

CCMN: A General Framework for Learning with Class-Conditional Multi-Label Noise

Ming-Kun Xie and Sheng-Jun Huang

Abstract—Class-conditional noise commonly exists in machine learning tasks, where the class label is corrupted with a probability depending on its ground-truth. Many research efforts have been made to improve the model robustness against the class-conditional noise. However, they typically focus on the single label case by assuming that only one label is corrupted. In real applications, an instance is usually associated with multiple labels, which could be corrupted simultaneously with their respective conditional probabilities. In this paper, we formalize this problem as a general framework of learning with Class-Conditional Multi-label Noise (CCMN for short). We establish two unbiased estimators with error bounds for solving the CCMN problems, and further prove that they are consistent with commonly used multi-label loss functions. Finally, a new method for partial multi-label learning is implemented with the unbiased estimator under the CCMN framework. Empirical studies on multiple datasets and various evaluation metrics validate the effectiveness of the proposed method.

Index Terms—Class-conditional noise, class-conditional multi-label noise, unbiased estimator, partial multi-label learning.



1 INTRODUCTION

In ordinary supervised classification problems, a common assumption is that the class labels of training data are always correct. However, in practice, the assumption hardly holds since the training examples are usually corrupted due to unavoidable reasons, such as measurement error, subjective labeling bias or human labeling error. Learning in presence of label noise has been a problem of theoretical as well as practical interest in machine learning communities [1]. In various applications, such as the image annotation [2] and text classification [3], many successful methods have been applied to train robust models against label noise.

In general, to deal with noise-corrupted data, it is crucial to discover the causes of label noise. A natural and simple formulation of label noise is that the labels are corrupted by a random noise process [4], which can be described by the random classification noise (RCN) framework [5]. RCN framework assumes that each label is flipped independently with a specific probability $\rho \in [0, \frac{1}{2})$. Despite the great impact that RCN framework has made, it is too simple to deal with practical tasks, where the cause of label noise may not follow a random process. In order to deal with such problems, in [6], [7], authors propose a class-conditional random label noise (CCN) framework, where the probability of a label flipping depends on its true label. Unfortunately, CCN framework only considers single label case and fail to deal with multiple noisy labels, where multiple labels assigned to one instance may be corrupted simultaneously.

In this paper, we extend CCN to a more general framework of Class Conditional Multi-label Noise (CCMN) to learn with multi-class and multi-label classification tasks.

In CCMN framework, each instance is associated with multiple labels and each of class labels may be flipped with a probability ρ_{+1}^j or ρ_{-1}^j , which depends on its true label $y_j \in \{+1, -1\}$. It is worth noting that a great deal of real-world applications in weak-supervised setting [8] can be regarded as special cases of this framework. Examples include learning from partial-labeled data [9], [10], [11] or learning with missing labels [12], which we will discuss detailedly in following content. To the best of our knowledge, general theoretical results in this setting have not been developed before.

To tackle corrupted data with class-conditional multiple noisy labels, we derive a modified loss function for learning a multi-label classifier with risk minimization. Theoretically, we show that the empirical risk minimization with the modified loss functions can be in an unbiased fashion from independent and dependent perspectives. We then provide the estimation error bound for the two unbiased estimators and further show that learning with class-conditional multiple noisy labels can be multi-label consistent to two commonly used losses, i.e., hamming loss and ranking loss, respectively. Finally, CCMN can be regarded as a general framework of various weakly-supervised learning scenarios, such as partial label learning [9], partial multi-label learning [10] and weak label learning [12]. Among them, we take PML as an examples, and propose a new approach for partial multi-label learning with unbiased estimator. Comprehensive empirical studies demonstrate the effectiveness of the proposed methods.

Our main contributions are summarized as follows:

- Ming-Kun Xie and Sheng-Jun Huang are with the MITT Key Laboratory of Pattern Analysis and Machine Intelligence, Nanjing University of Aeronautics and Astronautics, Collaborative Innovation Center of Novel Software Technology and Industrialization, Nanjing 211106, China. E-mail: {mkxie, huangsj}@nuaa.edu.cn
- Sheng-Jun Huang is the corresponding author.
- A general framework of learning from class-conditional multi-label noise is proposed. Varied weakly-supervised learning scenarios can be cast under CCMN framework.
- We propose two unbiased estimators for solving CCMN problems in both independent and dependent fashion. These two estimators are proven to be multi-label con-

sistent to hamming loss and ranking loss, respectively.

- A novel approach for Partial Multi-label Learning with unbiased estimator (uPML) is proposed under the CCMN framework.

2 RELATED WORK

There are a great deal of previous works aiming to learn a robust classification model in presence of label noise. Most of these methods only consider single label case and cannot tackle class-conditional multi-label noise.

