# GM-PLL: Graph Matching based Partial Label Learning

Gengyu Lyu, Songhe Feng, Tao Wang, Congyan Lang, Yidong Li

**Abstract**—Partial Label Learning (PLL) aims to learn from the data where each training example is associated with a set of candidate labels, among which only one is correct. The key to deal with such problem is to disambiguate the candidate label sets and obtain the correct assignments between instances and their candidate labels. In this paper, we interpret such assignments as instance-to-label matchings, and reformulate the task of PLL as a matching selection problem. To model such problem, we propose a novel *Graph Matching based Partial Label Learning* (GM-PLL) framework, where Graph Matching (GM) scheme is incorporated owing to its excellent capability of exploiting the instance and label relationship. Meanwhile, since conventional *one-to-one* GM algorithm does not satisfy the constraint of PLL problem that multiple instances may correspond to the same label, we extend a traditional *one-to-one* probabilistic matching algorithm to the *many-to-one* constraint, and make the proposed framework accommodate to the PLL problem. Moreover, we also propose a relaxed matching prediction model, which can improve the prediction accuracy via GM strategy. Extensive experiments on both artificial and real-world data sets demonstrate that the proposed method can achieve superior or comparable performance against the state-of-the-art methods.

**Index Terms**—Partial Label Learning, Matching Selection, Graph Matching, Many-to-one Constraint, Relaxed GM Predicted Model

✦

## 1 INTRODUCTION

As a weakly-supervised machine learning framework, partial label learning [1] learns from ambiguous labeling information where each training example corresponds to a candidate label set, among which only one is the ground-truth label [4] [5] [6]. During the training process, the correct label of each training example is concealed in its candidate label set and not directly accessible to the learning algorithm.

In many real-world scenarios, data with explicit labeling information (unique and correct label) is too scarce to obtain than that with implicit labeling information (redundant labels). Thus, when faced with such ambiguous data, conventional supervised learning framework based on *one instance one label* is out of its capability to learn from it accurately. Recently, Partial Label Learning (PLL) provides an effective solution to cope with it and has been widely used in many real-world scenarios. For example, in online annotation (Figure 1 (A)), users with varying knowledge and cultural backgrounds tend to annotate the same image with different labels. In order to learn from such ambiguous annotated collection, it is necessary to find the correspondence between each image and its ground-truth label. In naming faces (Figure 1 (B)), given a multi-figure image and its corresponding text description, the resulting set of images is ambiguously labeled if more than one name appear in the description. In other words, the specific correspondences between the faces and their names are unknown. In addition to the common scenarios mentioned above, PLL has also achieved competitive performance in many other applications, such as multimedia content analysis [7] [8] [9] [10], facial age estimation [11], web mining [12], ecoinformatics [13], etc.

The key to accomplish the task of learning from Partial-



Fig. 1. Examplar applications of partial-label learning.

Label (PL) data is disambiguation, which needs to fully explore the valuable information from ambiguous PL training data and obtain the correct assignments between the training INStances and their CandiDate Labels (INS-CDL). Recently, an Identification-based Disambiguation Strategy (IDS) is widely used in many PLL framework owing to its competitive performance on alleviating the interference of false positive labels [13] [14] [15] [16] [17] [18]. Among existing PLL methods based on IDS, some are often combined with the off-of-shelf learning schemes to identify the ground-truth label in an iterative manner, such as maximum like-lihood [13] [14] [15], maximum margin [16] [17] [19], etc. Others often try to explore the instance relationship from the ambiguous training data and directly disambiguate the candidate label sets [18]. Although the two kinds of PLL methods have obtained desirable performance in many real-world scenarios, they still suffer from some common defects. For example, for the instance relationship, they only consider the $k$-nearest-neighbor instances' similarity while simultaneously ignore the similarity among other instances and the dissimilarity among all instances, which makes the modeling output from unseen instance be overwhelmed by those from the negative nearest instances. And for the instance-label assignments, they usually utilize an iterative propagation procedure to implicitly obtain the objective labels, but neither explicitly describe the existing INS-CDL assignments relationship nor take the co-occurrence possibility of varying instance-label

- *School of Computer and Information Technology, Beijing Jiaotong University, No.3 Shangyuancun, Haidian District, Beijing 100044, China. E-mail: {lvgengyu, shfeng, twang, cylang, ydli}@bjtu.edu.cn*

1. In some literature, partial-label learning is also called as *superset label learning* [1], *ambiguous label learning* [2] or *soft label learning* [3].
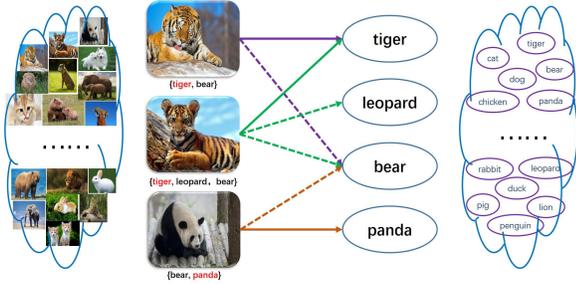
Fig. 2. Illustration of formulating PLL as a matching selection problem.

assignments into consideration to directly identify the optimal assignments, which may make the algorithm lose sight of direct instance-label assignments and result in its excessive attention to the instance relationship.

In order to overcome the above shortcomings, in this paper, we reinterpret the task of PLL as a matching selection problem, and simultaneously incorporate the instance relationship and the co-occurrence possibility of varying instance-label assignments into the same framework, then provide a novel solution for PLL problem. Specifically, we regard the INS-CDL correspondences as the instance-label matchings, and the task of PLL can be further reformulated as an instance-label matching selection problem (Figure 2), i.e. identifying the correct matching relationship between INSTances and their Ground-Truth Labels (INS-GTL). Afterwards, the goal of the PLL problem is transformed into how to solve the matching selection problem and obtain the optimal instance-label assignments. Recently, Graph Matching (GM) provides an effective solution for such problem, and owing to its excellent performance on utilizing structural information of training data, it has been widely used in many real-world applications [20] [21] [22] [23] [24]. Inspired by this, we incorporate the GM scheme into the PLL matching selection problem and propose a novel PLL learning framework named *Graph Matching based Partial Label Learning* (GM-PLL). Note that, existing graph matching algorithms are formulated with *one-to-one* constraint, which is not fully in accordance with the original task of PLL problem that one label can correspond to varying instances. Thus, we extend such *one-to-one* constraint to *many-to-one* constraint and propose a many-to-one probabilistic matching algorithm to make our method accommodate to the original PLL problem. Furthermore, during the establishment of the proposed framework, an affinity matrix is predetermined to describe the consistency relationship between varying INS-CDL assignments, where the similarity and dissimilarity of instances are simultaneously incorporated into the matrix. And these predetermined knowledge contributes the subsequent learning process and leads the algorithm to obtain the optimal solution. Moreover, to improve the predicted accuracy of test instances, we integrate the minimum error reconstruction scheme and graph matching scheme into a unified framework, and propose a relaxed GM predicted algorithm, where each unseen instance is first assigned with a candidate label set via minimum error reconstruction from its neighbor instances and then the predicted label is selected from $r$-maximum confidence candidate labels via graph matching strategy. Experimental results demonstrate that it can obtain higher classification accuracy than other predicted algorithms.

In summary, our main contributions lie in the following three aspects:

- Firstly, we reinterpret the conventional PLL problem and formulate the task of PLL as a matching selection problem.

To the best of our knowledge, it is the first time to regard PLL problem as a matching selection problem, and accordingly we propose a novel GM-based PLL framework (GM-PLL), where instance relationship and the co-occurrence possibility of varying instance-label assignments are simultaneously taken into consideration.

- Secondly, we extend conventional graph-matching algorithm with *one-to-one* constraint to a probabilistic matching algorithm with *many-to-one* constraint, which can guarantee that the proposed method fit the original task of PLL.

- Finally, we propose a relaxed GM prediction algorithm, which simultaneously incorporate the graph matching scheme and minimum error reconstruction scheme into the same framework to improve the classification accuracy.

We start the rest of the paper by giving a brief introduction about PLL, and then present technical details of the proposed GM-PLL algorithm and the comparative experiments with existing state-of-the-art methods. Finally, we conduct experimental analysis and conclude the whole paper.

## 2 RELATED WORK

Partial label learning, as a weakly supervised learning framework, focuses on solving the problem where data labeling information is excessively redundant. An intuitive strategy to cope with this issue is disambiguation, and existing disambiguation-based strategy are roughly grouped into three categories: *Averaging Disambiguation Strategy* (ADS), *Identification Disambiguation Strategy* (IDS) and *Disambiguation-Free Strategy* (DFS).

### 2.1 Averaging Disambiguation Strategy (ADS)

ADS-based methods usually assume that each candidate label has equal contribution to the learning model and they make prediction for unseen instances by averaging the outputs from all candidate labels. Following such strategy, Hullermeier et al. and Chen et al. adopt an instance-based model and disambiguate the ground-truth label by averaging the outputs of $k$-nearest neighbors following $\arg\max_{y \in \mathcal{Y}} \sum_{i \in \mathcal{N}_{(x^*)}} \mathbb{I}(y \in S_i)$ [25] [26]. Yu et al. utilize minimum error reconstruction criterion and obtain the predicted label via maximizing the confidence of $k$-nearest neighbors weighted-voting result [18]. Similarly, Tang et al. incorporate the boosting learning technique into its framework and improve the disambiguation classifier by adapting the weights of training examples and the ground-truth confidence of candidate labels [27]. Moreover, to further improve the disambiguation effectiveness, Zhang et al. facilitate its training process by taking the local topological information from feature space into consideration [11]. Obviously, the above PLL methods are clear and easy to implement, but they share a critical shortcoming that the output of the ground-truth label is overwhelmed by the outputs of the other false positive labels, which will enforce negative influence on the disambiguation of ground-truth label.

