

Journal Pre-proof

Two-stage label distribution learning with label-independent prediction based on label-specific features

Gui-Lin Li, Heng-Ru Zhang, Fan Min, Yu-Nan Lu

PII: S0950-7051(23)00176-4
DOI: <https://doi.org/10.1016/j.knosys.2023.110426>
Reference: KNOSYS 110426

To appear in: *Knowledge-Based Systems*

Received date: 12 September 2022

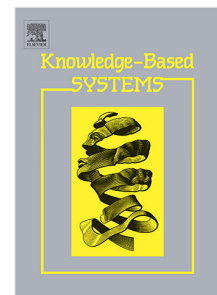
Revised date: 17 February 2023

Accepted date: 23 February 2023

Please cite this article as: G.-L. Li, H.-R. Zhang, F. Min et al., Two-stage label distribution learning with label-independent prediction based on label-specific features, *Knowledge-Based Systems* (2023), doi: <https://doi.org/10.1016/j.knosys.2023.110426>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2023 Published by Elsevier B.V.



Two-Stage Label Distribution Learning with Label-Independent Prediction Based on Label-Specific Features

Gui-Lin Li^a, Heng-Ru Zhang^{a,*}, Fan Min^{a,b}, Yu-Nan Lu^c

^a*School of Computer Science, Southwest Petroleum University, Chengdu 610500, China*

^b*Institute for Artificial Intelligence, Southwest Petroleum University, Chengdu 610500, China*

^c*School of Computer Science and Engineering, Nanjing University of Science and Technology, Nanjing, 210094, China*

Abstract

Label distribution learning (LDL) explicitly models label ambiguity based on the degree to which each label describes an instance. Existing LDL algorithms typically consider all features shared by all class labels. However, irrelevant or redundant information in the feature space can degrade the performance of the model. In this study, a novel two-stage LDL algorithm with label-independent prediction based on label-specific features (TSLIP) is proposed. The first stage combines sparse feature selection and instance correlation constraints. The original high-dimensional input feature space is mapped to a derived low-dimensional label-independent prediction space. The mapping operation not only supports the extraction of label-specific features but also brings the derived space closer to the label distribution space. The second stage employs the Pearson coefficient to explore the label correlation to enhance the performance of the model. Subsequently, the label-independent predictions are normalized to obtain the final label distribution. Extensive experiments are conducted on twelve real-world datasets. The results of the seven measures demonstrate the superiority of our algorithm compared with several state-of-the-art methods.

Keywords: Label correlations, Label distribution learning, Label-independent prediction, Label-specific features.

1. Introduction

Label ambiguity [1] is a popular topic in the fields of machine learning and artificial intelligence. Early label ambiguity is dominated by single-label learning (SLL), that is, assigning a unique label to an instance, for example, assigning “dog” to a picture or rating a paper “7.” However, in several scenarios, an instance may contain numerous labels, e.g., a photo is both a landscape and portrait. Thus, multilabel learning (MLL) [2] has been designed to handle such situations. MLL can further address label ambiguity in applications such as image classification [3, 4, 5, 6] and text recognition [7, 8, 9] by assigning several labels for each instance.

In addition to whether a label belongs to an instance, the importance of the label to the instance is crucial. Label distribution learning (LDL) [10] encodes labels in the form of histograms or probability distributions, essentially aiming to answer the question, “how much does each label describe the instance?” LDL uses more general semantics compared with MLL to address label ambiguity more effectively. Geng et al. [11] introduced the concept of label distribution in the age-estimation problem and proposed an IIS-LLD algorithm. Zhang et al. [12] introduced LDL in crowd-counting problems to reduce the impact of data imbalance in the training set. Ling et al.

[13] proposed a bidirectional sliding window method to generate label distributions for individual video segments.

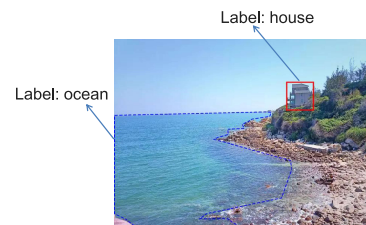


Figure 1: Description of label-specific features. A specific label is typically derived from a subset of all features. The label “house” is associated only with the features in the red box, whereas “ocean” is associated only with the features in the blue dashed frame.

However, the aforementioned studies do not consider label-specific features [14, 15]; This implies that a label is typically associated with only some features and has no clear relationship with others. Figure 1 shows an example of label-specific features; here, all the pixels constituting the image can be considered as the feature set. Intuitively, this instance could be assigned the label of “house” as the red box indicates that there is a house here. Similarly, the content of the image marked by the blue dashed frame indicates that the instance should have the label “ocean.” A specific label of the picture is often determined only by the part of the pixel value of the picture. The

*Corresponding author.

E-mail addresses: lgl@live.com (Gui-Lin Li), zhanghrswpu@163.com (Heng-Ru Zhang), minfan@swpu.edu.cn (Fan Min), luyun@njust.edu.cn (Yu-Nan Lu).

features that can infer specific labels are label-specific features.

In this regard, LDL can mine more information and significantly improve model performance by considering label-specific features. Ren et al. [16] attempted to learn specific labels for a single label and common features for all labels to enhance the model. To improve tolerance to noise, the ProLSFEO-LDL [17] method utilizes both label-specific features and prototype selection. However, in most real-world applications, even if all instances have the same values for the specific features of a given label, different instances are characterized by the label in variable degrees. Figure 2 presents an example.

Two persons may have the same feature values for the label “obesity” but different feature values for the label “cancer” (one person has cancer, whereas the other does not). When assessing their health, the obesity distribution in cancer patients is most likely to be markedly lower than that in others.

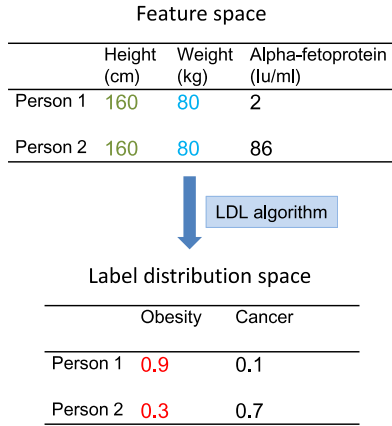


Figure 2: Common phenomenon in LDL. Different instances have the same feature values (height and weight) but different label distribution values (obesity).

To address this problem, this study proposed a two-stage LDL method with label-independent prediction based on label-specific features (TSLIP). Instead of directly generating distributions from features, a derived label-independent prediction space was injected between the feature space and label distribution space, and a two-stage learning model was built. Figure 3 shows the framework of this algorithm.

Label-independent prediction is a prediction value that considers only the specific feature of each label and foretells that label in the derived space. Instances should have the same independent prediction value on a label if they have the same specific feature value. Herein, a label-independent prediction space was introduced to effectively mine and preserve the data information in the feature space. Additionally, the label-independent prediction satisfied the instance correlation constraint. Essentially, instances with similar label-specific features should have similar label-independent prediction values. Furthermore, in the final label distribution space, the label correlation constraint

was adopted to improve the model performance.

The main contributions of this study are summarized as follows.

- 1) The phenomenon that different instances have specific feature values with the same label but different distribution values was highlighted. To the best of our knowledge, this problem has never been raised before.
- 2) The above problem was addressed by deriving a label-independent prediction space between the feature space and the label distribution space. Based on this space, a new two-stage learning model was proposed.
- 3) A sparse weighting matrix was designed to transform the feature space into a label-independent prediction space. This matrix not only selects label-specific features but also assigns different weights to each feature.
- 4) Instance correlation and label correlation were incorporated in these two stages, respectively, thereby effectively improving the generalization ability of the model.

The remainder of this paper is organized as follows. Section 2 provides a brief review of related studies. Section 3 presents the details of the proposed method. Section 4 analyzes the experimental results. Finally, Section 5 concludes the paper and provides suggestions for future work. The implementation of TSLIP is available at <https://github.com/zhanghrswpu/TSLIP>.

2. Related works

This section presents relevant work, including design strategies of LDL algorithms, introduction to label-specific features and presentation of label enhancement.

2.1. Label distribution learning

LDL generally involves three strategies in algorithm design. The first is the problem transformation (PT) strategy, which transforms the current problem into an existing problem. This is a strategy for fitting data into an algorithm. Typical algorithms include the PT-Bayes [10] and PT-SVM [18]. Both transform the LDL problem into an SLL problem. The second is the algorithm adaptation (AA) strategy, which directly uses an existing algorithm to deal with the LDL problem. This strategy adapts the algorithm to the data. Typical algorithms include AA- k NN [11] and AA-BP [19]. They obtain the label distribution value using the distribution of k neighbors and train a neural network. The third is the specialized algorithm (SA) strategy for modeling the input-to-output. Typical algorithms include the SA-IIS and SA-BFGS [10].

As an extension of the SA strategy, subsequent researchers conducted LDL through various ways. For example, Yang et al. [20] developed a multitask deep framework by jointly optimizing classification and distribution prediction. Zhao et al. [21] proposed an approach to simultaneously learn the label distribution and exploit label correlations based on optimal transport

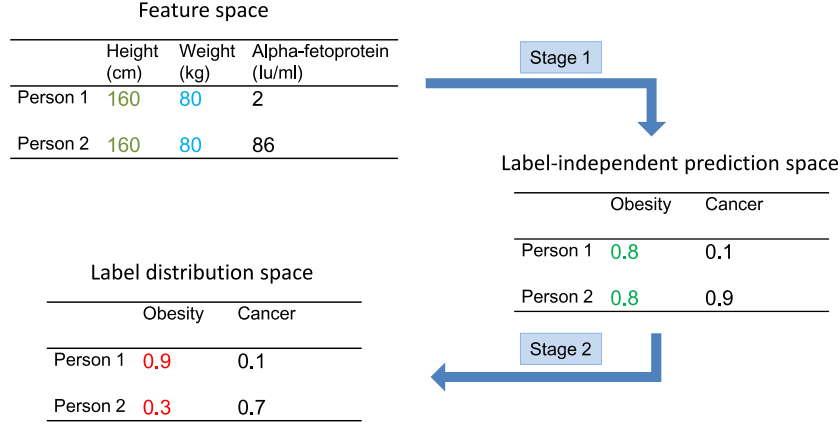


Figure 3: Illustration of the proposed label-independent prediction space and two-stage label distribution learning model. If the specific eigenvalues of a label are all the same, its label value in the label-independent prediction space is also the same. This implies that the features are the same but the distribution values are different in label distribution learning. The label-independent prediction space can be viewed as a missing space in label distribution learning.

theory. Jia et al. [22] modeled the correlation between different labels. Zhang et al. [23] based on a cosine augmentation method and label correlation for LDL. Jia et al. [24] proposed an LDL method based on a label-ranking exploration.

Numerous studies [25, 26, 27, 16] have confirmed the importance of correlation in LDL. In this study, both label correlation and instance correlation were exploited to improve the generalization ability of the model.

2.2. Label-specific features

Label-specific features is a processing strategy that customizes different features for each label instead of having all labels that share the same original features. It originates from multi-label learning and aims to extract features that are significantly related to each label to reduce noise from irrelevant features. Label-specific features can be divided into two categories based on their creation. The first is the transformation of the original features into a derived feature space. For example, LIFT [14] performs cluster analysis on positive and negative instances of each label, and then uses the Euclidean distance of each instance to the cluster center as a new feature. The LSDM [28] sets different ratio parameter values to cluster positive and negative instances of the same label. It reconstructs the label-specific feature space containing distance and spatial topology information. The second is to assign each label a specific subset of the original features, e.g., [29, 30] used $l_{2,1}$ -norm and l_1 -norm for feature sparse selection, respectively. Numerous studies [31, 32, 33] have demonstrated that label-specific features are an efficient processing strategy in multi-label learning.

