

Label correlation for partial label learning

*
GE Lingchi, FANG Min , LI Haikun, and CHEN Bo

School of Computer Science and Technology, Xidian University, Xi'an 710071, China

Abstract: Partial label learning aims to learn a multi-class classifier, where each training example corresponds to a set of candidate labels among which only one is correct. Most studies in the label space have only focused on the difference between candidate labels and non-candidate labels. So far, however, there has been little discussion about the label correlation in the partial label learning. This paper begins with a research on the label correlation, followed by the establishment of a unified framework that integrates the label correlation, the adaptive graph, and the semantic difference maximization criterion. This work generates fresh insight into the acquisition of the learning information from the label space. Specifically, the label correlation is calculated from the candidate label set and is utilized to obtain the similarity of each pair of instances in the label space. After that, the labeling confidence for each instance is updated by the smoothness assumption that two instances should be similar outputs in the label space if they are close in the feature space. At last, an effective optimization program is utilized to solve the unified framework. Extensive experiments on artificial and real-world data sets indicate the superiority of our proposed method to state-of-art partial label learning methods.

Keywords: pattern recognition, partial label learning, label correlation, disambiguation.

DOI: [10.23919/JSEE.2022.000102](https://doi.org/10.23919/JSEE.2022.000102)

1. Introduction

Although learning from the training examples associated with accurate labels is effective, collecting such labeled data is expensive in many real-world classification tasks. The aim of partial label learning is to learn a multi-class classifier from ambiguously labeled examples which can be easily obtained. Recently, partial label learning has arisen in many real-world applications, such as ecoinformatics [1], natural language processing [2], and automatic image annotation [3,4].

The difficulty of partial label learning is that the

ground-truth label hidden in the candidate label set cannot be accessed by the learning algorithms directly. The candidate label disambiguation is a straightforward method to deal with partial label learning. The current disambiguation-based approaches can be divided into two categories: the averaging-based strategy and the identification-based strategy. For the averaging-based strategy, the candidate labels of each training example are treated equally and the predictive model is made by averaging their modeling output [5,6]. For the identification-based strategy, the ground-truth label is regarded as a latent variable which can be determined by an iterative refining procedure [7–12].

It is well-known that the label correlation has a pivotal role in multi-label learning [13]. So far, little attention has been paid to the role of label correlation in the multi-class classification because of a lack of the label correlation information in the label space. However, partial label learning could be a contributing factor to build a multi-class classifier by the label correlation. Specifically, if a pair of classes is hard to distinguish, they are easy to appear in the candidate label set of the same sample in the partial label learning. In terms of this agreement, label correlation can be built from the label space. After that, label correlation is utilized to disambiguate candidate label sets by the smoothness assumption that two instances should be the similar output in the label space if they are close in the feature space.

In the label space, the global label information can be extracted by label correlation. The label relationship at the instance level can be extracted by the semantic difference maximization criterion [12]. At last, to overcome the influence of noise and outliers in the feature space, the adaptive graph [10] can be integrated into the unified framework. Then, a novel approach named label correlation, semantic difference maximization, and adaptive graph for partial label learning (PL-LCSA) is proposed in this paper.

There are two innovations in this paper: (i) To the best of our knowledge, this paper is the first one that label cor-

Manuscript received May 10, 2021.

*Corresponding author.

This work was supported by the National Natural Science Foundation of China (62176197;61806155) and the National Natural Science Foundation of Shaanxi Province (2020GY-062).

relation is introduced to deal with partial label learning, and the information from the label space can be further expanded. (ii) Label correlation, semantic difference maximization, and the adaptive graph are integrated into a learning framework, in which the information from the label space and the feature space can be learned simultaneously. Comprehensive experiments show that PL-LCSA achieves competitive performance against the state-of-the-art partial label learning approaches in the artificial and real-world data sets.

2. Related work

Formally, let $\mathcal{X}=\mathbf{R}^d$ be the d -dimensional feature space and $\mathcal{Y}=\{0,1\}^l$ be the label space with l class labels. The training set can be denoted by $\mathcal{D}=\{(\mathbf{x}_i, Z_i)|1\leq i\leq m\}$, where $\mathbf{x}_i\in\mathcal{X}$ is a d -dimensional feature vector and $Z_i\in\mathcal{Y}$ is the candidate label set associated with \mathbf{x}_i . The ground-truth label \mathbf{y}_i^* associated with \mathbf{x}_i is assumed to reside in the candidate label set Z_i , i.e., $\mathbf{y}_i^*\in Z_i$. The aim of the partial label learning is to learn a multi-class classifier $f: \mathcal{X}\rightarrow\mathcal{Y}$ from the training set \mathcal{D} .

The ground-truth label which resides in the candidate label set is not directly accessible to the learning algorithm. Disambiguation is a major approach to recover the ground-truth labeling information. There are currently two main disambiguation strategies, i.e., the averaging-based strategy and the identification-based strategy.

The averaging-based strategy is to treat the candidate labels in an identical manner and predict unseen instances by averaging the output of candidate labels. The k -nearest neighbor technique for partial label learning was proposed in [5], and the candidate label set of neighbor instances are integrated by weighted voting to make the prediction of unseen instances, i.e., $\arg\max_{\mathbf{y}\in\mathcal{Y}}\sum_{\mathbf{x}_i\in N_k(\mathbf{x}_i)}\omega_i\Pi(\mathbf{y}\in Z_i)$ ($N_k(\mathbf{x}_i)$ is the set of k -nearest neighbor for \mathbf{x}_i). The deficiency of the averaging-based strategy is that the output of the ground-truth label will be overwhelmed by false positive labels.

To overcome the weakness of the averaging-based strategy, the identification-based strategy is to make the ground-truth label a latent variable which can be identified by an iterative process. There are several techniques to iteratively refine the ground-truth labeling confidences, e.g., the maximum latent semantic differences criterion [14], the maximum-likelihood technique [15], the maximum margin criterion [7,16], the dictionary-based learning criterion [17], the boosting technique [18], the heterogeneous loss with sparse and low-rank regularization [19], artificial neural network technology [4,20].

Recently, the feature-aware disambiguation strategy aims to disambiguate the candidate label set by the infor-

mation from the feature space. The reconstructing information from the k -nearest neighbor is used to update the labeling confidence matrix iteratively [8]. The manifold structure of the feature space is propagated to the label space for disambiguation [9]. To overcome the noise and outliers in the feature space, the adaptive graph is utilized to guide disambiguation [11]. In the label space, the semantic difference maximization criterion aims to maximize the latent semantic differences of two instances which do not share any common candidate labels [14].

