

# Complementary to Multiple Labels: A Correlation-Aware Correction Approach

Yi Gao, Miao Xu, and Min-Ling Zhang, *Senior Member, IEEE*

**Abstract**—*Complementary label learning* (CLL) requires annotators to give *irrelevant labels* instead of relevant labels for instances. Currently, CLL has shown its promising performance on multi-class data by estimating a transition matrix. However, current multi-class CLL techniques cannot work well on multi-labeled data since they assume each instance is associated with one label while each multi-labeled instance is relevant to multiple labels. Here, we show theoretically how the estimated transition matrix in multi-class CLL could be distorted in multi-labeled cases as they ignore co-existing relevant labels. Moreover, theoretical findings reveal that calculating a transition matrix from label correlations in *multi-labeled CLL* (ML-CLL) needs multi-labeled data, while this is unavailable for ML-CLL. To solve this issue, we propose a two-step method to estimate the transition matrix from candidate labels. Specifically, we first estimate an initial transition matrix by decomposing the multi-label problem into a series of binary classification problems, then the initial transition matrix is corrected by label correlations to enforce the addition of relationships among labels. We further show that the proposal is classifier-consistent, and additionally introduce an MSE-based regularizer to alleviate the tendency of BCE loss overfitting to noises. Experimental results have demonstrated the effectiveness of the proposed method.

**Index Terms**—Complementary label learning, multi-label learning, transition matrix, label correlations.

## 1 INTRODUCTION

In *multi-label learning* (MLL), each instance is associated with a set of *relevant labels*, where the learned classifier aims to predict all relevant labels of unseen instances [1], [2], [3]. MLL is widely used in many real-world applications, such as text categorization [4], [5], image retrieval [6], [7], and medical domain [8], [9], etc. However, collecting precisely multi-labeled data is laborious because of the unknown number of relevant labels per instance and the existence of complex semantic labels [10]. For the example image in Fig. 1, besides the label *Architecture*, there exist other relevant labels whose accurate annotation needs one-by-one checking of the whole label space; in addition, annotators need special geographical and cultural domain knowledge to accurately label the image as *Paris*.

To release the laborious of annotating multi-labeled data, we explore the problem setting of *multi-labeled CLL* (ML-CLL), where each instance is associated with a *single complementary label* (an irrelevant label of the instance) instead of multiple relevant labels. Providing such weakly supervised information will ease the labeling process in large label space because selecting one complementary label is low-cost and requires less domain knowledge than selecting all relevant labels. One example of ML-CLL is given in Fig. 1 when selecting *desert* as the complementary label. Given



**The relevant label set**  
people    architecture  
sky        plant  
**Paris**  
**Complementary label**  
Not “desert”

Fig. 1. An example of ML-CLL. The relevant labels of the image include *people*, *architecture*, *sky*, *plant*, and *Paris*, indicating that these elements belong to the image. On the other hand, *desert* serves as the complementary label for this image, signifying the absence of a desert in the image. Notably, *Paris* is considered a complex semantic label, as it is difficult to be directly identified without domain knowledge.

the complementary label, the goal of ML-CLL is still the same as fully supervised MLL, i.e., learning a model that can accurately predict multiple relevant labels for unseen instances.

The setting of CLL was initially applied in the multi-class learning task [11], [12], [13], [14], [15], [16], [17]. Previous multi-class CLL approaches are based on an estimated transition matrix that summarizes the probability of a label being selected as a complementary label [11], [12], [13]. Although they have achieved a promising performance on multi-class data, they are restricted to the case where an instance is associated with only one relevant label. In this case, multi-class CLL approaches only consider the exclusive relationship among labels, while these approaches ignore that labels can bear other relationships in the multi-labeled case, especially the co-occurrence of labels. In fact, relationships among labels are crucial to solving ML-CLL problems since the selection of a complementary label of an instance in MLL is the combined result against multiple

- Yi Gao is with the School of Cyber Science and Engineering, Southeast University, Nanjing 210096, China and the Key Laboratory of Computer Network and Information Integration (Southeast University), Ministry of Education, China. E-mail: gao\_yi@seu.edu.cn
- Miao Xu is with The University of Queensland, Australia. E-mail: miao.xu@uq.edu.au
- Min-Ling Zhang (corresponding author) is with the School of Computer Science and Engineering, Southeast University, Nanjing 210096, China and the Key Laboratory of Computer Network and Information Integration (Southeast University), Ministry of Education, China. E-mail: zhangml@seu.edu.cn

relevant labels rather than against only a single relevant label. Misusing a technique targeting against a single relevant label to the multiple relevant labels case will result in a wrongly estimated transition matrix.

In this paper, we first theoretically analyze how the estimation of the transition matrix using the current multi-class CLL techniques could be distorted in multi-labeled cases. According to these findings, we observe that estimating the transition matrix in ML-CLL from label correlations needs to know relevant labels of instances, while these are unavailable. To remove this pain, we propose a two-step method to estimate the transition matrix in ML-CLL from candidate labels which are the complement of complementary labels. Our strategy includes: (1) estimating an initial transition matrix by decomposing the multi-label problem into binary classification problems; (2) using label correlations to correct the initial transition matrix by enforcing the addition of relationships among labels. The fast convergence of *Cross-Entropy* (CE) loss benefits from focusing on instances that are difficult to classify, which may result in CE loss overfitting to noisy labeled data. As a type of CE loss, *Binary CE* (BCE) loss has the same problem. The study of [18] indicates that *Mean Squared Error* (MSE) loss is less sensitive to noisy labels than CE loss. As *Binary CE* (BCE) loss is a benchmark of our approach, an MSE-based regularizer is further introduced to alleviate the tendency of it overfitting to noises.

In addition, we show that our proposed ML-CLL can be easily combined with learning from relevant labels, which significantly extends the application scenario of the proposed algorithm. This combination is particularly useful, e.g. when labels are collected via crowd-sourcing [19] where crowd-workers are asked to randomly select a complementary label and one or more relevant labels for an instance. Experimental results on various datasets demonstrate the effectiveness of the proposed approach. Especially in situation when each instance is only equipped with a complementary label and a relevant label, our proposal has superior performance, even comparable with the performance on fully supervised data. Our main contributions are summarized as follows:

- We theoretically analyze the distortion of the transition matrix estimated by multi-class CLL in multi-labeled cases, because multi-class CLL techniques ignore the co-existence of relevant labels. Theoretical findings reveal that multi-labeled data is indispensable for calculating the transition matrix from label correlations.
- To solve the problem of unavailable multi-labeled data, we propose a two-step method to estimate the transition matrix from candidate labels. Moreover, we show theoretically that the proposed approach is classifier-consistent under a mild assumption.
- We introduce a practical strategy – MSE-based regularization – to alleviate the overfitting tendency of BCE loss. Our empirical study shows that the proposal obtains comparable performance with state-of-the-art baselines, which proves the effectiveness of our approach.

The rest of this paper are organized as follows. Section 2 briefly reviews related work of ML-CLL. Then we formalize

the ML-CLL problem in Section 3, analyze it theoretically and describe our approach in Section 4. In Section 5, we introduce an MSE-based regularization and show how to adapt our method to bear an additional small amount of relevant labels. The experimental results are given in Section 6 and we conclude in Section 7.

## 2 RELATED WORK

In this section, we will give a brief review of related work of ML-CLL, including MLL, *partial multi-label learning* (PML) and multi-class CLL.

### 2.1 Multi-Label Learning

MLL problems aim to train a classifier that can predict a set of relevant labels for an unseen instance, where each training instance is associated with multiple relevant labels simultaneously. With the complexity of label correlation, the previous studies can be grouped into three categories [20], [21], [22], [23]: *first-order approach* [24], [25], [26], *second-order approach* [27], [28] and *high-order approach* [29], [30]. To solve MLL problems, the first-order approach decomposes MLL problems into a set of binary classification problems [24], [25]. However, these approaches ignore label correlations among labels, which play a crucial role in MLL [20]. After realizing the importance of label correlation, more and more studies attempt to exploit it to improve MLL performance. Among them, the second-order approach considers the pairwise label correlations that refer to the relationship between two labels. The kind of these approaches generally transform MLL problems into bipartite ranking problems by enforcing that relevant labels should be ranked higher than irrelevant labels [28], [31], [32]. Beyond second-order relationship, there exists more complex relationship between labels in many real-world scenarios. Therefore, many approaches begin to exploit high-order label correlations to handle the MLL problems recently [29], [33], [34], [35]. For example, Zhao et al. [35] leverage variational autoencoder to facilitate the learning process via exploiting high-order correlations among labels, while Wang et al. and Xun et al. [36], [37] both design special neural network blocks to automatically extract label correlations to improve the label prediction performance. Although high-order approaches have the ability of stronger label correlation-modeling, they may suffer from high computational cost comparing to first and second-orders approaches [38].

### 2.2 Partial Multi-Label Learning

Due to that the fully supervised data is difficult to collect, many researchers tend to explore the weakly supervision data form to alleviate the heavy load of labeled data collection [39]. PML is a recently emerging weakly supervised approach firstly proposed by Xie et al. [40]. In PML, each training instance is associated with a set of candidate labels that consist of *relevant* labels and *irrelevant* (noisy) labels and the goal is to learn a classifier assigning a set of labels accurately for unseen instances.

At the first glance, it seems that ML-CLL is an extreme case of PML, such that all PML methods are also applicable to ML-CLL. However, existing PML methods assume that

TABLE 1  
Summary of major mathematical notations.

Notations	Mathematical meanings
$\mathcal{X}$	the feature space of instances
$\mathcal{Y}$	the label space with $K$ possible labels $\{l_1, l_2, \dots, l_K\}$
$\mathbf{x}$	a multi-label instance ( $\mathbf{x} \in \mathcal{X}$ )
$Y$	the relevant label set of $\mathbf{x}$ ( $Y \subset \mathcal{Y}, Y \neq \emptyset$ nor $\mathcal{Y}$ )
$\mathbf{y}$	the binary vector of $Y$ ( $\mathbf{y} = [y^1, y^2, \dots, y^K]$ ), where $y^k = 1$ indicates that $l_k \in Y$ and 0 otherwise
$D$	multi-label dataset $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, Y)$
$f(\cdot)$	real-valued decision function $f : \mathcal{X} \rightarrow \mathbb{R}^K$ , where $f^k(\mathbf{x})$ predicts label $l_k$ being the relevant label of $\mathbf{x}$
$L(\cdot, \cdot)$	proper MLL loss function
$R_L(f)$	the expected risk of MLL
$f^*$	the minimizer of $R_L(f)$ ( $f^* = \operatorname{argmin}_f R_L(f)$ )
$\bar{\mathbf{y}}$	the complementary label of $\mathbf{x}$ ( $\bar{\mathbf{y}} \in \mathcal{Y} \setminus Y$ )
$\bar{\mathbf{y}}$	$K$ -dimensional vector of $\bar{\mathbf{y}}$ ( $\bar{\mathbf{y}} = [\bar{y}^1, \bar{y}^2, \dots, \bar{y}^K]$ ), where $\bar{y}^j = 1$ indicates $\bar{y} = l_j$ and 0 otherwise
$\hat{Y}$	the candidate label set of $\mathbf{x}$ , where $\hat{Y} = \mathcal{Y} \setminus \bar{Y}$
$\hat{\mathbf{y}}$	$K$ -dimensional vector of $\hat{Y}$ , where $\hat{\mathbf{y}} = \mathbf{1} - \bar{\mathbf{y}}$
$\bar{D}$	ML-CLL training set $\{(\mathbf{x}_i, \bar{\mathbf{y}}_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, \bar{Y})$
$\bar{L}(\cdot, \cdot)$	ML-CLL loss function
$R_{\bar{L}}(f)$	the expected risk of ML-CLL
$\mathcal{Y}'$	the label space excludes $\emptyset$ and $\mathcal{Y}$ , where $\mathcal{Y}' = \{2^{\mathcal{Y}} - \emptyset - \mathcal{Y}\}$
$C$	a subset of $\mathcal{Y}'$ ( $C \in \mathcal{Y}'$ ), where $C_k$ is the $k$ -th label subset in $\mathcal{Y}'$
$\tilde{\mathbf{T}}$	high-dimension transition matrix, $\tilde{\mathbf{T}} \in \mathbb{R}^{(2^K - 2) \times K}$
$\mathbf{T}$	low-dimension transition matrix $\mathbf{T} \in [0, 1]^{K \times K}$ used to replace $\tilde{\mathbf{T}}$
$\mathbf{Q}$	the transition matrix estimated in multi-class CLL, where $\mathbf{Q} \in [0, 1]^{K \times K}$
$\ell_j$	the difference between $\mathbf{T}$ and $\mathbf{Q}$ on label $l_j$ being the complementary label of $\mathbf{x}$
$\mathbf{S}$	initial transition matrix, where $\mathbf{S} \in [0, 1]^{K \times K}$
$A_k$	the subset of $\mathbf{x}$ in $\bar{D}$ with $\hat{y}^k = 1$
$\mathbf{C}$	the label correlation matrix ( $\mathbf{C} \in [0, 1]^{K \times K}$ ), which presents the relationship among labels
$\hat{\mathbf{T}}$	the estimation of $\mathbf{T}$ , where $\hat{\mathbf{T}}$ is calculated by the proposed method in the paper
$f_{CL}^*$	the minimizer of $R_{\bar{L}}(f)$ ( $f_{CL}^* = \operatorname{argmin}_f R_{\bar{L}}(f)$ )
$\bar{L}_{mse}(\cdot, \cdot)$	the MSE-based regularizer
$\bar{L}(\cdot, \cdot)$	the target loss function in this paper
$\tilde{Y}$	the subset of $Y$ , where $\tilde{Y} \subseteq Y$ and $\tilde{Y} \neq \emptyset$
$\tilde{\mathbf{y}}$	binary vector of $\tilde{Y}$ , where $\tilde{y}^1 = 1$ when label $l_k \in \tilde{Y}$ and 0 otherwise

noisy labels only compose a small portion in the candidate labels [38], [41], [42], [43], such that many approaches [38], [42], [43] adopt matrix factorization to tackle PML problems, which decompose the candidate label matrix into the low-rank multi-label matrix and the sparse noisy label matrix. Compared to PML, the studied ML-CLL problem in this paper are targeted at the problem with only one complementary label, resulting in a high-noise PML problem on which the existing approaches can not be applicable. We will demonstrate the performance difference in the experimental part.

### 2.3 Multi-Class Complementary Label Learning

Currently, CLL problem is only considered in multi-class learning, whose goal is to predict a single relevant label

per instance precisely from complementary labeled data. Previous approaches can be roughly grouped into two categories: (1) modeling the generative relationship between the complementary label and the relevant label [11], [12], [13], [17], [44]; (2) modeling the probability of complementary labels from the learned discriminative classifier directly [14], [15], [16].

The first multi-class CLL method belongs to category one. It models the generative relationship between complementary labels and relevant labels, and uses a such generative process to rewrite one-versus-all and pairwise comparison loss functions to derive an unbiased risk estimator [11]. Ishida et al. [12] realize that the method of [11] is restricted to loss functions and propose a new method which can use arbitrary losses and models. A typical way to make use of the modeled generative process is through a transition matrix, which summarizes the probabilities of a label being complementary labels when relevant labels are given. Then, approaches apply a transition matrix to recover relevant labels from complementary labels [12], [13], [44]. Compared with [11], [12], transition matrix-based methods can map more complex generative relationship rather than uniform one only. Therefore, we tend to design a transition matrix-based method to solve ML-CLL problem with a different estimating way.

Differ from category one, approaches residing in category two directly model the probabilities of complementary labels from the learned classifier without the generative relationship [14], [15], [16]. Chou et al. propose a surrogate complementary loss framework based on complementary labels providing negative feedback during the training process [14]. Although its losses fail to derive an unbiased risk estimator, it achieves good performance on the multi-class CLL. In light of the property of the complementary label that the predictive probability of the complementary label is expected to approach zero, [15] and [16] propose a discriminative solution by directly modeling the probabilities of complementary labels from learned classifier to avoid the generative assumption. Due to that multi-class CLL approaches are designed for a single relevant label case, which are not suitable for the ML-CLL case that an instance is associated with multiple labels simultaneously. We will demonstrate that in the experimental part.