Label-specific features have also been employed in several approaches for LDL. LDLSF [16] completes label-specific feature selection using two sparse weight matrices. ProLSFEO-LDL [17] uses the CHC framework, in which the search space is represented by a chromosome (or person), where the prototype choice and label-specific features are coded. In FSFL [34],

a subset of optimal and relevant features is selected from the top-ranked features based on the feature weights derived from two different strategies.

Herein, based on the correlation theory of label-specific features, a label-independent prediction space was derived and a two-stage LDL model was constructed. This is a new attempt at label-specific features in LDL. The experiments demonstrated the effectiveness of the proposed method.

2.3. Label enhancement

LDL can handle label ambiguity better by labeling an instance with a set of label distributions. However, the training sets required for LDL are limited because labeling is difficult and expensive. To address this issue, label enhancement (LE) [35] has been proposed to convert commonly used logical label datasets into distributed datasets. Numerous studies [36, 37, 38, 39] have demonstrated that label enhancement is helpful and advantageous for LDL. LE provides additional methods for data labeling as a supplement and extension of LDL.

3. TSLIP approach

This section proposes the novel TSLIP. First, the problem is defined. Second, a detailed description of the proposed algorithm is provided. Third, an optimization method for this algorithm is introduced. Finally, the time complexity of TSLIP is analyzed and a running example is provided.

3.1. Problem formulation

Let $\mathbf{Z} = \{z_1, z_2, \dots, z_c\}$ be the complete set of labels, where c denotes the number of labels and z_i denotes the i -th label. An LDL system is a tuple (\mathbf{X}, \mathbf{D}) where $\mathbf{X} = [\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_n] = [\mathbf{x}_{ij}]_{n \times m} \in \mathbb{R}^{n \times m}$ is the feature matrix, and

Table 1: Notations

Notation	Meaning
\mathbf{X}	The feature matrix of instances
\mathbf{D}	The ground-truth label distribution matrix
$\hat{\mathbf{D}}$	Predicted label distribution matrix
n	Total number of instances
m	Total number of features
c	Total number of labels
\mathbf{W}	The coefficient matrix of the label distribution for the features
\mathbf{Y}	The label-independent prediction matrix
\mathbf{W}_1	The coefficient matrix of the label-independent prediction for the features
\mathbf{W}_2	The coefficient matrix of the label distribution for the label-independent prediction
\mathbf{W}_3	Equivalent matrix of \mathbf{W}_1
\mathbf{x}_i	The feature vector of i -th instance
\mathbf{d}_i	The label distribution vector of i -th instance
\mathbf{y}_i	The label-independent prediction vector of i -th instance
d_{ij}	The j -th ground-truth label distribution value of instance \mathbf{x}_i
\hat{d}_{ij}	The j -th predicted label distribution value of instance \mathbf{x}_i
\hat{y}_{ij}	The j -th label-independent prediction value of instance \mathbf{x}_i
\mathbf{S}	Correlation matrix between labels
\mathbf{S}'	Correlation matrix between instances

$\mathbf{D} = [\mathbf{d}_1; \mathbf{d}_2; \dots; \mathbf{d}_n] = [d_{ik}]_{n \times c} \in [0, 1]^{n \times c}$ is the label distribution matrix that satisfies $\sum_{k=1}^c d_{ik} = 1$. In addition, n denotes the number of instances, m denotes the number of features, $\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{im}]$ is the feature vector of the i -th instance, and x_{ij} represents the j -th feature value of \mathbf{x}_i . $\mathbf{d}_i = [d_{i1}, d_{i2}, \dots, d_{ic}]$ denotes the label distribution associated with \mathbf{x}_i , d_{ik} represents the description degree of the label z_k to \mathbf{x}_i . The goal of LDL is to learn model $\Theta: \mathbf{X} \rightarrow \mathbf{D}$.

Table 1 lists a summary of the main notations used in this study.

3.2. Proposed method

Herein, $\mathbf{XW} = \hat{\mathbf{D}}$ was employed as the fundamental model, where $\mathbf{W} \in \mathbb{R}^{m \times c}$ denotes the regression coefficient matrix. Label-independent prediction serves as a transitional state between the feature and the distribution, which divides the model into two parts. Let $\mathbf{Y} \in \mathbb{R}^{n \times c}$ be the label-independent prediction matrix, then our model can be defined as follows.

$$\mathbf{XW}_1\mathbf{W}_2 = \mathbf{YW}_2 = \hat{\mathbf{D}}, \quad (1)$$

where $\mathbf{W}_1 \in \mathbb{R}^{m \times c}$ and $\mathbf{W}_2 \in \mathbb{R}^{c \times c}$ are the coefficient matrices of each part.

Subsequently, the formation process of the objective function is described.

Label-independent prediction was obtained using a sparse feature selection model. The general objective function can be defined as:

$$\min_{\mathbf{W}} \sum_{i=1}^n \|\mathbf{x}_i\mathbf{W} - \mathbf{d}_i\|_2^2 + \lambda \mathbf{R}(\mathbf{W}), \quad (2)$$

where $\mathbf{R}(\mathbf{W})$ is a regularizer that imposes a sparsity constraint on \mathbf{W} and λ is the regularization parameter.

As $l_{2,1}$ -norm is often used as a sparse regularizer [40, 41], our objective function can be written as:

$$\begin{aligned} \min_{\mathbf{W}_1, \mathbf{W}_2} & \|\mathbf{XW}_1\mathbf{W}_2 - \mathbf{D}\|_F^2 + \lambda_1 \|\mathbf{W}_1\|_{2,1} \\ \text{s.t.} & \quad \mathbf{XW}_1\mathbf{W}_2 \times \mathbf{1}^{c \times 1} = \mathbf{1}^{n \times 1} \\ & \quad \mathbf{XW}_1\mathbf{W}_2 \geq \mathbf{0}^{n \times c}, \end{aligned} \quad (3)$$

where $\|\cdot\|_F$ denotes the Frobenius norm of a matrix, $\|\cdot\|_{2,1}$ denotes the $l_{2,1}$ norm of a matrix, $\mathbf{1}^{c \times 1}$ and $\mathbf{1}^{n \times 1}$ represent matrices with element values of 1, and $\mathbf{0}^{n \times c}$ represents a matrix in which all elements are 0. Further, $\mathbf{XW}_1\mathbf{W}_2 \times \mathbf{1}^{c \times 1} = \mathbf{1}^{n \times 1}$ is the fundamental constraint of LDL that the sum of all the distribution values for an instance should equal 1.

Meanwhile, the label correlation constraint has been shown to yield good performance gains for LDL [16, 22, 25]. With the help of the labels associated with a given label, it can adjust the distribution value of that label. If a strong association exists between the two labels, their distribution values should be similar. In our algorithm, a label correlation constraint was imposed directly on the predicted distribution values. Therefore, the objective function can be rewritten as:

$$\begin{aligned} \min_{\mathbf{W}_1, \mathbf{W}_2} & \|\mathbf{XW}_1\mathbf{W}_2 - \mathbf{D}\|_F^2 + \lambda_1 \|\mathbf{W}_1\|_{2,1} \\ & + \lambda_2 \sum_{i=1}^n \sum_{p=1}^c \sum_{q=1}^c S_{pq} (\hat{d}_{ip} - \hat{d}_{iq})^2 \\ \text{s.t.} & \quad \mathbf{XW}_1\mathbf{W}_2 \times \mathbf{1}^{c \times 1} = \mathbf{1}^{n \times 1} \\ & \quad \mathbf{XW}_1\mathbf{W}_2 \geq \mathbf{0}^{n \times c}, \end{aligned} \quad (4)$$

where S_{pq} represents the correlation coefficient between the distributions of the p -th label z_p and q -th label z_q , λ_1 and λ_2 are the regularization parameters for the different regularization terms. S_{pq} can be calculated using the Pearson correlation function [42]:

$$S_{pq} = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{d_{ip} - \mu_p}{\sigma_p} \right) \left(\frac{d_{iq} - \mu_q}{\sigma_q} \right), \quad (5)$$

where μ_p and σ_p denote the mean and standard deviation of the p -th label distribution for all instances, respectively. μ_q and σ_q are similar to μ_p and σ_p .

Furthermore, instances with similar label-specific features should exhibit similar label-independent prediction values. By adding this constraint, our objective function can be rewritten as:

$$\begin{aligned} \min_{\mathbf{W}_1, \mathbf{W}_2} & \|\mathbf{XW}_1\mathbf{W}_2 - \mathbf{D}\|_F^2 + \lambda_1 \|\mathbf{W}_1\|_{2,1} \\ & + \lambda_2 \sum_{i=1}^n \sum_{p=1}^c \sum_{q=1}^c S_{pq} (\hat{d}_{ip} - \hat{d}_{iq})^2 \\ & + \lambda_3 \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^c S'_{ij} (y_{ik} - y_{jk})^2 \\ \text{s.t.} & \quad \mathbf{XW}_1\mathbf{W}_2 \times \mathbf{1}^{c \times 1} = \mathbf{1}^{n \times 1} \\ & \quad \mathbf{XW}_1\mathbf{W}_2 \geq \mathbf{0}^{n \times c}, \end{aligned} \quad (6)$$

where S'_{ij} represents the correlation coefficient between the i -th and j -th instances, which is still calculated using the Pearson correlation function, y_{ik} represents the k -th label-independent prediction of the i -th instance, and λ_3 is a new regularization parameter.

Eq. (6) can be derived as follows:

$$\begin{aligned} \min_{\mathbf{W}_1, \mathbf{W}_2} & \|\mathbf{XW}_1\mathbf{W}_2 - \mathbf{D}\|_F^2 + \lambda_1 \|\mathbf{W}_1\|_{2,1} \\ & + \lambda_2 \text{tr}(\widehat{\mathbf{D}}(\mathbf{V} - \mathbf{S})\widehat{\mathbf{D}}^T) \\ & + \lambda_3 \text{tr}(\mathbf{Y}^T(\mathbf{V}' - \mathbf{S}')\mathbf{Y}) \\ \text{s.t.} & \mathbf{XW}_1\mathbf{W}_2 \times \mathbf{1}^{c \times 1} = \mathbf{1}^{n \times 1} \\ & \mathbf{XW}_1\mathbf{W}_2 \geq \mathbf{0}^{n \times c}, \end{aligned} \quad (7)$$

where $\mathbf{S} = [S_{pq}]_{c \times c} \in [-1, 1]^{c \times c}$ denotes the correlation matrix of the labels, and $\mathbf{V} = \text{diag}(\mathbf{S} \times \mathbf{1}^{c \times 1}) \in \mathbb{R}^{c \times c}$ represents the diagonal matrix of the column vector transformation formed by row accumulation of \mathbf{S} . Similarly, $\mathbf{S}' = [S'_{ij}]_{n \times n} \in [-1, 1]^{n \times n}$ and $\mathbf{V}' = \text{diag}(\mathbf{S}' \times \mathbf{1}^{n \times 1}) \in \mathbb{R}^{n \times n}$ represent the correlation and diagonal matrices obtained by exploring instance correlations, respectively.