Different from the method of disambiguation, some algorithms work by binary decomposition for partial label learning. The disambiguation-free approach learns the predictive model by the error-correcting output codes (ECOC) technique [21]. The One-vs-One decomposition strategy is adapted to solving the partial label learning problem [22].

Most studies in the field of the label space have only focused on the difference between candidate labels and non-candidate labels, and ignored the label correlation. In the next section, we present a novel approach which introduces the label correlation to partial label learning and simultaneously utilizes the information from the feature space and the label space to disambiguate the candidate label set.

3. Approach

In this section, we introduce the proposed approach PL-LCSA. Firstly, we present the label correlation and the uniform learning framework. After that, an effective optimization procedure is utilized to deal with this framework.

3.1 Label correlation

We denote $\mathbf{X}=[\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m]^T\in\mathbf{R}^{m\times d}$ as the feature matrix, and $\mathbf{Y}=[\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_m]^T\in\mathbf{R}^{m\times l}$ as the partial label matrix where $y_{ij}=1$ means that the j -label is one of the candidate labels of \mathbf{x}_i , i.e., $y_{ij}\in Z_i$. Otherwise, $y_{ij}=0$. Let $\mathbf{P}=[p_{ij}]_{m\times l}$ be the labeling confidence matrix, where m denotes the number of training examples, l denotes the number of class labels, and p_{ij} indicates the probability of the j th label as the ground-truth label associated with instance \mathbf{x}_i . For each training example (\mathbf{x}_i, Z_i) , we aim to generate the labeling confidence vector $\mathbf{p}_i=[p_{i1}, p_{i2}, \dots, p_{il}]$. \mathbf{P} should be constrained to the following condition:

$$\sum_{j=1}^l p_{i,j} = 1, \quad \forall i \in \{1, 2, \dots, m\}, \quad (1)$$

$$0 \leq p_i \leq \mathbf{y}_i, \quad \forall i \in \{1, 2, \dots, m\}. \quad (2)$$

The aim of (1) is to normalize \mathbf{p}_i , and make the sum of

labeling confidence vector \mathbf{p}_i be 1 for each instance. The aim of (2) is to make sure that the labeling confidence of non-candidate labels must be 0, and the ground-truth label resides in the candidate label set. Once the labeling confidence vectors are normalized, the partial label training set \mathcal{D} can be transformed into its disambiguation counterpart $\mathcal{D}=\{(\mathbf{x}_i, \mathbf{p}_i)|1 \leq i \leq m\}$. Then the predictive model can be learned by the disambiguation results \mathcal{D} .

The label correlation matrix can be denoted as $\mathbf{B} \in \mathbf{R}^{l \times l}$, where $b_{i,j}$ is the similarity between the i th and the j th labels. The label correlation matrix can be calculated by the cosine similarity:

$$b_{i,j} = \cos(\mathbf{q}_i, \mathbf{q}_j) = \frac{\langle \mathbf{q}_i, \mathbf{q}_j \rangle}{\|\mathbf{q}_i\| \|\mathbf{q}_j\|} = \frac{\sum_{s=1}^m \mathbf{q}_{i,s} \mathbf{q}_{j,s}}{\sqrt{\sum_{s=1}^m \mathbf{q}_{i,s}^2} \sqrt{\sum_{s=1}^m \mathbf{q}_{j,s}^2}} \quad (3)$$

where \mathbf{q}_i is the i th column of label matrix $\mathbf{Y}=[\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_m]^T$, i.e., $\mathbf{Y}=[\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_l]$.

To disambiguate the candidate label set by the label correlation, the labeling confidence matrix will be generated by the smoothness assumption that two points should be the similar output in the label space if they are close in the feature space.

According to this assumption, we can calculate the similarity of instances in the label space by

$$\text{sim}L_{i,j} = \mathbf{p}_i^T \mathbf{B} \mathbf{p}_j = \sum_{k_1=1}^l \sum_{k_2=1}^l p_{i,k_1} b_{k_1,k_2} p_{j,k_2} \quad (4)$$

where $\text{sim}L_{i,j}$ is the similarity between instance \mathbf{x}_i and instance \mathbf{x}_j in the label space. The larger the similarity is, the bigger $\text{sim}L_{i,j}$ is.

Then the similarity of instances in the feature space can be obtained according to

$$h_{i,j} = \begin{cases} \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right), & \mathbf{x}_i \in N_k(\mathbf{x}_j) \\ 0, & \mathbf{x}_i \notin N_k(\mathbf{x}_j) \end{cases} \quad (5)$$

where σ is the average Euclidean distance among each pair of instances, i.e., $\sigma = \left(\sum_{i=1}^m \sum_{j=1}^m \|\mathbf{x}_i - \mathbf{x}_j\|_2\right) \cdot m^{-2}$. $h_{i,j}$ is the similarity between instance \mathbf{x}_i and instance \mathbf{x}_j in the feature space.

Then we should model the smoothness assumption by minimizing the gap of instances similarity between the label space and the feature space, and the label confidence matrix can be calculated by

$$\begin{aligned} \min_{\mathbf{P}} \sum_{i=1}^m \sum_{j=1}^m (h_{i,j} - \mathbf{p}_i^T \mathbf{B} \mathbf{p}_j)^2, \\ \text{s.t. } \mathbf{P} \mathbf{1}_l = \mathbf{1}_m, \mathbf{0}_{m \times l} \leq \mathbf{P} \leq \mathbf{Y}, \end{aligned} \quad (6)$$

where $\mathbf{1}_l=[1,1,\dots,1]^T \in \mathbf{R}^l$ is an all 1 vector with size l , $\mathbf{0}_{m \times l}$ is an all 0 matrix with size $m \times l$, and constraints are the matrix form of (1) and (2) respectively.

3.2 Uniform learning framework

The feature and label information are simultaneously utilized to disambiguate the candidate label sets in our uniform learning framework.