## 3 PROBLEM SETUP

In MLL, let  $\mathcal{X}$  be the feature space and  $\mathcal{Y} = \{l_1, l_2, \dots, l_K\}$  be the finite label space with  $K$  possible class labels ( $K > 2$ ). A multi-label instance  $\mathbf{x} \in \mathcal{X}$  is equipped with a set of relevant labels  $Y \subset \mathcal{Y}$ .  $(\mathbf{x}, Y)$  is independently sampled from an unknown joint probability distribution  $p(\mathbf{x}, Y)$ . Here we exclude the special cases of  $Y = \emptyset$  nor  $\mathcal{Y}$  to ensure relevant labels and complementary labels both exist. For convenience, we use a binary vector  $\mathbf{y} = [y^1, y^2, \dots, y^K] \in \{0, 1\}^K$  to denote  $Y$ , where  $y^k = 1$  indicates that  $l_k \in Y$  is relevant to  $\mathbf{x}$  and 0 otherwise. Suppose  $D = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, Y)$  is the training set with  $n$  instances. The goal of MLL is to learn a multi-label classifier  $h : \mathcal{X} \rightarrow 2^{\mathcal{Y}}$ , which can predict a set of relevant labels for any unseen instance. Instead of learning  $h$  directly, most MLL methods tend to learn a real-

valued decision function  $f : \mathcal{X} \rightarrow \mathbb{R}^K$  via minimizing the expected risk

$$R_L(f) = \mathbb{E}_{p(x, Y)}[L(f(x), y)], \quad (1)$$

where  $L$  is a proper MLL loss function [35], such as BCE loss.  $f(x)$  is usually interpreted as a probability vector:  $f^k(x)$  is the  $k$ -th entry of  $f(x)$  and predicts the confidence score that label  $l_k$  is relevant to  $x$ , i.e., if properly normalized then  $p(y^k = 1|x)$ . Denoting the optimal classifier learned from the expected risk as  $f^*$ , i.e.,  $f^* = \operatorname{argmin}_f R_L(f)$ .

In ML-CLL studied in this paper, each training instance is equipped with a *single* complementary label. The complementary labeled instance  $(x, \bar{y}) \in (\mathcal{X}, \mathcal{Y})$  is drawn from an unknown joint probability distribution  $p(x, \bar{y})$ , where  $\bar{y} \in \mathcal{Y} \setminus Y$  is a complementary label of  $x$ .  $\bar{y}$  can be presented as a  $K$ -dimensional vector  $\bar{y} = [\bar{y}^1, \bar{y}^2, \dots, \bar{y}^K]$ . If label  $l_j$  is selected as the complementary label to  $x$  ( $\bar{y} = l_j$ ), then  $\bar{y}^j$  is one and all other elements are zero in  $\bar{y}$ . We utilize  $\hat{Y} = \mathcal{Y} \setminus \bar{y}$  to denote the candidate label set of  $x$ . Let  $\hat{y} = [\hat{y}^1, \hat{y}^2, \dots, \hat{y}^K]$  to be the corresponding vector representation of subset  $\hat{Y}$ , where all elements are one except for the one corresponding to the complementary label, which is set to zero ( $\hat{y} = \mathbf{1} - \bar{y}$ ).

Let  $\bar{D} = \{(x_i, \bar{y}_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(x, \bar{y})$  be the ML-CLL training set with  $n$  instances. The expected risk of multi-labeled CLL is defined over  $p(x, \bar{y})$ :

$$R_{\bar{L}}(f) = \mathbb{E}_{p(x, \bar{y})}[\bar{L}(f(x), \bar{y})], \quad (2)$$

where  $\bar{L}$  denotes a ML-CLL loss, which will be proposed later this paper. For convenient reference, Table 1 provides a summary of the main notations used in this paper, along with their corresponding mathematical interpretations.

## 4 THE PROPOSED APPROACH

In this section, we first introduce the definition of the transition matrix in MLL and analyze why the estimated transition matrix using multi-class techniques is unsuitable for ML-CLL. Then, we describe an advanced two-step way to estimate the transition matrix in the MLL case. Finally, we prove our approach is classifier-consistent with a mild assumption.

### 4.1 Transition Matrix for ML-CLL

In ML-CLL, we start by introducing a transition matrix  $\tilde{\mathbf{T}}$  that summarizes the probabilities for a complementary label given a set of relevant labels. More specifically, the transition matrix  $\tilde{\mathbf{T}}$  is defined as  $\tilde{\mathbf{T}}_{kj} = p(\bar{y}^j = 1|Y = C_k)$  where  $C_k \in \mathcal{Y}' = \{2^{\mathcal{Y}} - \emptyset - \mathcal{Y}\}$  ( $k \in [2^K - 2]$ ) is the  $k$ -th label subset. If  $l_j \in C_k$ , then  $\tilde{\mathbf{T}}_{kj} = 0$  because the label  $l_j$  has no chance to be selected as the complementary label. In this paper, we employ the same class-dependent assumption as the multi-class CLL approach [13]:  $p(\bar{y}|Y, x) = p(\bar{y}|Y)$  as  $x$  are conditionally independent given  $Y$ . Then we can obtain the following equation:

$$p(\bar{y}^j = 1|x) = \sum_{C \in \mathcal{Y}', l_j \notin C} p(\bar{y}^j = 1|Y = C)p(Y = C|x), \quad (3)$$

where we assume the label  $l_j$  is a complementary label of  $x$ . Then, according to Eq. (3),  $p(\bar{y}|x)$  can be approximated by  $p(Y|x)$  when the transition matrix  $\tilde{\mathbf{T}}$  is known. If considering all possible label subsets of  $\mathcal{Y}'$  as  $C$ , we have  $\tilde{\mathbf{T}} \in \mathbb{R}^{(2^K - 2) \times K}$ , i.e., the size of  $\tilde{\mathbf{T}}$  depends on the size of the power set of  $\mathcal{Y}'$ . computing and storing the power set of  $\mathcal{Y}'$  would be infeasible due to its exponential growth, especially when  $K$  is large. For example, with a large number of possible labels,  $2^K - 2$  becomes excessively large. To solve this combinatorial explosion problem, we explore a more practical way to use an alternative lower-dimensional transition matrix to replace the higher-dimensional one. We start investigating the feasibility of the alternative lower-dimensional matrix from Theorem 1.

**Theorem 1.** *Given an instance  $x$ , suppose  $Y$  is the relevant label set and the label  $l_j$  is the complementary label which is randomly selected. Then the following equality holds:*

$$\begin{aligned} p(\bar{y}^j = 1|x) &= \sum_{C \in \mathcal{Y}', l_j \notin C} p(\bar{y}^j = 1|Y = C)p(Y = C|x) \\ &\geq \sum_{k=1, k \neq j}^K p(\bar{y}^j = 1|y^k = 1)p(y^k = 1|x). \end{aligned}$$

The second inequality holds because of addition rule of probability. The detailed proof is in Appendix A. Theorem 1 shows that using  $\mathbf{T}$  to approximate  $p(\bar{y}|x)$  is a lower bound of using  $\tilde{\mathbf{T}}$  to approximate  $p(\bar{y}|x)$ . Observed by Eq. (3), we find that our main goal transforms from precisely predicting the relevant label set  $Y$  of  $x$  to precisely predicting its complementary label  $\bar{y}$  via the transition matrix  $\tilde{\mathbf{T}}$ . This means that we need to maximize the predictive probability of the complementary label of  $x$ , i.e., maximizing  $p(\bar{y}|x)$ . From this point of view, Theorem 1 theoretically shows the feasibility of using a low-dimension transition matrix to replace the high-dimension  $\tilde{\mathbf{T}}$ , because we optimize by maximizing the lower bound of Eq. (3). Let  $\mathbf{T} \in [0, 1]^{K \times K}$  denote the lower-dimensional transition matrix, where the  $(k, j)$ -th element of  $\mathbf{T}$  is  $\mathbf{T}_{kj} = p(\bar{y}^j = 1|y^k = 1)$ , and  $\mathbf{T}_{kj} = 0$  when  $k = j$ . Thus, we adopt the  $K \times K$  matrix  $\mathbf{T}$  as the transition matrix in the following of the paper to avoid the pain in computation and storage brought up by the  $(2^K - 2) \times K$  matrix  $\tilde{\mathbf{T}}$ .

### 4.2 Distortion in Estimating the Transition Matrix

Before exploring how the transition matrix estimated by multi-class CLL is distorted from that of ML-CLL, we first introduce the transition matrix estimated by multi-class CLL techniques. Suppose  $\mathbf{Q} \in [0, 1]^{K \times K}$  be the transition matrix estimated in multi-class CLL. Recalling the approach [13], it estimates the transition matrix under a special assumption: for each label  $l_k$ , existing an anchor set  $\mathcal{S}_{x|l_k} \subset \mathcal{X}$  such that  $p(y^k = 1|x) = 1$  and  $p(y^{k'} = 1|x) = 0$  ( $l_{k'} \in \mathcal{Y} \setminus \{l_k\}$ ). With this assumption and regardless of label correlations, the estimation of  $\mathbf{Q}_{kj}$  is  $p(\bar{y}^j = 1|y^k = 1) = p(\bar{y}^j = 1|x)$  iff  $x$  is sampled from  $\mathcal{S}_{x|l_k}$ , where  $\mathbf{Q}_{kj}$  is the  $k$ -th row and  $j$ -th column element of  $\mathbf{Q}$ .

To measure the distortion between  $\mathbf{T}$  calculated in ML-CLL and the estimated  $\mathbf{Q}$ , we define their difference on the complementary label  $l_j$  of  $\mathbf{x}$  as follows

$$\ell_j = \sum_{k=1}^K |\mathbf{T}_{kj} - \mathbf{Q}_{kj}|. \quad (4)$$

The larger value of  $\sum_{j=1}^K \ell_j$  indicates that  $\mathbf{T}$  deviates further from  $\mathbf{Q}$ . As we know, label correlations and co-occurred multiple labels are key properties of MLL. Due to that the correlations among labels are intricate, directly calculating  $\mathbf{T}$  from all label correlations will bring high computational cost. For convenience, we give a simple case of MLL including label correlations – at most two labels can co-occur for an instance, and the rest of labels are mutually exclusive – to facilitate us calculating  $\mathbf{T}$  from label correlations and explore the distortion of  $\mathbf{T}$  and  $\mathbf{Q}$ . We start to study the above contents from the definition of mutually exclusive.

**Definition 2.** A label set  $\mathcal{Y}$  is mutually exclusive if, for every  $\mathbf{x} \in \mathcal{X}$ , only one element of  $\mathcal{Y}$  is relevant to  $\mathbf{x}$ . In other words, if  $y_1 \in \mathcal{Y}$  is the relevant label for  $\mathbf{x}$ , then no other label in  $\mathcal{Y} \setminus \{y_1\}$  is relevant to  $\mathbf{x}$ .

Under the simple case in MLL, in Theorem 3, we state how to estimate  $\mathbf{T}$  directly from label correlations, and the distortion of  $\mathbf{T}$  and  $\mathbf{Q}$ .

**Theorem 3.** Under a MLL scenario: suppose the labels  $l_{z_1}, l_{z_2} \in \mathcal{Y}$  ( $z_1, z_2 \in [K]$ ,  $z_1 \neq z_2$ ) are dependent, and the labels belonging to  $\mathcal{Y} \setminus \{l_{z_1}, l_{z_2}\}$  are mutually exclusive. For any  $\mathbf{x}$ , its label set  $Y \subseteq \{l_{z_1}, l_{z_2}\}$  and  $Y \neq \emptyset$ . Let the label  $l_j$  ( $j \in [K], j \neq z_1, z_2$ ) be the complementary label of  $\mathbf{x} \in \mathcal{X}$ .  $\mathbf{T}_{z_1j}$  and  $\mathbf{T}_{z_2j}$  calculated from label correlations satisfy

$$\begin{aligned} \mathbf{T}_{z_1j} &= \frac{p(\bar{y}^j = 1 | \mathbf{x})}{p(y^{z_2} = 1 | \bar{y}^j = 1, y^{z_1} = 1, \mathbf{x})p(y^{z_1} = 1 | \mathbf{x})}, \\ \mathbf{T}_{z_2j} &= \frac{p(\bar{y}^j = 1 | \mathbf{x})}{p(y^{z_1} = 1 | \bar{y}^j = 1, y^{z_2} = 1, \mathbf{x})p(y^{z_2} = 1 | \mathbf{x})}, \end{aligned}$$

where  $[K]$  denotes the integer set  $\{1, 2, \dots, K\}$ . The difference of  $\mathbf{T}$  and  $\mathbf{Q}$  on the complementary label  $l_j$  is

$$\ell_j \geq 2\left(\frac{1}{\xi^2} - 1\right)p(\bar{y}^j = 1 | \mathbf{x}),$$

where  $\xi = \max\{p(y^{z_2} = 1 | \bar{y}^j = 1, y^{z_1} = 1, \mathbf{x}), p(y^{z_1} = 1 | \bar{y}^j = 1, y^{z_2} = 1, \mathbf{x})\}$ .

The proof is provided in Appendix B. From Theorem 3, we can see that calculating the transition matrix from label correlations is more complex than estimating one without label correlations, and the relevant label sets of instances need to be known. Moreover, Theorem 3 shows that there is a distortion between  $\mathbf{T}$  and  $\mathbf{Q}$ , which widely exists in multi-labeled cases since each multi-label instance is relevant to multiple labels. The above learning scenario only considers the pairwise label correlations, while there exists a more complex relationship among labels. Similarly, under a realizable computational cost, we construct another simple MLL scenario with more complex label relationships to explore factors that affect  $\ell_j$  in Corollary 4.

**Corollary 4.** Under a MLL scenario: there are  $m$  ( $m \geq 2$ ) labels  $l_{z_1}, l_{z_2}, \dots, l_{z_m} \in \mathcal{Y}$  ( $z_1, \dots, z_m \in [K]$ ) that are dependent,

while the labels belong to  $\mathcal{Y} \setminus \{l_{z_1}, l_{z_2}, \dots, l_{z_m}\}$  are mutually exclusive. For any  $\mathbf{x} \in \mathcal{X}$ , its relevant set  $Y \subseteq \{l_{z_1}, l_{z_2}, \dots, l_{z_m}\}$  and  $Y \neq \emptyset$ . Suppose the label  $l_j$  is the complementary label of  $\mathbf{x}$ . The difference  $\ell_j$  between  $\mathbf{T}$  and  $\mathbf{Q}$  has

$$\ell_j \geq m\left(\frac{1}{\xi^m} - 1\right)p(\bar{y}^j = 1 | \mathbf{x}),$$

where  $\xi = \max\{p(y^{z_m} = 1 | \bar{y}^j = 1, y^{z_1} = 1, \dots, y^{z_{m-1}} = 1, \mathbf{x}), p(y^{z_{m-1}} = 1 | \bar{y}^j = 1, y^{z_1} = 1, \dots, y^{z_{m-2}} = 1, y^{z_m} = 1, \mathbf{x}), \dots, p(y^{z_1} = 1 | \bar{y}^j = 1, y^{z_2} = 1, \dots, y^{z_m} = 1, \mathbf{x})\}$  ( $\xi \in (0, 1]$ ).

The proof is shown in Appendix C. According to Corollary 4, when label correlations are more complex, the distortion of the transition matrix estimated by the multi-class CLL approach is more serious as  $m$  increases. Meanwhile, it demonstrates that the ML-CLL problem cannot be solved by current techniques in multi-class CLL.

### 4.3 Estimation $\mathbf{T}$ with Label Correlations

As discussed above, calculating the transition matrix  $\mathbf{T}$  from label correlations needs instances whose relevant label sets are known. Moreover, calculating  $\mathbf{T}$  is more and more difficult as relationships among labels become more complex by observing the results of  $\mathbf{T}$  in Theorem 3 and Corollary 4. Due to that multi-labeled data are unavailable for our setting, we propose a two-step method to estimate  $\mathbf{T}$  from candidate labels, and it can reduce the complexities in calculating  $\mathbf{T}$  from label correlations. This two-step method includes: (1) computing an initial transition matrix  $\mathbf{S} \in [0, 1]^{K \times K}$  from candidate labels by decomposing the multi-label problem into a series of binary classification problem; (2) obtaining the final estimation of  $\mathbf{T}$  by using label correlations to correct  $\mathbf{S}$ .