Let $\mathbf{L} = \mathbf{V} - \mathbf{S}$, and $\mathbf{L}' = \mathbf{V}' - \mathbf{S}'$; thus, the final objective function is

$$\begin{aligned} \min_{\mathbf{W}_1, \mathbf{W}_2} & \|\mathbf{XW}_1\mathbf{W}_2 - \mathbf{D}\|_F^2 + \lambda_1 \|\mathbf{W}_1\|_{2,1} \\ & + \lambda_2 \text{tr}((\mathbf{XW}_1\mathbf{W}_2)\mathbf{L}(\mathbf{XW}_1\mathbf{W}_2)^T) \\ & + \lambda_3 \text{tr}((\mathbf{XW}_1)^T \mathbf{L}' (\mathbf{XW}_1)) \\ \text{s.t.} & \mathbf{XW}_1\mathbf{W}_2 \times \mathbf{1}^{c \times 1} = \mathbf{1}^{n \times 1} \\ & \mathbf{XW}_1\mathbf{W}_2 \geq \mathbf{0}^{n \times c}. \end{aligned} \quad (8)$$

Algorithm 1: The TSLIP algorithm

Input: The feature matrix \mathbf{X} and the distribution matrix \mathbf{D}

Output: The coefficient matrix \mathbf{W}_1 and \mathbf{W}_2

- 1 Initialization: $\mathbf{W}_1^0, \mathbf{W}_2^0, \mathbf{W}_3^0, \mathbf{W}_1^0, \mathbf{W}_2^0, \mathbf{W}_3^0, \mathbf{I}_1^0, \mathbf{I}_2^0, \mathbf{I}_3^0, \lambda_1, \lambda_2, \lambda_3, \rho, t = 1$;
 - 2 Calculate the label correlation matrix \mathbf{S} ;
 - 3 Calculate the instance correlation matrix \mathbf{S}' ;
 - 4 Compute diagonal matrix \mathbf{V} and matrix \mathbf{L} based on \mathbf{S} ;
 - 5 Compute the diagonal matrix \mathbf{V}' and matrix \mathbf{L}' based on \mathbf{S}' ;
 - 6 **while** stop condition not met **do**
 - 7 update \mathbf{W}_1^{t+1} using Eq. (18) and Eq. (24) using BFGS;
 - 8 update \mathbf{W}_2^{t+1} using Eq. (19) and Eq. (25) using BFGS;
 - 9 update \mathbf{W}_3^{t+1} using Eq. (22);
 - 10 update \mathbf{W}_1^{t+1} using Eq. (23);
 - 11 update \mathbf{I}_1^{t+1} using Eq. (15);
 - 12 update \mathbf{I}_2^{t+1} using Eq. (16);
 - 13 update \mathbf{I}_3^{t+1} using Eq. (17);
 - 14 $t = t + 1$;
 - 15 **end**
-

3.3. Optimization

The alternating direction method of multipliers (ADMM) [43] has received increasing attention owing to its suitability for solving objective functions with multiple unknown coefficient

matrices. In this study, ADMM was used to solve the objective function.

First, Eq. (8) was transformed into the following equivalent form.

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3} & \|\mathbf{XW} - \mathbf{D}\|_F^2 + \lambda_1 \|\mathbf{W}_3\|_{2,1} \\ & + \lambda_2 \text{tr}((\mathbf{XW})\mathbf{L}(\mathbf{XW})^T) \\ & + \lambda_3 \text{tr}((\mathbf{XW}_1)^T \mathbf{L}' (\mathbf{XW}_1)) \\ \text{s.t.} & \mathbf{W} = \mathbf{W}_1\mathbf{W}_2 \\ & \mathbf{W}_1 = \mathbf{W}_3 \\ & \mathbf{XW} \times \mathbf{1}^{c \times 1} = \mathbf{1}^{n \times 1} \\ & \mathbf{XW} \geq \mathbf{0}^{n \times c}. \end{aligned} \quad (9)$$

The augmented Lagrangian for Eq. (9) is:

$$\begin{aligned} L_\rho(\mathbf{W}, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3, \mathbf{I}_1, \mathbf{I}_2, \mathbf{I}_3) & \\ = & \|\mathbf{XW} - \mathbf{D}\|_F^2 + \lambda_1 \|\mathbf{W}_3\|_{2,1} \\ & + \lambda_2 \text{tr}((\mathbf{XW})\mathbf{L}(\mathbf{XW})^T) + \lambda_3 \text{tr}((\mathbf{XW}_1)^T \mathbf{L}' (\mathbf{XW}_1)) \\ & + \frac{\rho}{2} \|\mathbf{W}^{t+1} - \mathbf{W}_1^t \mathbf{W}_2^t\|_F^2 + \langle \mathbf{I}_1^t, \mathbf{W}^{t+1} - \mathbf{W}_1^t \mathbf{W}_2^t \rangle \\ & + \frac{\rho}{2} \|\mathbf{W}_1^t - \mathbf{W}_3^t\|_F^2 + \langle \mathbf{I}_2^t, \mathbf{W}_1^t - \mathbf{W}_3^t \rangle \\ & + \frac{\rho}{2} \|\mathbf{XW}^t \times \mathbf{1}^{c \times 1} - \mathbf{1}^{n \times 1}\|_F^2 \\ & + \langle \mathbf{I}_3^t, \mathbf{XW}^t \times \mathbf{1}^{c \times 1} - \mathbf{1}^{n \times 1} \rangle, \end{aligned} \quad (10)$$

where ρ is the penalty factor and $\langle \cdot, \cdot \rangle$ is the Frobenius dot product. $\mathbf{I}_1 \in \mathbb{R}^{m \times c}$, $\mathbf{I}_2 \in \mathbb{R}^{m \times c}$ and $\mathbf{I}_3 \in \mathbb{R}^{n \times 1}$ are the Lagrange multipliers.

Eq. (10) can be solved using the following equations.

$$\begin{aligned} \mathbf{W}^{t+1} = \arg \min_{\mathbf{W}} & \|\mathbf{XW}^t - \mathbf{D}\|_F^2 + \lambda_2 \text{tr}[(\mathbf{XW}^t)\mathbf{L}(\mathbf{XW}^t)^T] \\ & + \frac{\rho}{2} \|\mathbf{W}^t - \mathbf{W}_1^t \mathbf{W}_2^t\|_F^2 + \langle \mathbf{I}_1^t, \mathbf{W}^t - \mathbf{W}_1^t \mathbf{W}_2^t \rangle \\ & + \frac{\rho}{2} \|\mathbf{XW}^t \times \mathbf{1}^{c \times 1} - \mathbf{1}^{n \times 1}\|_F^2 \\ & + \langle \mathbf{I}_3^t, \mathbf{XW}^t \times \mathbf{1}^{c \times 1} - \mathbf{1}^{n \times 1} \rangle, \end{aligned} \quad (11)$$

$$\begin{aligned} \mathbf{W}_1^{t+1} = \arg \min_{\mathbf{W}_1} & \lambda_3 \text{tr}[(\mathbf{XW}_1^t)^T \mathbf{L}' (\mathbf{XW}_1^t)] \\ & + \frac{\rho}{2} \|\mathbf{W}^{t+1} - \mathbf{W}_1^t \mathbf{W}_2^t\|_F^2 + \langle \mathbf{I}_1^t, \mathbf{W}^{t+1} - \mathbf{W}_1^t \mathbf{W}_2^t \rangle \\ & + \frac{\rho}{2} \|\mathbf{W}_1^t - \mathbf{W}_3^t\|_F^2 + \langle \mathbf{I}_2^t, \mathbf{W}_1^t - \mathbf{W}_3^t \rangle, \end{aligned} \quad (12)$$

$$\begin{aligned} \mathbf{W}_2^{t+1} = \arg \min_{\mathbf{W}_2} & \frac{\rho}{2} \|\mathbf{W}^{t+1} - \mathbf{W}_1^{t+1} \mathbf{W}_2^t\|_F^2 \\ & + \langle \mathbf{I}_1^t, \mathbf{W}^{t+1} - \mathbf{W}_1^{t+1} \mathbf{W}_2^t \rangle, \end{aligned} \quad (13)$$

$$\begin{aligned} \mathbf{W}_3^{t+1} = \arg \min_{\mathbf{W}_3} & \lambda_1 \|\mathbf{W}_3^t\|_{2,1} + \frac{\rho}{2} \|\mathbf{W}_1^{t+1} - \mathbf{W}_3^t\|_F^2 \\ & + \langle \mathbf{I}_2^t, \mathbf{W}_1^{t+1} - \mathbf{W}_3^t \rangle, \end{aligned} \quad (14)$$

$$\mathbf{I}_1^{t+1} = \mathbf{I}_1^t + \rho(\mathbf{W}^{t+1} - \mathbf{W}_1^{t+1} \mathbf{W}_2^{t+1}), \quad (15)$$

$$\mathbf{I}_2^{t+1} = \mathbf{I}_2^t + \rho(\mathbf{W}_1^{t+1} - \mathbf{W}_3^{t+1}), \quad (16)$$

$$\mathbf{I}_3^{t+1} = \mathbf{I}_3^t + \rho(\mathbf{XW}^{t+1} \times \mathbf{1}^{c \times 1} - \mathbf{1}^{n \times 1}). \quad (17)$$

Eqs. (15)– (17) can be directly calculated in iterations. Eqs. (11)– (14) can be simplified to the following form using the quadratic term absorption principle [44].

$$\begin{aligned} \mathbf{W}^{t+1} = \arg \min_{\mathbf{W}} & \|\mathbf{XW}^t - \mathbf{D}\|_F^2 + \lambda_2 \text{tr}[(\mathbf{XW}^t)\mathbf{L}(\mathbf{XW}^t)^T] \\ & + \frac{\rho}{2} \|\mathbf{W}^t - \mathbf{W}_1^t \mathbf{W}_2^t + \frac{1}{\rho} \mathbf{\Gamma}_1^t\|_F^2 \\ & + \frac{\rho}{2} \|\mathbf{XW}^t \times \mathbf{1}^{c \times 1} - \mathbf{1}^{n \times 1} + \frac{1}{\rho} \mathbf{\Gamma}_3^t\|_F^2, \end{aligned} \quad (18)$$

$$\begin{aligned} \mathbf{W}_1^{t+1} = \arg \min_{\mathbf{W}_1} & \lambda_3 \text{tr}[(\mathbf{XW}_1^t)^T \mathbf{L}' (\mathbf{XW}_1^t)] \\ & + \frac{\rho}{2} \|\mathbf{W}_1^t - \mathbf{W}_1^t \mathbf{W}_2^t + \frac{1}{\rho} \mathbf{\Gamma}_1^t\|_F^2 \\ & + \frac{\rho}{2} \|\mathbf{W}_1^t - \mathbf{W}_3^t + \frac{1}{\rho} \mathbf{\Gamma}_2^t\|_F^2, \end{aligned} \quad (19)$$

$$\mathbf{W}_2^{t+1} = \arg \min_{\mathbf{W}_2} \frac{\rho}{2} \|\mathbf{W}_1^{t+1} - \mathbf{W}_1^{t+1} \mathbf{W}_2^t + \frac{1}{\rho} \mathbf{\Gamma}_1^t\|_F^2, \quad (20)$$

$$\mathbf{W}_3^{t+1} = \arg \min_{\mathbf{W}_3} \lambda_1 \|\mathbf{W}_3^t\|_{2,1} + \frac{\rho}{2} \|\mathbf{W}_3^t - (\mathbf{W}_1^{t+1} + \frac{1}{\rho} \mathbf{\Gamma}_2^t)\|_F^2. \quad (21)$$

For Eq. (20), let its derivative be 0. Thus, its iterative solution is:

$$\mathbf{W}_2^{t+1} = (\mathbf{W}_1^{t+1T} \mathbf{W}_1^{t+1})^{-1} (\mathbf{W}_1^{t+1T} \mathbf{W}_1^{t+1} + \frac{1}{\rho} \mathbf{W}_1^{t+1T} \mathbf{\Gamma}_1^t). \quad (22)$$

For Eq. (21), the method proposed in [45] can be used to solve it. The optimal formulation can be expressed as:

$$\mathbf{W}_3^{t+1}(:, i) = \begin{cases} \frac{\|q_i\| - \frac{\lambda_1}{\rho}}{\|q_i\|} q_i, & \text{if } \frac{\lambda_1}{\rho} \leq \|q_i\|, \\ 0, & \text{otherwise,} \end{cases} \quad (23)$$

where $\mathbf{W}_3^{t+1}(:, i)$ denotes the i -th column of data in \mathbf{W}_3^{t+1} , q_i is the i -th column of data in matrix $\mathbf{Q} = \mathbf{W}_1^{t+1} + \frac{1}{\rho} \mathbf{\Gamma}_2^t$.