There has been a long line of studies in machine learning community on random label noise. The earliest work by Angluin and Laird [5] first propose the *random classification noise* (RCN) model. In [13], authors propose a boosting-based method and empirically show its robustness to random label noise. A theoretical analysis proposed in [14] proves that any method based on convex surrogate losses is inherently ill-suited to random label noise. In [4], authors consider random and adversarial label noise and try to deal with label noise by utilizing a robust SVM. Besides, there are many other heuristic methods without theoretical justification, such as confidence weighted learning [15], AROW [16] and the NHERD algorithm [17]. Noise-tolerance property is proposed in [18] to examine whether a loss function is robust to label noise.

The first attempt to tackle class-conditional label noise is in [6] where authors propose a variant of SVM to deal with noisy labels with theoretical guarantee. In [7], authors propose unbiased estimators to solve class-conditional label noise in the binary classification setting. Based on the assumption that the class-conditional distributions may overlap, a method called mixture proportion estimation is proposed in [19] to estimate the maximal proportion of one distribution that is present in another. Furthermore, authors extend the method into multi-class setting in [20].

Thanks to the great development of deep learning, there are various methods raised for utilizing deep neural networks to handle noisy labels, such as label correction methods [21], [22], [23], loss correction methods [24], [25], sample reweighting methods [26], [27], [28] and robust loss methods [29], [30].

In this paper, we also employ the CCMN framework to solve partial multi-label learning problems. In order to deal with partial-labeled data, the most commonly used strategy is disambiguation, which recovers ground-truth labeling information for candidate labels. Some methods perform disambiguation strategy by estimating a confidence for each candidate label [10], [31], [32], [33]. Other methods utilize decomposition scheme [34] or adversarial training [35]. Besides, label compression technique is utilized to deal with the large label space in PML tasks [36]. However, aforementioned methods never consider the generation process of noisy labels in candidate label set, which is a essential information for solving PML problems. In [37], authors first consider modeling the relationship between noisy labels and feature representations. Some studies aim to extend the PML framework to novel settings, such as semi-supervised learning [38] and multi-view learning [39]. Nevertheless, all of these methods failed to prove the consistency of the proposed method.

3 FORMULATION OF CCMN

Let $\mathbf{x} \in \mathcal{X}$ be a feature vector and $\mathbf{y} \in \mathcal{Y}$ be its corresponding label vector, where $\mathcal{X} \subset \mathbb{R}^d$ is the feature space and $\mathcal{Y} \subset \{-1, 1\}^q$ is the target space with q possible class labels. For notational convenience, the label y_j can be denoted by its index j . In the setting, $y_j = 1$ indicates the j -th label is a true label for instance \mathbf{x} ; $y_j = -1$, otherwise. In this paper, we focus on the multi-label learning problem, where each instance may be assigned with more than one label, i.e., $\sum_{j=1}^q I(y_j = 1) \geq 1$ holds, where I is the indicator function. Let $S = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$ be the given training data set, drawn *i.i.d.* according to the true distribution \mathcal{D} . We also use $[q]$ to denote the integer set $\{1, \dots, q\}$.

In this paper, for instance \mathbf{x} , its corresponding label \mathbf{y} can be corrupted and may be flipped into $\tilde{\mathbf{y}}$ following a class-conditional multi-label noise model as follows:

$$\begin{aligned} p(\tilde{y}_j = -1 | y_j = +1) &= \rho_{+1}^j, \\ p(\tilde{y}_j = +1 | y_j = -1) &= \rho_{-1}^j, \\ \forall j \in [q], \rho_{+1}^j + \rho_{-1}^j &< 1. \end{aligned}$$

where ρ_{+1}^j and ρ_{-1}^j are noise rates for each class label and are assumed to be known to the learner. In section 6, we will discuss how to estimate true noise rates in detail.

After injecting random noise into original samples S , the observed data set $S_\rho = \{(\mathbf{x}_1, \tilde{\mathbf{y}}_1), \dots, (\mathbf{x}_n, \tilde{\mathbf{y}}_n)\}$ are drawn *i.i.d.*, according to distribution \mathcal{D}_ρ . Our goal is to learn a prediction function $\mathbf{h} : \mathcal{X} \rightarrow \mathcal{Y}$ can accurately predict labels for any unseen instance. In general, it is not easy to learn \mathbf{h} directly, and alternatively, one usually learns a real-valued decision function $\mathbf{f} : \mathcal{X} \rightarrow \mathbb{R}^K$. Note that, for each instance \mathbf{x} , even though its final prediction depends on $\mathbf{h}(\mathbf{x})$, we also call \mathbf{f} itself the classifier. As mentioned before, CCMN can be used to solve multi-class and multi-label learning problems. In this paper, we focus on the multi-label classification task without loss of generality.

In general, CCMN framework has implications in varied real-world applications. In the following content, we discuss on two popular weakly-supervised learning scenarios, i.e., partial multi-label learning and weak label learning, which can be regarded as special cases of CCMN framework.

