, Xiwen Qu, Guorong Li, Feng Qin, Xiao Zheng, Qingming Huang, Multi-view multi-label learning with view-label-specific features, *IEEE Access* 7 (2019) 100979–100992.
- [31] Qiaoyu Tan, Guoxian Yu, Carlotta Domeniconi, Jun Wang, Zili Zhang, Multi-view weak-label learning based on matrix completion, in: *Proceedings of the 2018 SIAM International Conference on Data Mining, SIAM*, 2018, pp. 450–458.
- [32] Qiaoyu Tan, Guoxian Yu, Carlotta Domeniconi, Jun Wang, Zili Zhang, Incomplete multi-view weak-label learning, in: *IJCAI*, 2018, pp. 2703–2709.
- [33] Changming Zhu, Panhong Wang, Lin Ma, Rigui Zhou, Lai Wei, Global and local multi-view multi-label learning with incomplete views and labels, *Neural Comput. Appl.* 32 (18) (2020) 15007–15028.
- [34] Zhige Xie, Kai Xu, Wen Shan, Ligang Liu, Yueshan Xiong, Hui Huang, Projective feature learning for 3D shapes with multi-view depth images, in: *Computer Graphics Forum*, vol. 34, Wiley Online Library, 2015, pp. 1–11.
- [35] Q. Tan, G. Yu, J. Wang, C. Domeniconi, X. Zhang, Individuality- and commonality-based multiview multilabel learning, *IEEE Trans. Cybern.* (2019) 1–12.
- [36] Xuan Wu, Qingguo Chen, Yao Hu, Dengbao Wang, Xiaodong Chang, Xiaobo Wang, Minling Zhang, Multi-view multi-label learning with view-specific information extraction, in: *IJCAI*, 2019, pp. 3884–3890.
- [37] Fangwen Zhang, Xiuyi Jia, Weiwei Li, Tensor based multi-view label enhancement for multi-label learning, in: *IJCAI*, 2020.
- [38] Jia Zhang, Candong Li, Zhenqiang Sun, Zhiming Luo, Changen Zhou, Shaozi Li, Towards a unified multi-source-based optimization framework for multi-label learning, *Appl. Soft Comput.* 76 (2019) 425–435.
- [39] Lu Sun, Mineichi Kudo, Keigo Kimura, Multi-label classification with meta-label-specific features, in: *2016 23rd International Conference on Pattern Recognition (ICPR)*, IEEE, 2016, pp. 1612–1617.
- [40] Clara Pizzuti, A multi-objective genetic algorithm for community detection in networks, in: *2009 21st IEEE International Conference on Tools with Artificial Intelligence*, IEEE, 2009, pp. 379–386.
- [41] Jie Lu, Junyu Xuan, Guangquan Zhang, Xiangfeng Luo, Structural property-aware multilayer network embedding for latent factor analysis, *Pattern Recognit.* 76 (2018) 228–241.
- [42] Gao Huang, Guangbin Huang, Shiji Song, Keyou You, Trends in extreme learning machines: A review, *Neural Netw.* 61 (2015) 32–48.
- [43] Wouter F. Schmidt, Martin A. Kraaijveld, Robert P.W. Duin, et al., Feed forward neural networks with random weights, in: *International Conference on Pattern Recognition*, IEEE Computer Society Press, 1992, p. 1.
- [44] Yoh-Han Pao, Gwang-Hoon Park, Dejan J. Sobajic, Learning and generalization characteristics of the random vector functional-link net, *Neurocomputing* 6 (2) (1994) 163–180.

- [45] Yusheng Cheng, Dawei Zhao, Yibin Wang, Gensheng Pei, Multi-label learning with kernel extreme learning machine autoencoder, *Knowl.-Based Syst.* 178 (2019) 1–10.
- [46] Donald W. Marquardt, Generalized inverses, ridge regression, biased linear estimation, and nonlinear estimation, *Technometrics* 12 (3) (1970) 591–612.
- [47] Minling Zhang, Zhihua Zhou, A review on multi-label learning algorithms, *IEEE Trans. Knowl. Data Eng.* 26 (8) (2013) 1819–1837.
- [48] Janez Demšar, Statistical comparisons of classifiers over multiple data sets, *J. Mach. Learn. Res.* 7 (Jan) (2006) 1–30.