Herein, the quasi-Newton method (BFGS) [46, 47] was used to solve Eqs. (18) and (19). For optimization, BFGS is related primarily to the first-order gradient, which can be obtained using:

$$\begin{aligned} \nabla \mathbf{W}^{t+1} = & 2\mathbf{X}^T(\mathbf{XW}^t - \mathbf{D}) + \lambda_2 \mathbf{X}^T \mathbf{XW}^t (\mathbf{L}^T + \mathbf{L}) \\ & + \rho(\mathbf{W}^t - \mathbf{W}_1^t \mathbf{W}_2^t + \frac{1}{\rho} \mathbf{\Gamma}_1^t) \\ & + \rho \mathbf{X}^T (\mathbf{XW}^t \times \mathbf{1}^{c \times 1} - \mathbf{1}^{n \times 1} + \frac{1}{\rho} \mathbf{\Gamma}_3^t) \mathbf{1}^{1 \times c}, \end{aligned} \quad (24)$$

$$\begin{aligned} \nabla \mathbf{W}_1^{t+1} = & \lambda_3 \mathbf{X}^T (\mathbf{L}' + \mathbf{L}'^T) \mathbf{XW}_1^t \\ & - \rho(\mathbf{W}_1^{t+1} - \mathbf{W}_1^t \mathbf{W}_2^t + \frac{1}{\rho} \mathbf{\Gamma}_1^t) \mathbf{W}_2^T \\ & + \rho(\mathbf{W}_1^t - \mathbf{W}_3^t + \frac{1}{\rho} \mathbf{\Gamma}_2^t). \end{aligned} \quad (25)$$

Algorithm 1 shows the overall process of the proposed TSLIP algorithm.

Table 2: Computational complexity of Algorithm 1

Lines	Complexity	Description
Line 1	$O(1)$	Initialize the parameters
Line 2	$O(cn)$	Calculate the label correlation matrix
Line 3	$O(nm)$	Calculate the instance correlation matrix
Line 4	$O(c^2)$	Compute matrix \mathbf{V} and matrix \mathbf{L}
Line 5	$O(n^2)$	Compute matrix \mathbf{V}' and matrix \mathbf{L}'
Lines 6-15	$O(\tau_1 \tau_2 cn^2)$	Search the optimum values
Total	$O(1) + O(cn) + O(nm) + O(c^2) + O(n^2) + O(\tau_1 \tau_2 cn^2)$ $= O(\tau_1 \tau_2 cn^2)$	

Table 3: Raw medical data of several personnel

Personnel	Height (m)	SBP (mmHg)	TC (mmol/L)	LDC (mmol/L)	Weight (kg)	DP (mmHg)	Triglyceride (mmol/L)
Person 1	1.5	144	6.0	3.5	72	93	3.0
Person 2	1.5	120	1.2	3.0	72	85	2.0
Person 3	1.7	160	6.8	2.7	65	100	3.8

3.4. Time complexity analysis

This section elaborates on the time complexity of Algorithm 1 as follows.

Proposition 1. *The time complexity of the proposed algorithm is $O(\tau_1 \tau_2 cn^2)$.*

Proof 1. Suppose that our algorithm iterates τ_1 times, and the BFGS algorithm iterates τ_2 times. Note that the number of instances n is significantly larger than the number of features m and labels c . That is, $m \ll n$ and $c \ll n$. These conditions are considered prerequisites for eliminating small terms during the analysis process. Therefore, Line 1 initializes the parameters in $O(1)$. Line 2 calculates the label correlation matrix \mathbf{S} in $O(cn)$. Line 3 calculates the instance correlation matrix \mathbf{S}' in $O(nm)$. Line 4 computes diagonal matrix \mathbf{V} and matrix \mathbf{L} based on \mathbf{S} in $O(c^2)$. Line 5 computes the diagonal matrix \mathbf{V}' and matrix \mathbf{L}' based on \mathbf{S}' with $O(n^2)$. Lines 6 through 15 search for the optimum values in $O(\tau_1 \tau_2 cn^2)$.

The time complexity of Algorithm 1 is:

$$\begin{aligned} & O(1) + O(cn) + O(nm) + O(c^2) + O(n^2) + O(\tau_1 \tau_2 cn^2) \\ & = O(\tau_1 \tau_2 cn^2). \end{aligned} \quad (26)$$

This completes the proof. Table 2 describes the complexity of Algorithm 1 briefly.

3.5. Running example

Let us assume that a medical examination focuses primarily on obesity, high blood pressure, and high blood lipids. The

Table 4: Normalized medical data

Personnel	Height	SBP	TC	LDC	Weight	DP	Triglyceride
Person 1	0.75	0.72	0.60	0.35	0.72	0.465	0.30
Person 2	0.75	0.60	0.12	0.30	0.72	0.425	0.20
Person 3	0.85	0.80	0.68	0.27	0.65	0.500	0.38

Table 5: Label-independent prediction of three health problems

Personnel	Obesity	High Blood Pressure	High Blood Lipids
Person 1	0.852	0.636	0.640
Person 2	0.852	0.467	0.470
Person 3	0.620	0.735	0.880

Table 6: Label distribution of three health problems

Personnel	Obesity	High Blood Pressure	High Blood Lipids
Person 1	0.65	0.175	0.175
Person 2	0.90	0.050	0.050
Person 3	0.15	0.410	0.440

medical data of three persons are presented in Table 3, where SBP, TC, LDC, and DP represent the systolic blood pressure, total cholesterol, low-density cholesterol, and diastolic pressure, respectively. Table 4 presents the normalized data.

According to medical understanding, the two most important factors for evaluating whether a person is obese are their height and weight. Only systolic and diastolic blood pressure should be considered when evaluating a person's blood pressure. High blood lipid levels are correlated predominantly with levels of triglycerides, low-density cholesterol, and total cholesterol. Thus, for Person 1, the following conclusions can be drawn.

- 1) Height is 1.5 meters, weight is 72 kg, and BMI is 32, which is beyond the acceptable range of 18–24; hence, Person 1 is significantly obese.
- 2) Systolic/diastolic blood pressure is 144/93, which is slightly higher than the normal range of 140/90 and indicates moderate hypertension.
- 3) The cholesterol level is 6/3.5/3, which is a little higher than the usual level and indicates moderate hyperlipidemia.

The same analysis can also be performed for Persons 2 and 3.

Considering the severity of the three health problems independently and constraining them in the range [0,1], the health data are presented in Table 5. The degree of label representation in the domain can be regarded as a label-independent prediction. A label-independent prediction is based on label-specific features. If two instances have the same value for a particular label feature, their label-independent predictions for that label should also be the same.

Label-independent predictions are merely individual evaluations of a person's three health issues. When evaluating someone's health, all of their medical issues must be listed and a comprehensive evaluation must be provided. Thus, in our combined evaluation, the problem is marked as having a different severity if two persons have the same severity for one health issue but a different severity for the other. Therefore, a label-independent prediction should be transformed into a label distribution. Table 6 lists the transformed label distribution data.

Table 7: Characteristics of all 12 datasets

No.	Dataset	Instances	Features	Labels
1	s-JAFFE	213	243	6
2	M2B	1,240	250	5
3	Emotion6	1,980	168	7
4	Movie	7,755	1,869	5
5	Twitter.LDL	10,045	168	8
6	Flickr.LDL	11,150	168	8
7	Natural.Scene	2,000	294	9
8	FBP5500	5,500	512	5
9	Yeast.cdc	2,465	24	15
10	Yeast.cold	2,465	24	4
11	Yeast.spoem	2,465	24	2
12	Human.Gene	30,542	36	68

4. Experiments

This section presents the experiments conducted to analyze the effectiveness of the TSLIP algorithm.

4.1. Datasets

Herein, 12 real-world datasets from different domains were used for the experiments. Table 7 lists the statistics for all the datasets.

s-JAFFE [10] is an extension of the well-known facial expression image dataset, JAFFE [48]. In particular, each image was scored by 60 individuals on six basic emotion labels (happiness, sadness, surprise, fear, anger, and disgust) based on emotional conformity. The average score for each emotion was used to represent emotional intensity.

M2B [49] was processed by Ren et al. [16] based on the original dataset [50] for facial beauty perception. The label distribution was transformed using a k-wise model [49].

Emotion6 [51] contains 1,980 images assembled from Flickr for a sentiment prediction benchmark, which is annotated with votes for seven emotional categories (i.e., anger, disgust, joy, fear, sadness, surprise, and neutral). Herein, a 168-dimensional feature vector was extracted for each image using the method proposed by [16].

Movie [10] is a dataset that counts user ratings of movies. It uses 54,242,292 ratings of 7,755 movies from 478,656 users as the data source. The ratings come from Netflix and are on a scale of one–five (five labels). The rating label distribution was calculated for each movie as a percentage of the rating level. The features of the movie are extracted from metadata such as genre, director, actor, country, budget, etc. The feature vector of each movie is 1, 869-dimensional.

Flickr LDL and Twitter LDL [52] contain 11,150 and 10,045 images, respectively, whose labels fall within the typical eight-emotional space (anger, amusement, awe, contentment, disgust, excitement, fear, and sadness). Following the approach proposed in [16], feature vector extraction was performed on the two datasets using HOG[53] and Color Moment[54]. Simultaneously, PCA was used to reduce the dimension of high-dimensional features to 168 dimensions.