**Computing an initial transition matrix  $\mathbf{S}$ .** Let  $S_{kj} = p(\bar{y}^j = 1 | \hat{y}^k = 1)$  be an initial transition probability, which is a  $(k, j)$ -th element of  $\mathbf{S}$ . We calculate  $\mathbf{S}$  from candidate labels of instances. Multiplication theorem of probability<sup>1</sup> is applied to calculate  $S_{kj}$  and ensure that the following equation holds:

$$\begin{aligned} S_{kj} &= p(\bar{y}^j = 1 | \hat{y}^k = 1) \\ &= p(\bar{y}^j = 1 | \hat{y}^k = 1) \int p(\mathbf{x} | \bar{y}^j = 1, \hat{y}^k = 1) d\mathbf{x} \\ &= \int p(\bar{y}^j = 1 | \hat{y}^k = 1, \mathbf{x}) p(\mathbf{x} | \hat{y}^k = 1) d\mathbf{x} \\ &= \mathbb{E}_{p(\mathbf{x} | \hat{y}^k = 1)} [p(\bar{y}^j = 1 | \hat{y}^k = 1, \mathbf{x})], \end{aligned} \quad (5)$$

where  $j, k \in [K]$  and  $j \neq k$ . In practice,  $\mathbb{E}_{p(\mathbf{x} | \hat{y}^k = 1)} [p(\bar{y}^j = 1 | \hat{y}^k = 1, \mathbf{x})]$  can be approximated by the expectation of  $p(\bar{y}^j = 1 | \hat{y}^k = 1, \mathbf{x})$  over the conditional distribution  $p(\mathbf{x} | \hat{y}^k = 1)$ . Assuming  $\bar{y}$  and  $\hat{y}$  are conditionally independent given  $\mathbf{x}$ , so  $p(\bar{y}^j = 1 | \hat{y}^k = 1, \mathbf{x}) = p(\bar{y}^j = 1 | \mathbf{x})$ . Intuitively,  $p(\bar{y}^j = 1 | \mathbf{x})$  can be approximated by the classifier learned from  $\bar{D}$  to predict the probability of complementary labels. Let  $A_k$  denote the subset of  $\mathbf{x}$  in  $\bar{D}$  with  $\hat{y}^k = 1$ ,

$$\begin{aligned} 1 \cdot p(\mathbf{x}, \bar{y}^j = 1, \hat{y}^k = 1) &= p(\bar{y}^j = 1 | \hat{y}^k = 1, \mathbf{x})p(\mathbf{x} | \hat{y}^k = 1)p(\hat{y}^k = 1) \\ &= p(\bar{y}^j = 1 | \hat{y}^k = 1)p(\mathbf{x} | \bar{y}^j = 1, \hat{y}^k = 1)p(\hat{y}^k = 1) \Rightarrow p(\bar{y}^j = 1 | \hat{y}^k = 1, \mathbf{x})p(\mathbf{x} | \hat{y}^k = 1) = p(\bar{y}^j = 1 | \hat{y}^k = 1)p(\mathbf{x} | \bar{y}^j = 1, \hat{y}^k = 1) \end{aligned}$$

$S$		$C^T$		$\hat{T} = SC^T$			
$\bar{y}$	$l_1$	$l_2$	$l_3$	$\bar{y}$	$l_1$	$l_2$	$l_3$
$l_1$	0	0.7	0.3	$l_1$	0	0.9	0.1
$l_2$	0.4	0	0.6	$l_2$	0.9	0	0.6
$l_3$	0.9	0.1	0	$l_3$	0.1	0.6	0

  

$\bar{y}$	$l_1$	$l_2$	$l_3$
$l_1$	0	0.18	0.42
$l_2$	0.06	0	0.04
$l_3$	0.09	0.81	0

Fig. 2. An example of correcting  $S$  with label correlations.

which satisfies the conditional distribution  $p(x|\bar{y}^k = 1)$ . Thus,  $\mathbf{S}_{kj}$  can be estimated by

$$\begin{aligned}\mathbf{S}_{kj} &= \frac{1}{|A_k|} \sum_{\mathbf{x} \in A_k} p(\bar{y}^j = 1 | \bar{y}^k = 1, \mathbf{x}) \\ &= \frac{1}{|A_k|} \sum_{\mathbf{x} \in A_k} p(\bar{y}^j = 1 | \mathbf{x}).\end{aligned}\quad (6)$$

**Estimating  $\mathbf{T}$  with label correlations.** The calculating procedure of  $\mathbf{S}$  lacks exactly supervised data. Observed by the transition probabilities of  $\mathbf{T}$  calculated from label correlations in subsection 4.2, we can find that they are affected by label correlations. Moreover, a label that is low-co-occurred to the relevant labels could be preferentially selected as the complementary label from the view of label correlations. For example, considering *water* as the relevant label; in this case, *desert* (low-co-occurred label) will have a larger chance to be selected as the complementary label compared to *fish* (high-co-occurred label). Motivated by these findings, we use label correlations to correct the initial matrix  $\mathbf{S}$  to estimate  $\mathbf{T}$  by enforcing the addition of relationships among labels.

Suppose  $\mathbf{C} \in [0, 1]^{K \times K}$  be a label correlation matrix, where the element  $\mathbf{C}_{kj}$  represents the correlation between labels  $l_k$  and  $l_j$ . The value of  $\mathbf{C}_{kj}$  is larger when the correlation of labels  $l_k$  and  $l_j$  is stronger. Following [40], [45], we adopt the co-occurrence rate of two candidate labels as their correlations. Finally, the transition matrix  $\mathbf{T}$  can be estimated by  $\hat{\mathbf{T}} = \mathbf{SC}^T$ , where  $\hat{\mathbf{T}}_{kj} = 0$  if  $k = j$ , and normalizing  $\mathbf{T}$  by row.

Fig. 2 is an example of refining procedure. As can be seen from the Fig. 2, though the estimated initial probability of  $p(\bar{y}^2 = 1 | \bar{y}^1 = 1)$  is higher than  $p(\bar{y}^3 = 1 | \bar{y}^1 = 1)$  in  $\mathbf{S}$ , the value of  $p(\bar{y}^2 = 1 | \bar{y}^1 = 1)$  is lower than  $p(\bar{y}^3 = 1 | \bar{y}^1 = 1)$  in  $\hat{\mathbf{T}}$ . This is because the labels  $l_1$  and  $l_2$  have a strong correlation as shown in  $\mathbf{C}$ , so the label  $l_2$  has a lower chance to be selected as the complementary label for the label  $l_1$ . The corrected initial transition matrix  $\mathbf{S}$  agrees with our expectation on the low-co-occurred labels that tend to be selected as complementary labels preferentially. In practice, the estimation of  $\mathbf{T}$  depends on  $p(\bar{y}|\mathbf{x})$ , where the classifier should perfectly model the probability of complementary labels. When data equipped with complementary labels is sufficiently, the perfect model is capable of modeling  $p(\bar{y}|\mathbf{x})$ .

#### 4.4 A Classifier-Consistent Approach

According to the transition matrix  $\mathbf{T}$ , we can derive the probability of complementary labels from multi-label classifier. Let  $\bar{\mathbf{f}}(\mathbf{x}) \in \mathbb{R}^K$  be a complementary label classifier, which is defined as

$$\bar{\mathbf{f}}(\mathbf{x}) = \mathbf{T}^T \mathbf{f}(\mathbf{x}), \quad (7)$$

#### Algorithm 1: MLCL Algorithm

```

Input:
 $\bar{D}$ : the complementary-label training set
 $\{(\mathbf{x}_i, \bar{y}_i)\}_{i=1}^n$ ;
 $E$ : the number of epochs;
 $\mathcal{A}$ : an external stochastic optimization algorithm;
Output:
 $\theta$ : model parameter for  $\bar{\mathbf{f}}(\mathbf{x}; \theta)$ ;
1 if  $\mathbf{T}$  is unknown then
2   Train a classifier  $\bar{\mathbf{f}}(\mathbf{x})$  with the softmax output layer and Cross-Entropy loss on  $\bar{D}$ ;
3   Fill  $\mathbf{S} \in [0, 1]^{K \times K}$  with zeros;
4   for  $k = 1$  to  $K$  do
5     num = 0;
6     for  $(\mathbf{x}_i, \bar{y}_i) \in \bar{D}$  such that  $\bar{y}_i^k = 0$  do
7       num += 1;
8        $\mathbf{S}_{k.} += \bar{\mathbf{f}}(\mathbf{x}_i)$ ; //add  $\bar{\mathbf{f}}(\mathbf{x}_i)$  to  $k$ -th row of  $\mathbf{S}$ 
9     end
10     $\mathbf{S}_{k.} /= \text{num}$ ;
11  end
12   $\hat{\mathbf{T}} = \mathbf{S} \mathbf{C}^T$ ;
13 end
14 for  $t = 1$  to  $E$  do
15   Let  $\mathcal{L}$  be the risk,  $\mathcal{L} = \frac{1}{n} \sum_{i=1}^n \bar{\mathcal{L}}(\mathbf{f}(\mathbf{x}_i), \bar{y}_i) = \frac{1}{n} \sum_{i=1}^n (L(\hat{\mathbf{T}}^T \mathbf{f}(\mathbf{x}_i), \bar{y}_i) + \|\bar{y}_i - \hat{\mathbf{T}}^T \mathbf{f}(\mathbf{x}_i)\|_F^2)$ ;
16   Set gradient  $-\nabla_{\theta} \mathcal{L}$ ;
17   Update  $\theta$  by  $\mathcal{A}$ ;
18 end

```

where  $\bar{\mathbf{f}}(\mathbf{x})$  is applied to approximate  $p(\bar{y}|\mathbf{x})$ ,  $\bar{f}^j(\mathbf{x})$  refers to the  $j$ -th element of  $\bar{\mathbf{f}}(\mathbf{x})$ . ML-CLL problems aim to recover a set of relevant labels per instance from a complementary label. Since training instances are associated with complementary labels, the common loss functions of MLL are unsuitable for ML-CLL. Therefore, we define a complementary loss function  $\bar{\mathcal{L}}$  as

$$\bar{\mathcal{L}}(\mathbf{f}(\mathbf{x}), \bar{y}) = L(\bar{\mathbf{f}}(\mathbf{x}), \bar{y}) = L(\mathbf{T}^T \mathbf{f}(\mathbf{x}), \bar{y}). \quad (8)$$

Denote by  $\mathbf{f}_{CL}^*$  the minimizer of  $R_{\bar{\mathcal{L}}}(\mathbf{f})$ . Recalling the definition of classifier-consistent, if a classifier learned by an approach finally converges to the optimal classifier  $\mathbf{f}^*$  learned in MLL as the number of instances increases, then this approach is classifier-consistent [46], [47], [48]. We derive that our proposal is classifier-consistent based on a mild assumption:

**Assumption 5.** Suppose the transition matrix  $\mathbf{T}$  is invertible and can perfectly recover the relationship between relevant labels of  $\mathbf{x}$  and its complementary label. Then, we have  $\bar{y} = \mathbf{T}^T \mathbf{y}$ .

With Assumption 5, our approach trained on  $\bar{\mathcal{L}}$  can be inferred to be classifier-consistent, which is stated in Theorem 6. Naturally, Theorem 6 guarantees that the optimal classifier learned from complementary labeled data converges to the optimal one learned from fully supervised MLL.

**Theorem 6.** With Assumption 5, suppose the transition matrix  $\mathbf{T}$  is invertible, then the ML-CLL optimal classifier  $\mathbf{f}_{CL}^*$  converges to the MLL optimal classifier  $\mathbf{f}^*$ .

The proof is represented in Appendix D. Thanks to BCE loss is a popular loss function in MLL, we adopt BCE loss as the base in this paper, then  $\bar{L}$  is expressed as

$$\bar{L}(\mathbf{f}(\mathbf{x}), \bar{\mathbf{y}}) = -\bar{\mathbf{y}} \log(\mathbf{T}^T \mathbf{f}(\mathbf{x})) - (1 - \bar{\mathbf{y}}) \log(1 - \mathbf{T}^T \mathbf{f}(\mathbf{x})), \quad (9)$$

where  $\mathbf{1}$  denotes a  $K$ -dimensional vector with 1 for all elements.

## 5 REGULARIZATION-BASED ENHANCEMENT

In this section, an MSE-based regularization of our approach is described. And we attempt to combine a small amount of relevant labels to explore more possibilities of our proposal.

### 5.1 An MSE-Based Regularization

Previous works indicate that CE loss always makes the model focus on hard instances that are difficult to be classified precisely, while MSE loss and Mean Absolute Error (MAE) loss are less sensitive to hard instances since they treat per instance coequally [18], [49]. As this property, the convergence rate of CE loss is superior to MSE loss and MAE loss, whereas this property makes CE loss more prone to the overfitting problem than MSE loss and MAE loss when noisy labels present at training data [18], [49]. Actually, an excellent approach can converge quickly during the training process, and shows good generalization ability and robustness for unseen instances [16].

Obviously, BCE loss has a similar property to CE loss, which results in an excellent convergence rate of approaches. Meanwhile, approaches based on BCE loss are easy to suffer from the overfitting problem when using noisy labeled data to learn. In fact, ML-CLL is a problem setting with dense noisy labels, BCE loss may cause the overfitting problem of a model in ML-CLL. To cope with this problem, we introduce an MSE-based regularizer based on MSE loss (i.e.  $\ell_2$ -norm regularization) to balance the robustness and convergence requirement of the proposed approach. Hence, the MSE-based regularizer is defined as:

$$\bar{L}_{mse}(\mathbf{f}(\mathbf{x}), \bar{\mathbf{y}}) = \left\| \bar{\mathbf{y}} - \mathbf{T}^T \mathbf{f}(\mathbf{x}) \right\|_F^2. \quad (10)$$

Finally, we combine the complementary loss and the MSE-based regularizer term, which leads to our target loss:

$$\tilde{L}(\mathbf{f}(\mathbf{x}), \bar{\mathbf{y}}) = \bar{L}(\mathbf{f}(\mathbf{x}), \bar{\mathbf{y}}) + \beta \bar{L}_{mse}(\mathbf{f}(\mathbf{x}), \bar{\mathbf{y}}), \quad (11)$$

where  $\beta$  is the trade-off parameter and set as 1 (the selection shown in Section 6). The all procedure of the proposed approach (called MLCL) is shown in Algorithm 1.

### 5.2 Incorporation of Relevant Labels

In many practical situations, we can use complementary labels and relevant labels to learn more accurate classifiers, which is highly practical implementation. To this end, motivated by [11], [50], let us design a reasonable combination of the loss derived from complementary labeled data and relevant labeled data:

$$\tilde{\mathcal{L}}(\mathbf{f}(\mathbf{x}), \bar{\mathbf{y}}, \tilde{\mathbf{y}}) = \bar{L}(\mathbf{f}(\mathbf{x}), \bar{\mathbf{y}}) + \|\tilde{\mathbf{y}} - \mathbf{f}(\mathbf{x})\|_F^2, \quad (12)$$

TABLE 2  
Statistics of datasets.

Datasets	$ \mathcal{S} $	$dim(\mathcal{S})$	$L(\mathcal{S})$	$LCard(\mathcal{S})$
scene	2407	294	6	1.07
yeast	2417	103	14	4.23
eurlex_dc	8636	5000	15	1.02
eurlex_sm	13270	5000	15	1.74
corel5k	4194	499	15	1.70
corel16k	11103	120	15	1.77
bookmark	38912	2150	15	1.25
delicious	14784	500	15	4.32

where  $\tilde{\mathbf{y}} = [\tilde{y}^1, \dots, \tilde{y}^K] \in \{0, 1\}^K$  denotes a binary vector of relevant labels  $\tilde{Y}$  of  $\mathbf{x}$ , in which  $\tilde{y}^1 = 1$  when the label  $l_k \in \tilde{Y}$ . To provide more practicability, we do not restrict given relevant labels  $\tilde{Y}$  to must be equal to the set of relevant labels  $Y$ , which means  $\tilde{Y} \subseteq Y$  and  $\tilde{Y} \neq \emptyset$ .

As explained in the instruction, we can naturally collect data associated with complementary labels and relevant labels via crowdsourcing [19]. Our loss function Eq. (12) can leverage both kinds of labeled data to learn better classifiers. We will experimentally show the usefulness of this combination method in Section 6.

## 6 EXPERIMENTS

In this section, we will evaluate the effectiveness of MLCL, where five common MLL criteria, including *ranking loss*, *hamming loss*, *one error*, *coverage* and *average precision*, are employed in this paper. Smaller values for the first four criteria indicate better performance, while a higher value of *average precision* indicates better performance. The label set of  $\mathbf{x}$  is predicted by  $Y = \{l_k | f^k(\mathbf{x}) > 0.5, 1 \leq k \leq K\}$ . All experiments were conducted using PyTorch [51] and NVIDIA TESLA K80 GPU to implement. The code will be released after this paper has been accepted.

### 6.1 Experimental Settings

**Datasets.** We use eight widely-used MLL datasets, namely *corel5k*, *corel16k*, *delicious*, *eurlex\_dc*, *eurlex\_sm*, *yeast*, *bookmarks* and *scene*, to our experiments<sup>2</sup>. Following [40], [41], we adopt the same pre-processing to deal with the datasets. More specifically, rare class labels are filtered out for datasets with more than 15 class labels, whose class labels are kept under 15. Accordingly, instances that are relevant with removed class labels are filtered out as well. Detailed characteristics of these datasets are shown in Table 2.