In the feature space, we adopt the adaptive graph to recover the intrinsic manifold structure within the data more robustly and accurately. Let $\mathcal{G}=(\mathcal{N}, \mathcal{E}, \mathcal{S})$ be a weighted graph, where $\mathcal{N}=\{\mathbf{x}_i|1 \leq i \leq m\}$ is the node of the graph, and $\mathcal{E}=\{\{\mathbf{x}_i, \mathbf{x}_j\}|\mathbf{x}_j \in N_k(\mathbf{x}_i), 1 \leq i \leq m\}$ is a set of edges from \mathbf{x}_i to \mathbf{x}_j if and only if \mathbf{x}_j is among the k -nearest neighbors of \mathbf{x}_i . $\mathcal{S} \in \mathbf{R}^{m \times m}$ corresponds to the information of the manifold structure in the feature space. $\mathbf{L} \in \mathbf{R}^{m \times m}$ is an index matrix where $l_{i,j}=1$ if $(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{E}$; otherwise, $l_{i,j}=0$. The labeling confidence matrix can be calculated [11] by solving the following problem:

$$\begin{aligned} \min_{\mathbf{P}, \mathcal{S}, \mathbf{W}} \sum_{i=1}^m \|f(\mathbf{x}_i, \mathbf{W}) - \mathbf{p}_i\|_2^2 + \alpha \sum_{i=1}^m \|\mathbf{p}_i - \sum_{l_{i,j}=1} s_{i,j} \mathbf{p}_j\|_2^2 + \\ \beta \sum_{i=1}^m \|\mathbf{x}_i - \sum_{l_{i,j}=1} s_{i,j} \mathbf{x}_j\|_2^2 + \lambda \|\mathbf{W}\|_F^2, \\ \text{s.t. } \begin{cases} \mathbf{P} \mathbf{1}_l = \mathbf{1}_m \\ \mathbf{0}_{m \times l} \leq \mathbf{P} \leq \mathbf{Y} \\ \mathbf{S} \mathbf{1}_m = \mathbf{1}_m \\ \mathbf{0}_{m \times m} \leq \mathcal{S} \leq \mathbf{L} \end{cases}, \end{aligned} \quad (7)$$

where $f(\mathbf{x}_i, \mathbf{W})$ is a predictive model, and \mathbf{W} is the model parameter. The first term is to learn \mathbf{W} . The second term is to learn the labeling confidence matrix, the third term is to determine the adaptive graph weight matrix, and the fourth term is a regularization term.

In the label space, PL-LCSA aims to disambiguate candidate label sets by the label correlation and the semantic difference maximization. The aim of semantic difference maximization is to maximize the semantic differences of the two instances which do not share any common candidate labels. It can be expressed [14] by the following problem:

$$\max_{\mathbf{P}} \sum_{i=1}^m \sum_{j=1}^m r_{ij} \|\mathbf{p}_i - \mathbf{p}_j\|_2^2 \quad (8)$$

where $\mathbf{R}=[r_{i,j}]_{m \times m}$ is an index matrix; $r_{i,j}=1$ if $\mathbf{Y}_i^T \mathbf{Y}_j=0$, otherwise $r_{i,j}=0$.

By integrating (6), (7), and (8), the final objective

function is shown as follows:

$$\begin{aligned} & \min_{\mathbf{P}, \mathbf{S}, \mathbf{W}} \sum_{i=1}^m \|f(\mathbf{x}_i, \mathbf{W}) - \mathbf{p}_i\|_2^2 + \delta \sum_{i=1}^m \sum_{j=1}^m (h_{i,j} - \mathbf{p}_i^T \mathbf{B} \mathbf{p}_j)^2 + \\ & \alpha \sum_{i=1}^m \left\| \mathbf{p}_i - \sum_{l_{i,j}=1}^m s_{i,j} \mathbf{p}_j \right\|_2^2 + \beta \sum_{i=1}^m \left\| \mathbf{x}_i - \sum_{l_{i,j}=1}^m s_{i,j} \mathbf{x}_j \right\|_2^2 - \\ & \gamma \sum_{i=1}^m \sum_{j=1}^m r_{i,j} \|\mathbf{p}_i - \mathbf{p}_j\|_2^2 + \lambda \|\mathbf{W}\|_{\mathbb{F}}^2, \\ & \text{s.t. } \begin{cases} \mathbf{P} \mathbf{1}_l = \mathbf{1}_m \\ \mathbf{0}_{m \times l} \leq \mathbf{P} \leq \mathbf{Y} \\ \mathbf{S} \mathbf{1}_m = \mathbf{1}_m \\ \mathbf{0}_{m \times m} \leq \mathbf{S} \leq \mathbf{L} \end{cases}, \end{aligned} \quad (9)$$

where λ is the regularization parameter, δ , α , β , and γ are the trade-off parameters.

3.3 Alternative optimization

Before the optimization procedure, the matrices \mathbf{P} and \mathbf{S} should be initialized by solving the standard quadratic programming (QP) problems [11] as follows:

$$\begin{aligned} & \min_{\mathbf{S}} \sum_{i=1}^m \left\| \mathbf{x}_i - \sum_{l_{i,j}=1}^m s_{i,j} \mathbf{x}_j \right\|_2^2, \\ & \text{s.t. } \begin{cases} \mathbf{S} \mathbf{1}_m = \mathbf{1}_m \\ \mathbf{0}_{m \times m} \leq \mathbf{S} \leq \mathbf{L} \end{cases}, \end{aligned} \quad (10)$$

$$\begin{aligned} & \min_{\mathbf{P}} \sum_{i=1}^m \left\| \mathbf{p}_i - \sum_{l_{i,j}=1}^m s_{i,j} \mathbf{p}_j \right\|_2^2, \\ & \text{s.t. } \begin{cases} \mathbf{P} \mathbf{1}_l = \mathbf{1}_m \\ \mathbf{0}_{m \times l} \leq \mathbf{P} \leq \mathbf{Y} \end{cases}. \end{aligned} \quad (11)$$

Then, we utilize the alternative optimization procedure to solve the problem (9).

Update \mathbf{S} , with fixed \mathbf{P} and \mathbf{W} , the problem (9) can be stated as follows:

$$\begin{aligned} & \min_{\mathbf{S}} \alpha \sum_{i=1}^m \left\| \mathbf{p}_i - \sum_{l_{i,j}=1}^m s_{i,j} \mathbf{p}_j \right\|_2^2 + \\ & \beta \sum_{i=1}^m \left\| \mathbf{x}_i - \sum_{l_{i,j}=1}^m s_{i,j} \mathbf{x}_j \right\|_2^2 \\ & \text{s.t. } \begin{cases} \mathbf{S} \mathbf{1}_m = \mathbf{1}_m \\ \mathbf{0}_{m \times m} \leq \mathbf{S} \leq \mathbf{L} \end{cases}. \end{aligned} \quad (12)$$

The optimization problem (12) can be rewritten [11] as follows:

$$\begin{aligned} & \min_{\tilde{\mathbf{s}}_i} \tilde{\mathbf{s}}_i^T (\alpha (\mathbf{D}_i^p)^T (\mathbf{D}_i^p) + \beta (\mathbf{D}_i^x)^T (\mathbf{D}_i^x)) \tilde{\mathbf{s}}_i, \\ & \text{s.t. } \begin{cases} \mathbf{1}_k^T \tilde{\mathbf{s}}_i = 1 \\ \mathbf{0}_k \leq \tilde{\mathbf{s}}_i \leq \mathbf{1}_k \end{cases}, \end{aligned} \quad (13)$$

where $\tilde{\mathbf{s}}_i \in \mathbf{R}^k$ denotes the weight vector which shows the importance of the neighbor sample in reconstructing \mathbf{x}_i .

