

Label Distribution Learning by Exploiting Sample Correlations Locally

Xiang Zheng, Xiuyi Jia*

School of Computer Science and Engineering
Nanjing University of
Science and Technology
Nanjing 210014, China

Weiwei Li

College of Astronautics
Nanjing University of
Aeronautics and Astronautics
Nanjing 210016, China

Abstract

Label distribution learning (LDL) is a novel multi-label learning paradigm proposed in recent years for solving label ambiguity. Existing approaches typically exploit label correlations globally to improve the effectiveness of label distribution learning, by assuming that the label correlations are shared by all instances. However, different instances may share different label correlations, and few correlations are globally applicable in real-world applications. In this paper, we propose a new label distribution learning algorithm by exploiting sample correlations locally (LDL-SCL). To encode the influence of local samples, we design a local correlation vector for each instance based on the clustered local samples. Then we predict the label distribution for an unseen instance based on the original features and the local correlation vector simultaneously. Experimental results demonstrate that LDL-SCL can effectively deal with the label distribution problems and perform remarkably better than the state-of-the-art LDL methods.

Introduction

Learning with ambiguity is a hot topic in machine learning and data mining areas. There are mainly two sophisticated paradigms for solving label ambiguity at present, namely single-label learning and multi-label learning (Tsoumakas, Katakis, and Taniar 2007) respectively. In single-label learning framework, an instance is associated with a single class label, whereas in multi-label learning, an instance may have multiple class labels simultaneously. A great deal of research (Read, Pfahringer, and Holmes 2008; Read et al. 2011) have shown that multi-label learning is an effective and widely applied paradigm, which can be seen as an extension of single-label learning. However, multi-label learning can only model the ambiguity of “*what describes the instance*”, i.e., it usually outputs a set of labels which are associated with the instance. Sometimes we need to deal with the more general ambiguity of “*how to describe the instance*”. For example, in some cases, we need not only know which labels are associated with an instance, but also the extent to which each label describes the instance. To solve such problems, Geng (2010) proposed label distribution learning

(LDL), which is a further extension of multi-label learning. Different from traditional multi-label learning to output a label set, LDL outputs a label distribution, and each component in the distribution represents the degree of description of the corresponding label to the instance which is called the *description degree*.

After LDL has been put forward, more and more researchers tried to solve label ambiguity problems by using LDL and achieved good results. To name a few, Geng et al. (2010) proposed the IIS-LLD algorithm and applied it into the facial age estimation problem, in which the original single-label dataset was transformed into a label distribution dataset. In their further works, to improve the efficiency of IIS-LLD, they proposed conditional probability neural network (CPNN) algorithm (Geng, Yin, and Zhou 2013) and BFGS-LLD algorithm (Geng and Ji 2014). Yang et al. (2015) attempted to combine LDL with deep learning and proposed the deep label distribution learning (DLDL) algorithm.

However, the above research does not take into account the information of label correlations, i.e., one label may be related to another one under certain conditions. For example, when a movie has a high description degree on label *action*, it is more likely to have a higher description degree on label *adventure* than label *romance*. By considering the label correlations, LDL can exploit more data information to improve the performance. Certain studies used label correlations in different ways, e.g., in (Zhou et al. 2016), the relations of emotions were captured based on the Plutchik's wheel of emotions (Plutchik 1980); Zhou et al. (2015) considered the label correlations by seeking Pearson's correlation coefficients between two labels. Although above algorithms attempted to exploit label correlations, they usually considered label correlations in a global way by assuming that the correlations were shared by all instances. However, in most real-world applications, label correlations are usually appeared in a local way, where the correlations may be shared by only a subset of instances instead of all instances. For example, as shown in Figure 1, we assume *flower* and *butterfly* have a correlation. In image (b), *butterfly* is less prominent and thus could be difficult to predict; in this case, the correlation between *flower* and *butterfly* can be helpful in learning label *butterfly* since the label *flower* is relatively easier to predict in this image. However, for image (c), with

*Corresponding author: Xiuyi Jia (jiaxy@njust.edu.cn).

Copyright © 2018, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

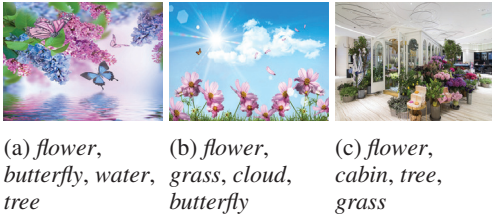


Figure 1: Illustration of non-global label correlations.

the *flower* clearly presented, this correlation turns to be misleading since *butterfly* is usually not appeared in a *cabin*. Therefore, exploiting such correlations at the global level will enforce unnecessary or even misleading constraints on some samples, and may hurt the performance by predicting some irrelevant labels.