In PML problems, each instance is associated with a candidate label set Y_c , which contains multiple relevant labels and some other noisy labels. One intuitive strategy to solve the task is that treats all labels in Y_c as relevant labels and transforms Y_c into \tilde{Y} . Here, we use \tilde{Y} , since there may exist irrelevant labels in \tilde{Y} . Then, a new training set is obtained, where each instance is associated with the label set \tilde{Y} . Besides relevant labels, \tilde{Y} is also injected into some irrelevant labels, which can be regarded as class-conditional multi-label noise, i.e., irrelevant labels are flipped into relevant labels with $\rho_{-1}^j > 0$ while $\rho_{+1}^j = 0$.

In weak label learning, also known as multi-label learning with missing labels, relevant labels of each instance are partially known. Specifically, each instance is associated with a relevant label set Y_{+1} while \tilde{Y}_{-1} is the irrelevant label set. Here, we use \tilde{Y}_{-1} , since there may exist missing labels in \tilde{Y}_{-1} , i.e., relevant labels are missed. Similarly, one can treat all labels in \tilde{Y}_{-1} as irrelevant labels. Therefore, besides irrelevant labels, \tilde{Y}_{-1} is also injected into some relevant labels which

can be regarded as class-conditional multi-label noise with $\rho_{+1}^j > 0$ while $\rho_{-1}^j = 0$.

4 LEARNING WITH CCMN

In this section, we first provide some necessary preliminaries, and then derive two CCMN solvers for independent and dependent cases. In the independent case, we solve each binary classification task of the CCMN problem independently while in the dependent case, we solve the CCMN problem by considering label correlations.

4.1 Preliminaries

To derive our results for solving CCMN problems, we introduce some notations and the property of multi-label consistency.

There are many multi-label loss functions (also called evaluation metrics), such as *hamming loss*, *ranking loss*, *coverage* and *average precision* [40], etc. In this paper, we focus on two well known loss functions, i.e., hamming loss and ranking loss, and leave the discussion on other loss functions in the future work.

The hamming loss considers how many instance-label pairs are misclassified. Given the decision function \mathbf{f} and prediction function F , the hamming loss can be defined by:

$$L_h(F(\mathbf{f}(\mathbf{x})), \mathbf{y}) = \frac{1}{q} \sum_{j=1}^q I(\hat{y}_j \neq y_j), \quad (1)$$

where $\hat{\mathbf{y}} = F(\mathbf{f}(\mathbf{x})) = [\hat{y}_1, \dots, \hat{y}_q]$.

The ranking loss considers label pairs that are ordered reversely for an instance. Given a real-value decision function $\mathbf{f} = [f_1, f_2, \dots, f_q]$, the ranking loss can be defined by:

$$L_r(\mathbf{f}, \mathbf{y}) = \sum_{1 \leq j < k \leq q} I(y_j < y_k) \ell(j, k) + I(y_j > y_k) \ell(k, j), \quad (2)$$

where

$$\ell(j, k) = I(f_j > f_k) + \frac{1}{2} I(f_j = f_k).$$

The risk of \mathbf{f} with respect to loss L is given by $R(\mathbf{f}) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}}[L(\mathbf{f}(\mathbf{x}), \mathbf{y})]$ and the minimal risk (also called the Bayes risk) can be defined by $R^* = \inf_{\mathbf{f}} R(\mathbf{f})$. For an instance \mathbf{x} , the conditional risk of \mathbf{f} can be defined as

$$l(\mathbf{p}, \mathbf{f}) = \sum_{\mathbf{y} \in \mathcal{Y}} p_{\mathbf{y}} L(\mathbf{f}, \mathbf{y}),$$

where $p_{\mathbf{y}} = [p(\mathbf{y}|\mathbf{x})]_{\mathbf{y} \in \mathcal{Y}}$ is a vector of conditional probability over $\mathbf{y} \in \mathcal{Y}$.

Note that the above mentioned two loss functions are discontinuous and computationally NP-hard, which makes the corresponding optimization problems hard to solve. In practice, a feasible solution is to consider alternatively a surrogate loss function \mathcal{L} which can be optimized efficiently. We will give the specific definition of \mathcal{L} in the next section. The \mathcal{L} -risk and Bayes \mathcal{L} -risk can be defined as $R_{\mathcal{L}}(\mathbf{f}) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}}[\mathcal{L}(\mathbf{f}(\mathbf{x}), \mathbf{y})]$ and $R_{\mathcal{L}}^* = \inf_{\mathbf{f}} R_{\mathcal{L}}(\mathbf{f})$, respectively. Accordingly, we define the empirical \mathcal{L} -risk as $\hat{R}_{\mathcal{L}}(\mathbf{f}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\mathbf{f}(\mathbf{x}_i), \mathbf{y}_i)$.