Denote matrix $\mathbf{D}_i^p = [\mathbf{p}_i - \mathbf{p}_i^1, \mathbf{p}_i - \mathbf{p}_i^2, \dots, \mathbf{p}_i - \mathbf{p}_i^k]^T \in \mathbf{R}^{k \times l}$ and $\mathbf{D}_i^x = [\mathbf{x}_i - \mathbf{x}_i^1, \mathbf{x}_i - \mathbf{x}_i^2, \dots, \mathbf{x}_i - \mathbf{x}_i^k]^T \in \mathbf{R}^{k \times d}$, where \mathbf{p}_i^j is the labeling confidence vector associated with the j th nearest neighbors of \mathbf{x}_i , and \mathbf{x}_i^j is the j th nearest neighbors of \mathbf{x}_i . The optimization problem (13) is a standard QP problem which can be solved by existing QP tools. We adopt interior point methods to solve (13) by the quadprog function in Matlab.

Updating \mathbf{P} , with fixed \mathbf{S} and \mathbf{W} , the problem (9) can be stated as follows:

$$\begin{aligned} & \min_{\mathbf{P}} \sum_{i=1}^m \|f(\mathbf{x}_i, \mathbf{W}) - \mathbf{p}_i\|_2^2 + \delta \sum_{i=1}^m \sum_{j=1}^m (h_{i,j} - \mathbf{p}_i^T \mathbf{B} \mathbf{p}_j)^2 + \\ & \alpha \sum_{i=1}^m \left\| \mathbf{p}_i - \sum_{l_{i,j}=1}^m s_{i,j} \mathbf{p}_j \right\|_2^2 - \gamma \sum_{i=1}^m \sum_{j=1}^m r_{i,j} \|\mathbf{p}_i - \mathbf{p}_j\|_2^2, \\ & \text{s.t. } \begin{cases} \mathbf{P} \mathbf{1}_l = \mathbf{1}_m \\ \mathbf{0}_{m \times l} \leq \mathbf{P} \leq \mathbf{Y} \end{cases}. \end{aligned} \quad (14)$$

In the optimization problem (14), the first three terms are convex, and the last term is concave. To optimize the second term, the first order Taylor approximation at \mathbf{P}_0 can be utilized to replace it.

$$\begin{aligned} & \sum_{i=1}^m \sum_{j=1}^m (h_{i,j} - \mathbf{p}_i^T \mathbf{B} \mathbf{p}_j)^2 = \|\mathbf{H} - \mathbf{P} \mathbf{B} \mathbf{P}^T\|_{\mathbb{F}}^2 \approx \\ & \text{tr}((\mathbf{H} - \mathbf{P}_0^T \mathbf{B} \mathbf{P}_0)^T (\mathbf{H} - \mathbf{P}_0^T \mathbf{B} \mathbf{P}_0)) + \\ & \text{tr}((-2\mathbf{H} \mathbf{P}_0 \mathbf{B}^T - 2\mathbf{H}^T \mathbf{P}_0 \mathbf{B} + 2\mathbf{P}_0 \mathbf{B}^T \mathbf{P}_0^T \mathbf{P}_0 \mathbf{B} + \\ & 2\mathbf{P}_0 \mathbf{B} \mathbf{P}_0^T \mathbf{P}_0 \mathbf{B}^T)^T (\mathbf{P} - \mathbf{P}_0)) \end{aligned}$$

where \mathbf{P}_0 is the updated value at the previous iteration of \mathbf{P} . After Taylor approximation of the second term, the first three terms are still convex, and the last term is concave. The problem (14) is a constrained convex-concave problem. Fortunately, the concave-convex procedure (CCCP) can be used to solve this optimization problem [23]. A rigorous analysis of the convergence of CCCP is provided by [24]. The idea of CCCP is to linearize the concave part of the objective function. Therefore, the last term can be linearized by its first order Taylor approximation:

$$\begin{aligned} & - \sum_{i=1}^m \sum_{j=1}^m r_{i,j} \|\mathbf{p}_i - \mathbf{p}_j\|_2^2 = -\text{tr}(\mathbf{P}^T \mathbf{L}_r \mathbf{P}) \approx \\ & -\text{tr}(\mathbf{P}_0^T \mathbf{L}_r \mathbf{P}_0 + (\mathbf{L}_r \mathbf{P}_0 + \mathbf{L}_r^T \mathbf{P}_0)^T (\mathbf{P} - \mathbf{P}_0)) \end{aligned}$$

where $\mathbf{L}_r = \text{diag}(\mathbf{R}_d) - \mathbf{R}$ is a Laplacian matrix where vector \mathbf{R}_d is the sum of row for matrix \mathbf{R} .

We define the modeling output matrix $\mathbf{F} = [f(\mathbf{x}_1, \mathbf{W}), f(\mathbf{x}_2, \mathbf{W}), \dots, f(\mathbf{x}_m, \mathbf{W})]^T$. In order to simplify notation, we define $\mathbf{V} = \mathbf{L}_r \mathbf{P}_0 + \mathbf{L}_r^T \mathbf{P}_0$ and $\mathbf{U} = -2\mathbf{H} \mathbf{P}_0 \mathbf{B}^T - 2\mathbf{H}^T \mathbf{P}_0 \mathbf{B} + 2\mathbf{P}_0 \mathbf{B}^T \mathbf{P}_0^T \mathbf{P}_0 \mathbf{B} + 2\mathbf{P}_0 \mathbf{B} \mathbf{P}_0^T \mathbf{P}_0 \mathbf{B}^T$. Then the optimization prob-

lem (14) can be rewritten as follows:

$$\min_{\tilde{\mathbf{p}}} \frac{1}{2} \tilde{\mathbf{p}}^T (\mathbf{\Theta} + \mathbf{\Theta}^T) \tilde{\mathbf{p}} + (\delta \tilde{\mathbf{u}} - 2\tilde{\mathbf{f}} - \tilde{\mathbf{v}})^T \tilde{\mathbf{p}},$$

$$\text{s.t.} \begin{cases} \mathbf{0} \leq \tilde{\mathbf{p}} \leq \tilde{\mathbf{y}} \\ \sum_{i=1, i \neq m}^m \tilde{p}_i = 1 \\ 0 \leq j \leq m-1 \end{cases}, \quad (15)$$

where $\tilde{\mathbf{p}}$ is the vectorization of matrix \mathbf{P} , i.e., $\tilde{\mathbf{p}} = \text{vec}(\mathbf{P}) \in \mathbf{R}^{ml}$, $\tilde{\mathbf{f}} = \text{vec}(\mathbf{F})$, $\tilde{\mathbf{u}} = \text{vec}(\mathbf{U})$, and $\tilde{\mathbf{v}} = \text{vec}(\mathbf{V})$. We define matrix $\mathbf{O} = \alpha(\mathbf{S}^T \mathbf{S} + \mathbf{E}_{m \times m} - \mathbf{S} - \mathbf{S}^T) + \mathbf{E}_{m \times m}$ and $\mathbf{\Theta} = \mathbf{E}_{l \times l} \otimes \mathbf{O} \in \mathbf{R}^{ml \times ml}$, \otimes is Kronecker product and \mathbf{E} is an identity matrix. This optimization problem (13) is a standard QP problem which can be solved by off-the-shelf QP tools. We adopt interior point methods to solve (13) by the quadprag function in Matlab.