### 2.2 Identification Disambiguation Strategy (IDS)

In order to overcome the shortcomings of ADS, the IDS based PLL methods are proposed to directly disambiguate the candidate label set. This strategy aims to build a direct mapping from instance space to label space, and accurately identify the ground-truth label for each training instance. Existing PLL algorithms following this strategy often view the ground-truth label as a latent variable first,

identified as $\arg\max_{y \in S_i} F(\mathbf{x}, \boldsymbol{\Theta}, y)$, and then refine the model parameter $\boldsymbol{\Theta}$ iteratively by utilizing Expectation-Maximization (EM) procedure [14]. Among these methods, some usually incorporate the maximum likelihood criterion and obtain the optimal label via maximizing the outputs of candidate labels, following $\sum_{i=1}^{n} \log(\sum_{y \in S_i} F(\mathbf{x}, \boldsymbol{\Theta}, y))$ [2] [13] [14] [28] [29] [30]. Others often utilize the maximum margin criterion and identify the ground-truth label according to maximizing the margin between the outputs of candidate labels and that of the non-candidate labels, following $\sum_{i=1}^{n}(\max_{y \in S_i} F(\mathbf{x}, \boldsymbol{\Theta}, y) - \max_{y \notin S_i} F(\mathbf{x}, \boldsymbol{\Theta}, y))$ [16] [17]. Experimental results demonstrate that IDS-based method has achieved superior and comparable performance than ADS-based methods.

## 2.3 Disambiguation-Free Strategy (DFS)

Recently, different from the two disambiguation-based PLL strategies mentioned above, some attempts have been made to learn from PL data by fitting the PL data to off-the-shelf learning techniques, where they can directly make prediction for the unseen instances without conduct the disambiguation on the candidate label set corresponding to the training instances. Following such strategy, Zhang et al. propose a disambiguation-free algorithm named PL-ECOC [31], which utilizes *Error-Correcting Output Codes* (ECOC) coding matrix [32] and transfers the PLL problem into binary learning problem. Wu et al. propose another disambiguation-free algorithm called *PALOC* [33], which enables binary decomposition for PLL data in a more concise manner without relying on extra manipulations such as coding matrix. Experimental results empirically demonstrate that FDS-based algorithms can achieve comparable performance with the other disambiguation based PLL methods.

Although the above methods have achieved good performance on solving the PLL problem, they still suffer from some common shortcomings, i.e. they neither consider non $k$-nearest neighbor instance-similarity nor take the instance-dissimilarity into consideration. Therefore, in this paper, we utilize the GM scheme and propose a novel partial label learning framework called GM-PLL, where the instance similarity and dissimilarity are simultaneously incorporated into the framework to improve the performance of disambiguation. The details of the framework is introduced in the following section.

## 3 THE GM-PLL METHOD

Formally speaking, we denote the $d$-dimensional input space as $\mathcal{X} = \mathbb{R}^d$, and the output space as $\mathcal{Y} = \{1, 2, \ldots, q\}$ with $q$ class labels. PLL aims to learn a classifier $f : \mathcal{X} \mapsto \mathcal{Y}$ from the PL training data $\mathcal{D} = \{(\mathbf{x}_i, S_i)\}(1 \le i \le m)$, where the instance $\mathbf{x}_i \in \mathcal{X}$ is described as a $d$-dimensional feature vector, the candidate label set $S_i = \{y_{i_1}, y_{i_2}, \ldots, y_{i_{|S_i|}}\} \subseteq \mathcal{Y}$ is associated with the instance $\mathbf{x}_i$ and $|S_i|$ represents the number of candidate labels for instance $\mathbf{x}_i$. Meanwhile, we denote $\mathbf{y} = \{y_1, y_2, \ldots, y_m\}$ as the ground-truth label assignments for training instances, where each $y_i \in S_i$ corresponding to $\mathbf{x}_i$ is not directly accessible to the algorithm.

### 3.1 Formulation

GM-PLL is a novel PLL framework based on GM scheme, which aims to explore valuable information from ambiguous PL data and establish an accurate assignment relationship between the instance

space $\mathcal{X}$ and the label space $\mathcal{Y}$. To make the proposed method easily understanding, we illustrate the GM-PLL method as a GM structure (Figure 3) before the following detailed introduction.
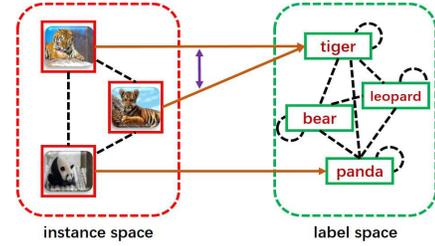


Fig. 3. The GM structure of GM-PLL. The GM structure originates from Figure 2.

As depicted in Figure 3, both the instance space and label space are formulated as two different undirected graphs $\mathbb{G}^i = (\mathbb{V}^i, \mathbb{E}^i)$ of size $n_i$, where $i \in \{1, 2\}$, and $n_1 = m, n_2 = q$. The nodes $\mathbb{V}^i$ in the two graphs represent the instances and labels respectively, while the edges $\mathbb{E}^i$ encode their similarities. The goal of GM-PLL is to establish the graph nodes correspondence between $\mathbb{G}^1$ and $\mathbb{G}^2$.

Here, we first denote $\mathbf{A}^i$ as the adjacent matrix for each graph $\mathbb{G}^i$, where $i = \{1, 2\}$. $\mathbf{A}^1 \in \mathbb{R}^{m \times m}$ encodes the instance-similarity, which is calculated by normalizing the popular Cosine Metric,

$$\mathbf{A}_{ij}^1 = \frac{\mathbf{x}_i^\top \cdot \mathbf{x}_j}{||\mathbf{x}_i||_2 \cdot ||\mathbf{x}_j||_2} \tag{1}$$

and $\mathbf{A}^2 \in \mathbb{R}^{q \times q}$ encodes the label-similarity,

$$\mathbf{A}_{i'j'}^2 = \begin{cases} 1, & \text{where } i' = j' \\ 0, & \text{where } i' \ne j' \end{cases} \tag{2}$$

where the similarity of different labels is set to 0 owing to the inherent characteristics of PLL problem that the prior pairwise-label relationship is always missing. Note that once the label relationship as prior knowledge can be obtained, the proposed GM-PLL can still be easily extended to satisfy the problem.

Then, we define $\mathbf{P} \in \{0, 1\}^{m \times q}$ to describe the graph node correspondences between $\mathbb{G}^1$ and $\mathbb{G}^2$, where $\mathbf{P}_{ij} = 1$ represents label $j$ is assigned to instance $\mathbf{x}_i$, and $\mathbf{P}_{ij} = 0$ otherwise. Among these correspondences that $\mathbf{P}_{ij} = 0$, a large number of them are invaluable to be considered since label $j$ is not contained in the candidate label set of instance $\mathbf{x}_i$. Accordingly, we exclude the assignments between instances and their non-candidate labels, and obtain the row-wise vectorized replica $\mathbf{p} = [p_1, p_2, \ldots, p_u]^\top \in \mathbb{R}^{u \times 1}$, where each element of $\mathbf{p}$ is defined as:

$$p_k = <\mathbf{x}_{i_k}, y_{l_k}> \tag{3}$$

here $i_k \in \{1, 2, \ldots, m\}$, $l_k \in \{1, 2, \ldots, |S_i|\}$, $u = \sum_{i=1}^{m} |S_i|$, $k \in \{1, 2, \ldots, u\}$ and the value of $<\mathbf{x}_{i_k}, y_{l_k}>$ represents the confidence of instance $\mathbf{x}_{i_k}$ assigned with its $l_k$-th candidate label.

Afterwards, the correspondence of INS-CDL can be obtained by solving the optimization problem **OP (1)**

$$\mathbf{P}^* = \arg\max_{\mathbf{P}} \sum_{i_a, l_a, i_b, l_b} d_{i_a, l_a, i_b, l_b} \mathbf{P}_{i_a, l_a} \mathbf{P}_{i_b, l_b}$$

$$s.t. \quad \mathbf{P}\underline{\mathbf{1}} = \underline{\mathbf{1}}.$$

where $d_{i_a,l_a,i_b,l_b}$ measures the pairwise consistency between instance edge $(i_a, i_b)$ and label edge $(l_a, l_b)$, which can also be regarded as the pairwise consistency between assignment $<\mathbf{x}_{i_a}, y_{l_a}>$ and assignment $<\mathbf{x}_{i_b}, y_{l_b}>$. Motivated by recent studies [24] [34] [35], we further formulate the **OP (1)** in a more general pairwise compatibility form **OP (2)**:

$$\mathbf{p}^* = \arg\max_{\mathbf{p}} \mathbf{p}^\top \mathbf{K}\mathbf{p}$$

$$s.t. \quad \mathbf{p} \in \{0,1\}^{u \times 1}$$

$$\mathbf{P}\mathbf{1} = \mathbf{1}.$$

where $\mathbf{K} \in \mathbb{R}^{u \times u}$ is the affinity matrix that will be introduced in the following subsection **Generation of Affinity Matrix K**. And the optimization details of **OP (2)** will also be exhibited in the following Section 3.2.

### 3.1.1 Generation of Affinity Matrix **K**

Affinity Matrix $\mathbf{K} \in \mathbb{R}^{u \times u}$ is defined to describe the matching consistency, and each element $\mathbf{K}_{ab}$ represents the INS-CDL correspondence between $\mathbf{p}_a$ and $\mathbf{p}_b$, i.e.

$$\mathbf{K}_{ab} = <\mathbf{p}_a, \mathbf{p}_b>$$
$$= <<\mathbf{x}_{i_a}, y_{l_a}>, <\mathbf{x}_{i_b}, y_{l_b}>> \quad (4)$$

here $a, b \in \{1, 2, \dots, u\}$, $<\mathbf{x}_{i_a}, y_{l_a}>$ represents the value of $s$-th element of $\mathbf{p}$ as the INS-CDL correspondence between the $i_a$-th instance $\mathbf{x}_{i_a}$ and its $l_a$-th candidate label $y_{l_a}$.