Furthermore, we define the conditional \mathcal{L} -risk of \mathbf{f}

$$W(\mathbf{p}, \mathbf{f}) = \sum_{\mathbf{y} \in \mathcal{Y}} p_{\mathbf{y}} \mathcal{L}(\mathbf{f}, \mathbf{y}),$$

and the conditional Bayes \mathcal{L} -risk

$$W^*(\mathbf{p}) = \inf_{\mathbf{f}} W(\mathbf{p}, \mathbf{f}).$$

Our goal is to learn a good classification model with the modified loss function $\tilde{\mathcal{L}}(\mathbf{f}(\mathbf{x}), \tilde{\mathbf{y}})$ from noise-corrupted data by minimizing empirical $\tilde{\mathcal{L}}$ -risk:

$$\hat{\mathbf{f}} = \arg \min_{\mathbf{f} \in \mathcal{F}} \hat{R}_{\tilde{\mathcal{L}}}(\mathbf{f}),$$

where \mathcal{F} is a function class.

4.2 Independent Case

In order to solve multi-label learning problems, the most straightforward method is to decompose the task into q independent binary classification problems [41], where the goal is to learn q functions, $\mathbf{f} = (f_1, f_2, \dots, f_q)$. However, as mentioned before, it is difficult to directly optimize the hamming loss due to its discontinuity. Alternatively, we consider the following surrogate loss

$$\mathcal{L}_h(\mathbf{f}(\mathbf{x}), \mathbf{y}) = \sum_{j=1}^q \phi(y_j f_j(\mathbf{x})), \quad (3)$$

where ϕ is a convex loss function. The common choices are least square loss $\phi(t) = (1-t)^2$ and hinge loss $\phi(t) = (1-t)_+$ in [42].

The modified loss function under class-conditional multi-label noise in the independent case can be defined as follows:

$$\tilde{\mathcal{L}}_h(\mathbf{f}(\mathbf{x}), \tilde{\mathbf{y}}) = \sum_{j=1}^q \tilde{\phi}(y_j f_j(\mathbf{x})), \quad (4)$$

where,

$$\tilde{\phi}(y_j f_j) = \kappa_j [(1 - \rho_{-y_j}) \phi(y_j f_j) - \rho_{y_j} \phi(-y_j f_j)].$$

Here, $\kappa_j = \frac{1}{1 - \rho_{+1}^j - \rho_{-1}^j}$ is a constant and we omit the superscript j of ρ_{y_j} .

We extend the results in [7] to have the following lemma, which shows unbiasedness of the estimator defined by Eq.(4).

Lemma 1. For any $y_j, \forall j \in [q]$, let $\phi(y_j f_j(\mathbf{x}))$ be any bounded loss function. Then, if $\tilde{\mathcal{L}}_h(\mathbf{f}, \tilde{\mathbf{y}})$ can be defined by Eq.(4), we have $\mathbb{E}_{\tilde{\mathbf{y}}} [\tilde{\mathcal{L}}_h(\mathbf{f}, \tilde{\mathbf{y}})] = \mathcal{L}_h(\mathbf{f}, \mathbf{y})$.

Let $\sigma = \{\sigma_1, \dots, \sigma_n\}$ be n Rademacher variables with σ_i independently uniform random variable taking value in $\{-1, +1\}$. Then, the Rademacher complexity with respect to function class \mathcal{F} and unbiased estimator $\tilde{\mathcal{L}}$ can be formulated as follows:

$$\mathcal{R}_n(\tilde{\mathcal{L}} \circ \mathcal{F}) = \mathbb{E}_{\mathbf{x}, \tilde{\mathbf{y}}, \sigma} \left[\sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \tilde{\mathcal{L}}(\mathbf{f}(\mathbf{x}_i), \tilde{\mathbf{y}}_i) \right].$$

Accordingly, the performance guarantee for the unbiased estimator can be derived as following theorems.

Theorem 1. Let $\mu = \max_j \frac{1}{1 - \rho_{-1}^j - \rho_{+1}^j}, \forall j \in [q]$. Then, for the loss function $\phi(\cdot)$ bounded by Θ , with probability at least $1 - \delta$, we have

$$R_{\mathcal{L}_h}(\hat{\mathbf{f}}) - \min_{\mathbf{f} \in \mathcal{F}} R_{\mathcal{L}_h}(\mathbf{f}) \leq 4qK_{\rho} \mathcal{R}_n(\mathcal{F}) + 2q\mu\Theta \sqrt{\frac{\ln \frac{1}{\delta}}{2n}},$$

where $\hat{\mathbf{f}}$ is trained by minimizing $\hat{R}_{\tilde{\mathcal{L}}_h}(\mathbf{f})$ and K_ρ is the Lipschitz constant of $\tilde{\mathcal{L}}_h$.