Update $\overline{\mathbf{W}}$, with fixed \mathbf{P} and \mathbf{S} . The linear model $f(\mathbf{x}_i, \mathbf{W}) = \mathbf{W}^T \mathbf{x}_i$ is utilized to predict the label of training examples. The problem (9) can be stated as follows:

$$\min_{\mathbf{W}} \sum_{i=1}^m \|\mathbf{W}^T \mathbf{x}_i - \mathbf{p}_i\|_2^2 + \lambda \|\mathbf{W}\|_F^2. \quad (16)$$

This is an unconstrained optimization problem. The predictive performance can be improved by the kernel technology. Let $\Phi(\cdot): \mathbf{R}^d \rightarrow \mathbf{R}^h$ be a feature mapping which maps the feature space to a Hilbert space with h -dimensions. The linear model is rewritten as $g(\mathbf{x}_i) = \overline{\mathbf{W}}^T \Phi(\mathbf{x}_i)$. We convert the problem (16) into an equality constrained minimization as follows:

$$\min_{\mathbf{W}} \sum_{i=1}^m \|\xi_i\|_2^2 + \lambda \|\overline{\mathbf{W}}\|_F^2,$$

$$\text{s.t.} \begin{cases} \overline{\mathbf{W}}^T \Phi(\mathbf{x}_i) - \mathbf{p}_i = \xi_i \\ 1 \leq i \leq m \end{cases}. \quad (17)$$

The method of Lagrange multipliers can be used to solve this problem. The Lagrange function is stated as follows:

$$\mathcal{L}(\overline{\mathbf{W}}, \mathbf{I}, \mathbf{A}) = \text{tr}(\mathbf{I}^T \mathbf{I}) + \lambda \text{tr}(\overline{\mathbf{W}}^T \overline{\mathbf{W}}) - \text{tr}(\mathbf{A}^T (\Phi \overline{\mathbf{W}} - \mathbf{P} - \mathbf{I}))$$

where $\Phi = [\Phi(\mathbf{x}_1), \Phi(\mathbf{x}_2), \dots, \Phi(\mathbf{x}_m)]^T$, $\mathbf{I} = [\xi_1, \xi_2, \dots, \xi_m]^T$, and $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m]^T$ is the Lagrange multiplier matrix. The optimal conditions of (17) are

$$\nabla \mathcal{L}(\overline{\mathbf{W}}, \mathbf{I}, \mathbf{A}) = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \overline{\mathbf{W}}} \\ \frac{\partial \mathcal{L}}{\partial \mathbf{I}} \\ \frac{\partial \mathcal{L}}{\partial \mathbf{A}} \end{bmatrix} = \begin{bmatrix} 2\lambda \overline{\mathbf{W}} - \Phi^T \mathbf{A} \\ 2\mathbf{I} + \mathbf{A} \\ \Phi \overline{\mathbf{W}} - \mathbf{P} - \mathbf{I} \end{bmatrix} = \mathbf{0}.$$

The optimal solution $\overline{\mathbf{W}}^*$ can be obtained as follows:

$$\begin{cases} \overline{\mathbf{W}}^* = \frac{\Phi^T \mathbf{A}}{2\lambda} \\ \mathbf{A} = \left(\frac{1}{2\lambda} \mathbf{K} + \frac{1}{2} \mathbf{E} \right)^{-1} \mathbf{P} \end{cases} \quad (18)$$

where $\mathbf{K} = \Phi \Phi^T$ is the kernel matrix with its element $k_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$, and $\kappa(\cdot, \cdot)$ is a kernel function. For PL-LCSA, we use Gaussian kernel to calculate \mathbf{K} . The modeling output matrix \mathbf{F} can be calculated by $\mathbf{F} = \Phi \overline{\mathbf{W}}^* = \mathbf{K} \mathbf{A} / 2\lambda$. Unseen instance \mathbf{x} can be predicted [11] by

$$\mathbf{y}^* = \arg \max_i \sum_{j=1}^m a_{ji} \kappa(\mathbf{x}, \mathbf{x}_j) \quad (19)$$

where a_{ij} is the element of matrix \mathbf{A} , and \mathbf{y}^* is the predicted label for \mathbf{x} . The pseudo-code of PL-LCSA is summarized as Algorithm 1.

Algorithm 1 Pseudo-code of PL-LCSA

Input:

\mathcal{D} : training label set $\mathcal{D} = \{(\mathbf{x}_i, Z_i) | 1 \leq i \leq m\}$

\mathbf{x} : the unseen instance

Parameter:

k : the number of nearest neighbors

$\delta, \alpha, \beta, \gamma, \lambda$: the trade-off coefficients

T : the number of maximum iterations

Output:

\mathbf{y}^* : the predicted label of \mathbf{x}

Process:

1: Calculate label correlation matrix \mathbf{B} according to (3)

2: Calculate \mathbf{H} by (5)

3: Calculate kernel matrix \mathbf{K}

4: Initialize \mathbf{S} according to (10)

5: Initialize \mathbf{P} according to (11)

6: **Repeat**

7: Update matrix \mathbf{A} according to (18)

8: Update modeling output matrix $\mathbf{F} = \mathbf{K} \mathbf{A} / 2\lambda$

9: Update matrix \mathbf{S} according to (13)

10: Update matrix \mathbf{P} by solving problem (15)

11: **Until** convergence or the maximum number of iterations

12: **Return** \mathbf{y}^* according to (19)

4. Experiments

4.1 Experimental setup

In this subsection, two series of comparative experiments based on synthetic data sets and real-world data sets are conducted to evaluate the performance of PL-LCSA. Table 1 summarizes the characteristics of six multi-class

data sets. Following the widely-used controlling protocol [6,8,11,14,25–27], artificial data sets are generated from multi-class data sets by three controlling parameters p , ϵ , and r . Here, p is the proportion of example which is ambiguous. r is the number of the false positive labels in the candidate label set (i.e., $|Z_i|=r+1$). ϵ is the probability that a false positive label co-occurs with the ground-truth label. The choice of false positive label is important, because it is the key factor of the label relationship in the synthetic data sets. In this paper, the reason why we choose multi-label data sets as synthetic data sets is that the multi-label data can be used to generate label correlation information. Firstly, (3) is used to calculate label cor-

relation in multi-label data sets. Then the most relevant label corresponding to the ground-truth label is taken as the false positive label. Table 2 summarizes the characteristics of six real-world data sets, where Avg. CLs is the average number of candidate labels.