To solve the above problem, we need to consider the label correlations at a local level. In this paper, we propose a label distribution learning algorithm by exploiting sample correlations locally (LDL-SCL). We assume that the instances can be separated into different clusters, and the instances in each cluster should be as similar as possible, thus the difference of their label distributions is as small as possible, i.e. the instances in each cluster should share the same label correlations. To reflect the influence of local samples, inspired by (Huang and Zhou 2012), we construct a local correlation vector as the additional features for each instance. Each item in the local correlation vector represents the impact of each cluster (local sample) on this instance. With the local correlation vector combined original features, we formulate LDL into a new optimization problem. Besides, we develop an alternating minimization method based on gradient descent for the optimization. Experiments on eight datasets show that our proposed algorithm exhibits a promising performance when considering local sample correlations.

The main contributions of this study can be summarized as follows: 1) local sample correlations are encoded as a local correlation vector; 2) a novel objective function and a novel output model of LDL are proposed. The rest of the paper is organized as follows. First, we discuss existing works related to our proposed approach. Second, we introduce the details of the proposed LDL-SCL algorithm. Finally, we report the experimental results and conclude the paper.

Related Works

There are three strategies for LDL at present. The first strategy is problem transformation, i.e., transforming the LDL problem into existing learning paradigms, such as traditional multi-label learning. The second strategy is algorithm adaptation, i.e., extending existing learning algorithms to deal with label distributions. The third strategy is to design specialized algorithms according to the characteristics of LDL. Different from the indirect strategies of problem transformation and algorithm adaptation, the specialized algorithms directly match the LDL problem. The related research (Geng 2016) have shown that the third strategy is more effective than other two strategies.

Generally speaking, a specialized algorithm consists of three parts, i.e., an output model, an objective function and an optimization method. For the output model, the specialized LDL algorithms mainly use the maximum entropy model (Berger, Pietra, and Pietra 1996), such as the IIS-LLD algorithm and the BFGS-LLD algorithm (Geng 2016). For the objective function, it usually chooses some functions that can measure the similarity between two distributions. Many divergence functions are adopted in existing research, such as the Kullback-Leibler (KL) divergence in the IIS-ALDL algorithm (Geng, Wang, and Xia 2014), the balanced divergence function defined in the EDL algorithm (Zhou et al. 2016), Jeffery divergence (Geng and Xia 2014), and weighted Jeffery divergence (Zhou, Xue, and Geng 2015). To solve the optimization problem, many methods were also applied in LDL. IIS-LLD algorithm (Geng 2016) used improved iterative scaling (IIS) method as the optimization method, while BFGS-LLD (Geng 2016) used BFGS and EDL (Zhou et al. 2016) used L-BFGS.

Similarly, our proposed LDL-SCL algorithm is also constructed from above three aspects. In this paper, the output model and the objective function are improved by taking into account the influence of local sample and an alternating solution is utilized for the optimization.

The LDL-SCL Approach

The Framework

Let $\mathcal{X} = \mathcal{R}^q$ denote the q -dimensional input space and $\mathcal{Y} = \{y_1, y_2, \dots, y_L\}$ denote the finite set of labels, where L is the number of labels. Given a training set $S = \{(x_1, D_1), (x_2, D_2), \dots, (x_n, D_n)\}$, where $x_i \in \mathcal{X}$ is an instance and $D_i = \{d_i^1, d_i^2, \dots, d_i^L\}$ is the label distribution associated with x_i . Such d_i^j is called the description degree of y_j to x_i and satisfies $\sum_{j=1}^L d_i^j = 1$. The goal of LDL is to learn a mapping function $f : \mathcal{X} \rightarrow D$ which can predict the label distributions for unseen instances.

As discussed in the introduction, we implement local sample correlations by adding additional features to each instance. To be specific, we construct a local correlation vector c_i for each instance x_i to reflect the influence of local samples, and expand the original feature representation with the vector. With the previous discussion, it is obvious that the length of the local correlation vector c is equal to the number of clusters. Generally, similar instances share the similar label correlations, which means similar instances will have similar c . Notice that we measure the similarity between instances in the label space instead of the feature space, because (1) instances with similar label distributions usually share similar label correlations; (2) the label space is usually much smaller than the feature space, resulting in reducing the running time. Considering both original features and local sample correlations simultaneously, we obtain the following LDL output model based on the maximum entropy model:

$$p(y_l|x_i; \theta, w, c) = \frac{1}{Z_i} \exp\left(\sum_k \theta_{l,k} x_i^k + \sum_t w_{l,t} c_i^t\right), \quad (1)$$

where $Z_i = \sum_l \exp(\sum_k \theta_{l,k} x_i^k + \sum_t w_{l,t} c_i^t)$ is a normalization term to satisfy the sum of all label description degrees of an instance equals 1. To simplify the formula, we use $p(y_l|x_i)$ to represent $p(y_l|x_i; \theta, w, c)$. The first term in the exponent part illustrates the information of original features while the second term illustrates the information of additional features, i.e., the local sample correlations. θ is the original feature coefficient matrix, in which $\theta_{l,k}$ is the element with respect to the k -th original feature and the l -th label. x_i^k is the k -th original feature of instance x_i . w is the additional feature coefficient matrix, in which $w_{l,t}$ is the element with respect to the t -th additional feature and the l -th label. c_i^t is the t -th element of local correlation vector c_i . Aiming to incorporate global discrimination fitting and local correlation influence into a unified framework, we optimize both θ, w, c by minimizing the following objective function:

$$\min_{\theta, w, c} V(\theta, w, c) + \lambda_1 \Omega_1(\theta) + \lambda_2 \Omega_2(w) + \lambda_3 Z(\theta, w, c). \quad (2)$$

V is the loss function defined on the training data. Ω_1 and Ω_2 are two penalty terms to control the complexity of the output model. Z is a penalty term to enhance the local characteristic of the correlation vector and λ_1, λ_2 and λ_3 are regularization parameters.

With the previous discussion, the goal of LDL is to make the predicted distribution and the true distribution as similar as possible, therefore the loss function should be a function which can measure the similarity of two distributions. Cha (2007) analyzed many functions that measure the similarity between two distributions, such as K-L divergence, Jeffery divergence, K divergence, etc. Here, we use K-L divergence as the loss function defined by

$$D_J(Q_a \| Q_b) = \sum_j Q_a^j \ln \frac{Q_a^j}{Q_b^j}, \quad (3)$$

where Q_a^j and Q_b^j are the j -th element of the two distributions Q_a and Q_b , respectively. Specifically, in this paper, the expression for V based on K-L divergence is defined as follows:

$$V(\theta, w, c) = \sum_{i=1}^n \sum_{l=1}^L (d_i^l \ln(\frac{d_i^l}{p(y_l|x_i)})), \quad (4)$$

where d_i^l and $p(y_l|x_i)$ denote the true description degree and the predicted description degree of label y_l to instance x_i , respectively. For the second and the third terms of Eq. (2), we utilize the F-norm of matrix to implement as follows:

$$\Omega_1(\theta) = \|\theta\|_F^2, \quad (5)$$

$$\Omega_2(w) = \|w\|_F^2. \quad (6)$$

The fourth term of Eq. (2) is utilized to enhance the local characteristic of the correlation vector. We assume that the training data can be divided into m clusters $\{G_1, G_2, \dots, G_m\}$. Each cluster denotes a local sample, in which all instances share a same label correlation. For easy implementation, we use k-means (Kanungo et al. 2002) as the clustering method. In addition, the Euclidean distance of

a cluster center and an instance is used to measure the correlation between the local sample and the instance. We use the following formula to obtain the label distribution of a cluster center:

$$p_j = \frac{1}{|G_j|} \sum_{x_k \in G_j} D_k, \quad (7)$$

where $|G_j|$ is the number of instances in the local sample G_j . Then, given an instance x_i , we define c_i^j to measure the local influence of G_j on x_i . The more similar D_i and p_j , the more likely that x_i shares the same correlation with instances in G_j , suggesting that the smaller the value of c_i^j , the smaller the impact of the local sample G_j on instance x_i . We construct a local correlation vector $c_i = [c_i^1, c_i^2, \dots, c_i^m]$ for each instance x_i , and define the penalty term on the vector as follows:

$$Z(\theta, w, c) = \sum_{i=1}^n \sum_{j=1}^m c_i^j \|p(y|x_i) - p_j\|_2^2, \quad (8)$$

where $p(y|x_i)$ is the predicted label distribution of instance x_i . By substituting Eqs. (4), (5), (6) and (8) into Eq. (2), we can obtain the following optimization problem:

$$\begin{aligned} \min_{\theta, w, c} & \sum_{i=1}^n \sum_{l=1}^L (d_i^l \ln(\frac{d_i^l}{p(y_l|x_i)})) + \lambda_1 \|\theta\|_F^2 + \lambda_2 \|w\|_F^2 \\ & + \lambda_3 \sum_{i=1}^n \sum_{j=1}^m c_i^j \|p(y|x_i) - p_j\|_2^2 \quad s.t. \quad c_i^j \geq 0, \end{aligned} \quad (9)$$

c_i^j measures the local influence of G_j on x_i , thus, it is constrained to be not less than 0.