By predetermining the prior knowledge into the learning framework, affinity matrix can imply valuable information exploited from PL training data, including both the similarity and dissimilarity between instances, and the INS-CDL mapping relationship as well. Thus, we initialize the affinity matrix $\mathbf{K}$ as follows

$$\mathbf{K}_{ab} = \begin{cases} \mathbf{A}^1_{ij}, & \mathbf{A}^2_{i'j'} = 1 \\ 1 - \mathbf{A}^1_{ij}, & \mathbf{A}^2_{i'j'} = 0. \end{cases} \quad (5)$$

It is worth noting that, compared with the conventional PLL methods based on $k$-nearest neighbor scheme, the proposed framework contributes more prior knowledge to the learning process:

A) It utilizes the similarity information from more training instances instead of only from the $k$-nearest neighbors.

B) It not only utilizes the instance similarity but also takes the dissimilarity between instances into consideration. Particularly, as shown in Eq (5), with a higher similarity degree between two instances ($\mathbf{x}_{i_a}$ and $\mathbf{x}_{i_b}$), the $\mathbf{K}_{ab}$ will get a higher value, i.e., the ground-truth labels ($y_a$ and $y_b$) of the two instances have higher probability to locate in the intersection of their candidate labels. On the contrary, if with a lower similarity degree between $\mathbf{x}_{i_a}$ and $\mathbf{x}_{i_b}$, $y_a$ and $y_b$ will have higher probability to belong to non-intersection of their candidate labels.

After initializing the affinity matrix $\mathbf{K}$, we take the issue of class imbalance with respect to training data into consideration, and incorporate the number of instance candidate labels as a bias into the generation of affinity matrix:

$$\mathbf{K}_{ab} = \mathbf{K}_{ab} \cdot [1 + \alpha \cdot \log_2(\sum_{a=1}^{u} h(\mathbf{K}_{ab} > 0) + \sum_{b=1}^{u} h(\mathbf{K}_{ab} > 0))], \quad (6)$$

here $\alpha$ is the weight parameter, $h(\cdot)$ is the indicator function such that $h(\cdot) = 1$ iff $(\cdot)$ is true, and $h(\cdot) = 0$ otherwise. To reduce

---

**Algorithm 1** The Training Algorithm of **GM-PLL**

**Inputs:**
  $\mathcal{D}$: the partial label training set $\{(\mathbf{x}_i, S_i)\}$;
**Process:**
**1.** Calculate the cosine distances between each instance and derive the instance similarity matrix $\mathbf{A}$ by Eq (1);
**2.** Calculate the affinity matrix $\mathbf{K}$ by Eq (5) and Eq (6);
**3.** Standardize the affinity matrix $\mathbf{K}$ and remove low-confidence assignment by $\mathbf{K}(\mathbf{K} < \beta) = \mathbf{0}$;
**4.** Set $\mathbf{K}^{(0)} = \mathbf{K}$ and $\mathbf{p}^{(0)} = \frac{1}{|S_i|}\mathbf{1}$ where $\mathbf{p}^{(0)} \in \mathbb{R}^{u \times 1}$;
**5. for** $t = 0$ **to** $iter$
**6.**   $\mathbf{q}^{(t)} = \mathbf{K}^{(t)}\mathbf{p}^{(t)}$;
**7.**   $\mathbf{p}^{(t+1)} = \text{Normalize}(\mathbf{q}^{(t)})$;
**8.**   $\mathbf{K}^{(t+1)}(a,b) = \mathbf{K}^{(t)}(a,b) \cdot (\mathbf{p}_a^{(t+1)}/\mathbf{p}_a^{(t)})$;
**9.**   **if** $(||\mathbf{p}^{(t+1)} - \mathbf{p}^{(t)}||_2) < \delta$;
**10.**     break;
**11.**   **end if**
**12. end for**
**13.** Discretize $\mathbf{p}^{(t+1)}$, and derive the assignment $(\mathbf{x}_i, y_i)$;
**Output:**
  $y_i$: the assigned label for $\mathbf{x}_i$;

---

noise and alleviate the computational complexity, we increase the sparsity of the affinity matrix $\mathbf{K}$ and set $\mathbf{K}_{ab} = 0$ if $\mathbf{K}_{ab} < \beta$, where $\beta$ is the threshold parameter and it will be analyzed in Section 5.1.

At this point, the prior knowledge has been encoded into the affinity matrix, and it can provide good guidance for the subsequence learning process.

### 3.2 Optimization

In this section, we extend the probabilistic graph matching scheme from [36] and derive a probabilistic graph matching partial label learning algorithm. The core of the proposed algorithm is based on the observation that we can use the solution of the spectral matching algorithm [37] to refine the estimate of the affinity matrix $\mathbf{K}$ and then solve a new assignment problem based on the refined matrix $\mathbf{K}$. Namely, we can attenuate the affinities corresponding to matches with small matching probabilities and thus prune the affinity matrix $\mathbf{K}$. In the same vein, we aim to adaptively increase the entries in $\mathbf{K}$ corresponding to assignments with high matching probabilities.

Concretely, we relax the first constraint of **OP (2)** to $\mathbf{p} \in [0,1]^{u \times 1}$ and interpret $\mathbf{p}$ as matching probabilities $P(<\mathbf{x}_i, y_l>)$. Then, the affinity matrix $\mathbf{K}$ can be further interpreted as a joint matching probabilities $P(<\mathbf{x}_{i_a}, y_{l_a}>, <\mathbf{x}_{i_b}, y_{l_b}>)$. Afterwards, we refine $\mathbf{K}$ and $\mathbf{p}$ in an iterative manner where each iteration can be partitioned into two steps: estimating the mapping confidence of $\mathbf{p}$ and refining the affinity matrix $\mathbf{K}$. In the former step, we relax the *one-to-one* constraints of [37] as a *many-to-one* constrain to accommodate that multiple instances may correspond to the same label. In the latter step, we follow [36] to make the refinement of $\mathbf{K}$ allow analytic interpretation and provable convergence.

Hence, we minimize the objective function **OP (3)**

$$[\mathbf{p}_a^*, (\mathbf{p}_a|\mathbf{p}_b)^*] = \arg\min_{a,b} \sum_a ((\sum_b (\mathbf{p}_a|\mathbf{p}_b) \cdot \mathbf{p}_b) - \mathbf{p}_a)^2$$

where $\mathbf{p}_a$ is the assignment probability $P(<\mathbf{x}_{i_a}, y_{l_a}>)$ and $(\mathbf{p}_a|\mathbf{p}_b)$ represents the conditional assignment probability $P(<$

$\mathbf{x}_{i_a}, y_{l_a} > | < \mathbf{x}_{i_b}, y_{l_b} >)$ that is the probability of assignment $< \mathbf{x}_{i_a}, y_{l_a} >$ when $< \mathbf{x}_{i_b}, y_{l_b} >$ is valid. In our scheme, the $\mathbf{p}_a$ and $(\mathbf{p}_a | \mathbf{p}_b)$ need to be updated simultaneously.

Specifically, in iteration $t$, we denote the estimation of $P^{(t)}(< \mathbf{x}_{i_a}, y_{l_a} > | < \mathbf{x}_{i_b}, y_{l_b} >)$ by $(\mathbf{p}_a^{(t)} | \mathbf{p}_b^{(t)})$ and $P^{(t)}(< \mathbf{x}_{i_a}, y_{l_a} >)$ by $\mathbf{p}_a^{(t)}$, respectively. Then, we update $\mathbf{p}_a^{(t)}$ by

$$\mathbf{p}_a^{(t+1)} = \sum_b (\mathbf{p}_a^{(t)}, \mathbf{p}_b^{(t)}) = \sum_b (\mathbf{p}_a^{(t)} | \mathbf{p}_b^{(t)}) \cdot \mathbf{p}_b^{(t)} \qquad (7)$$

where $(\mathbf{p}_a^{(t)}, \mathbf{p}_b^{(t)})$ represents the joint probability $P(< \mathbf{x}_{i_a}, y_{l_a} >, < \mathbf{x}_{i_b}, y_{l_b} >)$ which is the joint probability of assignment $< \mathbf{x}_{i_a}, y_{l_a} >$ and assignment $< \mathbf{x}_{i_b}, y_{l_b} >$.

Different from the *one-to-one* constraint of conventional GM problem, the framework of GM-PLL is formulated with *many-to-one* constraint. Thus, we induce the constraint $\sum_{l_a=1}^{|S_i|} P(< \mathbf{x}_{i_a}, y_{l_a} >) = 1$. And $\mathbf{p}_a^{(t+1)} = [\mathbf{p}_{a_1}^{(t+1)}, \mathbf{p}_{a_2}^{(t+1)}, \ldots, \mathbf{p}_{a_{S_i}}^{(t+1)}]$ can be normalized as:

$$\mathbf{p}_{a_i}^{(t+1)} = \frac{\mathbf{p}_{a_i}^{(t+1)}}{\sum_1^{|S_i|} \mathbf{p}_{a_i}^{(t+1)}} \qquad (8)$$

Next, we refine the conditional assignment probability by

$$(\mathbf{p}_a | \mathbf{p}_b)^{(t+1)} = (\mathbf{p}_a | \mathbf{p}_b)^{(t)} \cdot \frac{\mathbf{p}_a^{(t+1)}}{\mathbf{p}_a^{(t)}}. \qquad (9)$$

During the entire process of optimization, we first initialize the required variables, and then repeat the above steps until the algorithm converges. Finally, we get the assigned label for each training example. The whole training algorithm of GM-PLL is summarized in Algorithm 1.