Table 1 Characteristics of the multi-class data sets

Data set	Instance	Feature	PCA	Class
Enron	196	623	–	7
Medical	752	1287	–	30
Scene	2230	294	–	6
Bibtex	2744	1894	800	89
Mediamill	2854	120	–	10
Tmc2007	3130	15418	1000	21

Table 2 Characteristics of real-world partial label data sets

Data set	Example	Feature	Class	Avg. CLs	Task domain
FG-NET	1002	262	78	7.48	Facial age estimation [28]
Lost	1122	108	16	2.23	Automatic face naming [6]
MSRCv2	1758	48	23	3.16	Object classification [1]
BirdSong	4998	38	13	2.18	Bird song classification [29]
Soccer Player	17472	279	171	2.09	Automatic face naming [3]
Yahoo! News	22991	163	219	1.91	Automatic face naming [30]

The data sets in Table 1 are derived from multi-label benchmark data set by retaining instances with only one relevant label. They can be collected from Mulan (<http://mulan.sourceforge.net/index.html>). Because the features of Bibtex and Tmc2007 are relatively sparse, we make dimensionality reduction by principal component analysis (PCA), and the feature dimensions after dimensionality reduction are shown in the fourth column of Table 1.

In this paper, five partial label learning algorithms are utilized for comparative studies. Each algorithm is configured with the following literature:

(i) PL-KNN [5]: a k -nearest neighbor approach which makes the prediction for unseen instances by averaging the labeling information of its k -nearest neighbor (suggested configuration: $k=10$).

(ii) IPAL [8]: an instance-based approach which disambiguates the candidate label set via an iterative label propagation procedure (suggested configuration: $k=10$ and $\alpha=0.95$).

(iii) PL-ECOC [21]: a disambiguation-free approach which learns the predictive model by the error-correcting output codes (ECOC) technique (suggested configuration: $\tau=0.1 \cdot |D|$, $L=\lceil 40 \cdot \log_2(l) \rceil$).

(iv) PL-AGGD [11]: an approach which disambiguates the candidate label sets by adaptive graph to overcome the noise and outliers in the feature space (suggested configuration: $k=10$, $\lambda=1$, $\mu=1$, $\gamma=0.05$, and $T=10$).

(v) SDIM [14]: an approach which aims to maximize the latent semantic differences of the two instances whose ground-truth labels are definitely different (suggested configuration: $\lambda=0.05$, $\beta=0.001$).

The parameters of PL-LCSA are set as $\delta=0.5$, $\alpha=0.5$, $\beta=0.05$, $\gamma=0.5$, $\lambda=1$, $k=10$, and $T=10$. Ten-fold cross-validation is executed in each algorithm, and mean prediction accuracy and standard deviations will be recorded.

4.2 Results and discussion

4.2.1 Real-world data sets

Table 3 is the summary classification accuracy (mean \pm standard deviation (std)) of each comparing algorithm on the real-world data sets. • indicates the PL-LCSA is statistically superior/inferior to the comparing algorithms on each data set (pairwise t-test at 0.05 significance level). The performance of each algorithm is poor on the face and gesture recognition network (FG-NET) aging data set, because its Avg. CLs is extremely large.

As shown in Table 3, it can be seen that:

(i) On the Lost, MSRCv2 and Soccer Player data sets, the performance of PL-LCSA is superior to all the comparing algorithms;

(ii) On the BirdSong data set, the performance of PL-LCSA is comparable to the PL-ECOC and superior to the other comparing algorithms;

(iii) On the FG-NET data set, the performance of PL-

LCSA is comparable to the SDIM, PL-AGGD, PL-KNN and superior to the PL-ECOC and IPAL.

(iv) PL-LCSA is never outperformed by any comparing algorithms.

Table 3 Classification accuracy (mean±std) of each comparing algorithm on real-world partial label data sets

Data set	PL-LCSA	SDIM	PL-AGGD	PL-ECOC	IPAL	PL-KNN
Lost	0.789±0.030	0.797±0.030	0.744±0.020●	0.706±0.043●	0.645±0.034●	0.459±0.039●
MSRCv2	0.562±0.021	0.500±0.023●	0.509±0.028●	0.427±0.024●	0.531±0.037●	0.418±0.046●
BirdSong	0.756±0.008	0.734±0.012●	0.734±0.009●	0.751±0.013	0.712±0.015●	0.603±0.013●
Soccer Player	0.596±0.013	0.577±0.016●	0.539±0.016●	0.169±0.005●	0.544±0.014●	0.494±0.012●
Yahoo! News	0.670±0.008	0.663±0.013	0.647±0.009●	0.561±0.011●	0.607±0.012●	0.471±0.005●
FG-NET	0.079±0.020	0.076±0.037	0.076±0.027	0.005±0.007●	0.061±0.018●	0.066±0.018

From Table 3, PL-LCSA is capable of better performance compared with other comparison methods in the real-world data sets. Because three items are integrated in the PL-LCSA, we set three sets of comparative experiments in real-world data sets to determine the effect of these modules by parameters δ and γ . The first set of comparative experiments shows the performance of the adaptive graph. The second set shows the performance of the adaptive graph and the label correlation. The third set of comparative experiments shows the performance of PL-LCSA.

Table 4 shows the classification accuracy (mean±std) of these three sets of comparative experiments. A significantly increased performance was observed in the second set experiment compared with the first set experiment on MSRCv2, BirdSong, Soccer Player, and Yahoo! News. It is because the label correlation has contributed to the partial label learning. The performance of the third set of experiment achieves the better performance to the second, and the reason is that the label correlation and the

semantic difference maximization can jointly promote the performance.