Learning by Alternating Minimization

Eq. (9) is an optimization problem with inequality constraints, therefore we utilize an interior point method to convert it into an unconstrained problem:

$$\begin{aligned} \min_{\theta, w, c} & \sum_{i=1}^n \sum_{l=1}^L (d_i^l \ln(\frac{d_i^l}{p(y_l|x_i)})) + \lambda_1 \|\theta\|_F^2 + \lambda_2 \|w\|_F^2 \\ & + \lambda_3 \sum_{i=1}^n \sum_{j=1}^m c_i^j \|p(y|x_i) - p_j\|_2^2 + \mu \sum_{i=1}^n \sum_{j=1}^m \frac{1}{c_i^j}, \end{aligned} \quad (10)$$

where μ is the penalty factor and $\sum_{i=1}^n \sum_{j=1}^m \frac{1}{c_i^j}$ is the barrier function.

Eq. (10) can be solved by alternating minimization, i.e., optimizing one of the variables with the others fixed in each iteration. Specifically, we update one of the variables in $\{\theta, w, c\}$ with gradient descent, and leave the others fixed.

When we fix w and c to solve θ , Eq. (10) can be reduced to

$$\begin{aligned} \min_{\theta} & \sum_{i=1}^n \sum_{l=1}^L (d_i^l \ln(\frac{d_i^l}{p(y_l|x_i)})) + \lambda_1 \|\theta\|_F^2 \\ & + \lambda_3 \sum_{i=1}^n \sum_{j=1}^m c_i^j \|p(y|x_i) - p_j\|_2^2. \end{aligned} \quad (11)$$

We solve it with gradient decent and the gradient of the objective w.r.t. θ is

$$\begin{aligned} \nabla \theta_{l,k} = & - \sum_{i=1}^n d_i^l x_i^k + \sum_{i=1}^n \frac{\exp(x_i \cdot \theta_{\cdot l} + c_i \cdot w_{\cdot l}) x_i^k}{\sum_{i'} \exp(x_i \cdot \theta_{\cdot l'} + c_i \cdot w_{\cdot l'})} \\ & + 2\lambda_3 \sum_{i=1}^n \sum_{j=1}^m c_i^j (p(y_l | x_i) - p_{j,l}) \frac{\partial p(y_l | x_i)}{\partial \theta_{l,k}} \\ & + 2\lambda_1 \theta_{l,k}, \end{aligned} \quad (12)$$

where

$$\frac{\partial p(y_l | x_i)}{\partial \theta_{l,k}} = \frac{\exp(x_i \cdot \theta_{\cdot l} + c_i \cdot w_{\cdot l}) x_i^k \sum_{i' \neq l} \exp(x_i \cdot \theta_{\cdot i'} + c_i \cdot w_{\cdot i'})}{(\sum_{i'} \exp(x_i \cdot \theta_{\cdot i'} + c_i \cdot w_{\cdot i'}))^2}, x_i \cdot c_i \cdot$$

are row vectors, $\theta_{\cdot l}$ and $w_{\cdot l}$ are column vectors, and $p_{j,l}$ is the l -th element of cluster center p_j .

When we fix θ and c to solve w , Eq. (10) can be reduced to

$$\begin{aligned} \min_w \sum_{i=1}^n \sum_{l=1}^L (d_i^l \ln(\frac{d_i^l}{p(y_l | x_i)})) + \lambda_2 \|w\|_F^2 \\ + \lambda_3 \sum_{i=1}^n \sum_{j=1}^m c_i^j \|p(y_l | x_i) - p_j\|_2^2. \end{aligned} \quad (13)$$