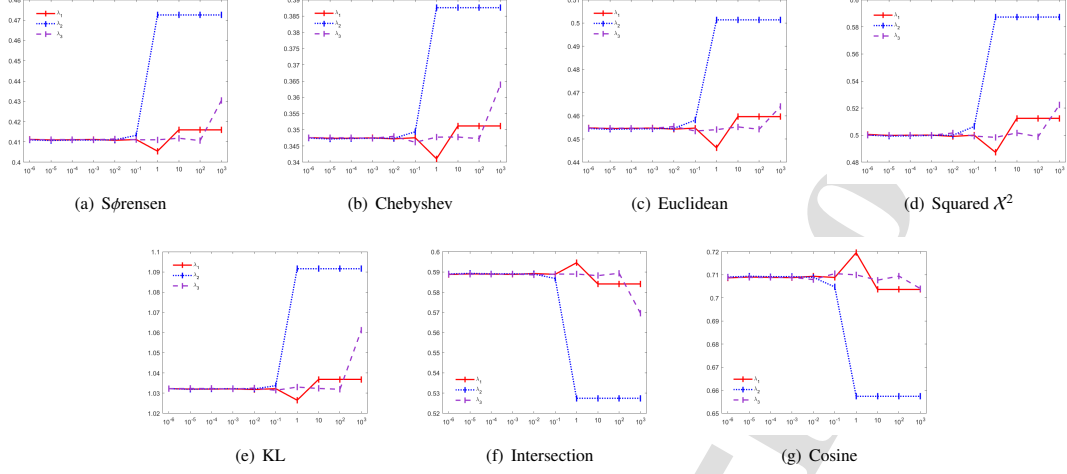
Figure 4: Influence of λ_1 , λ_2 , and λ_3 with seven metrics on dataset M2B

Table 8: Evaluation metrics for LDL algorithms

Name	Formula
Distance	Sørensen ↓ $Dis_1(d_i, \widehat{d}_i) = \frac{\sum_{j=1}^c d_{ij} - \widehat{d}_{ij} }{\sum_{j=1}^c d_{ij} + \widehat{d}_{ij} }$
	Chebyshev ↓ $Dis_2(d_i, \widehat{d}_i) = \max_{j=1}^c (d_{ij} - \widehat{d}_{ij})$
	Euclidean ↓ $Dis_3(d_i, \widehat{d}_i) = \sqrt{\sum_{j=1}^c (d_{ij} - \widehat{d}_{ij})^2}$
	Squared χ^2 ↓ $Dis_4(d_i, \widehat{d}_i) = \sum_{j=1}^c \frac{(d_{ij} - \widehat{d}_{ij})^2}{d_{ij} + \widehat{d}_{ij}}$
	Kullback-Leibler(KL) ↓ $Dis_5(d_i, \widehat{d}_i) = \sum_{j=1}^c d_{ij} \ln \frac{d_{ij}}{\widehat{d}_{ij}}$
Similarity	Intersection ↑ $Sim_1(d_i, \widehat{d}_i) = \sum_{j=1}^c \min(d_{ij}, \widehat{d}_{ij})$
	Cosine ↑ $Sim_2(d_i, \widehat{d}_i) = \frac{\sum_{j=1}^c d_{ij} \widehat{d}_{ij}}{\sqrt{\sum_{j=1}^c d_{ij}^2} \sqrt{\sum_{j=1}^c \widehat{d}_{ij}^2}}$

Finally, a Natural Scene [10] resulted from the consistent multilabel rankings of 2,000 natural scene images. FBP5500 [55] was obtained from the facial beauty perception of 5,500 color face photos. The next three datasets, Yeast_cdc, Yeast_cold, and Yeast_spoem [48], were obtained from biological experiments with yeast. The data of Human.Gene [10] were obtained from biological investigations of the connection between disease and human genes.

4.2. Evaluation metrics

To evaluate the performance of the proposed algorithm, seven widely used evaluation metrics [10, 23, 24], including five distance-based metrics (Sørensen, Chebyshev, Euclidean, Squared χ^2 , Kullback-Leibler(KL)), and two similarity-based metrics (Intersection and Cosine), were employed. Let $d_i = [d_{i1}, d_{i2}, \dots, d_{ic}]$ denote the true label distribution for the i th instance and $\widehat{d}_i = [\widehat{d}_{i1}, \widehat{d}_{i2}, \dots, \widehat{d}_{ic}]$ denotes the predicted distribution value for the same instance. Table 8 lists the formulae of these seven evaluation indicators. \downarrow indicates that the smaller the value, the better the algorithm's effects. For similarity-based metrics, \uparrow indicates that the larger the value, the better the effects of the algorithm.

4.3. Experiment settings

The proposed algorithm was compared with seven state-of-the-art algorithms, including PT-Bayes, SA-BFGS, LDL-HR [56], LDL-LDM [26], LDL-SF [16], EDL [57], and Adam-LDL-SCL [58]. PT-Bayes uses maximum likelihood estimation to estimate gaussian-like conditional probability density functions. For the remaining algorithms, the suggested parameters reported in the relevant literature were used. In particular, for SA-BFGS, the parameters of the strong Wolfe conditions were set to $c_1 = 10^{-4}$ and $c_2 = 0.9$. For LDL-HR, $\lambda_1 = 0.001$, λ_2 , and λ_3 were tuned from the candidate set $\{10^{-3}, \dots, 1\}$, and $\rho = 0.01$. For LDL-LDM, λ_1 , λ_2 , and λ_3 are tuned from $\{10^{-3}, \dots, 10^3\}$, and g was selected from 1 to 14. For LDL-SF, the parameters λ_1 , λ_2 , and λ_3 were selected from $10^{[-6, -5, \dots, -2, -1]}$, respectively, and ρ was set to 10^{-3} . For EDL, the ξ , ξ_1 , and ξ_2 were set as 0.25, 0.0001, and 0.1, respectively. For Adam-LDL-SCL, the parameters λ_1 , λ_2 and λ_3 were tuned from candidate set $10^{[-3, -2, \dots, 2, 3]}$, and the number of clusters obtained by k -means was tuned from 0 to 14. For TSLIP, the parameters λ_1 , λ_2 and λ_3 were selected from the same candidate set $10^{[-6, -5, \dots, 1, 2, 3]}$, and ρ was set to 10^{-2} .

Table 9: Comparison results on all datasets, shown as “mean \pm std.” $\uparrow(\downarrow)$ indicates the higher (lower), the better. The best results on each row are highlighted. “AVG.Rank” is the average ranking of each algorithm on the corresponding evaluation metric