### 3.3 Prediction

During the stage of label prediction for unseen instances, we propose a graph matching based PLL prediction algorithm, which simultaneously takes the similarity reconstruction scheme and the GM scheme into consideration. The details of the prediction algorithm is introduced as follows.

We first integrate both the training instances and test instances into a large instances set, and then calculate a new instance-similarity matrix following Eq (1). Afterwards, we assign the candidate label set for each test instance $\mathbf{x}^*$ according to the weighted-voting results of its $k$-nearest neighbor instances $\mathcal{N}(\cdot)$, where the weights $\mathbf{w} \in \mathbb{R}^{k \times 1}$ are calculated via minimum error reconstruction scheme **OP (4)**:

$$w_c^* = \min_{w_c} \left\| \mathbf{x}^* - \sum_{c=1}^k w_c \cdot \mathbf{x}_c \right\|^2$$

$$s.t. \quad w_c \geq 0, \quad \sum_{c=1}^k w_c = 1, \quad (\mathbf{x}_c \in \mathcal{N}(\mathbf{x}^*), 1 \leq c \leq k)$$

here, $w_c$ is an element of $\mathbf{w}$ and $c \in \{1, 2, \ldots, k\}$.

Based on the weighted-voting results, we obtain the confidence of each candidate label assigned to $\mathbf{x}^*$, and then we can rank these labels according to the confidence in a descending order. Afterwards, we select the $r$-maximum confidence labels to constitute the candidate label set for $\mathbf{x}^*$. Subsequently, the construction of candidate label set for each unseen instance has been completed.

Apparently, when the value of $r$ equals to the total number of candidate-label categories $q$, the predicted model will degenerate into disambiguation from all candidate labels, which is commonly

in existing methods. In contrast, if only one label is retained ($r = 1$), the ground-truth label will be assigned with the maximum probability label, which is the same as [18]. The larger the value of $r$ is, the higher probability that the ground-truth label can be contained in the candidate label set, but meanwhile it would draw massive false labels that can decrease the effectiveness of the model. On the contrary, the smaller the value of $r$ is, the less false labels would be contained in the candidate label set, which would also result in the fact that the ground-truth label may be removed from the candidate label set.

Based on the above analysis, we can conclude that the total number of class labels (CL*) and the average number of class labels (AVG-CL*) for each instance have significant influence on the selecting of the number of assigned candidate labels $r$. Concretely, on one hand, more class labels means more noise class labels, thus we tend to assign $r$ with a smaller value to avoid the negative effect of these noise labels when CL* is larger. On the other hand, the average number of class labels can represent the average number of positive labels, thus we tend to choose larger $r$ when AVG-CL* is larger. At this point, we can calculate the $r$ by the following formula:

$$r = \left[ 1 + \frac{\text{AVG-CL*}}{lg(\text{CL*})} \right] \qquad (10)$$

here $[\triangle]$ is the integral function, which represents the rounding operation for $\triangle$.

Finally, once the above operations are completed, we follow the idea of Algorithm 1 to rebuild the affinity matrix and utilize the GM scheme to recover the correct mapping between test instances and their ground-truth labels.

## 4 EXPERIMENTS

### 4.1 Experimental Setup

To verify the effectiveness of the proposed GM-PLL method, we conduct experiments on **nine** controlled UCI data sets and **six** real-world data sets:

**(1) Controlled UCI data sets**. Under specified configuration of two controlling parameters (i.e. $p$ and $r$), the nine UCI data sets generate 189 ($7 \times 3 \times 9$) artificial partial-label data sets [2] [38]. Here, $p \in \{0.1, 0.2, \ldots, 0.7\}$ is the proportion of instances with partial labeling and $r \in \{1, 2, 3\}$ is the number of candidate labels except the ground-truth label. Table 1 summarizes the characteristics of the nine UCI data sets, including the number of examples (**EXP***), the number of the features (**FEA***), the whole number of class labels (**CL***) and their common configurations (**CONFIGURATIONS**).

TABLE 1
Characteristics of the controlled data sets

| UCI data sets | EXP* | FEA* | CL* | CONFIGURATIONS |
|---|---|---|---|---|
| Glass | 214 | 10 | 7 | |
| Ecoli | 336 | 7 | 8 | |
| Dermatology | 364 | 23 | 6 | $r = 1, p \in \{0.1, 0.2, \ldots, 0.7\}$ |
| Vehicle | 846 | 18 | 4 | |
| Segment | 2310 | 18 | 7 | $r = 2, p \in \{0.1, 0.2, \ldots, 0.7\}$ |
| Abalone | 4177 | 7 | 29 | |
| Letter | 5000 | 16 | 26 | $r = 3, p \in \{0.1, 0.2, \ldots, 0.7\}$ |
| Satimage | 6345 | 36 | 7 | |
| Pendigits | 10992 | 16 | 10 | |

**(2) Real-World (RW) data sets** . These data sets are collected from the four following task domains: (A) ***Facial Age Estimation***
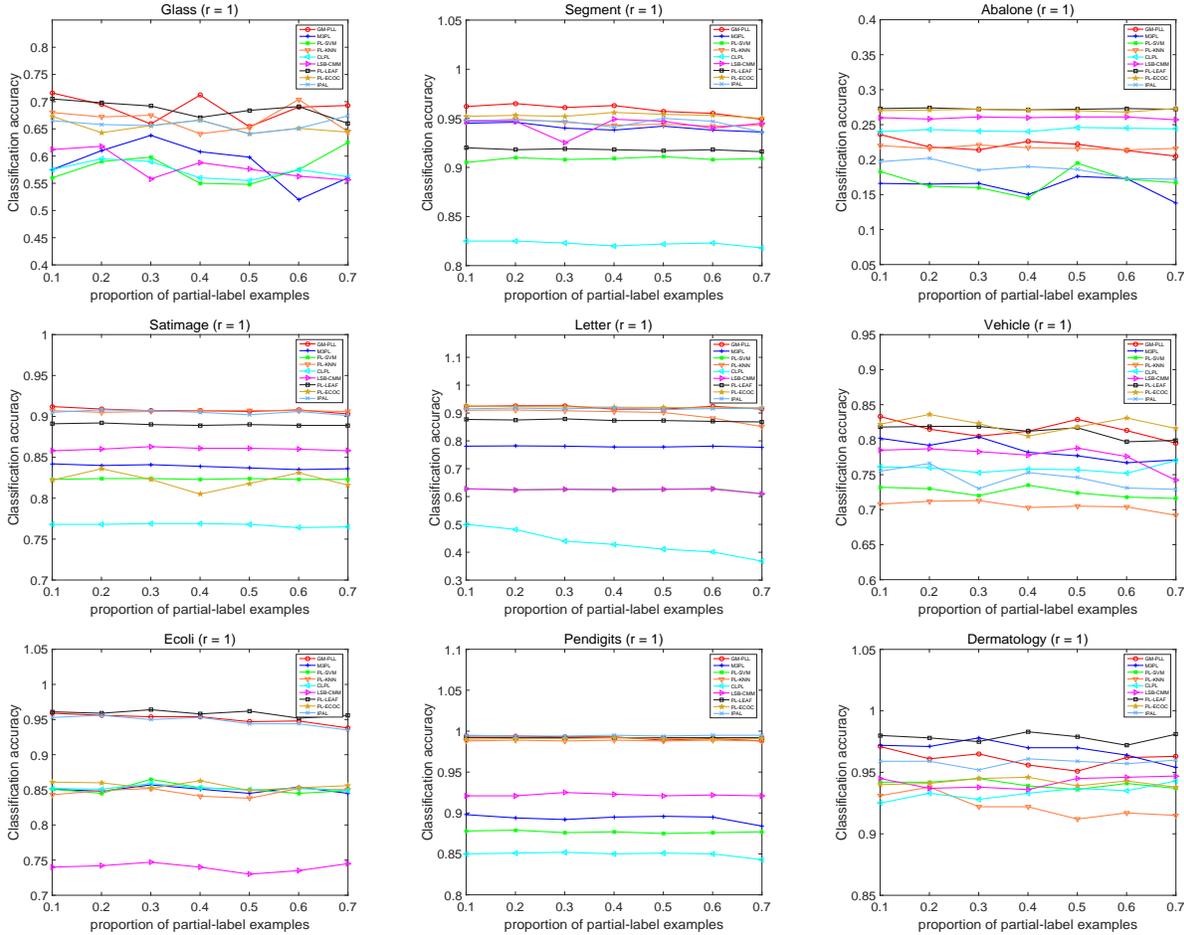
Fig. 4. The classification accuracy of each comparing method on nine controlled UCI data sets with one false positive candidate label (r = 1)

Human faces are represented as instances and the ages annotated by ten crowd-sourced labelers together with the ground-truth ages are regarded as candidate labels; (B) *Automatic Face Naming* Human faces copped from images or videos are represented as instances and each candidate label set is composed of the names extracted from the corresponding captions or subtitles; (C) *Object Classification* Image segmentations constitute the instance space and the objects appearing within the same image constitute the candidate label sets; (D) *Bird Song Classification* Singing syllables of the birds are represented as instances while bird species jointly singing during a 10-seconds period are regarded as candidate labels; Table 2 summarizes the characteristics of the above real world data sets, including not only the number of examples (**EXP\***), the number of the feature (**FEA\***) and the whole number of class labels (**CL\***), but also the average number of class labels (**AVG-CL\***) and their task domains (**TASK DOMAIN**).