Again, we use gradient descent and the gradient of the objective w.r.t. w is

$$\begin{aligned} \nabla w_{l,t} = & - \sum_{i=1}^n d_i^l c_i^t + \sum_{i=1}^n \frac{\exp(x_i \cdot \theta_{\cdot l} + c_i \cdot w_{\cdot l}) c_i^t}{\sum_{i'} \exp(x_i \cdot \theta_{\cdot l'} + c_i \cdot w_{\cdot l'})} \\ & + 2\lambda_3 \sum_{i=1}^n \sum_{j=1}^m c_i^j (p(y_l | x_i) - p_{j,l}) \frac{\partial p(y_l | x_i)}{\partial w_{l,t}} \\ & + 2\lambda_2 w_{l,t}, \end{aligned} \quad (14)$$

where

$$\frac{\partial p(y_l | x_i)}{\partial w_{l,t}} = \frac{\exp(x_i \cdot \theta_{\cdot l} + c_i \cdot w_{\cdot l}) c_i^t \sum_{i' \neq l} \exp(x_i \cdot \theta_{\cdot i'} + c_i \cdot w_{\cdot i'})}{(\sum_{i'} \exp(x_i \cdot \theta_{\cdot i'} + c_i \cdot w_{\cdot i'}))^2}.$$

When we fix θ and w to solve c , Eq. (10) can be reduced to

$$\begin{aligned} \min_c \sum_{i=1}^n \sum_{l=1}^L (d_i^l \ln(\frac{d_i^l}{p(y_l | x_i)})) + \mu \sum_{i=1}^n \sum_{j=1}^m \frac{1}{c_i^j} \\ + \lambda_3 \sum_{i=1}^n \sum_{j=1}^m c_i^j \|p(y_l | x_i) - p_j\|_2^2. \end{aligned} \quad (15)$$

Similarly, we use gradient descent and the gradient of the objective w.r.t. c is

$$\begin{aligned} \nabla c_i^j = & - \sum_{l=1}^L d_i^l w_{l,j} + \frac{\sum_l (\exp(x_i \cdot \theta_{\cdot l} + c_i \cdot w_{\cdot l}) w_{l,j})}{\sum_l \exp(x_i \cdot \theta_{\cdot l} + c_i \cdot w_{\cdot l})} \\ & + 2\lambda_3 c_i^j \sum_l (p(y_l | x_i) - p_{j,l}) \frac{\partial p(y_l | x_i)}{\partial c_i^j} \\ & + \lambda_3 \|p(y_l | x_i) - p_j\|_2^2 - \mu \frac{1}{c_i^j}, \end{aligned} \quad (16)$$

where

$$\frac{\partial p(y_l | x_i)}{\partial c_i^j} = - \frac{\exp(x_i \cdot \theta_{\cdot l} + c_i \cdot w_{\cdot l}) \sum_{i' \neq l} \exp((x_i \cdot \theta_{\cdot i'} + c_i \cdot w_{\cdot i'}) w_{l',j})}{(\sum_{i'} \exp(x_i \cdot \theta_{\cdot i'} + c_i \cdot w_{\cdot i'}))^2}.$$

Algorithm 1: The LDL-SCL algorithm

Input: training set $\{X, D\}$, parameters $\lambda_1, \lambda_2, \lambda_3$ and m .

Output: the label distribution D_t .

- 1 **Train:**
 - 2 initialize θ, w, c and μ ;
 - 3 **repeat**
 - 4 update θ according to Eq. (12);
 - 5 update w according to Eq. (14);
 - 6 update c according to Eq. (16);
 - 7 $\mu \leftarrow \beta \mu, 0 < \beta < 1$;
 - 8 **until convergence or maximum number of iterations;**
 - 9 **for** $j = 1$ **to** m **do**
 - 10 train a linear regression model R_j ;
 - 11 **end**
 - 12 **Test:**
 - 13 **for** $j = 1$ **to** m **do**
 - 14 predict c_t^j using R_j ;
 - 15 **end**
 - 16 return the label distribution D_t according to Eq. (1).
-

Table 1: Number of labels in eight datasets.

Dataset	alpha	cdc	elu	diau	heat	spo	cold	dtl
#Samples	2465	2465	2465	2465	2465	2465	2465	2465
#Features	24	24	24	24	24	24	24	24
#Labels	18	15	14	7	6	6	4	4

Notice that given a test instance x_t , its local correlation vector c_t is unknown. We thus train m regression models using the original features and the local correlation matrix of training instances. Then, in testing phase, the local correlation vector of the test instance can be obtained from the outputs of regression models.

The pseudo code of LDL-SCL is presented in Algorithm 1. The coefficient matrices θ and w are initialized with all elements equal 1. The local correlation matrix c is initialized with the result of k-means clustering on the label space. In detail, c_i^k is assigned 1 if x_i is in the k -th cluster and 0 otherwise. Then, we use an alternating solution to update the variables θ, w and c . Since Eq. (9) is a convex function, it will converge to a global minimum. After that, m regression models are trained with the original features as inputs and the local correlation matrix as outputs. Given a test instance x_t , the local correlation vector $c_t = [c_t^1, c_t^2, \dots, c_t^m]$ is obtained with the m regression models. Notice that the regression models can be trained with a low computational cost because usual the cluster number m is not large. Finally, the label distribution D_t is obtained by Eq. (1).