Metric	Dataset	PT-Bayes	SA-BFGS	LDL-HR	LDL-LDM	LDL-SF	EDL	Adam-LDL-SCL	TSLLIP (our)
Sørensen \downarrow	s-JAFFE	.1547 \pm .0025	.1236 \pm .0009	.1374 \pm .0074	.1487 \pm .0047	.1157 \pm .0009	.1589 \pm .0081	.6919 \pm .2750	.1134\pm.0011
	M2B	.7099 \pm .0005	.4442 \pm .0026	.4656 \pm .0085	.4326 \pm .0042	.4327 \pm .0022	.4797 \pm .0060	.4400 \pm .0058	.4175\pm.0065
	Emotion6	.6031 \pm .0048	.4279 \pm .0008	.4448 \pm .0070	.4273 \pm .0029	.4152 \pm .0034	.4740 \pm .0052	.4210 \pm .0041	.3049\pm.0004
	Movie	.2800 \pm .0008	.1826 \pm .0011	.2294 \pm .0024	.1725 \pm .0004	.1882 \pm .0004	.2550 \pm .0023	.1756 \pm .0031	.1679\pm.0027
	Twitter.LDL	.5436 \pm .0063	.3812 \pm .0020	.6425 \pm .0027	.4843 \pm .0018	.3772 \pm .0004	.6608 \pm .0011	.3961 \pm .0029	.3747\pm.0038
	Flickr.LDL	.6054 \pm .0031	.4083 \pm .0001	.5732 \pm .0021	.4598 \pm .0004	.4073 \pm .0001	.5921 \pm .0011	.4224 \pm .0022	.4060\pm.0045
	Natural.Scene	.6517 \pm .0008	.4631 \pm .0030	.5059 \pm .0171	.5021 \pm .0077	.4603 \pm .0006	.6373 \pm .0087	.5245 \pm .0055	.4588\pm.0105
	FBP5500	.4985 \pm .0003	.2010 \pm .0018	.3893 \pm .0168	.1783 \pm .0021	.1632 \pm .0005	.4982 \pm .0051	.1576 \pm .0037	.1636 \pm .0031
	Yeast_cdc	.2357 \pm .0035	.0429 \pm .0004	.0476 \pm .0008	.0432 \pm .0002	.0427 \pm .0002	.0429 \pm .0008	.0421\pm.0005	.0425 \pm .0011
	Yeast_cold	.2113 \pm .0035	.0595 \pm .0009	.0598 \pm .0009	.0596 \pm .0003	.0593 \pm .0002	.0663 \pm .0020	.0593 \pm .0013	.0590\pm.0017
	Yeast_spoem	.1795 \pm .0109	.0872 \pm .0004	.0873 \pm .0037	.0871 \pm .0014	.0869 \pm .0002	.0916 \pm .0033	.0894 \pm .0035	.0868\pm.0032
	Human_Gene	.5330 \pm .0121	.2160 \pm .0003	.2439 \pm .0035	.2154 \pm .0000	.2156 \pm .0001	.2172 \pm .0015	.2163 \pm .0013	.2146\pm.0048
	AVG.Rank	7.7	3.9	5.9	3.9	2.5	6.8	4.0	1.3
	AVG.Rank	7.7	3.9	5.9	3.9	2.5	6.8	4.0	1.3
Chebyshev \downarrow	s-JAFFE	.1208 \pm .0031	.0923 \pm .0001	.1068 \pm .0099	.1182 \pm .0035	.0875 \pm .0014	.9226 \pm .0488	4.7385 \pm .19567	.0865\pm.0010
	M2B	.7063 \pm .0004	.3878 \pm .0011	.3840 \pm .0093	.3617 \pm .0052	.3608 \pm .0015	.27408 \pm .0364	2.6191 \pm .0299	.3482\pm.0007
	Emotion6	.5278 \pm .0075	.3168 \pm .0036	.3295 \pm .0092	.3188 \pm .0053	.3053 \pm .0001	.38232 \pm .0362	3.7315 \pm .0361	.3042\pm.0002
	Movie	.2030 \pm .0012	.1311 \pm .0009	.1604 \pm .0026	.1209 \pm .0005	.1317 \pm .0001	1.4102 \pm .0148	1.0672 \pm .0174	.1175\pm.0017
	Twitter.LDL	.5207 \pm .0068	.3004 \pm .0021	.5276 \pm .0034	.3875 \pm .0024	.2984 \pm .0006	6.2116 \pm .0104	6.2511 \pm .0129	.2969\pm.0046
	Flickr.LDL	.5622 \pm .0037	.3065 \pm .0010	.4358 \pm .0041	.3514 \pm .0009	.3053 \pm .0001	5.5473 \pm .0111	5.4434 \pm .0147	.3046\pm.0046
	Natural.Scene	.4092 \pm .0013	.3381 \pm .0011	.4427 \pm .0200	.3276 \pm .0070	.3208 \pm .0007	7.0062 \pm .0662	6.8069 \pm .0490	.3188\pm.0145
	FBP5500	.4082 \pm .0003	.1764 \pm .0017	.3887 \pm .0171	.1585 \pm .0013	.1452\pm.0004	.29587 \pm .0128	2.1896 \pm .0184	.1453 \pm .0033
	Yeast_cdc	.1143 \pm .0026	.0163 \pm .0001	.0181 \pm .0003	.0164 \pm .0002	.0162\pm.0002	.6518 \pm .0135	.6405 \pm .0081	.0162 \pm .0006
	Yeast_cold	.1905 \pm .0035	.0513 \pm .0006	.0515 \pm .0009	.0515 \pm .0001	.0511 \pm .0002	.2688 \pm .0083	.2407 \pm .0056	.0509\pm.0017
	Yeast_spoem	.1795 \pm .0109	.0872 \pm .0004	.0873 \pm .0037	.0871 \pm .0014	.0869 \pm .0002	.1888 \pm .0070	.1851 \pm .0077	.0868\pm.0032
	Human_Gene	.2104 \pm .0111	.0534 \pm .0002	.0536 \pm .0027	.0533 \pm .0000	.0533 \pm .0000	14.5609 \pm .0963	14.4879 \pm .0837	.0529\pm.0030
	AVG.Rank	5.8	3.5	5.0	3.4	2.0	7.8	7.2	1.2
	AVG.Rank	5.8	3.5	5.0	3.4	2.0	7.8	7.2	1.2
Euclidean \downarrow	s-JAFFE	.1560 \pm .0029	.1251 \pm .0003	.1388 \pm .0090	.1510 \pm .0043	.1168 \pm .0012	.1607 \pm .0074	.7643 \pm .3095	.1150\pm.0011
	M2B	.8912 \pm .0010	.5074 \pm .0012	.4892 \pm .0089	.4604 \pm .0050	.4753 \pm .0022	.5011 \pm .0049	.4672 \pm .0059	.4593\pm.0047
	Emotion6	.6555 \pm .0067	.4186 \pm .0012	.4279 \pm .0086	.4117 \pm .0043	.4043 \pm .0015	.4504 \pm .0065	.4101 \pm .0051	.4043\pm.0002
	Movie	.2887 \pm .0012	.1875 \pm .0012	.2334 \pm .0026	.1747 \pm .0004	.1912 \pm .0003	.2585 \pm .0022	.1784 \pm .0033	.1784\pm.0033
	Twitter.LDL	.6704 \pm .0083	.3823 \pm .0025	.6226 \pm .0029	.4706 \pm .0021	.3783 \pm .0005	.6412 \pm .0012	.3903 \pm .0032	.3772\pm.0048
	Flickr.LDL	.7101 \pm .0042	.3947 \pm .0004	.5339 \pm .0033	.4373 \pm .0006	.3940 \pm .0001	.5507 \pm .0012	.4052 \pm .0024	.3933\pm.0050
	Natural.Scene	.5658 \pm .0013	.4495 \pm .0029	.5603 \pm .0202	.4428 \pm .0082	.4340 \pm .0007	.5251 \pm .0131	.4691 \pm .0077	.4303\pm.0132
	FBP5500	.5304 \pm .0004	.2274 \pm .0028	.4945 \pm .0219	.2058 \pm .0019	.1895 \pm .0005	.5166 \pm .0046	.1846 \pm .0044	.1897 \pm .0038
	Yeast_cdc	.1652 \pm .0028	.0281 \pm .0002	.0310 \pm .0005	.0282 \pm .0002	.0279 \pm .0002	.0283 \pm .0006	.0276\pm.0004	.0278 \pm .0008
	Yeast_cold	.2444 \pm .0042	.0685 \pm .0010	.0688 \pm .0010	.0686 \pm .0003	.0683 \pm .0002	.0766 \pm .0023	.0684 \pm .0015	.0680\pm.0012
	Yeast_spoem	.2538 \pm .0155	.1234 \pm .0006	.1235 \pm .0053	.1232 \pm .0019	.1229 \pm .0001	.1295 \pm .0047	.1265 \pm .0050	.1228\pm.0045
	Human_Gene	.2774 \pm .0109	.0866 \pm .0002	.0941 \pm .0028	.0863 \pm .0000	.0873 \pm .0001	.0868 \pm .0012	.0868 \pm .0011	.0859\pm.0033
	AVG.Rank	7.8	4.1	5.9	3.8	2.8	6.6	3.8	1.2
	AVG.Rank	7.8	4.1	5.9	3.8	2.8	6.6	3.8	1.2
Squared χ^2 \downarrow	s-JAFFE	.0708 \pm .0025	.0509 \pm .0006	.0577 \pm .0063	.0664 \pm .0027	.0455 \pm .0007	.0738 \pm .0054	1.1600 \pm .5481	.0431\pm.0007
	M2B	1.1917 \pm .0034	.5676 \pm .0047	.5552 \pm .0159	.5001\pm.0072	.5394 \pm .0043	.5815 \pm .0104	.5134 \pm .0089	.5082 \pm .0085
	Emotion6	.9319 \pm .0099	.5397 \pm .0002	.5643 \pm .0145	.5349 \pm .0063	.5188 \pm .0057	.6217 \pm .0116	.5243 \pm .0096	.5179\pm.0007
	Movie	.2355 \pm .0011	.1168 \pm .0010	.1598 \pm .0030	.1026\pm.0005	.1313 \pm .0006	.1899 \pm .0031	.1147 \pm .0051	.1030 \pm .0032
	Twitter.LDL	.8837 \pm .0161	.5162 \pm .0037	1.0333 \pm .0053	.6744 \pm .0035	.5101 \pm .0001	1.0846 \pm .0024	.5412 \pm .0041	.5067\pm.0059
	Flickr.LDL	.9805 \pm .0053	.5483 \pm .0006	.8543 \pm .0043	.6314 \pm .0013	.5447 \pm .0002	.8974 \pm .0023	.5741 \pm .0040	.5431\pm.0084
	Natural.Scene	1.0218 \pm .0018	.6694 \pm .0050	.7802 \pm .0352	.7326 \pm .0134	.6662 \pm .0009	.9880 \pm .0204	.7811 \pm .0094	.6571\pm.0202
	FBP5500	.6526 \pm .0006	.1641 \pm .0018	.4839 \pm .0310	.1321 \pm .0034	.1191 \pm .0006	.6458 \pm .0103	.1085\pm.0042	.1200 \pm .0044
	Yeast_cdc	.0620 \pm .0052	.0071 \pm .0001	.0083 \pm .0003	.0072 \pm .0000	.0071 \pm .0002	.0073 \pm .0003	.0069\pm.0001	.0070 \pm .0004
	Yeast_cold	.1631 \pm .0038	.0123 \pm .0005	.0126 \pm .0006	.0124 \pm .0000	.0123 \pm .0001	.0151 \pm .0010	.0124 \pm .0005	.0122\pm.0012
	Yeast_spoem	.1177 \pm .0116	.0250 \pm .0002	.0263 \pm .0029	.0264 \pm .0001	.0248\pm.0001	.0274 \pm .0021	.0262 \pm .0019	.0248 \pm .0021
	Human_Gene	.7535 \pm .0256	.1843 \pm .0007	.2125 \pm .0063	.1836 \pm .0000	.1864 \pm .0002	.1865 \pm .0029	.1848 \pm .0024	.1826\pm.0083
	AVG.Rank	7.7	3.7	5.8	3.8	2.8	6.9	3.8	1.3
	AVG.Rank	7.7	3.7	5.8	3.8	2.8	6.9	3.8	1.3
KL \downarrow	s-JAFFE	.0744 \pm .0026	.0525 \pm .0004	.0591 \pm .0071	.0676 \pm .0027	.0506 \pm .0010	.0784 \pm .0057	7.8077 \pm .3.8722	.0445\pm.0017
	M2B	12.4931 \pm .0145	.7659 \pm .0120	.5900 \pm .0182	.7820 \pm .0114	1.1540 \pm .0395	.6204 \pm .0115	.5519\pm.0120	1.0316 \pm .0065
	Emotion6	.39536 \pm .0475	.6344 \pm .0007	.6364 \pm .0185	8.7785 \pm .1085	.7578 \pm .0203	.7065 \pm .0154	.5946\pm.0135	.6953 \pm .0026
	Movie	.2430 \pm .0175	.1378 \pm .0373	.1577 \pm .0032	.1324 \pm .0026	.2420 \pm .0023	.1887 \pm .0032	.2175 \pm .0280	.1243\pm.0045
	Twitter.LDL	2.4935 \pm .0128	.6533 \pm .0006	1.2325 \pm .0076	14.6794 \pm .0708	.7112 \pm .0004	1.3232 \pm .0045	.6598 \pm .0066	.6988 \pm .0160
	Flickr.LDL	2.0693 \pm .0130	.6571\pm.0008	.9951 \pm .0061	13.2738 \pm .0295	.7310 \pm .0013	1.0645 \pm .0032	.6737 \pm .0055	.7276 \pm .0160
	Natural.Scene	2.3733 \pm .0335	.9067 \pm .0556	7.2284 \pm .6466	7.1245 \pm .0527	1.1245 \pm .0101	1.1729 \pm .0322	.9014\pm.0126	.9671 \pm .0572
	FBP5500	1.6059 \pm .0179	.1594 \pm .0105	5.0716 \pm .7691	1.2375 \pm .0641	.1760 \pm .0022	.7005 \pm .0135	.1149\pm.0047	.1786 \pm .0121
	Yeast_cdc	.3058 \pm .0063	.0071 \pm .0001	.0082 \pm .0003	.0070 \pm .0001	.0070 \pm .0001	.0072 \pm .0003	.0071 \pm .0004	.0069\pm.0004
	Yeast_cold	.2455 \pm .0081	.0123 \pm .0001	.0125 \pm .0006	.0127 \pm .0000	.0122 \pm .0001	.0151 \pm .0010	.0125 \pm .0005	.0121\pm.0012
	Yeast_spoem	.1971 \pm .0005	.0246 \pm .0001	.0260 \pm .0029	.0415 \pm .0003	.0247 \pm .0001	.0271 \pm .0020	.0259 \pm .0018	.0245\pm.0020
	Human_Gene	1.8418 \pm .0379	.2367 \pm .0010	.2636 \pm .0107	.2256\pm.0000	.2368 \pm .0001	.2394 \pm .0050	.2375 \pm .0043	.2347 \pm .0134
	AVG.Rank	7.3	2.4	5.2	5.8	4.0	5.5	3.2	

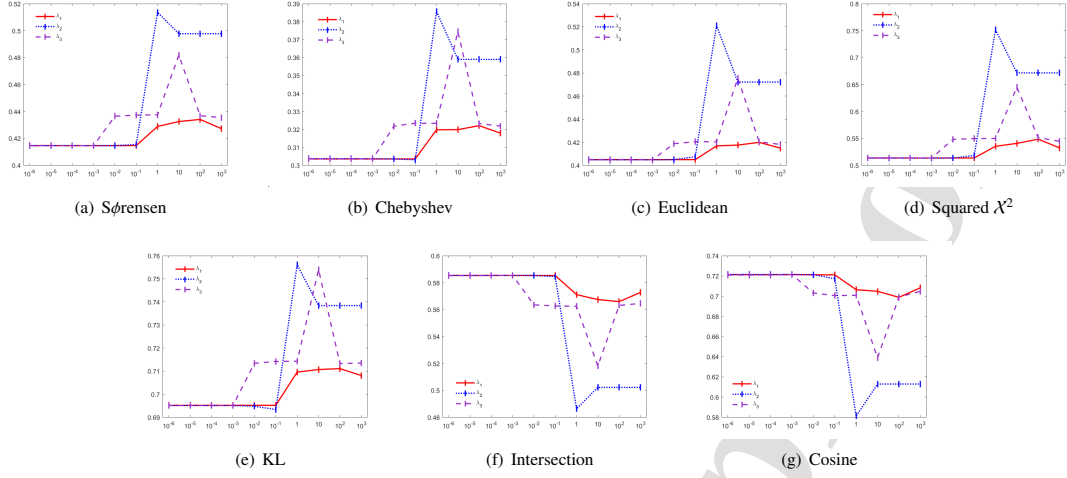
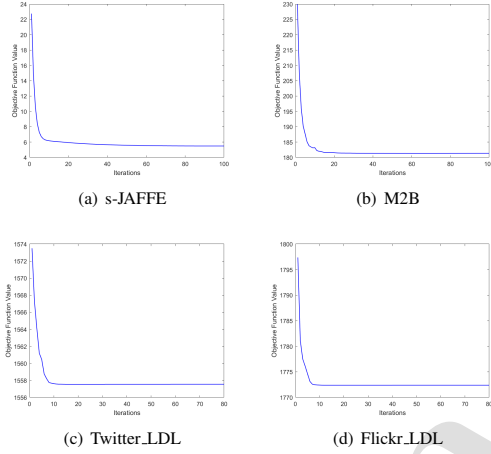
Figure 5: Influence of λ_1 , λ_2 , and λ_3 with seven metrics on dataset Emotion6

Figure 6: Convergence of TSLIP on s-JAFFE, M2B, Twitter.LDL, and Flickr.LDL.