Theorem 2. If ϕ is convex function with $\phi'(0) < 0$, then learning from CCMN examples with the modified surrogate loss $\tilde{\mathcal{L}}_h$ defined by Eq.(4) is consistent w.r.t hamming loss, i.e., there exists a non-negative concave function ξ with $\xi(0) = 0$, such that

$$R_{L_h}(\hat{\mathbf{f}}) - R_{L_h}^* \leq \xi(R_{\mathcal{L}_h}(\hat{\mathbf{f}}) - R_{\mathcal{L}_h}^*).$$

As shown in Theorem 1, the generalization error is dependent to the noise rates ρ_{-1} and ρ_{+1} . It is intuitive that smaller noise rates lead to a better generalization performance owing to a smaller μ . By combining Theorem 1 with Theorem 2, we obtain a performance guarantee for learning from class-conditional multiple noisy labels with respect to hamming loss. As $n \rightarrow \infty$, we have the consistency: if $R_{\mathcal{L}_h}(\hat{\mathbf{f}}_n) \rightarrow R_{\mathcal{L}_h}^*$ (as shown in Theorem 1), then $R_{L_h}(\hat{\mathbf{f}}_n) \rightarrow R_{L_h}^*$, since $\mathcal{R}_n(\mathcal{F}) \rightarrow 0$ for all parametric models with a bounded norm such as deep networks trained with weight decay [43].

4.3 Dependent Case

In multi-label learning problems, a common assumption is that there exist label correlations among labels [44], [45]. Therefore, it is fundamental to learn with class-conditional multi-label noise in a dependent fashion. The ranking loss considers the second-order label correlation and its surrogate loss is commonly defined as

$$\begin{aligned} \mathcal{L}_r(\mathbf{f}(\mathbf{x}), \mathbf{y}) &= \sum_{1 \leq j < k \leq q} I(y_j > y_k) \phi(f_{jk}) + I(y_j < y_k) \phi(f_{kj}) \\ &= \sum_{1 \leq j < k \leq q} \phi(y_{jk}(f_j - f_k)), \end{aligned} \quad (5)$$

where $y_{jk} = \frac{y_j - y_k}{2}$ and $f_{jk} = f_j - f_k$.
If $y_j \neq y_k$, let

$$a = (1 - \rho_{-y_{jk}}^j)(1 - \rho_{y_{jk}}^k), b = \rho_{y_{jk}}^j \rho_{-y_{jk}}^k,$$

and if $y_j = y_k$, let

$$c = \rho_{y_j}^j (1 - \rho_{-y_k}^k), d = \rho_{y_k}^k (1 - \rho_{-y_j}^j).$$

Then, the modified loss function under class-conditional multi-label noise in the dependent case can be defined as

$$\tilde{\mathcal{L}}_r(\mathbf{f}(\mathbf{x}), \mathbf{y}) = \sum_{1 \leq j < k \leq q} \tilde{\phi}((f_j, f_k), (y_j, y_k)), \quad (6)$$

where $\tilde{\phi}((f_j, f_k), (y_j, y_k))$ can be abbreviated by $\tilde{\phi}(j, k)$,

$$\tilde{\phi}(j, k) = \begin{cases} \kappa_{jk} [a\phi(y_{jk}f_{jk}) + b\phi(-y_{jk}f_{jk})] & \text{if } y_j \neq y_k \\ -\kappa_{jk} [c\phi(-y_jf_{jk}) + d\phi(y_jf_{jk})] & \text{if } y_j = y_k \end{cases}$$

Here, $\kappa_{jk} = \frac{1}{(1-\rho_{+1}^j-\rho_{-1}^j)(1-\rho_{+1}^k-\rho_{-1}^k)}$ is a constant.

The unbiasedness of the estimator defined by Eq.(6) can be shown as following lemma.

Lemma 2. For any $y_j, y_k, \forall j, k \in [q]$, let $\phi(\cdot)$ be any bounded loss function. Then, if $\tilde{\mathcal{L}}_r(\mathbf{f}, \tilde{\mathbf{y}})$ can be defined by Eq.(6), we have $\mathbb{E}_{\tilde{\mathbf{y}}} [\tilde{\mathcal{L}}_r(\mathbf{f}, \tilde{\mathbf{y}})] = \mathcal{L}_r(\mathbf{f}, \mathbf{y})$.

Accordingly, the performance guarantee for the unbiased estimator can be derived as following theorems.

Theorem 3. Let $\mu = \max_j \frac{1+|\rho_{-1}^j-\rho_{+1}^j|}{(1-\rho_{-1}^j-\rho_{+1}^j)^2}, \forall j \in [q]$. Then, for the loss function $\phi(\cdot)$ bounded by Θ , with probability at least $1 - \delta$, we have

$$\begin{aligned} R_{\mathcal{L}_r}(\hat{\mathbf{f}}) - \min_{\mathbf{f} \in \mathcal{F}} R_{\mathcal{L}_r}(\mathbf{f}) \\ \leq 4q(q-1)K_\rho \mathcal{R}_n(\mathcal{F}) + 2q(q-1)\mu\Theta \sqrt{\frac{\ln \frac{1}{\delta}}{2n}}, \end{aligned}$$

where $\hat{\mathbf{f}}$ is trained by minimizing $\hat{R}_{\tilde{\mathcal{L}}_r}(\mathbf{f})$ and K_ρ is the Lipschitz constant of $\tilde{\mathcal{L}}_r$.