Experiments

Datasets

The datasets used in the experiments were collected from biological experiments on the budding yeast *Saccharomyces cerevisiae* (Eisen et al. 1998). There are 2465 yeast genes in

Table 2: Evaluation measures for LDL algorithms.

	Name	Formula
Distance	Euclidean↓	$Dis_1 = \sqrt{\sum_{j=1}^c (P_j - Q_j)^2}$
	Sørensen↓	$Dis_2 = \frac{\sum_{j=1}^c P_j - Q_j }{\sum_{j=1}^c P_j + Q_j }$
	Squared χ^2 ↓	$Dis_3 = \sum_{j=1}^c \frac{(P_j - Q_j)^2}{P_j + Q_j}$
	Kullback-Leibler(KL)↓	$Dis_4 = \sum_{j=1}^c P_j \ln \frac{P_j}{Q_j}$
Similarity	Intersection↑	$Sim_1 = \sum_{j=1}^c \min(P_j, Q_j)$
	Fidelity↑	$Sim_2 = \sum_{j=1}^c \sqrt{P_j Q_j}$

total, each of that is represented by an associated phylogenetic profile vector of length 24. The labels correspond to the discrete time points in different biological experiments. The gene expression level (after normalization) at each time point provides a natural measure of the description degree of the corresponding label. There are ten datasets in total in this series, and we just choose the datasets with the number of labels greater than or equal to 4 since the datasets with fewer labels are lack of the information of label correlations. The details of the eight datasets are summarized in Table 1.

Evaluation Measures

In this paper, six measures are chosen as the evaluation measures for the LDL algorithms. The names and formulas of the six measures are presented in Table 2, where P_j and Q_j are the j -th element of the ground-truth label distribution and the predicted distribution, respectively. For the former four distance measures, “↓” indicates “the smaller the better”. For the latter two similarity measures, “↑” indicates “the larger the better”.

Experimental Setting

The LDL-SCL algorithm proposed in this paper is compared with seven state-of-the-art algorithms, i.e., PT-Bayes (Geng and Ji 2014), PT-SVM (Geng, Wang, and Xia 2014), AA-kNN (Geng, Smith-Miles, and Zhou 2010), AA-BP (Geng, Yin, and Zhou 2013), IIS-LLD (Geng, Smith-Miles, and Zhou 2010), BFGS-LLD (Geng, Yin, and Zhou 2013) and EDL (Zhou, Xue, and Geng 2015). The parameter settings of all algorithms are as follows. For PT-Bayes, maximum likelihood estimation is used to estimate the Gaussian class-conditional probability density functions. PT-SVM is implemented as the “C-SVC” type in LIBSVM using the RBF kernel with the parameters $C = 1.0$ and $Gamma = 0.01$. The number of neighbors k in AA-kNN is set to 5. The number of hidden-layer neurons for AA-BP is set to 60. For LDL-SCL, the parameters in Eq. (9) are set to: $\lambda_1 = 0.001$, $\lambda_2 = 0.001$ and $\lambda_3 = 0.001$.

Experimental Results

For each dataset, 10 times 10-fold cross-validation are employed and average results are recorded. The experimental

results are shown in Table 3.

The experimental results are presented in the form of “mean±std (rank)”. Ranking refers to the prediction effect of all LDL algorithms in metrics. The smaller the value the better the performance. Furthermore, the best performance among the 8 comparing algorithms is shown in boldface. The 8 LDL algorithms include two algorithms obtained from problem transformation (PT-Bayes and PT-SVM), two algorithms obtained from algorithm adaptation (AA-kNN and AA-BP) and four specialized algorithms (IIS-LLD, BFGS-LLD, EDL and LDL-SCL). The LDL-SCL algorithm is our proposed approach in this paper.

By analyzing the experimental results shown in Table 3, in summary, the specialized LDL algorithms generally perform better than those algorithms obtained from problem transformation and algorithm adaptation. Furthermore, the proposed LDL-SCL algorithm has the best performance on datasets *alpha*, *cdc*, *elu*, *heat* and *cold*, while has the sub-optimum performance on datasets *diau*, *spo* and *dtt*. The result demonstrates the effectiveness of LDL-SCL. It is worth mentioning that, on all datasets, the LDL-SCL algorithm performs better than the EDL algorithm which exploits label correlations globally.