Table 4 Classification accuracy (mean±std) of control variables for PL-LCSA on real-world data sets

Data set	$\delta = 0.5, \gamma = 0.5$	$\delta = 0.5, \gamma = 0$	$\delta = 0, \gamma = 0$
Lost	0.789±0.030	0.767±0.032	0.746±0.025
MSRCv2	0.562±0.021	0.544±0.030	0.507±0.028
BirdSong	0.756±0.008	0.751±0.010	0.733±0.011
Soccer Player	0.596±0.013	0.589±0.013	0.538±0.015
Yahoo! News	0.670±0.008	0.666±0.009	0.648±0.009
FG-NET	0.079±0.020	0.071±0.023	0.077±0.024

4.2.2 Synthetic data sets

Fig. 1 illustrates the classification accuracy of each comparing algorithm on the synthetic data sets as the co-occurring probability ϵ varies from 0.1 to 0.9 with step size 0.1, where $r=2$ and $p=1$.

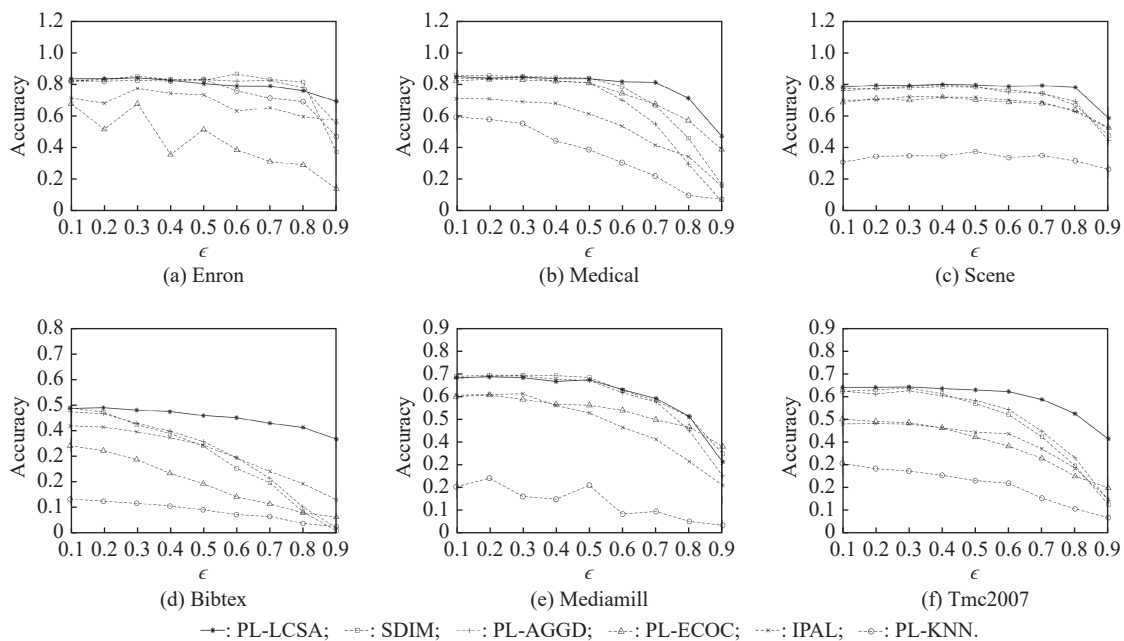


Fig. 1 Classification performance on synthetic data sets with ϵ ranging from 0.1 to 0.9 with step size 0.1 ($r=2, p=1$)

In general, PL-LCSA achieves competitive or better performance at synthetic data sets. Fig. 1 reveals that there has been a gradual decrease in the second half. The reason may be that the greater the similarity between the ground-truth label and the false positive label is, the more difficult it is to distinguish. Compared with the PL-AGGD, PL-LCSA achieves competitive performance at a smaller value of ϵ and better performance at a bigger value of ϵ on Enron, Scene, and Mediamill. The reason is that the size of label space on Enron, scene, and Mediamill is small, and the synthetic data does not have the information about the label correlation at a smaller value of ϵ and it is disadvantageous to PL-LCSA. On the contrary, when the value of ϵ is high, PL-LCSA achieves superior performance against PL-AGGD. It shows that the more the information about the label correlation is, the better the performance of PL-LCSA is.

4.2.3 Parameter sensitivity

Fig. 2 shows the accuracy of PL-LCSA under different configurations for parameters δ and γ on Lost and MSRCv2. As γ increases, PL-LCSA starts to take into consideration the semantic difference maximization criterion and the classification accuracy increases. For δ , as the weight of the label correlation increases, the classification accuracy decreases first, and then increases. In practice, we suggest users to choose δ and γ around 0.5 for all data sets.

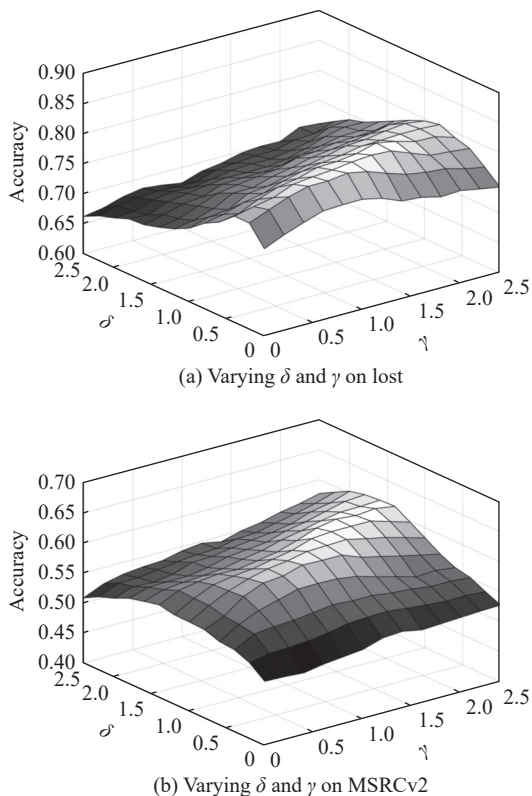


Fig. 2 Parameter sensitivity analysis for PL-LCSA

5. Conclusions

In this paper, we propose a unified framework that simultaneously focuses on the label and feature space. Meanwhile, this work generates fresh insight into the acquisition of the learning information from the label space, i.e., the label correlation. The framework integrates the label correlation, the adaptive graph, and the semantic difference maximization criterion. The relationship of instances can be learned by the adaptive graph in the feature space, the semantic difference analyzes the label relationship at the instance level, and the label correlation learns at the global label level. An effective optimization method is also proposed for this framework. Experiments on real-world and artificial data sets have demonstrated the superiority of PL-LCSA to the state-of-the-art partial label learning approaches.