Meanwhile, we employ **four** classical (PL-SVM, PL-KNN, CLPL, LSB-CMM) and **four** state-of-the-art (M3PL, PL-LEAF, PL-ECOC, IPAL) partial label learning algorithms that are based on different disambiguation strategies [2] for comparative studies, where the configured parameters of each method are utilized following the suggestions in respective literatures:

2. We partially use the open source codes from Zhang Minling's homepage: http://cse.seu.edu.cn/PersonalPage/zhangml/

- **PL-SVM** [16]: Based on IDS, it gets the predicted-label according to incorporating maximum margin scheme. [suggested configuration: $\lambda \in \{10^{-3}, 10^{-2}, \ldots, 10^3\}$] ;

- **PL-KNN** [25]: Based on ADS, it obtains the predicted-label according to averaging the outputs of the $k$-nearest neighbors. [suggested configuration: $k$=10];

- **CLPL** [38]: A convex optimization partial-label learning method based on ADS. [suggested configuration: SVM with hinge loss];

- **LSB-CMM** [13]: Based on IDS, it makes prediction according to calculating the maximum-likelihood value of the model with unseen instances input. [suggested configuration: $q$ mixture components];

- **M3PL** [17]: Originated from PL-SVM, it is also based on the maximum-margin strategy, and it gets the predicted-label via calculating the maximum values of model outputs. [suggested configuration: $C_{max} \in \{10^{-2}, 10^{-1}, \ldots, 10^2\}$] ;

- **PL-LEAF** [11]: A partial-label learning method via feature-aware disambiguation. [suggested configuration: $k$=10, $C_1 = 10$, $C_2 = 1$];

- **IPAL** [18]: it disambiguates the candidate label set by taking the instance similarity into consideration. [suggested configuration: $k$=10];

- **PL-ECOC** [31]: Based on a coding-decoding procedure, it learns from partial-label training examples in a disambiguation-free manner. [suggested configuration: the

TABLE 2
Characteristics of the real-world data sets

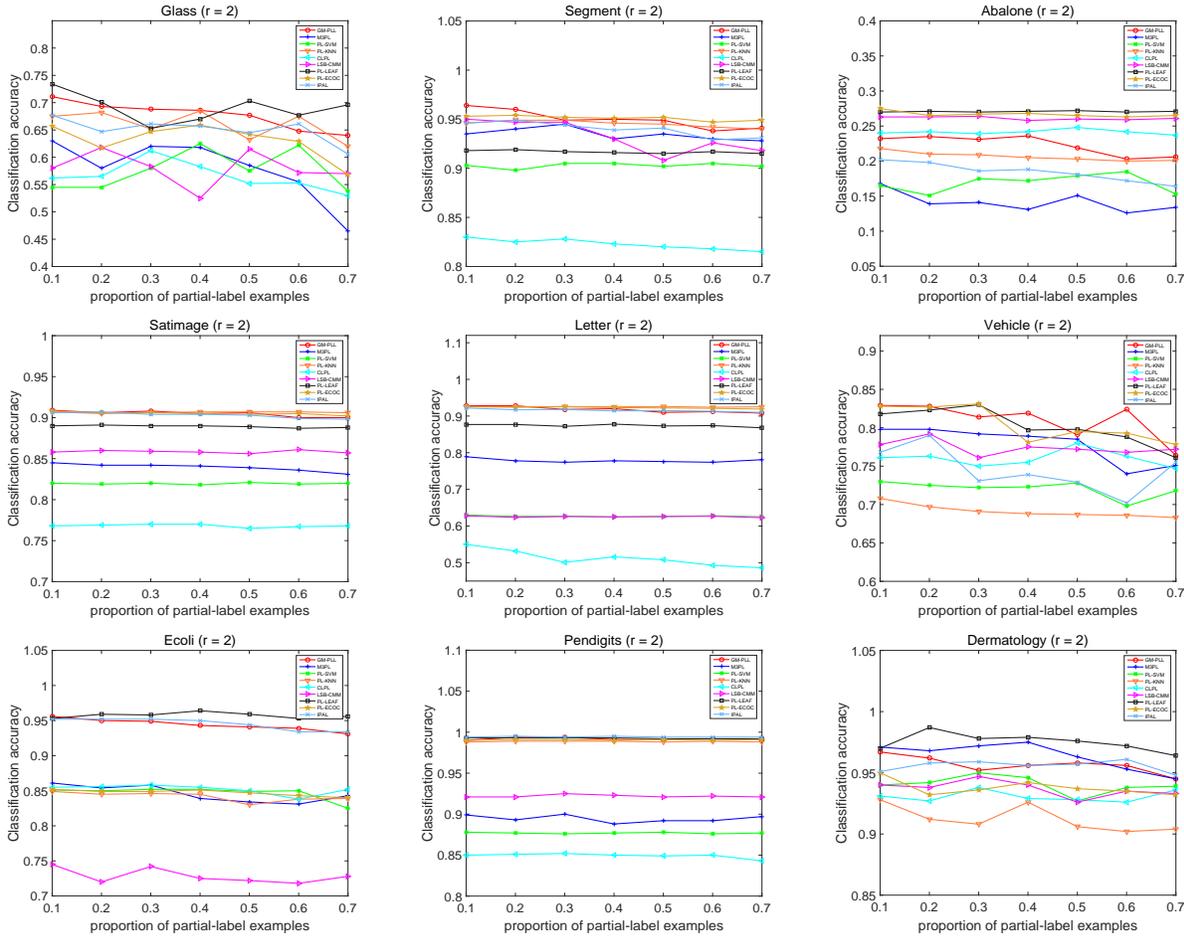| RW data sets | EXP* | FEA* | CL* | AVG-CL* | TASK DOMAIN |
|---|---|---|---|---|---|
| Lost | 1122 | 108 | 16 | 2.33 | *Automatic Face Naming* [38] |
| MSRCv2 | 1758 | 48 | 23 | 3.16 | *Image Classification* [39] |
| FG-NET | 1002 | 262 | 99 | 7.48 | *Facial Age Estimation* [40] |
| Soccer Player | 17472 | 279 | 171 | 2.09 | *Automatic Face Naming* [7] |
| Yahoo! News | 22991 | 163 | 219 | 1.91 | *Automatic Face Naming* [41] |



Fig. 5. The classification accuracy of each comparing method on nine controlled UCI data sets with two false positive candidate labels (r = 2)

TABLE 3
Win/tie/loss counts of the GM-PLL's classification performance against each comparing method on UCI data sets (pairwise $t$-test at 0.05 significance level)

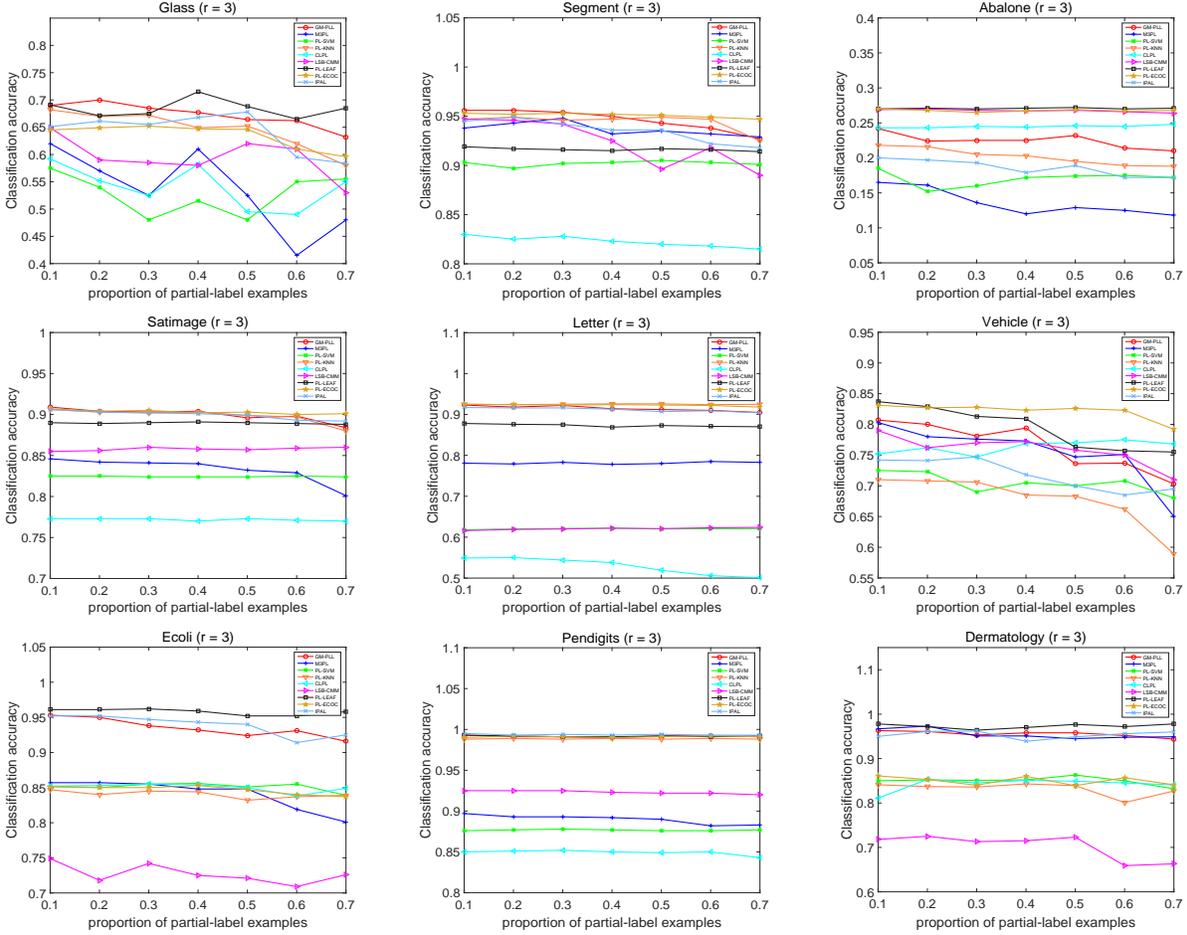| Data set | PL-KNN | PL-SVM | LSB-CMM | CLPL | M3PL | PL-LEAF | PL-ECOC | IPAL | sum |
|---|---|---|---|---|---|---|---|---|---|
| glass | 19/2/0 | 21/0/0 | 21/0/0 | 21/0/0 | 21/0/0 | 7/4/10 | 21/0/0 | 19/2/0 | 150/8/10 |
| segment | 21/0/0 | 21/0/0 | 21/0/0 | 21/0/0 | 21/0/0 | 21/0/0 | 16/5/0 | 21/0/0 | 163/5/0 |
| vehicle | 21/0/0 | 21/0/0 | 17/3/1 | 18/0/3 | 19/2/0 | 8/7/6 | 7/5/9 | 21/0/0 | 132/17/19 |
| letter | 14/7/0 | 21/0/0 | 21/0/0 | 21/0/0 | 21/0/0 | 21/0/0 | 5/16/0 | 15/6/0 | 139/29/0 |
| satimage | 19/2/0 | 21/0/0 | 21/0/0 | 21/0/0 | 21/0/0 | 20/1/0 | 15/6/0 | 19/2/0 | 157/11/0 |
| abalone | 21/0/0 | 21/0/0 | 0/0/21 | 0/10/11 | 21/0/0 | 0/0/21 | 0/0/21 | 21/0/0 | 84/10/74 |
| ecoli | 12/9/0 | 21/0/0 | 21/0/0 | 21/0/0 | 21/0/0 | 1/13/7 | 21/0/0 | 11/10/0 | 129/32/7 |
| dermatology | 14/7/0 | 21/0/0 | 21/0/0 | 21/0/0 | 6/14/1 | 0/14/7 | 21/0/0 | 13/8/0 | 117/43/8 |
| pendigits | 2/19/0 | 21/0/0 | 21/0/0 | 21/0/0 | 21/0/0 | 9/12/0 | 11/10/0 | 1/20/0 | 107/61/0 |
| sum | 163/22/4 | 178/6/5 | 142/0/42 | 148/14/27 | 153/15/21 | 79/44/66 | 110/37/42 | 125/55/9 | - |