Influence of Parameters

In order to examine the robustness of the proposed algorithm, we also analyze the influence of parameters in the experiment, including λ_1 , λ_2 , λ_3 in Eq. (9) and the number of clusters m . We run LDL-SCL with λ_1 which is set to 0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000 and so does with λ_2 and λ_3 . In addition, we run LDL-SCL with m varying from 0 to 14 with step size of 2. Due to the page limit, we only present the experimental results on dataset *cold* with 6 evaluation measures. The results are shown in Figure 2 and Figure 3. Notice that, for criteria *Euclidean*, *Sørensen*, *Squared χ^2* and *KL*, the smaller the value, the better the performance; but for criteria *Intersection* and *Fidelity*, the larger the value, the better the performance. As for parameters λ_1 , λ_2 and λ_3 shown in Figure 2, we can find that, the performance gets worse when three parameters take large values since the objective function is not dominated by the first term. Besides, we can see that the influence of λ_2 is smaller than the influence of λ_1 and λ_3 . As for the number of clusters m shown in Figure 3, it has the worst result when m is equal to 0, which means the influence of local samples is ignored and the result is similar to the results of IIS-LLD and BFGS-LLD. Furthermore, when m is less than 6, the performance of LDL-SCL rises steadily, otherwise, the performance rises slowly and tends to be stable.

Convergency

In order to investigate the convergence of the alternating minimization algorithm to solve LDL-SCL model, we plot the value of objective function Eq. (10) on two datasets (*alpha* and *cold*) in Figure 4. It can be observed that the objective function value decreases with respect to iterations, and the value approaches to be a fixed value after a few iterations (about 10 iterations for *alpha* and *cold*).

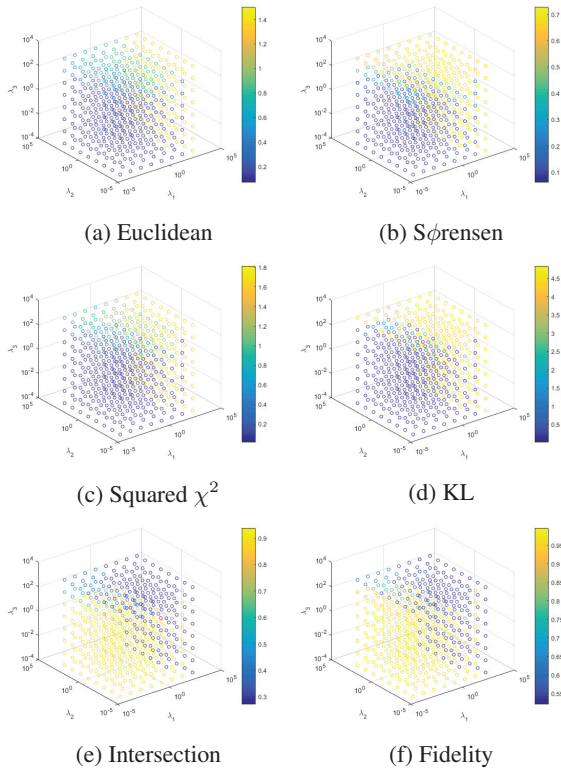


Figure 2: Influence of λ_1 , λ_2 and λ_3 with 6 measures on dataset *cold*.

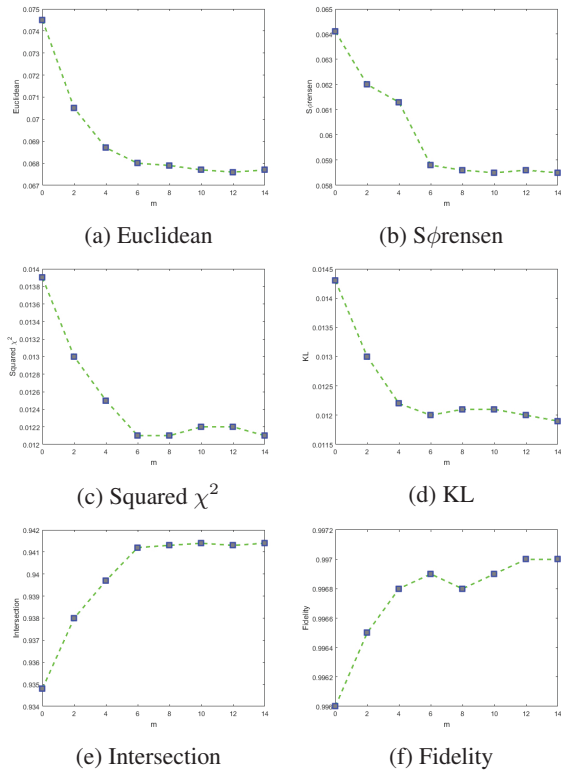


Figure 3: Influence of m with 6 measures on dataset *cold*.