**Base models.** The linear model is used as the base model.

**Baselines.** Two typical MLL approaches, ML-KNN [26] and LIFT [52], are utilized as baselines, which deal with ML-CLL via regarding all possible labels in the candidate label set as relevant labels for a training instance. Similarly, three recent PML approaches are employed as comparing approaches, including PML-lc [40], fpml [43] and PML-LRS [42], which learn from training instances associated with candidate labels. In addition, we employ a multi-class CLL approach, called L-UW [15], as a baseline, which uses BCE loss and *sigmoid* output layer instead of CE loss and *softmax* output layer respectively to make L-UW suit for multi-labeled data.

2. Publicly available at <http://mulan.sourceforge.net/datasets>.

TABLE 3

Experimental results (mean  $\pm$  std) on training data with uniform complementary labels. The best performance of each dataset is presented in **boldface**, where  $\bullet/\circ$  indicates whether MLCL is superior/inferior to baselines (with 5% t-test).

Methods	ML-KNN	LIFT	fpml	PML-lc	PML-LRS	L-UW	MLCL
Ranking loss $\downarrow$							
scene	.340 $\pm$ .032 $\bullet$	.289 $\pm$ .020 $\bullet$	.504 $\pm$ .025 $\bullet$	.490 $\pm$ .025 $\bullet$	<b>.258<math>\pm</math>.007<math>\circ</math></b>	.372 $\pm$ .028 $\bullet$	.259 $\pm$ .030
yeast	.247 $\pm$ .012 $\bullet$	.298 $\pm$ .012 $\bullet$	.233 $\pm$ .013 $\bullet$	.251 $\pm$ .015 $\bullet$	.464 $\pm$ .019 $\bullet$	.214 $\pm$ .011	<b>.211<math>\pm</math>.013</b>
eurlex_dc	.303 $\pm$ .016 $\bullet$	.286 $\pm$ .016 $\bullet$	.488 $\pm$ .033 $\bullet$	.347 $\pm$ .025 $\bullet$	.316 $\pm$ .011 $\bullet$	.598 $\pm$ .024 $\bullet$	<b>.229<math>\pm</math>.026</b>
eurlex_sm	.336 $\pm$ .010 $\bullet$	.346 $\pm$ .012 $\bullet$	.488 $\pm$ .006 $\bullet$	.436 $\pm$ .011 $\bullet$	.332 $\pm$ .009 $\bullet$	.646 $\pm$ .015 $\bullet$	<b>.312<math>\pm</math>.014</b>
corel5k	.379 $\pm$ .034	.433 $\pm$ .037 $\bullet$	.444 $\pm$ .026 $\bullet$	.406 $\pm$ .075 $\bullet$	<b>.334<math>\pm</math>.009<math>\circ</math></b>	.367 $\pm$ .031	.349 $\pm$ .035
corel16k	.328 $\pm$ .047	.392 $\pm$ .027 $\bullet$	.420 $\pm$ .033 $\bullet$	.457 $\pm$ .046 $\bullet$	.303 $\pm$ .005	.303 $\pm$ .035	<b>.289<math>\pm</math>.042</b>
bookmark	.384 $\pm$ .006 $\bullet$	.310 $\pm$ .007 $\bullet$	.469 $\pm$ .019 $\bullet$	.454 $\pm$ .036 $\bullet$	.260 $\pm$ .004	.303 $\pm$ .010 $\bullet$	<b>.252<math>\pm</math>.013</b>
delicious	.398 $\pm$ .004 $\bullet$	.383 $\pm$ .003 $\bullet$	.438 $\pm$ .008 $\bullet$	.445 $\pm$ .015 $\bullet$	.305 $\pm$ .002	.302 $\pm$ .006	.310 $\pm$ .003
One Error $\downarrow$							
scene	.692 $\pm$ .030 $\bullet$	.605 $\pm$ .023 $\bullet$	.815 $\pm$ .027 $\bullet$	.717 $\pm$ .021 $\bullet$	.540 $\pm$ .023 $\bullet$	.609 $\pm$ .041 $\bullet$	<b>.427<math>\pm</math>.018</b>
yeast	.297 $\pm$ .029 $\bullet$	.284 $\pm$ .028 $\bullet$	.251 $\pm$ .025	.583 $\pm$ .026 $\bullet$	.738 $\pm$ .102 $\bullet$	.251 $\pm$ .025	<b>.251<math>\pm</math>.023</b>
eurlex_dc	.776 $\pm$ .031 $\bullet$	.670 $\pm$ .013 $\bullet$	.925 $\pm$ .016 $\bullet$	.774 $\pm$ .015 $\bullet$	.847 $\pm$ .010 $\bullet$	.837 $\pm$ .034 $\bullet$	<b>.594<math>\pm</math>.035</b>
eurlex_sm	.689 $\pm$ .012 $\bullet$	.679 $\pm$ .009 $\bullet$	.872 $\pm$ .011 $\bullet$	.662 $\pm$ .012	.731 $\pm$ .005 $\bullet$	.696 $\pm$ .008 $\bullet$	<b>.656<math>\pm</math>.029</b>
corel5k	.815 $\pm$ .048 $\bullet$	.842 $\pm$ .056 $\bullet$	.854 $\pm$ .035 $\bullet$	.811 $\pm$ .062 $\bullet$	.756 $\pm$ .010	.769 $\pm$ .034	<b>.736<math>\pm</math>.065</b>
corel16k	.736 $\pm$ .056	.789 $\pm$ .046 $\bullet$	.816 $\pm$ .025 $\bullet$	.946 $\pm$ .028 $\bullet$	.730 $\pm$ .000 $\bullet$	.693 $\pm$ .057	<b>.690<math>\pm</math>.056</b>
bookmark	.801 $\pm$ .006 $\bullet$	.649 $\pm$ .016 $\bullet$	.885 $\pm$ .020 $\bullet$	.798 $\pm$ .005 $\bullet$	.584 $\pm$ .005 $\bullet$	.590 $\pm$ .022 $\bullet$	<b>.509<math>\pm</math>.012</b>
delicious	.592 $\pm$ .018 $\bullet$	.533 $\pm$ .015 $\bullet$	.618 $\pm$ .017 $\bullet$	.679 $\pm$ .011 $\bullet$	.452 $\pm$ .007	.467 $\pm$ .023 $\bullet$	<b>.448<math>\pm</math>.016</b>
Hamming loss $\downarrow$							
scene	.820 $\pm$ .002 $\bullet$	.820 $\pm$ .003 $\bullet$	.819 $\pm$ .002 $\bullet$	<b>.251<math>\pm</math>.007</b>	.814 $\pm$ .000 $\bullet$	.518 $\pm$ .042 $\bullet$	.264 $\pm$ .027
yeast	.697 $\pm$ .012 $\bullet$	.697 $\pm$ .013 $\bullet$	.697 $\pm$ .013 $\bullet$	.268 $\pm$ .010 $\bullet$	.316 $\pm$ .000 $\bullet$	.243 $\pm$ .010	<b>.235<math>\pm</math>.008</b>
eurlex_dc	.932 $\pm$ .000 $\bullet$	.932 $\pm$ .000 $\bullet$	.118 $\pm$ .006 $\bullet$	.104 $\pm$ .002 $\bullet$	.890 $\pm$ .039 $\bullet$	.806 $\pm$ .015 $\bullet$	<b>.092<math>\pm</math>.005</b>
eurlex_sm	.883 $\pm$ .001 $\bullet$	.883 $\pm$ .001 $\bullet$	.148 $\pm$ .005 $\bullet$	<b>.138<math>\pm</math>.002</b>	.825 $\pm$ .027 $\bullet$	.773 $\pm$ .008 $\bullet$	.139 $\pm$ .005
corel5k	.886 $\pm$ .007 $\bullet$	.887 $\pm$ .007 $\bullet$	.887 $\pm$ .007 $\bullet$	<b>.155<math>\pm</math>.004</b>	.869 $\pm$ .002 $\bullet$	.463 $\pm$ .018 $\bullet$	.229 $\pm$ .068
corel16k	.882 $\pm$ .009 $\bullet$	.882 $\pm$ .009 $\bullet$	.882 $\pm$ .009 $\bullet$	<b>.177<math>\pm</math>.011<math>\circ</math></b>	.862 $\pm$ .001 $\bullet$	.423 $\pm$ .033 $\bullet$	.202 $\pm$ .067
bookmark	.917 $\pm$ .001 $\bullet$	.916 $\pm$ .001 $\bullet$	.420 $\pm$ .009 $\bullet$	<b>.123<math>\pm</math>.001<math>\circ</math></b>	.813 $\pm$ .001 $\bullet$	.409 $\pm$ .014 $\bullet$	.140 $\pm$ .004
delicious	.711 $\pm$ .003 $\bullet$	.711 $\pm$ .003 $\bullet$	.711 $\pm$ .003 $\bullet$	.394 $\pm$ .011 $\bullet$	.459 $\pm$ .002 $\bullet$	.369 $\pm$ .027 $\bullet$	<b>.289<math>\pm</math>.004</b>
Coverage $\downarrow$							
scene	.299 $\pm$ .026 $\bullet$	.256 $\pm$ .017 $\bullet$	.434 $\pm$ .021 $\bullet$	.420 $\pm$ .021 $\bullet$	<b>.230<math>\pm</math>.006</b>	.328 $\pm$ .022 $\bullet$	.234 $\pm$ .025
yeast	.579 $\pm$ .018 $\bullet$	.649 $\pm$ .020 $\bullet$	.553 $\pm$ .033 $\bullet$	<b>.506<math>\pm</math>.023</b>	.742 $\pm$ .027 $\bullet$	.525 $\pm$ .017	.525 $\pm$ .021
eurlex_dc	.285 $\pm$ .014 $\bullet$	.269 $\pm$ .015 $\bullet$	.458 $\pm$ .031 $\bullet$	.326 $\pm$ .023 $\bullet$	.298 $\pm$ .010 $\bullet$	.334 $\pm$ .017 $\bullet$	<b>.204<math>\pm</math>.023</b>
eurlex_sm	.416 $\pm$ .010 $\bullet$	.427 $\pm$ .013 $\bullet$	.569 $\pm$ .010 $\bullet$	.509 $\pm$ .013 $\bullet$	.419 $\pm$ .010 $\bullet$	.519 $\pm$ .008 $\bullet$	<b>.365<math>\pm</math>.014</b>
corel5k	.473 $\pm$ .034	.516 $\pm$ .035 $\bullet$	.529 $\pm$ .028 $\bullet$	.492 $\pm$ .072	<b>.429<math>\pm</math>.008</b>	.457 $\pm$ .038	.445 $\pm$ .048
corel16k	.430 $\pm$ .044	.488 $\pm$ .027 $\bullet$	.513 $\pm$ .035 $\bullet$	.537 $\pm$ .051 $\bullet$	.405 $\pm$ .008	.407 $\pm$ .033	<b>.393<math>\pm</math>.042</b>
bookmark	.359 $\pm$ .007 $\bullet$	.328 $\pm$ .008 $\bullet$	.475 $\pm$ .019 $\bullet$	.458 $\pm$ .035 $\bullet$	.280 $\pm$ .004	.292 $\pm$ .011 $\bullet$	<b>.279<math>\pm</math>.011</b>
delicious	.712 $\pm$ .006 $\bullet$	.703 $\pm$ .004 $\bullet$	.726 $\pm$ .009 $\bullet$	.695 $\pm$ .009 $\bullet$	<b>.609<math>\pm</math>.003<math>\circ</math></b>	.613 $\pm$ .006 $\circ$	.632 $\pm$ .007
Average Precision $\uparrow$							
scene	.543 $\pm$ .024 $\bullet$	.600 $\pm$ .017 $\bullet$	.417 $\pm$ .021 $\bullet$	.465 $\pm$ .018 $\bullet$	.637 $\pm$ .011 $\bullet$	.568 $\pm$ .026 $\bullet$	<b>.699<math>\pm</math>.017</b>
yeast	.677 $\pm$ .019 $\bullet$	.636 $\pm$ .017 $\bullet$	.688 $\pm$ .017 $\bullet$	.610 $\pm$ .016 $\bullet$	.459 $\pm$ .032 $\bullet$	.712 $\pm$ .020	<b>.718<math>\pm</math>.019</b>
eurlex_dc	.412 $\pm$ .018 $\bullet$	.471 $\pm$ .012 $\bullet$	.232 $\pm$ .022 $\bullet$	.373 $\pm$ .015 $\bullet$	.346 $\pm$ .009 $\bullet$	.250 $\pm$ .031 $\bullet$	<b>.549<math>\pm</math>.025</b>
eurlex_sm	.419 $\pm$ .010 $\bullet$	.421 $\pm$ .010 $\bullet$	.273 $\pm$ .006 $\bullet$	.367 $\pm$ .009 $\bullet$	.402 $\pm$ .005 $\bullet$	.285 $\pm$ .009 $\bullet$	<b>.474<math>\pm</math>.017</b>
corel5k	.355 $\pm$ .035 $\bullet$	.307 $\pm$ .038 $\bullet$	.297 $\pm$ .023 $\bullet$	.330 $\pm$ .044 $\bullet$	<b>.397<math>\pm</math>.010</b>	.371 $\pm$ .028	.391 $\pm$ .037
corel16k	.405 $\pm$ .050	.350 $\pm$ .035 $\bullet$	.325 $\pm$ .022 $\bullet$	.248 $\pm$ .026 $\bullet$	.424 $\pm$ .006	.437 $\pm$ .044	<b>.449<math>\pm</math>.049</b>
bookmark	.383 $\pm$ .007 $\bullet$	.480 $\pm$ .010 $\bullet$	.267 $\pm$ .019 $\bullet$	.329 $\pm$ .016 $\bullet$	.534 $\pm$ .004 $\bullet$	.506 $\pm$ .014 $\bullet$	<b>.584<math>\pm</math>.013</b>
delicious	.487 $\pm$ .006 $\bullet$	.511 $\pm$ .004 $\bullet$	.457 $\pm$ .006 $\bullet$	.446 $\pm$ .010 $\bullet$	<b>.580<math>\pm</math>.002</b>	.570 $\pm$ .009	.572 $\pm$ .005

## 6.2 Comparison on Uniform Complementary Labels

**Setup.** Weight-decay is set as  $1e - 4$  and learning rate is selected from  $\{1e - 1, 1e - 2, 1e - 3\}$  for all data sets. We employ Adam [53] optimization method, and set the number of batch-size and epoch as 256 and 200 respectively. L-UW applies the same model and hyper-parameters as ours. Here, we estimate  $\mathbf{T}$  with a linear model. We use Ten-fold cross-validation to evaluate experiments, where training data is associated with complementary labels that are generated by randomly selecting one of possible labels excepting relevant labels (uniform complementary labels), and test data is equipped with the set of relevant labels.

The mean metrics value and *standard deviation* (std) will be reported as final experimental results for all approaches.

**Results.** Table 3 is utilized to report experimental results of various approaches on eight data sets equipped with uniform complementary labels.  $\uparrow / \downarrow$  indicates the larger/smaller the value, the better the performance.

According to reported results in Table 3, we can observe that results of MLCL are superior or comparable performance against baselines out of different data sets on five criteria. Our approach achieves the best performance in most cases. Specifically, the proposed approach outperforms LIFT on eight datasets across all metrics. This is because

TABLE 4

Experimental results (mean  $\pm$  std) on training data with biased complementary labels. The best performance of each dataset is presented in **boldface**, where  $\bullet/\circ$  represents whether MLCL is superior/inferior to baselines (with 5% t-test).