Fig. 6. The classification accuracy of each comparing method on nine controlled UCI data sets with three false positive candidate labels (r = 3)

codeword length $L = \lceil log_2(q) \rceil$;

Before conducting the experiments, we give the range of the required variables. In detail, during the training phase, the threshold variable $\beta$ is set among $\{0.3, 0.4, \ldots, 0.8\}$ to exploit the most valuable similarity information and dissimilarity information. And the coefficient parameter $\alpha$ is chosen from $\{0, 0.1, 0.2\}$ to balance the effect of the number of varying label categories. During the test phase, inspired by [18], we empirically set $k = 10$ for $k$-nearest neighbor instances to complete the candidate label set of each unseen instance, and meanwhile the size of the label set $r$ is empirically set to more than 1 to guarantee that the ground-truth label can be involved in the assigned candidate label set. After initializing the above variables, we adopt ten-fold cross-validation to train the model and get the average classification accuracy on each data set.

## 4.2 Experimental Results

Since the origins of the two kinds of data sets are different, nine UCI data sets are constructed manually while six RW data sets come from real world scenarios, we conduct two series of experiments to evaluate the proposed method and the experimental results are exhibited in the following two subsections separately. In our paper, the experimental results of the comparing algorithms originate from two aspects: one is from the results we implemented by utilizing the source codes provided by the authors; the other is from the results exhibited in the respective literatures.

### 4.2.1 Controlled UCI data sets

Figure 4-6 illustrate the classification accuracy of each comparing method on the nine controlled data sets as $p$ increases from $0.1$ to $0.7$ with the step-size $0.1$. Together with the ground-truth label, the $r$ class labels are randomly chosen from $\mathcal{Y}$ to constitute the rest of each candidate label set, where $r = 1, 2, 3$. Table 3 summaries the win/tie/loss counts between GM-PLL and other comparing methods. Out of 189 (9 data sets $\times$ 21 configurations) statistical comparisons show that GM-PLL achieves either superior or comparable performance against the eight comparing methods, which is embodied in the following aspects:

- Among the comparing methods, GM-PLL achieves superior performance against PL-KNN, PL-SVM, LSB-CMM, CLPL and M3PL in most cases. And compared with PL-LEAF, PL-ECOC and IPAL, it also achieves superior or comparable performance in 65.08%, 77.78%, 95.23% cases, respectively. These results demonstrate that the proposed method has superior capacity of disambiguation against other methods based on varying disambiguation strategies, as well as disambiguation-free strategy.

- Compared with the methods that directly establish INS-GTL assignments, GM-PLL achieves superior performance on most data sets. For example, the average classification accuracy of GM-PLL is 11.2% higher than M3PL on *Glass* data set and 29.5% higher than PL-SVM on *Satimage* data

TABLE 4
Inductive accuracy (mean ± std) of each comparing algorithm on real-world data sets. ●/○ indicates that GM-PLL is statistically superior / inferior to the comparing algorithm on each data set (pairwise $t$-text at 0.05 significate level).

|  | Lost | MSRCv2 | Yahoo! News | BirdSong | SoccerPlayer | FG-NET |
|---|---|---|---|---|---|---|
| GM-PLL | **0.737±0.043** | **0.530±0.019** | 0.629±0.007 | 0.663±0.010 | **0.549±0.009** | 0.065±0.021 |
| PL-SVM | 0.639±0.056 ● | 0.417±0.027 ● | 0.636±0.018 ○ | 0.662±0.032 ● | 0.430±0.004 ● | 0.058±0.010 ● |
| CLPL | 0.670±0.024 ● | 0.375±0.020 ● | 0.462±0.009 ● | 0.632±0.017 ● | 0.347±0.004 ● | 0.047±0.017 ● |
| PL-KNN | 0.332±0.030 ● | 0.417±0.012 ● | 0.457±0.009 ● | 0.614±0.024 ● | 0.494±0.004 ● | 0.037±0.008 ● |
| LSB-CMM | 0.591±0.019 ● | 0.431±0.008 ● | 0.648±0.015 ○ | 0.717±0.024 ○ | 0.506±0.006 ● | 0.056±0.008 ● |
| M3PL | 0.732±0.035 ● | 0.521±0.030 ● | 0.655±0.010 ○ | 0.709±0.010 ○ | 0.446±0.013 ● | 0.037±0.025 ● |
| PL-LEAF | 0.664±0.020 ● | 0.459±0.013 ● | 0.597±0.012 ● | 0.706±0.012 ○ | 0.515±0.004 ● | **0.072±0.010** ○ |
| IPAL | 0.726±0.041 ● | 0.523±0.025 ● | **0.667±0.014** ○ | 0.708±0.014 ○ | 0.547±0.014 ● | 0.057±0.023 ● |
| PL-ECOC | 0.703±0.052 ● | 0.505±0.027 ● | 0.662±0.010 ○ | 0.740±0.016 ○ | 0.537±0.020 ● | 0.040±0.018 ● |

TABLE 5
Transductive accuracy (mean ± std) of each comparing algorithm on real-world data sets. ●/○ indicates that GM-PLL is statistically superior / inferior to the comparing algorithm on each data set (pairwise $t$-text at 0.05 significate level).

|  | Lost | MSRCv2 | Yahoo! News | BirdSong | SoccerPlayer | FG-NET |
|---|---|---|---|---|---|---|
| GM-PLL | 0.881±0.005 | **0.770±0.013** | 0.705±0.612 | 0.834±0.010 | 0.668±0.003 | **0.186±0.021** |
| PL-SVM | 0.887±0.012 ○ | 0.653±0.024 ● | 0.871±0.002 ○ | 0.825±0.012 ● | 0.688±0.014 ○ | 0.136±0.021 ● |
| CLPL | **0.894±0.005** ○ | 0.656±0.010 ● | 0.834±0.002 ○ | 0.822±0.004 ● | 0.680±0.010 ● | 0.158±0.018 ● |
| PL-KNN | 0.615±0.036 ● | 0.616±0.006 ● | 0.692±0.010 ● | 0.772±0.021 ● | 0.492±0.015 ● | 0.173±0.017 ● |
| LSB-CMM | 0.721±0.010 ● | 0.524±0.007 ● | **0.872±0.001** ○ | 0.716±0.014 ● | 0.704±0.002 ○ | 0.138±0.019 ● |
| M3PL | 0.860±0.006 ● | 0.732±0.025 ● | 0.870±0.002 ○ | 0.855±0.030 ○ | **0.761±0.010** ○ | 0.127±0.013 ● |
| PL-LEAF | 0.809±0.022 ● | 0.645±0.015 ● | 0.827±0.002 ○ | 0.882±0.014 ○ | 0.702±0.003 ○ | 0.148±0.009 ● |
| IPAL | 0.840±0.041 ● | 0.714±0.015 ● | 0.823±0.008 ○ | 0.833±0.030 ● | 0.673±0.014 ● | 0.158±0.024 ● |
| PL-ECOC | 0.851±0.013 ● | 0.555±0.030 ● | 0.862±0.007 ○ | **0.886±0.014** ○ | 0.671±0.003 ● | 0.132±0.019 ● |

set. Meanwhile, GM-PLL also has higher or comparable classification accuracy against the comparing state-of-the-art methods on other controlled UCI data sets. We attribute such success to that it can utilize the co-occurrence possibility of varying instance-label assignments to obtain the accurate INS-GTL assignments.