References

- [1] LIU L P, DIETTERICH T G. A conditional multinomial mixture model for superset label learning. *Advances in Neural Information Processing Systems*, 2012, 1: 548–556.
- [2] ZHOU D, ZHANG Z, ZHANG M L, et al. Weakly supervised POS tagging without disambiguation. *ACM Transactions on Asian and Low-Resource Language Information Processing*, 2018, 17(4): 1–19.
- [3] ZENG Z N, XIAO S J, JIA K, et al. Learning by associating ambiguously labeled images. *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013: 708–715.
- [4] YAO Y, DENG J H, CHEN X H, et al. Deep discriminative CNN with temporal ensembling for ambiguously-labeled image classification. *Proc. of the AAAI Conference on Artificial Intelligence*, 2020, 34(7): 12669–12676.
- [5] HULLERMEIER E, BERINGER J. Learning from ambiguously labeled examples. *Intelligent Data Analysis*, 2006, 10(5): 419–439.
- [6] COUR T, SAPP B, TASKAR B. Learning from partial labels. *Journal of Machine Learning Research*, 2011, 12: 1501–1536.
- [7] NGUYEN N, CARUANA R. Classification with partial labels. *Proc. of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2008: 551–559.
- [8] ZHANG M L, YU F. Solving the partial label learning problem: an instance-based approach. *Proc. of the International Joint Conference on Artificial Intelligence*, 2015: 4048–4054.
- [9] ZHANG M L, ZHOU B B, LIU X Y. Partial label learning via feature-aware disambiguation. *Proc. of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016: 1335–1344.
- [10] GONG C, LIU T L, TANG Y Y, et al. A regularization approach for instance-based superset label learning. *IEEE Trans. on Cybernetics*, 2017, 48(3): 967–978.
- [11] WANG D B, LI L, ZHANG M L. Adaptive graph guided disambiguation for partial label learning. *Proc. of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2019: 83–91.
- [12] LYU G Y, FENG S H, WANG T, et al. GM-PLL: graph

- matching based partial label learning. *IEEE Trans. on Knowledge and Data Engineering*, 2021, 33(2): 521–535.
- [13] ZHU Y, KWOK J T, ZHOU Z H. Multi-label learning with global and local label correlation. *IEEE Trans. on Knowledge and Data Engineering*, 2017, 30(6): 1081–1094.
- [14] FENG L, AN B. Partial label learning by semantic difference maximization. *Proc. of the International Joint Conference on Artificial Intelligence*, 2019: 2294–2300.
- [15] JIN R, GHARAMANI Z. Learning with multiple labels. *Proc. of the Annual Conference on Neural Information Processing Systems*, 2002, 2: 897–904.
- [16] YU F, ZHANG M L. Maximum margin partial label learning. *Proc. of the Asian Conference on Machine Learning*, 2016: 96–111.
- [17] CHEN Y C, PATEL V M, CHELLAPPA R, et al. Ambiguously labeled learning using dictionaries. *IEEE Trans. on Information Forensics and Security*, 2014, 9(12): 2076–2088.
- [18] TANG C Z, ZHANG M L. Confidence-rated discriminative partial label learning. *Proc. of the AAAI Conference on Artificial Intelligence*, 2017, 31(1): 2611–1618.
- [19] LYU G Y, FENG S H, LI Y D, et al. Hera: partial label learning by combining heterogeneous loss with sparse and low-rank regularization. *ACM Trans. on Intelligent Systems and Technology*, 2020, 11(3): 1–19.
- [20] ZHANG Y B, YANG G, ZHAO S Y, et al. Partial label learning via generative adversarial nets. *Proc. of the European Conference on Artificial Intelligence*, 2020: 1674–1681.
- [21] ZHANG M L, YU F, TANG C Z. Disambiguation-free partial label learning. *IEEE Trans. on Knowledge and Data Engineering*, 2017, 29(10): 2155–2167.
- [22] WU X, ZHANG M L. Towards enabling binary decomposition for partial label learning. *Proc. of the International Joint Conference on Artificial Intelligence*, 2018: 2868–2874.
- [23] YUILLE A L, RANGARAJAN A. The concave-convex procedure (CCCP). *Advances in Neural Information Processing Systems*, 2002, 2: 1033–1040.
- [24] SRIPERUMBUDUR B K, LANCKRIET G R G. On the convergence of the concave-convex procedure. *Proc. of the Annual Conference on Neural Information Processing Systems*, 2009, 9: 1759–1767.
- [25] XU S, YANG M, ZHOU Y, et al. Partial label metric learning by collapsing classes. *International Journal of Machine Learning and Cybernetics*, 2020, 11: 2453–2460.
- [26] LYU J Q, XU M, FENG L, et al. Progressive identification of true labels for partial-label learning. *Proc. of the International Conference on Machine Learning*, 2020: 6500–6510.
- [27] YAN Y, GUO Y H. Partial label learning with batch label correction. *Proc. of the AAAI Conference on Artificial Intelligence*, 2020, 34(4): 6575–6582.
- [28] PANIS G, LANITIS A, TSAPATSOU LIS N, et al. Overview of research on facial ageing using the FG-NET ageing database. *IET Biometrics*, 2016, 5(2): 37–46.
- [29] BRIGGS F, FERN X Z, RAICH R. Rank-loss support instance machines for MIML instance annotation. *Proc. of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2012: 534–542.
- [30] GUILLAUMIN M, VERBEEK J, SCHMID C. Multiple instance metric learning from automatically labeled bags of faces. *European Conference on Computer Vision*, 2010: 634–647.

Biographies



GE Lingchi was born in 1997. He received his B.S. degree from Chang'an University, Xi'an, China, in 2018. He is currently pursuing his M.S. degree in computer science at Xidian University, Xi'an, China. His research interests are pattern recognition, partial label learning, and multi-label learning.
E-mail: 15109285306@163.com



FANG Min was born in 1965. She received her B.S. degree in computer control, M.S. degree in computer software engineering, and Ph.D. degree in computer application from Xidian University, Xi'an, China, in 1986, 1991, and 2004, respectively, where she is currently a professor. Her research interests include intelligent information process, multi-agent system, and network

technology.

E-mail: fanglabtg@163.com



LI Haikun was born in 1990. He received his M.S. degree from Yunnan University, Kunming, China, in 2017. He is currently a Ph.D. candidate in technology of computer application at Xidian University, Xi'an, China. His research interests include pattern recognition, machine learning, and partial-label learning.
E-mail: haikun1990@163.com



CHEN Bo was born in 1993. He received his B.S. degree in 2014 from Xi'an University of Post & Telecommunications, Xi'an, China. He is currently a Ph.D. candidate in computer science at Xidian University, Xi'an, China. His research interests include pattern recognition, machine learning, and time series.
E-mail: bchen_0314@stu.xidian.edu.cn