Methods	ML-KNN	LIFT	fpml	PML-lc	PML-LRS	L-UW	MLCL
Ranking loss↓							
scene	<b>.086±.015○</b>	.319±.025	.486±.027●	.492±.019●	.258±.013○	.368±.025●	.326±.050
yeast	.240±.014●	.297±.016●	.227±.013●	.248±.012●	.454±.024●	.202±.012	<b>.199±.012</b>
eurlex_dc	.668±.009●	.636±.021●	.537±.015●	.349±.028●	.326±.009	.586±.036●	<b>.308±.034</b>
eurlex_sm	.364±.020●	.392±.014●	.499±.019●	.447±.012●	.333±.009●	.641±.015●	<b>.316±.016</b>
corel5k	<b>.324±.038○</b>	.431±.030●	.474±.028●	.386±.047	.357±.012	.382±.033	.358±.039
corel16k	.413±.063●	.431±.041●	.454±.033●	.471±.068●	.375±.015	.373±.029	<b>.357±.040</b>
bookmark	.567±.007●	.449±.042●	.552±.018●	.491±.016●	.244±.003●	.326±.008●	<b>.211±.011</b>
delicious	.430±.005●	.413±.005●	.452±.008●	.433±.011●	<b>.314±.003○</b>	.349±.012○	.360±.008
One Error↓							
scene	<b>.228±.032○</b>	.669±.043●	.803±.038●	.720±.018●	.613±.017●	.696±.025●	.553±.054
yeast	.330±.032●	.280±.025●	.254±.028	.583±.027●	.546±.097●	.256±.025	<b>.254±.024</b>
eurlex_dc	.977±.005●	.959±.014●	.947±.008●	.774±.015●	.822±.004●	.822±.038●	<b>.695±.074</b>
eurlex_sm	.699±.016●	.753±.036●	.886±.024●	.664±.014	.737±.011●	.704±.012●	<b>.650±.045</b>
corel5k	<b>.738±.067</b>	.851±.038●	.861±.034●	.828±.059●	.747±.016	.792±.039●	.752±.037
corel16k	.780±.061●	.827±.049●	.837±.025●	.952±.021●	.730±.000	.731±.053	<b>.707±.063</b>
bookmark	.906±.007●	.804±.037●	.925±.008●	.792±.004●	.576±.003●	.635±.022●	<b>.502±.008</b>
delicious	.585±.012●	.557±.013●	.617±.025●	.681±.012●	<b>.434±.006○</b>	.485±.016●	.463±.017
Hamming loss ↓							
scene	<b>.088±.009○</b>	.819±.002●	.820±.002●	<b>.252±.006</b>	.814±.000●	.523±.048●	.290±.029
yeast	.697±.012●	.697±.013●	.697±.013●	.268±.010●	.316±.000●	.253±.017●	<b>.239±.008</b>
eurlex_dc	.932±.000●	.932±.000●	.118±.007●	<b>.104±.002</b>	.889±.039●	.799±.035●	.109±.011
eurlex_sm	.883±.001●	.883±.001●	.148±.005●	.139±.002	.825±.027●	.772±.009●	<b>.138±.007</b>
corel5k	<b>.114±.008○</b>	.887±.007●	.887±.007●	.157±.003●	.869±.002●	.498±.012●	.208±.033
corel16k	.882±.009●	.882±.009●	.882±.009●	<b>.178±.010</b>	.862±.001●	.481±.028●	.207±.086
bookmark	.917±.001●	.916±.001●	.419±.009●	<b>.122±.001○</b>	.813±.003●	.549±.046●	.146±.003
delicious	.711±.003●	.711±.003●	.711±.003●	.388±.013●	.459±.002●	.453±.015●	<b>.304±.005</b>
Coverage↓							
scene	<b>.086±.013○</b>	.280±.020	.420±.023●	.420±.016●	.229±.011○	.321±.021●	.286±.041
yeast	.551±.017●	.638±.028●	.533±.012●	<b>.493±.025</b>	.723±.040●	.500±.018	.498±.021
eurlex_dc	.626±.008●	.596±.019●	.504±.014●	.328±.026●	.306±.009●	.333±.018●	<b>.274±.030</b>
eurlex_sm	.432±.018●	.456±.014●	.579±.015●	.520±.015●	.418±.009●	.512±.009●	<b>.362±.016</b>
corel5k	<b>.419±.055</b>	.515±.024●	.555±.031●	.480±.041	.451±.013	.470±.036	.449±.038
corel16k	.498±.052●	.521±.038●	.542±.035●	.533±.066●	.454±.018	.468±.030	<b>.453±.039</b>
bookmark	.565±.006●	.455±.039●	.553±.017●	.492±.014●	.265±.003●	.308±.013●	<b>.231±.011</b>
delicious	.736±.004●	.723±.005●	.737±.008●	.691±.009	<b>.625±.003○</b>	.671±.012○	.688±.006
Average Precision ↑							
scene	<b>.860±.020○</b>	.559±.028●	.428±.026●	.462±.014●	.608±.013	.529±.020●	.618±.046
yeast	.670±.023●	.634±.016●	.691±.022●	.614±.015●	.500±.026●	.719±.020	<b>.726±.018</b>
eurlex_dc	.145±.005●	.166±.016●	.201±.009●	.371±.020●	.357±.005●	.266±.031●	<b>.456±.061</b>
eurlex_sm	.405±.013●	.373±.016●	.262±.016●	.366±.010●	.400±.007●	.282±.011●	<b>.482±.025</b>
corel5k	<b>.409±.040</b>	.300±.030●	.282±.017●	.325±.048●	.392±.017	.352±.032	.380±.037
corel16k	.355±.054●	.318±.033●	.301±.024●	.240±.030●	.393±.054	.384±.036	<b>.407±.047</b>
bookmark	.219±.004●	.320±.037●	.212±.007●	.320±.004●	.544±.003●	.469±.014●	<b>.599±.008</b>
delicious	.473±.006●	.490±.006●	.450±.008●	.449±.010●	<b>.581±.002○</b>	.544±.010	.544±.009

our approach is better at tackling the issue that training data is associated with relevant labels and irrelevant labels simultaneously than fully supervised MLL algorithms. Furthermore, experimental results of PML-lc and PML-LRS are inferior to ours in most cases, which demonstrates that PML approaches are indeed inferior to our approach in cases of dense noisy labels. Similarly, based on the results of L-UW shown in Table 3, we observe that our approach outperforms L-UW on almost all datasets and metrics other than *ranking loss* and *coverage* on the delicious dataset. This reflects that label correlations are important to solve ML-CLL problems, which leads to the proposed approach taking

label correlations into account surpasses L-UW that ignores label correlations.

### 6.3 Comparison on Biased Complementary Labels

**Setup.** To evaluate the effectiveness of our approach in different situations, we utilize training data with biased complementary labels that are generated via the co-occurrence rate of relevant labels. Specifically, we select a complementary label of an instance  $x$  from  $\mathcal{Y} \setminus Y$ , and the selecting rule follows: the class label with a lower co-occurrence rate has a higher probability to be selected as a complementary label. We adopt training data with biased complementary labels

TABLE 5

Ablation experimental results (mean  $\pm$  std) on training data with uniform complementary labels. The best performance is in **boldface**.

Methods	Uniform complementary labels				Biased complementary labels			
	scene	yeast	eurlex_dc	corel5k	scene	yeast	eurlex_dc	corel5k
	Hamming loss $\downarrow$							
MLCL	<b>.264<math>\pm</math>.027</b>	.235 $\pm$ .008	<b>.092<math>\pm</math>.005</b>	<b>.229<math>\pm</math>.068</b>	<b>.290<math>\pm</math>.029</b>	.239 $\pm$ .008	.109 $\pm$ .011	<b>.208<math>\pm</math>.033</b>
Without C	.290 $\pm$ .039	.421 $\pm$ .011	.109 $\pm$ .018	.466 $\pm$ .025	.294 $\pm$ .029	.409 $\pm$ .012	<b>.088<math>\pm</math>.004</b>	.444 $\pm$ .031
Without $\bar{L}_{mse}$	.510 $\pm$ .044	<b>.229<math>\pm</math>.007</b>	.509 $\pm$ .043	.461 $\pm$ .053	.481 $\pm$ .047	<b>.230<math>\pm</math>.009</b>	.512 $\pm$ .046	.489 $\pm$ .036
Ranking loss $\downarrow$								
MLCL	<b>.259<math>\pm</math>.030</b>	<b>.211<math>\pm</math>.013</b>	<b>.229<math>\pm</math>.026</b>	<b>.349<math>\pm</math>.035</b>	<b>.326<math>\pm</math>.050</b>	<b>.199<math>\pm</math>.012</b>	.308 $\pm$ .034	<b>.358<math>\pm</math>.039</b>
Without C	.282 $\pm$ .063	.419 $\pm$ .018	.277 $\pm$ .041	.487 $\pm$ .021	.348 $\pm$ .046	.406 $\pm$ .016	<b>.268<math>\pm</math>.024</b>	.467 $\pm$ .026
Without $\bar{L}_{mse}$	.379 $\pm$ .024	.216 $\pm$ .010	.303 $\pm$ .028	.362 $\pm$ .030	.353 $\pm$ .018	.204 $\pm$ .011	.320 $\pm$ .025	.387 $\pm$ .027
One error $\downarrow$								
MLCL	<b>.427<math>\pm</math>.018</b>	.251 $\pm$ .023	<b>.594<math>\pm</math>.035</b>	.736 $\pm$ .065	<b>.553<math>\pm</math>.054</b>	<b>.254<math>\pm</math>.024</b>	.695 $\pm$ .074	<b>.752<math>\pm</math>.037</b>
Without C	.474 $\pm$ .047	.633 $\pm$ .043	.708 $\pm$ .106	.866 $\pm$ .019	.560 $\pm$ .042	.612 $\pm$ .051	<b>.564<math>\pm</math>.029</b>	.855 $\pm$ .027
Without $\bar{L}_{mse}$	.607 $\pm$ .037	<b>.250<math>\pm</math>.025</b>	.740 $\pm$ .048	<b>.734<math>\pm</math>.058</b>	.686 $\pm$ .013	.256 $\pm$ .025	.753 $\pm$ .044	.773 $\pm$ .068
Coverage $\downarrow$								
MLCL	<b>.234<math>\pm</math>.025</b>	<b>.525<math>\pm</math>.021</b>	<b>.204<math>\pm</math>.023</b>	<b>.445<math>\pm</math>.048</b>	<b>.286<math>\pm</math>.041</b>	<b>.498<math>\pm</math>.021</b>	.274 $\pm$ .030	<b>.449<math>\pm</math>.038</b>
Without C	.255 $\pm$ .055	.683 $\pm$ .029	.247 $\pm$ .035	.565 $\pm$ .032	.306 $\pm$ .039	.660 $\pm$ .023	<b>.240<math>\pm</math>.023</b>	.547 $\pm$ .031
Without $\bar{L}_{mse}$	.334 $\pm$ .020	.527 $\pm$ .011	.249 $\pm$ .024	.451 $\pm$ .035	.310 $\pm$ .015	.501 $\pm$ .015	.265 $\pm$ .022	.473 $\pm$ .023
Average precision $\uparrow$								
MLCL	<b>.699<math>\pm</math>.017</b>	<b>.718<math>\pm</math>.019</b>	<b>.549<math>\pm</math>.025</b>	<b>.391<math>\pm</math>.037</b>	<b>.618<math>\pm</math>.046</b>	<b>.726<math>\pm</math>.018</b>	<b>.456<math>\pm</math>.061</b>	<b>.380<math>\pm</math>.037</b>
Without C	.671 $\pm$ .045	.472 $\pm$ .018	.469 $\pm$ .085	.274 $\pm$ .014	.611 $\pm$ .038	.489 $\pm$ .015	.447 $\pm$ .021	.289 $\pm$ .022
Without $\bar{L}_{mse}$	.566 $\pm$ .023	.711 $\pm$ .019	.426 $\pm$ .040	.389 $\pm$ .041	.541 $\pm$ .013	.717 $\pm$ .020	.411 $\pm$ .034	.359 $\pm$ .050

TABLE 6

Parameter sensitivity analysis on uniform complementary-label data, where metric is *average precision*. The best performance is in **boldface**.

$\beta$	scene	yeast	eurlex_dc	eurlex_sm	corel5k	corel16k	bookmark	delicious
0.1	.678 $\pm$ .017	.714 $\pm$ .019	.545 $\pm$ .019	.451 $\pm$ .025	.374 $\pm$ .033	.444 $\pm$ .046	.565 $\pm$ .007	.554 $\pm$ .005
0.3	.683 $\pm$ .015	.716 $\pm$ .018	.549 $\pm$ .021	.460 $\pm$ .021	.378 $\pm$ .032	.447 $\pm$ .047	.579 $\pm$ .011	.565 $\pm$ .005
0.5	.687 $\pm$ .016	.718 $\pm$ .018	.547 $\pm$ .022	.463 $\pm$ .016	.385 $\pm$ .031	.447 $\pm$ .048	.583 $\pm$ .008	<b>.575<math>\pm</math>.005</b>
0.8	.693 $\pm$ .016	.718 $\pm$ .018	.541 $\pm$ .022	.469 $\pm$ .018	.387 $\pm$ .037	.448 $\pm$ .048	.582 $\pm$ .007	.572 $\pm$ .006
1	<b>.699<math>\pm</math>.017</b>	<b>.718<math>\pm</math>.019</b>	<b>.549<math>\pm</math>.025</b>	<b>.474<math>\pm</math>.017</b>	<b>.391<math>\pm</math>.037</b>	<b>.449<math>\pm</math>.049</b>	<b>.584<math>\pm</math>.013</b>	.572 $\pm$ .005

to train the model, while test data is equipped with relevant label sets to evaluate the effectiveness of our approach. For other experimental settings, we apply same settings with Subsection 5.2.

**Results.** The mean and std of results on test data are shown in Table 4. According to results shown in Table 4, we can summarize the following impressive observations: (1) MLCL achieves superior or comparable performance to LIFT, fpml, PML-lc, PML-LRS and L-UW on different data sets, which proves that the proposed approach can predict the set of proper labels for unseen instances from complementary labeled data; (2) Although MLCL fails to achieve the best result on the scene dataset, our approach is better than other baselines in the rest of datasets, which indicates that our approach can effectively deal with ML-CLL problems than others. These observations demonstrate that the proposed method can both hold for the situation of data with uniform and biased complementary labels.

#### 6.4 Additional Experiments

**Ablation experiments.** We then explore the effect of different learning components on MLCL performance. Table 5 summarizes results of MLCL without the different components, which are trained on the data with uniform complementary labels. In Table 5, without C refers to MLCL directly using the estimated initial transition matrix S for

training, and without  $\bar{L}_{mse}$  indicates that MLCL only utilizes Eq. (9) for optimization.

From results reported in Table 5, the performance of MLCL surpasses that without different components in most cases, which shows that two components, including using label correlations to correct and an MSE-based regularizer, are beneficial for our approach to improve the performance. Especially, estimating T based on label correlations pushes the proposed approach performance forward significantly compared with that without C on most cases. Similarly, an MSE-based regularizer brings significant benefits for our approach, which demonstrates that an MSE-based regularizer balances the robustness and convergence rate of BCE loss. These indicate that using label correlations to estimate the transition matrix T and an MSE-based regularizer are effective strategies to alleviate ML-CLL problems.

**Trade-off parameter  $\beta$ .** Table 6 reports the performance of MLCL with varying  $\beta$  values that trade-off the complementary loss function  $\bar{L}$  and an MSE-based regularization  $\bar{L}_{mse}$ . Here, *average precision* is regarded as the criterion, and the training data is with uniform complementary labels.  $\beta$  is selected from the candidate value list {0.1, 0.3, 0.5, 0.8, 1}. We can observe the best results of most datasets is achieved at  $\beta = 1$  and the performance drops when  $\beta$  takes a smaller value. In general, a relatively large  $\beta$  ( $\beta \leq 1$ ) usually leads to better performance than a small value. Therefore, we set

TABLE 7

Experimental results (mean  $\pm$  std) of five criteria. “Fully supervised” is the linear model training with the fully supervised data (fully supervised MLL). “CL” denotes that each instance is associated with a **complementary label** sampled uniformly. “RL” presents that the model only using **relevant labels**  $\tilde{Y}$  to train. “CL & RL” uses the linear model with the loss function Eq. (12) to train, where each instance is equipped with a **complementary label** and a **relevant label**.