- Compared with the methods utilizing the instance similarity, GM-PLL also achieves competitive performance. From the perspective of the Average Classification Accuracy, GM-PLL gets 1.2% higher than IPAL on *Segment* data set and 1.4% higher than PL-LEAF on *Letter* data set, respectively; And from the perspective of the Max-Min of classification accuracy, GM-PLL is only 0.84% higher on *Glass* data set while all other methods are more than 1%. Moreover, the standard deviation of GM-PLL classification accuracy is lower than the other comparing methods on most data sets. These results clearly indicate the advantage of the proposed method against other instance-similarity based methods.

### 4.2.2 *Real-world (RW) data sets*

We compare the GM-PLL with all above comparing algorithms on the real-world data sets. The comparison results of inductive accuracy and transductive accuracy are separately reported in Table 4 and Table 5, where the recorded results are based on ten-fold cross-validation.

The transductive classification accuracy reflects the disambiguation capacity of PLL methods in recovering ground-truth labeling information from candidate label set, while the inductive classification accuracy reflects the prediction capacity of obtaining

the ground-truth label for unseen examples. According to Table 4 and Table 5, it is clear to observe that GM-PLL performs better than most comparing PLL algorithms on these RW data sets. The superiority of GM-PLL can be embodied in the following aspects:

- As shown in Table 4, GM-PLL significantly outperforms all comparing methods on *Lost*, *MSRCv2*, and *SoccerPlayer* data sets, respectively. Especially, compared with the classical methods, the classification accuracy of the proposed method is 40.5% higher than that of PL-KNN on *Lost* data set, and 20.2% higher than that of CLPL on *SoccerPlayer* data set. Even compared with the state-of-the-art methods, it also can achieve 2.5% higher than PL-ECOC on *MSRCv2* and 1.1% higher than IPAL on *Lost* data set.
- Meanwhile, GM-PLL also achieves competitive performance on other RW data sets. Specifically, for the *FG-NET* data set, GM-PLL outperforms all comparing methods except PL-LEAF, where it is only 0.7% lower than PL-LEAF. But on *Yahoo! News* data set, GM-PLL performs great superiority than PL-LEAF, where the classification accuracy is 3.4% higher than that of PL-LEAF. Besides, among all comparing methods, it is impressive that GM-PLL outperforms CLPL and PL-KNN on all six RW data sets. And, it also exceeds other comparing methods over four in six RW data sets. The experimental results demonstrate the superiority of GM-PLL.
- As shown in Table 5, GM-PLL shows significantly superior disambiguation ability on *Lost*, *MSRCv2* and *FG-NET* data set and competitive disambiguation ability on *BirdSong* and *SoccerPlayer* data sets, which demonstrates the superiority of the GM scheme on disambiguation. But for *Yahoo! News* data
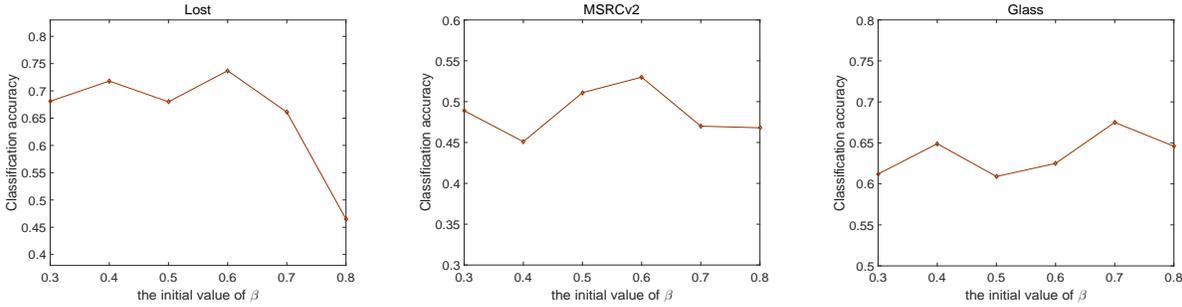
Fig. 7. The classification accuracy of the proposed methods on *Lost*, *MSRCv2* and *Glass* data sets with $r$ fixed ($r = 3$ on *Lost* data set, $r = 4$ on *MSRCv2* data set and $r = 4$ on *Glass* data set respectively)
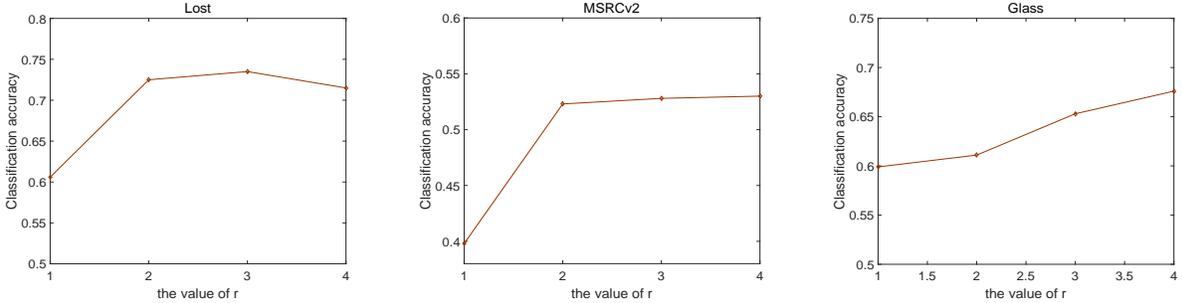


Fig. 8. The classification accuracy of SP-PLL on *Lost*, *MSRCv2* and *Glass* data sets with $\beta$ fixed ($\beta = 0.3$ on *Lost* data set, $\beta = 0.6$ on *MSRCv2* data set and $\beta = 0.7$ on *Glass* data set respectively)

set, GM-PLL is inferior to some comparing state-of-the-art methods. Even so, it can still achieve superior or comparable performance against other comparing methods on making prediction for unseen instances, which demonstrates the superiority of GM scheme on making prediction for unseen instances. In summary, the experimental results demonstrate the effectiveness of our proposed GM-PLL algorithm.

- We notice that the performance of GM-PLL is inferior to most comparing methods on *Yahoo! News* data set, which is attributed to the low intra-class instance similarity. Especially, over 8440 examples come from two categories, among which the intra-class instance similarity of over 65% examples is less than 0.60. Obviously, such low intra-class instance similarity may decrease the effectiveness of our proposed method.

#### 4.2.3 Summary

The two series of experiments mentioned above powerfully demonstrate the effectiveness of GM-PLL, and we attribute the success to the superiority of GM scheme, i.e. simultaneously taking the instance relationship and the co-occurrence possibility of varying instance-label assignments into the same framework. Concretely speaking, for the instance relationship, especially the instance dissimilarity, it can alleviate the effect of the similar instance with varying labels and avoid the outputs of instances be overwhelmed by that of its negative nearest instances. And for the instance-label assignments, the co-occurrence possibility can lead the algorithm to pay more attention to matching selection and reducing its dependence on instance relationship. The two schemes jointly improve the effectiveness and robustness of the proposed

method. And as expected, the experimental results demonstrate the effectiveness of our method.

## 5 FURTHER ANALYSIS

### 5.1 Parameter Sensitivity

The proposed method learns from the PL examples by utilizing two important parameters, i.e. $\beta$ (threshold parameter) and $r$ (the number of candidate labels assigned to unseen instances). Figure 7 and Figure 8 respectively illustrate how GM-PLL performs under different $\beta$ and $r$ configurations. We study the sensitivity analysis of GM-PLL in the following subsection.

#### 5.1.1 The threshold parameter $\beta$

The threshold parameter controls the percentage of prior knowledge incorporated into the learning framework. More prior knowledge can be added into the framework as $\beta$ is small, while less prior knowledge contributes to the learning process when $\beta$ becomes larger. On the other hand, small $\beta$ will draw more noise into the learning framework and large $\beta$ will lose more valuable information, two of which have negative effects on the learning model. Faced with varying data sets, we set the threshold parameter $\beta$ among $\{0.3, 0.4, \ldots, 0.8\}$ via cross-validation and the specific value is shown in Table 6.

TABLE 6
The optimal value of $\beta$ for GM-PLL

| Data set | Lost | MSRCv2 | FG-NET | BirdSong | SoccerPlayer | Yahoo! News |
|----------|------|--------|--------|----------|--------------|-------------|
| $\beta$ | 0.6 | 0.6 | 0.8 | 0.5 | 0.3 | 0.7 |

### 5.1.2 *The number $r$ of candidate label for unseen instances*

As mentioned above, the percentage of candidate labels assigned to unseen instances has great influence on making prediction for unseen instances. According to the analysis in section 3.3, we simultaneously take the total number of class labels (CL*) and the average number of class labels (AVG-CL*) into consideration, and then utilize Eq (10) to obtain the number of assigned labels $r$. To demonstrate the validness of Eq (10) empirically, we conduct the experiments under different $r$ configuration and express the comparing results in Figure 8.

As described in Figure 8, with the increasing of $r$, the classification accuracy of GM-PLL at first increases and later decreases. And such phenomenon is intuitive, i.e. algorithm with smaller $r$ indicates that less noisy labels need to be removed but the ground-truth label has lower possibility to be contained in the candidate label set; and larger $r$ indicates that the ground-truth label has higher possibility to be contained in the candidate label set but it tends to draw more noisy labels into the candidate label set. The number comparison of assigned candidate labels between empirically optimal value and calculation results of Eq (10) on each RW data set is exhibited in Table 7. As shown in Table 7, except the *FG-NET* data set, the empirically optimal number of candidate labels $r^{**}$ is basically identical to the calculation results $r^{*}$ of Eq (10).