4.4. Influence of parameters

To examine the robustness of the algorithm, the influence of the parameters in the experiment, i.e., λ_1 , λ_2 , and λ_3 in Eq. (8), was analyzed. For the analysis of λ_1 , λ_2 was fixed to 10^{-2} and λ_3 to 10^{-3} , and TSLIP was run with λ_1 set to $10^{-6}, 10^{-5}, 10^{-4}, \dots, 10^2, 10^3$. For the analysis of λ_2 , similar to λ_1 , that was set to $10^{-6}, 10^{-5}, 10^{-4}, \dots, 10^2, 10^3$. For the analysis of λ_3 , λ_1 was set to 10^{-4} and λ_2 to 10^{-2} , and TSLIP was run with λ_3 that was set to $10^{-6}, 10^{-5}, 10^{-4}, \dots, 10^2, 10^3$. Figures 4 and 5 show the experimental results for the seven metrics on the datasets M2B and Emotion6, respectively. Note that for the metrics Sørensen, Chebyshev, Euclidean, Squared χ^2 , and KL, the smaller the value, the better the performance. However, for metrics Intersection, and Cosine, the larger the value, the better the performance. The parameter λ_1 shown in Figure 4 reveals that when the value of λ_1 was 1, the performance was the

best, and when λ_1 took a small value (e.g., $10^{-3}, 10^{-2}, 10^{-1}$) or large value (e.g., $10^1, 10^2, 10^3$), the performance worsened for all seven evaluations. This is because when λ_1 is too small, the proposed label-independent prediction loss term plays a small role, whereas when λ_1 is large, the objective function shown in Eq. (8) is not dominated by the third or fourth terms. Similarly, the results for λ_2 and λ_3 exhibit the same trend as λ_1 . These phenomena provide additional proof that the proposed method is reliable.

4.5. Convergency

To examine the ADMM algorithm's effectiveness in solving the TSLIP model, the value of the objective function Eq. (8) on four datasets (s-JAFFE, M2B, Twitter.LDL, and Flickr.LDL) are plotted in Figure 6. Evidently, as the iterations progressed, the objective function value decreased. In detail, on s-JAFFE and M2B, the convergency occurred after 60 and 30 iterations, respectively, whereas on Twitter.LDL and Flickr.LDL, the convergency occurred after 20 iterations.

4.6. Results and discussion

Ten times ten-fold cross-validation were performed on each dataset, and the mean and standard deviation of each metric were recorded. Table 9 lists the experimental results of the eight comparison algorithms on the 12 datasets based on seven metrics. Each data is marked with "mean \pm std" and the best results on each row are highlighted. "AVG.Rank" is the average ranking of each algorithm on the corresponding evaluation metric. Evidently from the table, the proposed TSLIP algorithm ranked almost in the top two among the 84 sets of results, corresponding to 12 datasets and seven evaluation metrics. Among them, 58 ranked first, accounting for 69.1%, and 14 ranked second, accounting for 16.7%.

To conduct a more reliable comparative analysis, the Friedman test [59], which is a useful statistical test for comparing more than two algorithms over multiple datasets, was further examined. Table 10 summarizes the Friedman statistics F_F and

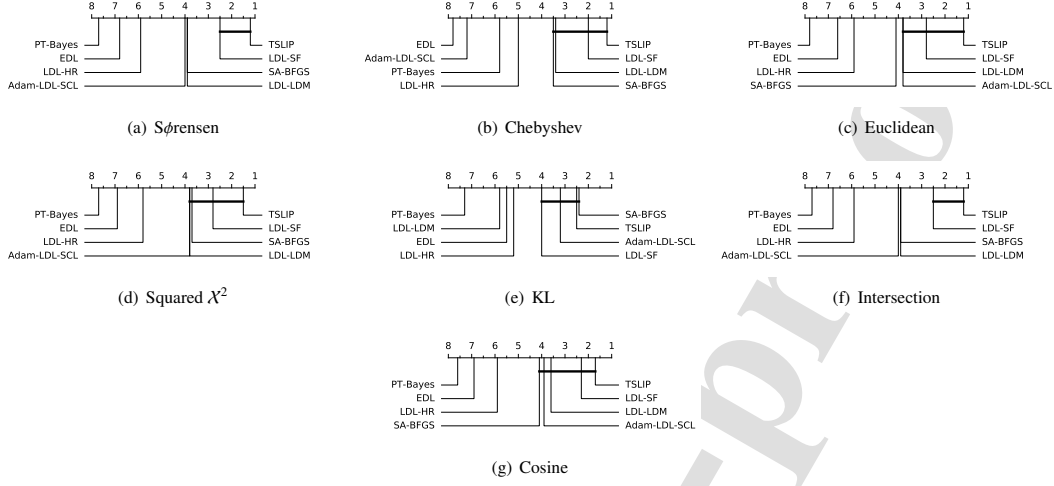


Figure 7: Comparison of our algorithm against seven comparison algorithms with the *Bonferroni-Dunn* test. Algorithms not connected with TSLIP are considered to have a significantly different performance from TSLIP (CD = 2.69 at 0.05 significance level).

Table 10: Friedman statistics F_F in terms of each evaluation metric and the critical value at 0.05 significance level (comparing algorithms $k = 8$ and datasets $N = 12$)

Metric	F_F	Critical Value
Sørensen	37.3770	2.1310
Chebyshev	123.6939	
Euclidean	36.0468	
Squared χ^2	31.7778	
KL	10.4584	
Intersection	37.3770	
Cosine	32.3396	

corresponding critical value for each measure. As presented in Table 10, for each metric, the null hypothesis of indistinguishable performance at the 0.05 significance level among the compared algorithms was rejected. Furthermore, the Bonferroni-Dunn test [59] at a 0.05 significance level was used to test whether the proposed TSLIP algorithm has competitive performance against other algorithms. The performances of the two algorithms are significantly different if their average ranks over all datasets differ by at least one critical difference (CD). Figure 7 shows CD diagrams for each metric. In each subgraph, any algorithm is connected to the TSLIP if its average rank and TSLIP's average rank are within a CD interval. Algorithms independent of TSLIP are considered to have significantly different performances. As shown in Figure 7, the proposed TSLIP method outperformed PT-Bayes, EDL, and LDL-HR on every metric. The proposed algorithm outperformed Adam-LDL-SCL on Sørensen, Chebyshev, and the intersection. Further, TSLIP performed better on Sørensen, KL, and intersections compared with LDL-LDM. In addition, the proposed approach outperformed SA-BFGS in terms of Sørensen and Euclidean.

Based on the experimental findings, we can conclude that the proposed TSLIP algorithm achieves competitive performance compared to other state-of-the-art algorithms.

5. Conclusions and future works

This study proposed a novel two-stage LDL algorithm with label-independent prediction based on label-specific features. Herein, feature sparse selection and weight coefficients were combined to map the original feature space into a derived label-independent prediction space; this not only supports the extraction of label-specific features but also brings the derived space closer to the label distribution space. Label-independent predictions were normalized to obtain the final label distribution. Experimental results on 12 real-world datasets based on seven metrics validated the effectiveness of the proposed algorithm.

Nonetheless, some issues still exist that require further investigation.

- 1) Modification of the model framework. To improve the effectiveness of model solving, techniques such as softmax can be used to remove the constraints of optimizing the objective function.
- 2) Exploring more effective ways of feature selection. In this study, the $l_{2,1}$ -norm was used for feature selection. In future work, we will try other constraint terms, such as l_1 -norm, Frobenius norm, $l_{1,2}$ -norm, $l_{1/2}$ -norm, and $l_{2/3}$ -norm, to complete the similar work.
- 3) Exploring more efficient ways to employ label-independent prediction on LDL. Label-independent prediction space is a missing space in LDL derived from inference. An important focus of our future work will be on how to more effectively

actualize the mining of this space and utilize it to anticipate label distribution.

- 4) Combination of matrix factorization and label-independent prediction for LDL. To improve computational efficiency, we can convert the initial high-dimensional sparse feature matrix into a low-rank space.
- 5) Application for LE. LE is dedicated to transforming the original multi-label database into an LDL database. To improve performance, we can add label-independent prediction constraints during label enhancement.

6. Acknowledgments

This study was supported by the National Natural Science Foundation of China (61902328), the Applied Basic Research Project of Science and Technology Bureau of Nanchong City (SXHZ040), and Central Government Funds of Guiding Local Scientific and Technological Development (2021ZYD0003).