Datasets	scene	yeast	eurlex_dc	eurlex_sm	corel5k	corel16k	bookmark	delicious
Hamming loss $\downarrow$								
Fully supervised	.120 $\pm$ .013	.208 $\pm$ .009	.004 $\pm$ .000	.033 $\pm$ .001	.198 $\pm$ .012	.196 $\pm$ .012	.098 $\pm$ .004	.276 $\pm$ .006
CL	.264 $\pm$ .027	.235 $\pm$ .008	.092 $\pm$ .005	.139 $\pm$ .005	.229 $\pm$ .068	.202 $\pm$ .067	.140 $\pm$ .004	.289 $\pm$ .004
RL	.128 $\pm$ .012	.241 $\pm$ .011	.005 $\pm$ .001	.055 $\pm$ .005	.184 $\pm$ .014	.152 $\pm$ .010	.076 $\pm$ .002	.285 $\pm$ .004
CL & RL	.124 $\pm$ .008	.225 $\pm$ .010	.005 $\pm$ .001	.053 $\pm$ .002	.178 $\pm$ .012	.172 $\pm$ .010	.085 $\pm$ .002	.285 $\pm$ .004
Ranking loss $\downarrow$								
Fully supervised	.075 $\pm$ .009	.169 $\pm$ .009	.003 $\pm$ .001	.019 $\pm$ .001	.258 $\pm$ .029	.222 $\pm$ .029	.090 $\pm$ .005	.226 $\pm$ .004
CL	.259 $\pm$ .030	.211 $\pm$ .013	.229 $\pm$ .026	.312 $\pm$ .014	.349 $\pm$ .035	.289 $\pm$ .042	.252 $\pm$ .013	.310 $\pm$ .003
RL	.081 $\pm$ .009	.23 $\pm$ .013	.005 $\pm$ .001	.052 $\pm$ .027	.269 $\pm$ .031	.250 $\pm$ .025	.116 $\pm$ .006	.295 $\pm$ .008
CL & RL	.082 $\pm$ .011	.191 $\pm$ .011	.005 $\pm$ .001	.044 $\pm$ .002	.268 $\pm$ .031	.227 $\pm$ .021	.102 $\pm$ .004	.267 $\pm$ .004
One Error $\downarrow$								
Fully supervised	.222 $\pm$ .032	.223 $\pm$ .023	.019 $\pm$ .004	.069 $\pm$ .005	.627 $\pm$ .038	.588 $\pm$ .056	.313 $\pm$ .009	.340 $\pm$ .012
CL	.427 $\pm$ .018	.251 $\pm$ .023	.594 $\pm$ .035	.656 $\pm$ .029	.736 $\pm$ .065	.690 $\pm$ .056	.509 $\pm$ .012	.448 $\pm$ .016
RL	.231 $\pm$ .032	.280 $\pm$ .024	.022 $\pm$ .005	.106 $\pm$ .02	.641 $\pm$ .042	.616 $\pm$ .042	.338 $\pm$ .007	.438 $\pm$ .019
CL & RL	.229 $\pm$ .033	.255 $\pm$ .032	.022 $\pm$ .005	.098 $\pm$ .007	.639 $\pm$ .040	.600 $\pm$ .044	.324 $\pm$ .007	.398 $\pm$ .017
Coverage $\downarrow$								
Fully supervised	.077 $\pm$ .009	.451 $\pm$ .019	.004 $\pm$ .000	.074 $\pm$ .002	.347 $\pm$ .044	.315 $\pm$ .024	.112 $\pm$ .005	.527 $\pm$ .007
CL	.234 $\pm$ .025	.525 $\pm$ .021	.204 $\pm$ .023	.365 $\pm$ .014	.445 $\pm$ .048	.393 $\pm$ .042	.279 $\pm$ .011	.632 $\pm$ .007
RL	.083 $\pm$ .008	.561 $\pm$ .024	.007 $\pm$ .002	.123 $\pm$ .035	.364 $\pm$ .048	.353 $\pm$ .025	.142 $\pm$ .006	.582 $\pm$ .008
CL & RL	.084 $\pm$ .010	.474 $\pm$ .021	.006 $\pm$ .001	.113 $\pm$ .004	.363 $\pm$ .048	.326 $\pm$ .020	.125 $\pm$ .004	.564 $\pm$ .006
Average Precision $\uparrow$								
Fully supervised	.868 $\pm$ .018	.760 $\pm$ .015	.988 $\pm$ .003	.943 $\pm$ .004	.494 $\pm$ .024	.530 $\pm$ .038	.766 $\pm$ .007	.662 $\pm$ .005
CL	.699 $\pm$ .017	.718 $\pm$ .019	.549 $\pm$ .025	.474 $\pm$ .017	.391 $\pm$ .037	.449 $\pm$ .049	.584 $\pm$ .013	.572 $\pm$ .005
RL	.860 $\pm$ .018	.704 $\pm$ .013	.984 $\pm$ .003	.888 $\pm$ .032	.484 $\pm$ .028	.505 $\pm$ .028	.737 $\pm$ .006	.590 $\pm$ .009
CL & RL	.860 $\pm$ .019	.734 $\pm$ .018	.985 $\pm$ .004	.899 $\pm$ .004	.485 $\pm$ .028	.523 $\pm$ .030	.753 $\pm$ .006	.618 $\pm$ .005

$\beta = 1$  for MLCL.

## 6.5 Combination of Complementary Labels and Relevant Labels

**Setup.** Finally, we demonstrate the effectiveness of combining relevant labeled data and complementary labeled one. The training data is associated with uniform complementary labels and relevant labels simultaneously. More specifically, an instance  $x$  is associated with a complementary label  $\bar{y}$  and relevant labels  $\tilde{Y}$ , where  $\bar{y}$  is uniformly selected and  $\tilde{Y}$  is randomly selected from the relevant label set  $Y$  of  $x$  (i.e.,  $\tilde{Y} \subseteq Y$ ). Here, we set  $|\tilde{Y}| = 1$  that means each instance only associated with a complementary label and a relevant label. The other experimental settings are the same with Subsection 5.2.

**Results.** We compare three methods: (1) the “Fully supervised” method uses the linear model to train with the fully supervised data, which is fully supervised MLL; (2) the “CL” method refers to MLCL training with the uniform complementary-label data; (3) the “RL” method uses the linear model to train with relevant labels  $\tilde{Y}$ ; (4) the combination (“CL & RL”) method adopts the linear model with the loss function Eq. (12) to train, where the training data is equipped with a combination of complementary labels and relevant labels. Table 7 reports the experimental results on five criteria. We can see that the performance of “CL& RL” method is much superior to “CL” method on all datasets over *hamming loss*, *ranking loss*, *one error*, *coverage* and *average precision*, such as “CL& RL” method

outperforms “CL” method by a large margin over *average precision* (+0.436 on eurlex\_dc and +0.425 on eurlex\_sm). This demonstrates that the ML-CLL is easily applied to fully supervised MLL scenarios, MLL with missing labels [54], [55] or other MLL scenarios. Moreover, the results in “CL & RL” achieve comparable performance to “RL” method, which illustrates that using complementary labels improves performance and the information encompassed by relevant labels does not overshadow or encompass the information from complementary labels. Finally, the results of “CL & RL” method close to “Fully supervised” method, which illustrates that ML-CLL can get excellent results with just a small amount of additional information. This is useful for application in the real world because ML-CLL can obtain good performance through less expensive labeled data.

## 7 CONCLUSION

In this paper, we theoretically analyze the reason why the estimated transition matrix in multi-class CLL is distorted in ML-CLL. To alleviate the difficulty in directly calculating the transition matrix from complex label correlations under multi-labeled data, we propose a two-step method to estimate the transition matrix  $T$  in ML-CLL, which uses label correlations to correct an initial transition matrix. Furthermore, we theoretically show that the proposed approach is classifier-consistent. Additionally, due to the robustness of MSE loss, an MSE-based regularizer is introduced to alleviate the tendency of the fast convergent BCE loss overfitting to noises. Finally, we show that our proposed ML-

CLL can be easily combined with relevant labels and the proposed method can achieve performance comparable to fully supervised MLL with just a small amount of additional information.

## ACKNOWLEDGMENTS

The authors wish to thank the anonymous reviewers for their helpful comments and suggestions. This work was supported by the National Science Foundation of China (62176055). MX is supported by the Australian Research Council (DE230101116). We thank the Big Data Center of Southeast University for providing the facility support on the numerical calculations in this paper.

## REFERENCES

- [1] M.-L. Zhang and Z.-H. Zhou, "A review on multi-label learning algorithms," *IEEE Trans. Knowl. Data Eng.*, vol. 26, no. 8, pp. 1819–1837, 2014.
- [2] M.-L. Zhang and L. Wu, "Lift: Multi-label learning with label-specific features," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 37, no. 1, pp. 107–120, 2015.
- [3] W. Tang, W. Zhang, and M.-L. Zhang, "Multi-instance partial-label learning: Towards exploiting dual inexact supervision," *Science China Information Sciences*, vol. 67, no. 3, pp. 1–14, 2024.
- [4] T. N. Rubin, A. Chambers, P. Smyth, and M. Steyvers, "Statistical topic models for multi-label document classification," *Mach. Learn.*, vol. 88, no. 1-2, pp. 157–208, 2012.
- [5] P.-J. Tang, M. Jiang, B. N. Xia, J. W. Pitera, J. Welser, and N. V. Chawla, "Multi-label patent categorization with non-local attention-based graph convolutional network," in *Proceedings of the 34th Conference on Artificial Intelligence*, York, NY, 2020, pp. 9024–9031.
- [6] A. Lambrecht and C. Tucker, "When does retargeting work? information specificity in online advertising," *Journal of Marketing research*, vol. 50, no. 5, pp. 561–576, 2013.
- [7] Y. Yang, J. Guo, G. Li, L. Li, W. Li, and J. Yang, "Alignment efficient image-sentence retrieval considering transferable cross-modal representation learning," *Frontiers of Computer Science*, vol. 18, no. 1, p. 181335, 2024.
- [8] W. Zhang, F. Liu, L. Luo, and J. Zhang, "Predicting drug side effects by multi-label learning and ensemble learning," *BMC bioinformatics*, vol. 16, pp. 1–11, 2015.
- [9] Y. Qiu, F. Lin, W. Chen, and M. Xu, "Pre-training in medical data: A survey," *Machine Intelligence Research*, vol. 20, no. 2, pp. 147–179, 2023.
- [10] Y. Gao, M. Xu, and M.-L. Zhang, "Unbiased risk estimator to multi-labeled complementary label learning," in *Proceedings of the 32nd International Joint Conference on Artificial Intelligence*, Macao, China, 2023, pp. 3732–3740.
- [11] T. Ishida, G. Niu, W.-H. Hu, and M. Sugiyama, "Learning from complementary labels," in *Advances in Neural Information Processing Systems 30*, Long Beach, CA, 2017, pp. 5639–5649.
- [12] T. Ishida, G. Niu, A. K. Menon, and M. Sugiyama, "Complementary-label learning for arbitrary losses and models," in *Proceedings of the 36th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, Long Beach, CA, 2019, pp. 2971–2980.
- [13] X.-Y. Yu, T.-L. Liu, M.-M. Gong, and D.-C. Tao, "Learning with biased complementary labels," in *Proceedings of the 15th European Conference on Computer Vision*, Munich, Germany, 2018, pp. 69–85.
- [14] Y.-T. Chou, G. Niu, H.-T. Lin, and M. Sugiyama, "Unbiased risk estimators can mislead: A case study of learning with complementary labels," in *Proceedings of the 37th International Conference on Machine Learning*, Virtual Event, 2020, pp. 1929–1938.
- [15] Y. Gao and M.-L. Zhang, "Discriminative complementary-label learning with weighted loss," in *Proceedings of the 38th International Conference on Machine Learning*, Virtual Event, 2021, pp. 3587–3597.
- [16] D.-B. Wang, L. Feng, and M.-L. Zhang, "Learning from complementary labels via partial-output consistency regularization," in *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence*, Virtual Event, 2021, pp. 3075–3081.
- [17] L. Feng, T. Kaneko, B. Han, G. Niu, B. An, and M. Sugiyama, "Learning with multiple complementary labels," in *Proceedings of the 37th International Conference on Machine Learning*, Virtual Event, 2020, pp. 3072–3081.
- [18] A. Ghosh, H. Kumar, and P. S. Sastry, "Robust loss functions under label noise for deep neural networks," in *Proceedings of the 31st AAAI Conference on Artificial Intelligence*, San Francisco, CA, 2017, pp. 1919–1925.
- [19] T. S. Sindlinger, "Crowdsourcing: why the power of the crowd is driving the future of business," 2010.
- [20] M.-L. Zhang and Z.-H. Zhou, "A review on multi-label learning algorithms," *IEEE Trans. Knowl. Data Eng.*, vol. 26, no. 8, pp. 1819–1837, 2014.
- [21] F. Wu, Z.-H. Wang, Z.-F. Zhang, Y. Yang, J.-B. Luo, W.-W. Zhu, and Y.-T. Zhuang, "Weakly semi-supervised deep learning for multi-label image annotation," *IEEE Trans. Big Data*, vol. 1, no. 3, pp. 109–122, 2015.
- [22] S. S. Bucak, R. Jin, and A. K. Jain, "Multi-label learning with incomplete class assignments," in *The 24th IEEE Conference on Computer Vision and Pattern Recognition*, Colorado Springs, CO, 2011, pp. 2801–2808.
- [23] W.-W. Liu, I. W. Tsang, and K. Müller, "An easy-to-hard learning paradigm for multiple classes and multiple labels," *J. Mach. Learn. Res.*, vol. 18, pp. 94:1–94:38, 2017.
- [24] M.-L. Zhang, Y.-K. Li, X.-Y. Liu, and X. Geng, "Binary relevance for multi-label learning: an overview," *Frontiers Comput. Sci.*, vol. 12, no. 2, pp. 191–202, 2018.
- [25] M. R. Boutell, J.-B. Luo, X.-P. Shen, and C. M. Brown, "Learning multi-label scene classification," *Pattern Recognit.*, vol. 37, no. 9, pp. 1757–1771, 2004.
- [26] M.-L. Zhang and Z.-H. Zhou, "ML-KNN: A lazy learning approach to multi-label learning," *Pattern Recognit.*, vol. 40, no. 7, pp. 2038–2048, 2007.
- [27] A. Elisseeff and J. Weston, "A kernel method for multi-labelled classification," in *Advances in Neural Information Processing Systems 14*, Vancouver, Canada, 2001, pp. 681–687.
- [28] J. Fürnkranz, E. Hüllermeier, E. Loza Mencía, and K. Brinker, "Multilabel classification via calibrated label ranking," *Machine learning*, vol. 73, no. 2, pp. 133–153, 2008.
- [29] J. Read, B. Pfahringer, G. Holmes, and E. Frank, "Classifier chains for multi-label classification," *Mach. Learn.*, vol. 85, no. 3, pp. 333–359, 2011.
- [30] G. Tsoumakas, I. Katakis, and I. P. Vlahavas, "Random k-labelsets for multilabel classification," *IEEE Trans. Knowl. Data Eng.*, vol. 23, no. 7, pp. 1079–1089, 2011.
- [31] M.-L. Zhang and Z.-H. Zhou, "Multilabel neural networks with applications to functional genomics and text categorization," *IEEE transactions on Knowledge and Data Engineering*, vol. 18, no. 10, pp. 1338–1351, 2006.
- [32] Y.-C. Li, Y. Song, and J.-B. Luo, "Improving pairwise ranking for multi-label image classification," in *Proceedings of 2017 IEEE conference on computer vision and pattern recognition*, Honolulu, HI, 2017, pp. 3617–3625.
- [33] S.-W. Ji, L. Tang, S.-P. Yu, and J.-P. Ye, "A shared-subspace learning framework for multi-label classification," *ACM Trans. Knowl. Discov. Data*, vol. 4, no. 2, pp. 8:1–8:29, 2010.
- [34] W. Gerych, T. Hartvigsen, L. Buquicchio, E. Agu, and E. A. Rundensteiner, "Recurrent bayesian classifier chains for exact multi-label classification," in *Advances in Neural Information Processing Systems 34*, virtual event, 2021, pp. 15 981–15 992.
- [35] W.-T. Zhao, S.-F. Kong, J.-W. Bai, D. Fink, and C. P. Gomes, "HOT-VAE: learning high-order label correlation for multi-label classification via attention-based variational autoencoders," in *Proceedings of 35th AAAI Conference on Artificial Intelligence*, Virtual Event, 2021, pp. 15 016–15 024.
- [36] L.-C. Wang, Z.-M. Ding, S.-J. Han, J.-J. Han, C. Choi, and Y. Fu, "Generative correlation discovery network for multi-label learning," in *Proceedings of 2019 IEEE International Conference on Data Mining*, Beijing, China, 2019, pp. 588–597.
- [37] G.-X. Xun, K. Jha, J.-H. Sun, and A.-D. Zhang, "Correlation networks for extreme multi-label text classification," in *Proceedings of 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, Virtual Event, 2020, pp. 1074–1082.
- [38] L. Sun, S. Feng, J. Liu, G. Lyu, and C. Lang, "Global-local label correlation for partial multi-label learning," *IEEE Transactions on Multimedia*, vol. PP, no. 99, pp. 1–1, 2021.