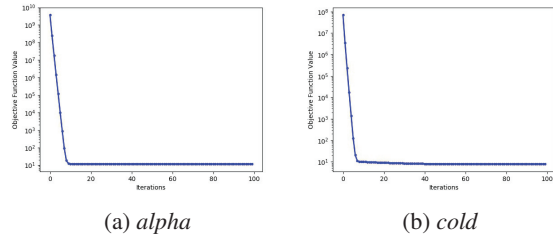


Figure 4: Convergence of LDL-SCL on *alpha* and *cold*.

and propose the LDL-SCL algorithm. The experimental results demonstrate that LDL-SCL can effectively deal with the label distribution problems and perform remarkable better than the state-of-the-art LDL methods. In future work, we will try other measures to reflect the influence of local samples and design related optimization methods with faster convergence speed.

Acknowledgements

This work was supported by the National Natural Science Foundation of China (61773208, 61403200), and the Natural Science Foundation of Jiangsu Province (BK20170809).

References

- Berger, A. L.; Pietra, V. J. D.; and Pietra, S. A. D. 1996. A maximum entropy approach to natural language processing. *Computational Linguistics* 22(1):39–71.
- Cha, S. H. 2007. Comprehensive survey on distance/similarity measures between probability density functions. *International Journal of Mathematical Models & Methods in Applied Sciences* 1(4):300–307.
- Eisen, M. B.; Spellman, P. T.; Brown, P. O.; and Botstein, D. 1998. Cluster analysis and display of genome-wide expression patterns. *Proceedings of the National Academy of Sciences of the United States of America* 95(25):14863–14868.
- Geng, X., and Ji, R. 2014. Label distribution learning. In *IEEE International Conference on Data Mining Workshops*, 377–383.
- Geng, X., and Xia, Y. 2014. Head pose estimation based on multivariate label distribution. In *Computer Vision and Pattern Recognition*, 1837–1842.
- Geng, X.; Smith-Miles, K.; and Zhou, Z. H. 2010. Facial age estimation by learning from label distributions. In *Twenty-Fourth AAAI Conference on Artificial Intelligence*, 451–456.
- Geng, X.; Wang, Q.; and Xia, Y. 2014. Facial age estimation by adaptive label distribution learning. In *IEEE International Conference on Pattern Recognition*, 4465–4470.
- Geng, X.; Yin, C.; and Zhou, Z. H. 2013. Facial age estimation by learning from label distributions. *IEEE Transactions on Pattern Analysis & Machine Intelligence* 35(10):2401–2412.
- Geng, X. 2016. Label distribution learning. *IEEE Transactions on Knowledge & Data Engineering* 28(7):1734–1748.

- Huang, S. J., and Zhou, Z. H. 2012. Multi-label learning by exploiting label correlations locally. In *Twenty-Sixth AAAI Conference on Artificial Intelligence*, 949–955.
- Kanungo, T.; Mount, D. M.; Netanyahu, N. S.; Piatko, C. D.; Silverman, R.; and Wu, A. Y. 2002. An efficient k-means clustering algorithm: Analysis and implementation. *IEEE Transactions on Pattern Analysis & Machine Intelligence* 24(7):881–892.
- Plutchik, R. 1980. *Chapter 1 - A General Psychoevolutionary Theory of Emotion*. Elsevier Inc.
- Read, J.; Pfahringer, B.; Holmes, G.; and Frank, E. 2011. Classifier chains for multi-label classification. *Machine Learning* 85(3):333–348.
- Read, J.; Pfahringer, B.; and Holmes, G. 2008. Multi-label classification using ensembles of pruned sets. In *Eighth IEEE International Conference on Data Mining*, 995–1000.
- Tsoumakas, G.; Katakis, I.; and Taniar, D. 2007. Multi-label classification: An overview. *International Journal of Data Warehousing & Mining* 3(3):1–13.
- Yang, X.; Gao, B. B.; Xing, C.; and Huo, Z. W. 2015. Deep label distribution learning for apparent age estimation. In *IEEE International Conference on Computer Vision Workshop*, 344–350.
- Zhou, D.; Zhang, X.; Zhou, Y.; Zhao, Q.; and Geng, X. 2016. Emotion distribution learning from texts. In *Conference on Empirical Methods in Natural Language Processing*, 638–647.
- Zhou, Y.; Xue, H.; and Geng, X. 2015. Emotion distribution recognition from facial expressions. In *ACM International Conference on Multimedia*, 1247–1250.