Theorem 4. If ϕ is a differential and non-increasing function with $\phi'(0) < 0$ and $\phi(t) + \phi(-t) = 2\phi(0)$, then learning from CCMN examples with the modified surrogate loss $\tilde{\mathcal{L}}_r$ defined by Eq.(6) is consistent w.r.t ranking loss, i.e., there exists a non-negative concave function ξ with $\xi(0) = 0$, such that

$$R_{L_r}(\hat{\mathbf{f}}) - R_{L_r}^* \leq \xi(R_{\mathcal{L}_r}(\hat{\mathbf{f}}) - R_{\mathcal{L}_r}^*).$$

From Theorem 3, it can be observed that besides the noise rates ρ_{-1} and ρ_{+1} , the generalization error is also dependent to the difference between noise rates, i.e., $|\rho_{-1} - \rho_{+1}|$. Generally, a smaller difference also leads to a better generalization performance owing to a smaller μ . By combining Theorem 3 with Theorem 4, we obtain a performance guarantee for learning from class-conditional multiple noisy labels with respect to ranking loss. As $n \rightarrow \infty$, we have the consistency: if $R_{\mathcal{L}_r}(\hat{\mathbf{f}}_n) \rightarrow R_{\mathcal{L}_r}^*$ (as shown in Theorem 3), then $R_{L_r}(\hat{\mathbf{f}}_n) \rightarrow R_{L_r}^*$, since $\mathcal{R}_n(\mathcal{F}) \rightarrow 0$ for all parametric models with a bounded norm such as deep networks trained with weight decay [43].

Testing Phase. Regarding the classifier $\hat{\mathbf{f}}$ trained by minimizing the loss function $\tilde{\mathcal{L}}_h$ or \mathcal{L}_h , for each testing instance \mathbf{x}_t , we use the $\text{sgn}(f_j(\mathbf{x}_t))$ to predict its labels, where $\text{sgn}(a)$ is a function which outputs +1 if $a \geq 0$; -1, otherwise.

However, regarding classifier \mathbf{f} trained by minimizing the loss function $\tilde{\mathcal{L}}_r$ or \mathcal{L}_r , it is unreasonable to use 0 as a threshold directly. Instead, to perform predictions, in training phase, we introduce a *dummy* label $y_0 = 0$ for each instance, and then, in testing phase, the output of the classifier for label y_0 is used as the threshold to decide the label assignment for each instance.

Specifically, for instance \mathbf{x} , suppose that its corresponding label vector can be represented by $\mathbf{y} = [y_0, y_1, \dots, y_q]$, where y_0 is indexed by 0. In the training phase, the cumulative loss $\tilde{\mathcal{L}}_0$ with respect to y_0 can be written as follows:

$$\tilde{\mathcal{L}}_0(\mathbf{f}(\mathbf{x}), \mathbf{y}) = \sum_{j=1}^q \tilde{\phi}(f_{j0}(\mathbf{x}), y_{j0}), \quad (7)$$

where $f_{j0}(\mathbf{x}) = f_j(\mathbf{x}) - f_0(\mathbf{x})$ and $y_{j0} = y_j - y_0$. Note that $y_{j0} = y_j$, since we have $y_0 = 0$. In the situation, the formulation of loss $\tilde{\mathcal{L}}_0$ is similar to Eq.(4) and can be solved efficiently.

5 PARTIAL MULTI-LABEL LEARNING WITH CCMN

As discussed above, partial multi-label learning (PML) is a recently proposed framework, and is a typical CCMN task. In this section, we take PML as an example to examine the effectiveness of the proposed framework. Specifically, we propose a new approach for partial multi-label learning with unbiased estimator (uPML for short). In the PML setting, the noise rates satisfy $\forall j \in [q], p(\tilde{y}_j = +1 | y_j = -1) = \rho^j, p(\tilde{y}_j = -1 | y_j = +1) = 0$. Here, we omit the subscript of ρ^j for notational simplicity.

Based on the results in Section 4, in the independent case, the objective function of uPML $\tilde{\mathcal{L}}_h^{\text{PML}}$ can be defined as follows:

$$\tilde{\mathcal{L}}_h^{\text{PML}}(\mathbf{f}(\mathbf{x}), \mathbf{y}) = \sum_{j=1}^q \tilde{\phi}(y_j f_j(\mathbf{x})), \quad (8)$$

where,

$$\tilde{\phi}(y_j f_j(\mathbf{x})) = \begin{cases} \frac{\phi(-f_j(\mathbf{x})) - \rho^j \phi(f_j(\mathbf{x}))}{1 - \rho^j} & \text{if } y_j = -1 \\ \phi(f_j(\mathbf{x})) & \text{otherwise} \end{cases}$$