- [39] Z.-H. Zhou, "A brief introduction to weakly supervised learning," *National science review*, vol. 5, no. 1, pp. 44–53, 2018.
- [40] M.-K. Xie and S.-J. Huang, "Partial multi-label learning," in *Proceedings of the 32nd AAAI Conference on Artificial Intelligence*, New Orleans, LA, 2018, pp. 4302–4309.
- [41] S.-J. H. Ming-Kun Xie, "Partial multi-label learning with noisy label identification," in *Proceedings of 34th AAAI Conference on Artificial Intelligence*, New York, NY, 2020, pp. 6454–6461.
- [42] L.-J. Sun, S.-H. Feng, T. Wang, C.-Y. Lang, and Y. Jin, "Partial multi-label learning by low-rank and sparse decomposition," in *Proceedings of the 33rd AAAI Conference on Artificial Intelligence*, Honolulu, HI, 2019, pp. 5016–5023.
- [43] G.-X. Yu, X. Chen, C. Domeniconi, J. Wang, Z. Li, Z.-L. Zhang, and X.-D. Wu, "Feature-induced partial multi-label learning," in *Proceedings of 2018 IEEE International Conference on Data Mining*, Singapore, 2018, pp. 1398–1403.
- [44] Y.-W. Xu, M.-M. Gong, J.-X. Chen, T.-L. Liu, K. Zhang, and K. Battanghelich, "Generative-discriminative complementary learning," in *Proceedings of the 34th AAAI Conference on Artificial Intelligence*, New York, NY, 2020, pp. 6526–6533.
- [45] S. Diplaris, G. Tsoumakas, P. A. Mitkas, and I. P. Vlahavas, "Protein classification with multiple algorithms," in *Advances in 10th Panhellenic Conference on Informatics*, vol. 3746, Volos, Greece, 2005, pp. 448–456.
- [46] G. Patrini, A. Rozza, A. K. Menon, R. Nock, and L.-Z. Qu, "Making deep neural networks robust to label noise: A loss correction approach," in *Proceedings of 2017 IEEE Conference on Computer Vision and Pattern Recognition*, Honolulu, HI, 2017, pp. 2233–2241.
- [47] X.-B. Xia, T.-L. Liu, N.-N. Wang, B. Han, C. Gong, G. Niu, and M. Sugiyama, "Are anchor points really indispensable in label-noise learning?" in *Advances in Neural Information Processing Systems* 32, Vancouver, Canada, 2019, pp. 6835–6846.
- [48] J.-Q. Lv, M. Xu, L. Feng, G. Niu, X. Geng, and M. Sugiyama, "Progressive identification of true labels for partial-label learning," in *Proceedings of the 37th International Conference on Machine Learning*, Virtual Event, 2020, pp. 6500–6510.
- [49] Z.-L. Zhang and M. R. Sabuncu, "Generalized cross entropy loss for training deep neural networks with noisy labels," in *Advances in Neural Information Processing Systems* 31, Montréal, Canada, 2018, pp. 8792–8802.
- [50] Y. Katsura and M. Uchida, "Bridging ordinary-label learning and complementary-label learning," in *Proceedings of the 12th Asian Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, Bangkok, Thailand, 2020, pp. 161–176.
- [51] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z.-M. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J.-J. Bai, and S. Chintala, "Pytorch: An imperative style, high-performance deep learning library," in *Advances in Neural Information Processing Systems* 32, Vancouver, Canada, 2019, pp. 8024–8035.
- [52] M.-L. Zhang and L. Wu, "Lift: Multi-label learning with label-specific features," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 37, no. 1, pp. 107–120, 2015.
- [53] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proceedings of the 3rd International Conference on Learning Representations*, San Diego, CA, 2015.
- [54] L. Feng, J. Huang, S.-L. Shu, and B. An, "Regularized matrix factorization for multilabel learning with missing labels," *IEEE Trans. Cybern.*, vol. 52, no. 5, pp. 3710–3721, 2022.
- [55] C.-X. Wang, Y.-J. Lin, and J.-H. Liu, "Feature selection for multi-label learning with missing labels," *Appl. Intell.*, vol. 49, no. 8, pp. 3027–3042, 2019.
- [56] M. D. Reid and R. C. Williamson, "Composite binary losses," *J. Mach. Learn. Res.*, vol. 11, pp. 2387–2422, 2010.



**Yi Gao** received the BSc and MSc degrees in computer science from Northwest University, China, in 2017 and 2020 respectively. She is currently pursuing the PhD degree in Southeast University, China. Her main research interests include machine learning and data mining, with a focus on learning from complementary labels.



**Miao Xu** is a senior lecturer in the School of Electrical Engineering and Computer Science at the University of Queensland, Australia. She was awarded the Australian Research Council Discovery Early Career Researcher Award (DECRA) in 2023. Dr Xu specializes in machine learning and data mining, particularly focusing on the challenges of learning from imperfect information. Dr Xu earned a PhD from Nanjing University, where research efforts led to notable recognitions including the CAAI Outstanding Doctoral Dissertation Award.



**Min-Ling Zhang** received the BSc, MSc, and PhD degrees in computer science from Nanjing University, China, in 2001, 2004 and 2007, respectively. Currently, he is a Professor at the School of Computer Science and Engineering, Southeast University, China. His main research interests include machine learning and data mining. In recent years, Dr. Zhang has served as the General Co-Chairs of ACML'18, Program Co-Chairs of PAKDD'19, CCF-ICAI'19, ACML'17, CCF-FAI'17, PRICAL'16, Senior PC member or Area Chair of KDD 2021-2022, AAAI 2017-2020, IJCAI 2017-2023, ICDM 2015-2022, etc. He is also on the editorial board of IEEE Transactions on Pattern Analysis and Machine Intelligence, ACM Transactions on Intelligent Systems and Technology, Science China Information Sciences, Frontiers of Computer Science, Machine Intelligence Research, etc. Dr. Zhang is the Steering Committee Member of ACML and PAKDD, Vice-Chair of the CAAI Machine Learning Society, standing committee member of the CCF Artificial Intelligence & Pattern Recognition Society. He is a Distinguished Member of CCF, CAAI, and Senior Member of AAAI, ACM, IEEE.

## APPENDIX A

### THE PROOF OF THEOREM 1

**Theorem 1.** Given an instance  $\mathbf{x}$ , suppose  $Y$  is the relevant label set and the label  $l_j$  is the complementary label which is randomly selected. Then the following equality holds:

$$p(\bar{y}^j = 1|\mathbf{x}) = \sum_{C \in \mathcal{Y}', l_j \notin C} p(\bar{y}^j = 1|Y = C)p(Y = C|\mathbf{x}) \geq \sum_{k=1, k \neq j}^K p(\bar{y}^j = 1|y^k = 1)p(y^k = 1|\mathbf{x}).$$

*Proof.* Firstly, we should introduce the addition rule of probability:  $p(AB) = p(A) + p(B) - p(A \cup B)$ , so we have  $p(AB) \geq p(A) + p(B)$ . We start to prove the above inequality. According to the assumption:  $p(\bar{y}|Y) = p(\bar{y}|Y, \mathbf{x})$ , we have

$$\begin{aligned} p(\bar{y}^j = 1|\mathbf{x}) &= \sum_{C \in \mathcal{Y}', l_j \notin C} p(\bar{y}^j = 1|Y = C)p(Y = C|\mathbf{x}) \\ &= \sum_{C \in \mathcal{Y}', l_j \notin C} p(\bar{y}^j = 1|Y = C, \mathbf{x})p(Y = C|\mathbf{x}) \\ &= \sum_{C \in \mathcal{Y}', l_j \notin C} p(\bar{y}^j = 1, Y = C|\mathbf{x}) \\ &= \sum_{C \in \mathcal{Y}', l_j \notin C} p(Y = C|\bar{y}^j = 1, \mathbf{x})p(\bar{y}^j = 1|\mathbf{x}). \end{aligned}$$

According to addition rule of probability, so we have

$$\begin{aligned} p(\bar{y}^j = 1|\mathbf{x}) &\geq \sum_{C \in \mathcal{Y}', l_j \notin C} \left[ \sum_{k=1, k \neq j, l_k \in C}^K p(y^k = 1|\bar{y}^j = 1, \mathbf{x}) + \sum_{k=1, l_k \notin C}^K p(y^k = 0|\bar{y}^j = 1, \mathbf{x}) \right] p(\bar{y}^j = 1|\mathbf{x}) \\ &\geq \sum_{C \in \mathcal{Y}', l_j \notin C} \sum_{k=1, k \neq j, l_k \in C}^K p(y^k = 1|\bar{y}^j = 1, \mathbf{x})p(\bar{y}^j = 1|\mathbf{x}) \quad \because \sum_{k=1, l_k \notin C}^K p(y^k = 0|\bar{y}^j = 1, \mathbf{x}) \geq 0 \\ &= \sum_{C \in \mathcal{Y}', l_j \notin C} \sum_{k=1, k \neq j, l_k \in C}^K p(\bar{y}^j = 1|y^k = 1, \mathbf{x})p(y^k = 1|\mathbf{x}) \\ &= \sum_{C \in \mathcal{Y}', l_j \notin C} \sum_{k=1, k \neq j}^K p(\bar{y}^j = 1|y^k = 1, \mathbf{x})p(y^k = 1|\mathbf{x}) \quad \because p(y^k = 1|\mathbf{x}) = 0 \text{ if } l_k \notin Y \\ &= \sum_{k=1, k \neq j}^K \sum_{C \in \mathcal{Y}', l_j \notin C} p(\bar{y}^j = 1|y^k = 1, \mathbf{x})p(y^k = 1|\mathbf{x}) \\ &= \sum_{k=1, k \neq j}^K (2^{K-1} - 1)p(\bar{y}^j = 1|y^k = 1, \mathbf{x})p(y^k = 1|\mathbf{x}) \\ &\geq \sum_{k=1, k \neq j}^K p(\bar{y}^j = 1|y^k = 1, \mathbf{x})p(y^k = 1|\mathbf{x}) \\ &= \sum_{k=1, k \neq j}^K p(\bar{y}^j = 1|y^k = 1)p(y^k = 1|\mathbf{x}). \end{aligned}$$

□

## APPENDIX B

### THE PROOF OF THEOREM 3

**Theorem 3.** Under a MLL scenario: suppose the labels  $l_{z_1}, l_{z_2} \in \mathcal{Y}$  ( $z_1, z_2 \in [K], z_1 \neq z_2$ ) are dependent, and the labels belonging to  $\mathcal{Y} \setminus \{l_{z_1}, l_{z_2}\}$  are mutually exclusive. For any  $\mathbf{x} \in \mathcal{X}$ , its label set  $Y \subseteq \{l_{z_1}, l_{z_2}\}$  and  $Y \neq \emptyset$ . Let the label  $l_j$  ( $j \in [K], j \neq z_1, z_2$ ) be the complementary label of  $\mathbf{x}$ .  $\mathbf{T}_{z_1j}$  and  $\mathbf{T}_{z_2j}$  calculated from label correlations satisfy

$$\begin{aligned} \mathbf{T}_{z_1j} &= \frac{p(\bar{y}^j = 1|\mathbf{x})}{p(y^{z_2} = 1|\bar{y}^j = 1, y^{z_1} = 1, \mathbf{x})p(y^{z_1} = 1|\mathbf{x})}, \\ \mathbf{T}_{z_2j} &= \frac{p(\bar{y}^j = 1|\mathbf{x})}{p(y^{z_1} = 1|\bar{y}^j = 1, y^{z_2} = 1, \mathbf{x})p(y^{z_2} = 1|\mathbf{x})}, \end{aligned}$$

where  $[K]$  denotes the integer set  $\{1, 2, \dots, K\}$ . The difference of  $\mathbf{T}$  and  $\mathbf{Q}$  on the complementary label  $l_j$  is

$$\ell_j \geq 2\left(\frac{1}{\xi^2} - 1\right)p(\bar{y}^j = 1|\mathbf{x}),$$

where  $\xi = \max\{p(y^{z_2} = 1|\bar{y}^j = 1, y^{z_1} = 1, \mathbf{x}), p(y^{z_1} = 1|\bar{y}^j = 1, y^{z_2} = 1, \mathbf{x})\}$ .

*Proof.* We start calculating the difference  $\ell_j$  from estimating the transition probabilities  $\mathbf{T}_{z_1j}$  and  $\mathbf{T}_{z_2j}$ . According to Definition 2 and the description of Theorem 3, we have

$$\begin{aligned} p(\bar{y}^j = 1|\mathbf{x}) &= \sum_{k=1, k \neq j, z_1, z_2}^K p(\bar{y}^j = 1|y^k = 1, \mathbf{x})p(y^k = 1|\mathbf{x}) + p(\bar{y}^j = 1|y^{z_1} = 1, y^{z_2} = 1, \mathbf{x})p(y^{z_1} = 1, y^{z_2} = 1|\mathbf{x}) \\ &\quad + p(\bar{y}^j = 1|y^{z_1} = 1, y^{z_2} = 0, \mathbf{x})p(y^{z_1} = 1, y^{z_2} = 0|\mathbf{x}) + p(\bar{y}^j = 1|y^{z_1} = 0, y^{z_2} = 1, \mathbf{x})p(y^{z_1} = 0, y^{z_2} = 1|\mathbf{x}) \\ &\quad + p(\bar{y}^j = 1|y^{z_1} = 0, y^{z_2} = 0, \mathbf{x})p(y^{z_1} = 0, y^{z_2} = 0|\mathbf{x}) \\ &= \sum_{k=1, k \neq j, z_1, z_2}^K p(\bar{y}^j = 1|y^k = 1, \mathbf{x})p(y^k = 1|\mathbf{x}) + p(y^{z_2} = 1|\bar{y}^j = 1, y^{z_1} = 1, \mathbf{x})p(\bar{y}^j = 1|y^{z_1} = 1, \mathbf{x})p(y^{z_1} = 1|\mathbf{x}) \\ &\quad + p(y^{z_1} = 1|\bar{y}^j = 1, y^{z_2} = 0, \mathbf{x})p(\bar{y}^j = 1|y^{z_2} = 0, \mathbf{x})p(y^{z_2} = 0|\mathbf{x}) \\ &\quad + p(y^{z_2} = 1|\bar{y}^j = 1, y^{z_1} = 0, \mathbf{x})p(\bar{y}^j = 1|y^{z_1} = 0, \mathbf{x})p(y^{z_1} = 0|\mathbf{x}) \\ &\quad + p(y^{z_2} = 0|\bar{y}^j = 1, y^{z_1} = 0, \mathbf{x})p(\bar{y}^j = 1|y^{z_1} = 0, \mathbf{x})p(y^{z_1} = 0|\mathbf{x}). \end{aligned}$$

Based on the assumption that  $\bar{y}$  and  $\mathbf{x}$  are conditionally independent given  $Y$ , then we can have

$$\begin{aligned} p(\bar{y}^j = 1|\mathbf{x}) &= \sum_{k=1, k \neq j, z_1, z_2}^K p(\bar{y}^j = 1|y^k = 1)p(y^k = 1|\mathbf{x}) + p(y^{z_2} = 1|\bar{y}^j = 1, y^{z_1} = 1, \mathbf{x})p(\bar{y}^j = 1|y^{z_1} = 1)p(y^{z_1} = 1|\mathbf{x}) \\ &\quad + p(y^{z_1} = 1|\bar{y}^j = 1, y^{z_2} = 0, \mathbf{x})p(\bar{y}^j = 1|y^{z_2} = 0)p(y^{z_2} = 0|\mathbf{x}) \\ &\quad + p(y^{z_2} = 1|\bar{y}^j = 1, y^{z_1} = 0, \mathbf{x})p(\bar{y}^j = 1|y^{z_1} = 0)p(y^{z_1} = 0|\mathbf{x}) \\ &\quad + p(y^{z_2} = 0|\bar{y}^j = 1, y^{z_1} = 0, \mathbf{x})p(\bar{y}^j = 1|y^{z_1} = 0)p(y^{z_1} = 0|\mathbf{x}). \end{aligned}$$

Since  $p(\bar{y}^j = 1|y^{z_1} = 0)$  and  $p(\bar{y}^j = 1|y^{z_2} = 0)$  do not hold according to the definition of the transition matrix, we can obtain

$$\begin{aligned} p(\bar{y}^j = 1|\mathbf{x}) &= \sum_{k=1, k \neq j, z_1, z_2}^K p(\bar{y}^j = 1|y^k = 1)p(y^k = 1|\mathbf{x}) + p(y^{z_2} = 1|\bar{y}^j = 1, y^{z_1} = 1, \mathbf{x})p(\bar{y}^j = 1|y^{z_1} = 1)p(y^{z_1} = 1|\mathbf{x}) \\ &= p(y^{z_2} = 1|\bar{y}^j = 1, y^{z_1} = 1, \mathbf{x})p(\bar{y}^j = 1|y^{z_1} = 1)p(y^{z_1} = 1|\mathbf{x}) \quad \text{because } p(y^k = 1|\mathbf{x}) = 0 \text{ if } l_k \notin Y \\ \Rightarrow \mathbf{T}_{z_1j} &= p(\bar{y}^j = 1|y^{z_1} = 1) = \frac{p(\bar{y}^j = 1|\mathbf{x})}{p(y^{z_2} = 1|\bar{y}^j = 1, y^{z_1} = 1, \mathbf{x})p(y^{z_1} = 1|\mathbf{x})}. \end{aligned}$$