TABLE 7
The number comparison of candidate labels between the optimal value of $r^{**}$ and the calculation results $r^{*}$ of Eq (10)

| Data set | Lost | MSRCv2 | FG-NET | BirdSong | SoccerPlayer | Yahoo! News |
|---|---|---|---|---|---|---|
| $r^{*}$ | 3 | 3 | 4 | 3 | 2 | 2 |
| $r^{**}$ | 3 | 4 | 1 | 4 | 2 | 1 |

## 5.2 Time Consumption

Although we have conducted corresponding strategies to reduce the computational complexity of the proposed algorithm, the time consumption of the proposed prediction model is still longer than some comparing methods on some large-scale data sets. Nonetheless, such time consumption is acceptable for the PLL problem. Specifically, on most UCI data sets, the time consumptions are no more than 30 seconds; meanwhile, on some small-scale or medium-scale RW data sets, it is also no more than 20 seconds. Moreover, although the time consumption of the prediction model is longer than some comparing methods, the total running time cost (combining training time and testing time) is appropriate and sometimes even less than some state-of-the-art PLL methods, such as PL-LEAF. According to our experimental results, the running time cost of our proposed methods is no more than 1.5h on all RW data sets, which is only 1/10 of that of PL-LEAF. Table 8 illustrates the total running time and testing time consumption of our proposed algorithm on both UCI and RW data sets, measured within Matlab environment equipped with Intel E5-2650 CPU.

## 6 CONCLUSION

In this paper, we have proposed a novel graph-matching based partial label learning method GM-PLL. To the best of our knowledge, it is the first time to reformulate the PLL problem into a graph matching structure. By incorporating much prior knowledge and establishing INS-CDL assignments, the proposed GM-PLL

TABLE 8
Total running time and testing time consumption of our proposed algorithm on UCI and RW data sets

| Data set | Lost | MSRCv2 | FG-NET | BirdSong | SoccerPlayer |
|---|---|---|---|---|---|
| running time | 37.046s | 127.818s | 198.160s | 281.765s | 3271.877s |
| testing time | 0.837s | 1.431s | 1.254s | 21.879s | 422.012s |
| Data set | Yahoo! News | glass | segment | satimage | vehicle |
| running time | 8612.220s | 2.080s | 80.095s | 236.724s | 7.138s |
| testing time | 1025.886s | 0.204s | 5.901s | 29.574s | 0.743s |
| Data set | letter | abalone | ecoli | dermatology | pendigits |
| running time | 312.502s | 268.547s | 1.916s | 2.924s | 116.202s |
| testing time | 28.344s | 21.380s | 0.287s | 0.334s | 11.538s |

algorithm can effectively contribute the valuable information to the learning model. Extensive experiments have demonstrated the effectiveness of our proposed method. In the future, we will further explore other knowledge from PL data and improve the denoising method to further improve the effectiveness and robustness of the model.

## REFERENCES

[1] L. Liu and T. Dietterich, "Learnability of the superset label learning problem," in *International Conference on Machine Learning*, 2014, pp. 1629–1637.

[2] Y. Chen, V. Patel, R. Chellappa, and P. Phillips, "Ambiguously labeled learning using dictionaries," *IEEE Transactions on Information Forensics and Security*, vol. 9, no. 12, pp. 2076–2088, 2014.

[3] L. Oukhellou, T. Denux, and P. Aknin, "Learning from partially supervised data using mixture models and belief functions," *Pattern Recognition*, vol. 42, no. 3, pp. 334–348, 2009.

[4] J. Wang and M.-L. Zhang, "Towards mitigating the class-imbalance problem for partial label learning," in *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2018, pp. 2427–2436.

[5] Y.-C. Chen, V. M. Patel, J. K. Pillai, R. Chellappa, and P. J. Phillips, "Dictionary learning from ambiguously labeled data," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 353–360.

[6] M.-L. Zhang, "Disambiguation-free partial label learning," in *Proceedings of the 2014 SIAM International Conference on Data Mining*, 2014, pp. 37–45.

[7] Z. Zeng, S. Xiao, K. Jia, T. Chan, S. Gao, D. Xu, and Y. Ma, "Learning by associating ambiguously labeled images," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 708–715.

[8] M. Xie and S. Huang, "Partial multi-label learning," in *AAAI Conference on Artificial Intelligence*, 2018, pp. 1–1.

[9] C. H. Chen, V. M. Patel, and R. Chellappa, "Learning from ambiguously labeled face images," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. PP, no. 99, pp. 1–1, 2018.

[10] T. Cour, B. Sapp, C. Jordan, and B. Taskar, "Learning from ambiguously labeled images," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 919–926.

[11] M. Zhang, B. Zhou, and X. Liu, "Partial label learning via feature-aware disambiguation," in *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016, pp. 1335–1344.

[12] J. Luo and F. Orabona, "Learning from candidate labeling sets," in *Advances in Neural Information Processing Systems*, 2010, pp. 1504–1512.

[13] L. Liu and T. G. Dietterich, "A conditional multinomial mixture model for superset label learning," in *Advances in Neural Information Processing Systems*, 2012, pp. 548–556.

[14] R. Jin and Z. Ghahramani, "Learning with multiple labels," in *Advances in Neural Information Processing Systems*, 2003, pp. 921–928.

[15] L. Feng and B. An, "Leveraging latent label distributions for partial label learning," in *International Joint Conference on Artificial Intelligence*, 2018, pp. 2107–2113.

[16] N. Nguyen and R. Caruana, "Classification with partial labels," in *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2008, pp. 551–559.

[17] F. Yu and M. Zhang, "Maximum margin partial label learning," *Machine Learning*, vol. 106, no. 4, pp. 573–593, 2017.

[18] M. Zhang and F. Yu, "Solving the partial label learning problem: an instance-based approach," in *International Joint Conference on Artificial Intelligence*, 2015, pp. 4048–4054.

[19] F. Yu and M.-L. Zhang, "Maximum margin partial label learning," in *Asian Conference on Machine Learning*, 2016, pp. 96–111.

[20] M. Chertok and Y. Keller, "Spectral symmetry analysis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, no. 7, pp. 1227–1238, 2009.

[21] A. Egozi, Y. Keller, and H. Guterman, "Improving shape retrieval by spectral matching and meta similarity," *IEEE Transactions on Image Processing*, vol. 19, no. 5, pp. 1319–1327, 2010.

[22] J. H. Hays, M. Leordeanu, A. A. Efros, and Y. Liu, "Discovering texture regularity via higher-order matching," in *European Conference on Computer Vision*, 2006, pp. 522–535.

[23] T. Wang and H. Ling, "Gracker: A graph-based planar object tracker," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 6, pp. 1494–1501, 2018.

[24] T. Wang, H. Ling, C. Lang, and S. Feng, "Graph matching with adaptive and branching path following," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, no. 1, pp. 1–1, 2017.

[25] E. Hullermeier and J. Beringer, "Learning from ambiguously labeled examples," *International Symposium on Intelligent Data Analysis*, vol. 10, no. 5, pp. 168–179, 2005.

[26] G. Chen, T. Liu, Y. Tang, Y. Jian, Y. Jie, and D. Tao, "A regularization approach for instance-based superset label learning," *IEEE Transactions on Cybernetics*, vol. 48, no. 3, pp. 967–978, 2017.

[27] C. Tang and M. Zhang, "Confidence-rated discriminative partial label learning," in *AAAI Conference on Artificial Intelligence*, 2017, pp. 2611–2617.

[28] Y. Grandvalet and Y. Bengio, "Learning from partial labels with minimum entropy," *Cirano Working Papers*, pp. 512–517, 2004.

[29] Y. Zhou, J. He, and H. Gu, "Partial label learning via gaussian processes," *IEEE Transactions on Cybernetics*, vol. 47, no. 12, pp. 4443–4450, 2016.

[30] P. Vannoorenberghe and P. Smets, "Partially supervised learning by a credal em approach," in *European Conference on Symbolic and Quantitative Approaches to Reasoning and Uncertainty*, 2005, pp. 956–967.

[31] M. Zhang, F. Yu, and C. Tang, "Disambiguation-free partial label learning," *IEEE Transactions on Knowledge and Data Engineering*, vol. 29, no. 10, pp. 2155–2167, 2017.

[32] T. G. Dietterich and G. Bakiri, "Solving multiclass learning problems via error-correcting output codes," *Journal of Artificial Intelligence Research*, vol. 48, pp. 263–286, 1994.

[33] X. Wu and M.-L. Zhang, "Towards enabling binary decomposition for partial label learning," in *International Joint Conference on Artificial Intelligence*, 2018, pp. 2868–2874.

[34] M. Cho, J. Lee, and K. M. Lee, "Reweighted random walks for graph matching," in *European Conference on Computer vision*, 2010, pp. 492–505.

[35] Z.-Y. Liu and H. Qiao, "Gnccp graduated non convexity and concavity procedure," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 6, pp. 1258–1267, 2014.

[36] A. Egozi, Y. Keller, and H. Guterman, "A probabilistic approach to spectral graph matching," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 1, pp. 18–27, 2013.

[37] M. Leordeanu and M. Hebert, "A spectral technique for correspondence problems using pairwise constraints," in *International Conference on Computer Vision*, 2005, pp. 1482–1489.

[38] T. Cour, B. Sapp, and B. Taskar, "Learning from partial labels," *IEEE Transactions on Knowledge and Data Engineering*, vol. 12, no. 5, pp. 1501–1536, 2011.

[39] F. Briggs, X. Fern, and R. Raich, "Rank-loss support instance machines for miml instance annotation," in *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2012, pp. 534–542.

[40] G. Panis and A. Lanitis, "An overview of research activities in facial age estimation using the fg-net aging database," *Journal of American History*, vol. 5, no. 2, pp. 37–46, 2016.

[41] M. Guillaumin, J. Verbeek, and C. Schmid, "Multiple instance metric learning from automatically labeled bags of faces," in *European Conference on Computer Vision*, 2010, pp. 634–647.