In the dependent case, the objective function of uPML $\tilde{\mathcal{L}}_r^{\text{PML}}$ can be defined as follows:

$$\tilde{\mathcal{L}}_r^{\text{PML}}(\mathbf{f}(\mathbf{x}), \mathbf{y}) = \sum_{1 \leq j < k \leq q} \tilde{\phi}(y_{jk} f_{jk}(\mathbf{x})), \quad (9)$$

where,

$$\tilde{\phi}(y_{jk} f_{jk}(\mathbf{x})) = \begin{cases} \frac{\phi(f_{jk})}{1 - \rho^k} & \text{if } y_{jk} = +1 \\ \frac{\phi(-f_{jk})}{1 - \rho^j} & \text{if } y_{jk} = -1 \\ \frac{-\rho^j \phi(f_{jk}) - \rho^k \phi(-f_{jk})}{(1 - \rho^j)(1 - \rho^k)} & \text{if } y_j = y_k = -1 \\ 0 & \text{if } y_j = y_k = +1 \end{cases}$$

We derive the generalization error bound for the proposed uPML method, which can be regarded as special cases of Theorem 1 and 3.

Corollary 1. For the least square loss $\phi(t) = (1 - t)^2$, with probability at least $1 - \delta$, we have

$$R_{\mathcal{L}_h}(\hat{\mathbf{f}}) - \min_{\mathbf{f} \in \mathcal{F}} R_{\mathcal{L}_h}(\mathbf{f}) \leq 4qK_\rho \mathcal{R}_n(\mathcal{F}) + \frac{8q}{1 - \rho_{\min}} \sqrt{\frac{\ln \frac{1}{\delta}}{2n}},$$

where $\rho_{\min} = \min_j \rho^j$ and $\hat{\mathbf{f}}$ is trained by minimizing $\hat{R}_{\tilde{\mathcal{L}}_h^{\text{PML}}}(\mathbf{f})$, and

$$\begin{aligned} & R_{\mathcal{L}_r}(\hat{\mathbf{f}}) - \min_{\mathbf{f} \in \mathcal{F}} R_{\mathcal{L}_r}(\mathbf{f}) \\ & \leq 4q(q-1)K_\rho \mathcal{R}_n(\mathcal{F}) + 8q(q-1)\mu \sqrt{\frac{\ln \frac{1}{\delta}}{2n}}, \end{aligned}$$

where $\mu = \min_j \frac{1 + \rho^j}{(1 - \rho^j)^2}$ and $\hat{\mathbf{f}}$ is trained by minimizing $\hat{R}_{\tilde{\mathcal{L}}_r^{\text{PML}}}(\mathbf{f})$.

It is easy to prove the consistency of the proposed uPML method based on the results derived in previous sections. To the best of our knowledge, uPML is the first consistent method for solving PML problems.

6 NOISE RATE ESTIMATION

In the aforementioned discussions, the noise rates ρ are assumed to be known. However, in many real-world scenarios, the true noise rates ρ would be unknown and need to be estimated. In this section, we provide an efficient method to estimate ρ , which is a multi-label extension of the recent noise estimator [24], [46].

When learning with class-conditional multiple noisy labels, the information of ground-truth labels is no longer accessible by the learner. Without any extra information, it is impractical to estimate the true noise rates. Based on the two assumptions, we propose to estimate the noise rates for CCMN as following theorem.

Theorem 5. Assume $p(\mathbf{x}, \mathbf{y})$ is such that:

(1) For each label $y_j, \forall j \in [q]$, there exist anchor points in the sense that

$$\exists \bar{\mathbf{x}}_l^j \in \mathcal{X}, p(\bar{\mathbf{x}}_l^j) > 0 \wedge p(y_j = l | \bar{\mathbf{x}}_l^j) = 1, l \in \{+1, -1\}.$$

(2) given sufficiently many corrupted training examples, \mathbf{f}^1 is rich enough to model $p(\tilde{y}_j = l | \mathbf{x})$ accurately.

It follows that $\forall j \in [q], l \in \{-1, +1\}, \rho_l^j = p(\tilde{y}_j = l | \bar{\mathbf{x}}_l^j)$.

Proof. For any $j \in [q]$ and any $\mathbf{x} \in \mathcal{X}$, we have:

$$p(\tilde{y}_j = l | \mathbf{x}) = \rho_{-l}^j p(y_j = -l | \mathbf{x}) + (1 - \rho_l^j) p(y_j = l | \mathbf{x}).$$

Given condition (1), when $\mathbf{x} = \bar{\mathbf{x}}_l^j$, we have $p(y_j = l | \bar{\mathbf{x}}_l^j) = 1$, which means $\rho_l^j = p(\tilde{y}_j = -l | \bar{\mathbf{x}}_l^j)$. \square

Theorem 5 tells us that the noise rates can be estimated based on noisy label probability estimates, i.e., the outputs of sigmoid of a network trained on noisy labels. In particular, as shown in [24], it even does not require these examples to have any clean labels at all. Specifically, $\forall j \in [q], l \in \{+1, -1\}$, the noise rate ρ_l^j can be approximated with two procedures:

$$\bar{\mathbf{x}}_j^l = \arg \max_{\mathbf{x} \in S_\rho} \hat{p}(\tilde{y}_j = l | \mathbf{x})$$

$$\hat{\rho}_l^j = \hat{p}(\tilde{y}_j = -l | \bar{\mathbf{x}}_j^l),$$

where $\hat{p}(\tilde{y}_j | \mathbf{x})$ is the sigmoid output and can be regarded as an estimation of class-conditional probabilities $p(\tilde{y}_j | \mathbf{x})$.