References

- [1] Rolf Jagerman, Julia Kiseleva, and Maarten de Rijke. Modeling label ambiguity for neural list-wise learning to rank. *arXiv preprint arXiv:1707.07493*, 2017.
- [2] Grigoris Tsoumakas and Ioannis Katakis. Multi-label classification: An overview. *International Journal of Data Warehousing and Mining*, 3(3):1–13, 2007.
- [3] Yang Liu, Kai-Wen Wen, Quan-Xue Gao, Xin-Bo Gao, and Fei-Ping Nie. Svm based multi-label learning with missing labels for image annotation. *Pattern Recognition*, 78:307–317, 2018.
- [4] Bao-Yuan Wu, Siwei Lyu, Bao-Gang Hu, and Qiang Ji. Multi-label learning with missing labels for image annotation and facial action unit recognition. *Pattern Recognition*, 48:2279–2289, 2015.
- [5] Dong Liang, Xin-Bo Gao, Wen Lu, and Li-Huo He. Deep multi-label learning for image distortion identification. *Signal Processing*, 172:107536, 2020.
- [6] Amin Hashemi, Mohammad Bagher Dowlatshahi, and Hossein Nezamabadi-Pour. A bipartite matching-based feature selection for multi-label learning. *International Journal of Machine Learning and Cybernetics*, 12(2):459–475, 2021.
- [7] Sheng-Jun Huang and Zhi-Hua Zhou. Multi-label learning by exploiting label correlations locally. In *AAAI*, pages 949–955, 2012.
- [8] Xuan Wu, Qing-Guo Chen, Yao Hu, Deng-Bao Wang, Xiao-Dong Chang, Xiao-Bo Wang, and Min-Ling Zhang. Multi-view multi-label learning with view-specific information extraction. In *IJCAI*, pages 3884–3890, 2019.
- [9] Jia-Qi Lv, Tian-Ran Wu, Cheng-Lun Peng, Yun-Peng Liu, Ning Xu, and Xin Geng. Compact learning for multi-label classification. *Pattern Recognition*, 113:107–833, 2021.
- [10] Xin Geng. Label distribution learning. *IEEE Transactions on Knowledge and Data Engineering*, 28(7):1734–1748, 2016.
- [11] Xin Geng, Kate Smith-Miles, and Zhi-Hua Zhou. Facial age estimation by learning from label distributions. In *AAAI*, pages 451–456, 2010.
- [12] Zhao-Xiang Zhang, Mo Wang, and Xin Geng. Crowd counting in public video surveillance by label distribution learning. *Neurocomputing*, 166:151–163, 2015.
- [13] Miao-Gen Ling and Xin Geng. Soft video parsing by label distribution learning. *Frontiers of Computer Science*, 13(2):302–317, 2019.
- [14] Min-Ling Zhang and Lei Wu. Lift: Multi-label learning with label-specific features. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(1):107–120, 2014.
- [15] Ze-Bang Yu and Min-Ling Zhang. Multi-label classification with label-specific feature generation: A wrapped approach. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9):5199–5210, 2021.
- [16] Ting-Ting Ren, Xiu-Yi Jia, Wei-Wei Li, Lei Chen, and Ze-Chao Li. Label distribution learning with label-specific features. In *IJCAI*, pages 3318–3324, 2019.
- [17] Manuel González, José-Ramón Cano, and Salvador García. Protsfo-ldl: Prototype selection and label-specific feature evolutionary optimization for label distribution learning. *Applied Sciences*, 10(9):3089–3104, 2020.
- [18] Xin Geng, Qin Wang, and Yu Xia. Facial age estimation by adaptive label distribution learning. In *ICPR*, pages 4465–4470, 2014.
- [19] Xin Geng, Chao Yin, and Zhi-Hua Zhou. Facial age estimation by learning from label distributions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 2401–2412, 2013.
- [20] Ju-Feng Yang, Dong-Yu She, and Ming Sun. Joint image emotion classification and distribution learning via deep convolutional neural network. In *IJCAI*, pages 3266–3272, 2017.
- [21] Peng Zhao and Zhi-Hua Zhou. Label distribution learning by optimal transport. In *AAAI*, pages 4506–4513, 2018.
- [22] Xiu-Yi Jia, Wei-Wei Li, Jun-Yu Liu, and Yu Zhang. Label distribution learning by exploiting label correlations. In *AAAI*, pages 3310–3317, 2018.
- [23] Heng-Ru Zhang, Yu-Ting Huang, Yuan-Yuan Xu, and Fan Min. Cos-ldl: Label distribution learning by cosine-based distance-mapping correlation. *IEEE Access*, pages 63961–63970, 2020.
- [24] Xiu-Yi Jia, Yu-Nan Lu, and Fang-Wen Zhang. Label enhancement by maintaining positive and negative label relation. *IEEE Transactions on Knowledge and Data Engineering*, pages 1–13, 2021.
- [25] Wen-Bin Qian, Yin-Song Xiong, Jun Yang, and Wen-Hao Shu. Feature selection for label distribution learning via feature similarity and label correlation. *Information Sciences*, pages 38–59, 2022.
- [26] Jing Wang and Xin Geng. Label distribution learning by exploiting label distribution manifold. *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–14, 2021.
- [27] Xiang Zheng, Xiu-Yi Jia, and Wei-Wei Li. Label distribution learning by exploiting sample correlations locally. In *AAAI*, pages 4556–4563, 2018.
- [28] Yu-Meng Guo, Fu-Lai Chung, Guo-Zheng Li, Jian-Cong Wang, and James C Gee. Leveraging label-specific discriminant mapping features for multi-label learning. *ACM Transactions on Knowledge Discovery from Data*, 13(2):1–23, 2019.
- [29] Jun-Long Li, Pei-Pei Li, Xue-Gang Hu, and Kui Yu. Learning common and label-specific features for multi-label classification with correlation information. *Pattern Recognition*, 121:1–15, 2022.
- [30] Xiu-Yi Jia, Sai-Sai Zhu, and Wei-Wei Li. Joint label-specific features and correlation information for multi-label learning. *Journal of Computer Science and Technology*, 35(2):247–258, 2020.
- [31] Jia Zhang, Can-Dong Li, Dong-Lin Cao, Yao-Jin Lin, Song-Zhi Su, Liang Dai, and Shao-Zi Li. Multi-label learning with label-specific features by resolving label correlations. *Knowledge-Based Systems*, 159:148–157, 2018.
- [32] Zi-Wei Cheng and Zi-Wei Zeng. Joint label-specific features and label correlation for multi-label learning with missing label. *Applied Intelligence*, 50(11):4029–4049, 2020.
- [33] Jun Huang, Feng Qin, Xiao Zheng, Ze-Kai Cheng, Zhi-Xiang Yuan, Wei-Gang Zhang, and Qing-Ming Huang. Improving multi-label classification with missing labels by learning label-specific features. *Information Sciences*, 492:124–146, 2019.
- [34] Wen-Bin Qian, Qian-Zhi Ye, Yi-Hui Li, and Shi-Ming Dai. Label distribution feature selection with feature weights fusion and local label correlations. *Knowledge-Based Systems*, 256:1–16, 2022.
- [35] Ning Xu, Yun-Peng Liu, and Xin Geng. Label enhancement for label distribution learning. *IEEE Transactions on Knowledge and Data Engineering*, 33(4):1632–1643, 2019.
- [36] Ning Xu, Jun Shu, Yun-Peng Liu, and Xin Geng. Variational label enhancement. In *ICML*, pages 10597–10606, 2020.
- [37] Ning Xu, Cong-Yu Qiao, Jia-Qi Lv, Xin Geng, and Min-Ling Zhang. One positive label is sufficient: Single-positive multi-label learning with label enhancement. In *NeurIPS*, pages 1–17, 2022.
- [38] Ning Xu, Cong-Yu Qiao, Xin Geng, and Min-Ling Zhang. Instance-dependent partial label learning. In *NeurIPS*, volume 34, pages 27119–27130, 2021.
- [39] Qing-Hai Zheng, Ji-Hua Zhu, Hao-Yu Tang, Xin-Yuan Liu, Zhong-Yu Li, and Hui-Min Lu. Generalized label enhancement with sample correlations. *IEEE Transactions on Knowledge and Data Engineering*, pages

- 1–14, 2021.
- [40] Yan Cui, Jie-Lin Jiang, Zhi-Hui Lai, Zuo-Jin Hu, Yu-Quan Jiang, and Wai-Keung Wong. New semi-supervised classification using a multi-modal feature joint ℓ_{21} -norm based sparse representation. *Signal Processing: Image Communication*, pages 94–106, 2018.
 - [41] Zheng Zhang, Yong Xu, Jian Yang, Xue-Long Li, and David Zhang. A survey of sparse representation: algorithms and applications. *IEEE Access*, pages 490–530, 2015.
 - [42] Jacob Benesty, Jing-Dong Chen, Yi-Teng Huang, and Israel Cohen. Pearson correlation coefficient. In *Noise Reduction in Speech Processing*, pages 1–4, 2009.
 - [43] Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, Jonathan Eckstein, et al. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning*, pages 1–122, 2011.
 - [44] Bin Liu, Ying-Ming Li, and Zeng-Lin Xu. Manifold regularized matrix completion for multi-label learning with admm. *Neural Networks*, pages 57–67, 2018.
 - [45] Guang-Can Liu, Zhou-Chen Lin, and Yong Yu. Robust subspace segmentation by low-rank representation. In *ICML*, pages 663–670, 2010.
 - [46] Jr John E Dennis and JJ Moré. Quasi-newton methods, motivation and theory. *SIAM Review*, 19(1):46–89, 1977.
 - [47] Ya-Xiang Yuan. A modified bfgs algorithm for unconstrained optimization. *IMA Journal of Numerical Analysis*, 11(3):325–332, 1991.
 - [48] Michael Lyons, Shigeru Akamatsu, Miyuki Kamachi, and Jiro Gyoba. Coding facial expressions with gabor wavelets. In *FG*, pages 200–205, 1998.
 - [49] Yi Ren and Xin Geng. Sense beauty by label distribution learning. In *IJCAI*, pages 2648–2654, 2017.
 - [50] Tam V Nguyen, Si Liu, Bing-Bing Ni, Jun Tan, Yong Rui, and Shui-Cheng Yan. Sense beauty via face, dressing, and/or voice. In *ACM MM*, pages 239–248, 2012.
 - [51] Kuan-Chuan Peng, Tsuhan Chen, Amir Sadovnik, and Andrew C Gallagher. A mixed bag of emotions: Model, predict, and transfer emotion distributions. In *CVPR*, pages 860–868, 2015.
 - [52] Ju-Feng Yang, Ming Sun, and Xiao-Xiao Sun. Learning visual sentiment distributions via augmented conditional probability neural network. In *AAAI*, pages 224–230, 2017.
 - [53] Navneet Dalal and Bill Triggs. Histograms of oriented gradients for human detection. In *CVPR*, pages 886–893, 2005.
 - [54] Markus Andreas Stricker and Markus Orengo. Similarity of color images. In *Storage and Retrieval for Image and Video Databases III*, pages 381–392, 1995.
 - [55] Ling-Yu Liang, Luo-Jun Lin, Lian-Wen Jin, Duo-Rui Xie, and Meng-Ru Li. Scut-fbp5500: A diverse benchmark dataset for multi-paradigm facial beauty prediction. In *Pattern Recognition*, pages 1598–1603, 2018.
 - [56] Jing Wang and Xin Geng. Learn the highest label and rest label description degrees. In *IJCAI*, pages 3097–3103, 2021.
 - [57] De-Yu Zhou, Xuan Zhang, Yin Zhou, Quan Zhao, and Xin Geng. Emotion distribution learning from texts. In *EMNLP*, pages 638–647, 2016.
 - [58] Xiu-Yi Jia, Ze-Chao Li, Xiang Zheng, Wei-Wei Li, and Sheng-Jun Huang. Label distribution learning with label correlations on local samples. *IEEE Transactions on Knowledge and Data Engineering*, 33(4):1619–1631, 2019.
 - [59] Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine Learning Research*, 7:1–30, 2006.

Highlights

- By embedding a derived label-independent prediction space in the feature space and the label distribution space, we propose a Two-Stage label distribution learning algorithm with Label-Independent Prediction based on label-specific features (TSLIP).
- Label-independent prediction can ensure that different instances have the same label value under the foundation of the same specific feature values.
- We design a sparse weighting matrix to transform the feature space into the label-independent prediction space for extracting label-specific features.
- We incorporate instance correlation and label correlation to improve the generalization ability of the model.

Credit Author Statement

Credit Author Statement

Gui-Lin Li: Methodology, Software, Writing- Original draft preparation

Heng-Ru Zhang: Conceptualization, Supervision, Funding acquisition, Resources

Fan Min: Conceptualization, Supervision, Writing - Review & Editing

Yu-Nan Lu: Software, Validation, Investigation

Conflict of interest statement

We declare that we have no financial and personal relationships with other people or organizations that can inappropriately influence our work. There is no professional or other personal interest of any nature or kind in any product, service and/or company that could be construed as influencing the position presented in, or the review of, the manuscript entitled.