Similarly, we can get

$$\mathbf{T}_{z_2j} = p(\bar{y}^j = 1|y^{z_2} = 1) = \frac{p(\bar{y}^j = 1|\mathbf{x})}{p(y^{z_1} = 1|\bar{y}^j = 1, y^{z_2} = 1, \mathbf{x})p(y^{z_2} = 1|\mathbf{x})}.$$

Next, we calculate the difference  $\ell_j$ . The rest of the elements of  $\mathbf{T}_{\cdot j}$  are same as those estimated by multi-class CLL. According the definition of  $\ell_j$ , we have

$$\begin{aligned} \ell_j &= \sum_{k=1}^K |\mathbf{T}_{kj} - \mathbf{Q}_{kj}| \\ &= |\mathbf{T}_{z_1j} + \mathbf{T}_{z_2j} - 2p(\bar{y}^j = 1|\mathbf{x})| \\ &= \left| \frac{p(\bar{y}^j = 1|\mathbf{x})}{p(y^{z_2} = 1|\bar{y}^j = 1, y^{z_1} = 1, \mathbf{x})p(y^{z_1} = 1|\mathbf{x})} + \frac{p(\bar{y}^j = 1|\mathbf{x})}{p(y^{z_1} = 1|\bar{y}^j = 1, y^{z_2} = 1, \mathbf{x})p(y^{z_2} = 1|\mathbf{x})} - 2p(\bar{y}^j = 1|\mathbf{x}) \right| \\ &\geq \left| 2\left(\frac{1}{\xi^2} - 1\right)p(\bar{y}^j = 1|\mathbf{x}) \right| \\ &= 2\left(\frac{1}{\xi^2} - 1\right)p(\bar{y}^j = 1|\mathbf{x}). \quad \because \frac{1}{\xi^2} \geq 1 \end{aligned}$$

Because  $0 \leq p(y^{z_1} = 1|\mathbf{x}) \leq p(y^{z_1} = 1|\bar{y}^j = 1, y^{z_2} = 1, \mathbf{x}) \leq 1$  and  $0 \leq p(y^{z_2} = 1|\mathbf{x}) \leq p(y^{z_2} = 1|\bar{y}^j = 1, y^{z_1} = 1, \mathbf{x}) \leq 1$ ,  $\xi$  is defined as  $\xi = \max\{p(y^{z_2} = 1|\bar{y}^j = 1, y^{z_1} = 1, \mathbf{x}), p(y^{z_1} = 1|\bar{y}^j = 1, y^{z_2} = 1, \mathbf{x})\}$ , the above inequation holds.  $\square$

## APPENDIX C

### THE PROOF OF COROLLARY 4

**Corollary 4.** Under a MLL scenario: there are  $m$  ( $m \geq 2$ ) labels  $l_{z_1}, l_{z_2}, \dots, l_{z_m} \in \mathcal{Y}$  ( $z_1, \dots, z_m \in [K]$ ) that are dependent, while the labels belong to  $\mathcal{Y} \setminus \{l_{z_1}, l_{z_2}, \dots, l_{z_m}\}$  are mutually exclusive. For any  $\mathbf{x} \in \mathcal{X}$ , its relevant set  $Y \subseteq \{l_{z_1}, l_{z_2}, \dots, l_{z_m}\}$  and  $Y \neq \emptyset$ . Suppose the label  $l_j$  is the complementary label of  $\mathbf{x}$ . The difference  $\ell_j$  between  $\mathbf{T}$  and  $\mathbf{Q}$  has

$$\ell_j \geq m\left(\frac{1}{\xi^m} - 1\right)p(\bar{y}^j = 1|\mathbf{x}),$$

where  $\xi = \max\{p(y^{z_m} = 1|\bar{y}^j = 1, y^{z_1} = 1, \dots, y^{z_{m-1}} = 1, \mathbf{x}), p(y^{z_{m-1}} = 1|\bar{y}^j = 1, y^{z_1} = 1, \dots, y^{z_{m-2}} = 1, y^{z_m} = 1, \mathbf{x}), \dots, p(y^{z_1} = 1|\bar{y}^j = 1, y^{z_2} = 1, \dots, y^{z_m} = 1, \mathbf{x})\}$  ( $\xi \in (0, 1]$ ).

*Proof.* Here, we apply induction to compute the difference as  $m$  increases. We start by computing the difference in the case of  $m = 3$ . Suppose class labels  $l_{z_1}, l_{z_2}, l_{z_3} \in \mathcal{Y}$  are dependent, while the rest of labels in the label space are mutually exclusive.  $\mathbf{x}$  is associated with  $Y \subseteq \{l_{z_1}, l_{z_2}, l_{z_3}\}$ , and  $Y \neq \emptyset$ . Then we calculate transition probabilities in  $\mathbf{T}$  from label correlations according to Theorem 3 as follows:

$$\begin{aligned} p(\bar{y}^j = 1|\mathbf{x}) &= \sum_{k=1, k \neq j, z_1, z_2, z_3}^K p(\bar{y}^j = 1|y^k = 1, \mathbf{x})p(y^k = 1|\mathbf{x}) + p(\bar{y}^j = 1, y^{z_1} = 1, y^{z_2} = 1, y^{z_3} = 1|\mathbf{x}) \\ &= p(\bar{y}^j = 1, y^{z_1} = 1, y^{z_2} = 1, y^{z_3} = 1|\mathbf{x}) \\ &= p(y^{z_3} = 1|\bar{y}^j = 1, y^{z_1} = 1, y^{z_2} = 1, \mathbf{x})p(y^{z_2} = 1|\bar{y}^j = 1, y^{z_1} = 1, \mathbf{x})p(\bar{y}^j = 1|y^{z_1} = 1, \mathbf{x})p(y^{z_1} = 1|\mathbf{x}) \\ &= p(y^{z_3} = 1|\bar{y}^j = 1, y^{z_1} = 1, y^{z_2} = 1, \mathbf{x})p(y^{z_2} = 1|\bar{y}^j = 1, y^{z_1} = 1, \mathbf{x})p(\bar{y}^j = 1|y^{z_1} = 1)p(y^{z_1} = 1|\mathbf{x}) \\ \Rightarrow \mathbf{T}_{z_1j} &= p(\bar{y}^j = 1|y^{z_1} = 1) = \frac{p(\bar{y}^j = 1|\mathbf{x})}{p(y^{z_3} = 1|\bar{y}^j = 1, y^{z_1} = 1, y^{z_2} = 1, \mathbf{x})p(y^{z_2} = 1|\bar{y}^j = 1, y^{z_1} = 1, \mathbf{x})p(y^{z_1} = 1|\mathbf{x})}. \end{aligned}$$

$\mathbf{T}_{z_2j}$  and  $\mathbf{T}_{z_3j}$  use the same way to estimate. Due to  $0 \leq p(y^{z_1} = 1|\mathbf{x}) \leq p(y^{z_1} = 1|\bar{y}^j = 1, y^{z_2} = 1, \mathbf{x}) \leq p(y^{z_1} = 1|\bar{y}^j = 1, y^{z_2} = 1, y^{z_3} = 1, \mathbf{x}) \leq 1$ , let  $\xi = \max\{p(y^{z_3} = 1|\bar{y}^j = 1, y^{z_1} = 1, y^{z_2} = 1, \mathbf{x}), p(y^{z_2} = 1|\bar{y}^j = 1, y^{z_1} = 1, y^{z_3} = 1, \mathbf{x}), p(y^{z_1} = 1|\bar{y}^j = 1, y^{z_2} = 1, y^{z_3} = 1, \mathbf{x})\}$ , we can obtain

$$\mathbf{T}_{z_1j} = p(\bar{y}^j = 1|y^{z_1} = 1) \geq \frac{1}{\xi^3}p(\bar{y}^j = 1|\mathbf{x}).$$

Similarly, we can compute  $\mathbf{T}_{z_2j}, \mathbf{T}_{z_3j} \geq \frac{1}{\xi^3}p(\bar{y}^j = 1|\mathbf{x})$ . Then the difference  $\ell_j$  is

$$\begin{aligned} \ell_j &= \sum_{k=1}^K |\mathbf{T}_{kj} - \mathbf{Q}_{kj}| \\ &= |\mathbf{T}_{z_1j} + \mathbf{T}_{z_2j} + \mathbf{T}_{z_3j} - 3p(\bar{y}^j = 1|\mathbf{x})| \\ &\geq 3\left(\frac{1}{\xi^3} - 1\right)p(\bar{y}^j = 1|\mathbf{x}). \end{aligned}$$

Similarly, for any  $m$  ( $0 < m < K$ ), suppose class labels  $l_{z_1}, l_{z_2}, \dots, l_{z_m} \in \mathcal{Y}$  are strongly dependent, while the rest of labels in the label space are mutually exclusive.  $\mathbf{x}$  is associated with  $Y \subseteq \{l_{z_1}, l_{z_2}, l_{z_3}\}$  and  $Y \neq \emptyset$ . Then we calculate transition probabilities from label correlations:

$$\begin{aligned} p(\bar{y}^j = 1|\mathbf{x}) &= \sum_{k=1, k \neq j, z_1, \dots, z_m}^K p(\bar{y}^j = 1|y^k = 1, \mathbf{x})p(y^k = 1|\mathbf{x}) + p(\bar{y}^j = 1, y^{z_1} = 1, \dots, y^{z_m} = 1|\mathbf{x}) \\ &= p(\bar{y}^j = 1, y^{z_1} = 1, \dots, y^{z_m} = 1|\mathbf{x}) \\ &= p(y^{z_m} = 1|\bar{y}^j = 1, y^{z_1} = 1, \mathbf{x})p(\bar{y}^j = 1|y^{z_1} = 1)p(y^{z_1} = 1|\mathbf{x})\Pi_{i=3}^m p(y^{z_i} = 1|\bar{y}^j = 1, y^{z_1} = 1, \dots, y^{z_{i-1}} = 1, \mathbf{x}) \\ \Rightarrow \mathbf{T}_{z_1j} &= p(\bar{y}^j = 1|y^{z_1} = 1) = \frac{p(\bar{y}^j = 1|\mathbf{x})}{p(y^{z_m} = 1|\bar{y}^j = 1, y^{z_1} = 1, \mathbf{x})p(y^{z_1} = 1|\mathbf{x})\Pi_{i=3}^m p(y^{z_i} = 1|\bar{y}^j = 1, y^{z_1} = 1, \dots, y^{z_{i-1}} = 1, \mathbf{x})}. \end{aligned}$$

As discussed above,  $\mathbf{T}_{z_1j} \geq \frac{1}{\xi^m}p(\bar{y}^j = 1|\mathbf{x})$  since  $\xi = \max\{p(y^{z_m} = 1|\bar{y}^j = 1, y^{z_1} = 1, \dots, y^{z_{m-1}} = 1, \mathbf{x}), p(y^{z_{m-1}} = 1|\bar{y}^j = 1, y^{z_1} = 1, \dots, y^{z_{m-2}} = 1, y^{z_m} = 1, \mathbf{x}), \dots, p(y^{z_1} = 1|\bar{y}^j = 1, y^{z_2} = 1, \dots, y^{z_m} = 1, \mathbf{x})\}$  ( $\xi \in (0, 1]$ ). Using the same calculation way, we can obtain  $\mathbf{T}_{z_2j}, \dots, \mathbf{T}_{z_mj} \geq \frac{1}{\xi^m}p(\bar{y}^j = 1|\mathbf{x})$ . Based on induction, we can summarize the difference  $\ell_j = \sum_{k=1}^K |\mathbf{T}_{kj} - \mathbf{Q}_{kj}| \geq m\left(\frac{1}{\xi^m} - 1\right)p(\bar{y}^j = 1|\mathbf{x})$ .  $\square$

## APPENDIX D

### THE PROOF OF THEOREM 6

**Theorem 6.** With Assumption 5, suppose the transition matrix  $\mathbf{T}$  is invertible, then the ML-CLL optimal classifier  $f_{CL}^*$  converges to the MLL optimal classifier  $f^*$ .

*Proof.* Before presenting the proof, we need to introduce the definition of *proper composite losses* as defined in [56]. According to the definition, many losses are considered *composite*, comprising a loss and a *link function* denoted as  $\psi$  (where  $\psi$  is invertible). Let  $L_\psi$  be a composite loss, which can be expressed with the assistance of a link function as follows:

$$L_\psi(\mathbf{g}(\mathbf{x}), \mathbf{y}) = L(\psi^{-1}(\mathbf{g}(\mathbf{x})), \mathbf{y}),$$

where the *inverse link function* ( $\psi^{-1}$ ) represents the *sigmoid* function in the case of *binary cross-entropy* (BCE) loss. In this context,  $\mathbf{g}(\cdot) \in \mathbb{R}^K$  denotes the output of a neural network before the application of the *sigmoid* function, and  $g^k(\cdot)$  denotes the  $k$ -th element of this output. Within neural networks, the real-valued function  $\mathbf{f}(\mathbf{x})$  denotes the output of  $\mathbf{g}(\mathbf{x})$  after passing through the *sigmoid* function. Moreover, [56] introduced the property of composite losses when the composite loss is considered *proper*, which means:

$$\operatorname{argmin}_{\mathbf{g}} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} [L_\psi(\mathbf{g}(\mathbf{x}), \mathbf{y})] = \psi(p(\mathbf{y}|\mathbf{x})).$$

In cases where composite losses are deemed *proper*, their minimizer exhibits a specific form corresponding to the link function applied to the class-conditional probabilities  $p(\mathbf{y}|\mathbf{x})$ . It's noteworthy that *binary cross-entropy* (BCE) loss and square loss are both examples of *proper composite* losses [46].

With Eq. (8) and the preceding equations, it is straightforward to have

$$\begin{aligned} \bar{L}(\mathbf{f}(\mathbf{x}), \bar{\mathbf{y}}) &= L(\mathbf{T}^T \mathbf{f}(\mathbf{x}), \bar{\mathbf{y}}) \\ &= L(\mathbf{T}^T \psi^{-1}(\mathbf{g}(\mathbf{x})), \bar{\mathbf{y}}) \\ &= L_\phi(\mathbf{g}(\mathbf{x}), \bar{\mathbf{y}}), \end{aligned}$$

where  $\phi^{-1}$  refers to the inverse link function, and we denote  $\phi^{-1} = \psi^{-1} \circ \mathbf{T}^T$ . Since  $\phi^{-1}$  is composed by invertible functions,  $\phi^{-1}$  is also invertible. Equivalently,  $\phi = (\mathbf{T}^{-1})^T \circ \psi$ . Leveraging the property of proper composite losses, the minimizer of the expected risk is:

$$\begin{aligned} \operatorname{argmin}_{\mathbf{g}} \mathbb{E}_{(\mathbf{x}, \bar{\mathbf{y}}) \sim \bar{\mathcal{D}}} [\bar{L}(\mathbf{f}(\mathbf{x}), \bar{\mathbf{y}})] &= \phi(p(\bar{\mathbf{y}}|\mathbf{x})) \\ &= \psi((\mathbf{T}^{-1})^T p(\bar{\mathbf{y}}|\mathbf{x})) \\ &= \psi((\mathbf{T}^{-1})^T \mathbf{T}^T p(\mathbf{y}|\mathbf{x})) \\ &= \psi(p(\mathbf{y}|\mathbf{x})) \\ &= \operatorname{argmin}_{\mathbf{g}} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} [L_\psi(\mathbf{g}(\mathbf{x}), \mathbf{y})]. \end{aligned}$$

Since the *sigmoid* function preserves the rank of its inputs, and  $\mathbf{f}(\mathbf{x})$  is the output of  $\mathbf{g}(\mathbf{x})$  after passing through the *sigmoid* function, it holds for that

$$\begin{aligned} \operatorname{argmin}_{\mathbf{f}} \mathbb{E}_{(\mathbf{x}, \bar{\mathbf{y}}) \sim \bar{\mathcal{D}}} [\bar{L}(\mathbf{f}(\mathbf{x}), \bar{\mathbf{y}})] &= \operatorname{argmin}_{\mathbf{f}} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} [L(\mathbf{f}(\mathbf{x}), \mathbf{y})] \\ \Leftrightarrow \operatorname{argmin}_{\mathbf{f}} R_{\bar{L}}(\mathbf{f}) &= \operatorname{argmin}_{\mathbf{f}} R_L(\mathbf{f}). \end{aligned}$$

Finally, it can be observed that  $f_{CL}^*$  converges to  $f^*$ , where  $f_{CL}^*$  and  $f^*$  are minimizers of  $R_{\bar{L}}(\mathbf{f})$  and  $R_L(\mathbf{f})$ , respectively.  $\square$