In practice, it would satisfy assumption (1) of Theorem 5 when the noisy training set S_ρ is large enough. However, it is more difficult to justify assumption (2) of Theorem 5, since the true class-conditional probability is unknown. In our experiments, it can be observed that the proposed methods obtain promising performance based on the estimated noise rates. This validates the effectiveness of the proposed estimator.

7 PROOF

In this section, we provide detailed proofs of the theorems derived in previous sections.

1. Note that $\mathbf{f}(\cdot)$ passes through the sigmoid function $1/(1 + e^{-f_j(\mathbf{x})})$, which can be interpreted as a vector approximating the class-conditional probabilities $p(\mathbf{y} | \mathbf{x})$.

7.1 Proof of Lemma 1

It is straightforward to show the assertion of Lemma 1 as follows. Recall that here we use ρ_j and ρ_{-j} to denote ρ_{y_j} and ρ_{-y_j} , respectively, then, we have

$$\begin{aligned} & \mathbb{E}_{\tilde{\mathbf{y}}} [\tilde{\mathcal{L}}_h(\mathbf{f}, \tilde{\mathbf{y}})] \\ &= \sum_{j=1}^q (1 - \rho_j) \tilde{\ell}(f_j, y_j) + \rho_j \tilde{\ell}(f_j, -y_j) \\ &= \sum_{j=1}^q \kappa_j \{ (1 - \rho_j) [(1 - \rho_{-j}) \ell(f_j, y_j) - \rho_j \ell(f_j, -y_j)] \\ &\quad + \rho_j [(1 - \rho_j) \ell(f_j, -y_j) - \rho_{-j} \ell(f_j, y_j)] \} \\ &= \sum_{j=1}^K \kappa_j (1 - \rho_j - \rho_{-j}) \ell(f_j, y_j) = \mathcal{L}_h(\mathbf{f}, \mathbf{y}), \end{aligned}$$

which completes proof. \square

7.2 Proof of Theorem 2

The proof is mainly composed of the following two lemmas.

Lemma 3. Let $\mathcal{R}_n(\tilde{\mathcal{L}}_h \circ \mathcal{F})$ be the Rademacher complexity of $\tilde{\mathcal{L}}_h$ and \mathcal{F} over S_ρ with n training points drawn from \mathcal{D}_ρ , which can be defined as

$$\mathcal{R}_n(\tilde{\mathcal{L}}_h \circ \mathcal{F}) = \mathbb{E}_{S_\rho} \mathbb{E}_\sigma \left[\sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \tilde{\mathcal{L}}_h(\mathbf{f}(\mathbf{x}_i), \tilde{\mathbf{y}}_i) \right].$$

Then,

$$\mathcal{R}_n(\tilde{\mathcal{L}}_h \circ \mathcal{F}) \leq q K_\rho \mathcal{R}_n(\mathcal{F}),$$

where K_ρ is the Lipschitz constant of $\tilde{\mathcal{L}}_h$.

Proof. Recall that $\tilde{\mathcal{L}}_h(\mathbf{f}(\mathbf{x}), \mathbf{y}) = \sum_{j=1}^q \tilde{\phi}(y_j f_j(\mathbf{x}))$, then, we have

$$\begin{aligned} & \mathcal{R}_n(\tilde{\mathcal{L}}_h \circ \mathcal{F}) \\ &= \mathbb{E}_{S_\rho} \mathbb{E}_\sigma \left[\sup_{f_1, \dots, f_q \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \sum_{j=1}^q \tilde{\phi}(y_j^{(i)} f_j(\mathbf{x}_i)) \right] \\ &= \mathbb{E}_{\mathcal{X}} \mathbb{E}_\sigma \left[\sup_{f_1, \dots, f_q \in \mathcal{F}} \frac{1}{n} \sum_{\mathbf{x}_i \in \mathcal{X}} \sigma_i \sum_{j=1}^q \tilde{\phi}(y_j^{(i)} f_j(\mathbf{x}_i)) \right] \\ &\leq \sum_{j=1}^q \mathbb{E}_{\mathcal{X}} \mathbb{E}_\sigma \left[\sup_{f_j \in \mathcal{F}} \frac{1}{n} \sum_{\mathbf{x}_i \in \mathcal{X}} \sigma_i \tilde{\phi}(y_j^{(i)} f_j(\mathbf{x}_i)) \right] \\ &= q \mathcal{R}_n(\tilde{\phi} \circ \mathcal{F}). \end{aligned}$$

Sequentially, according to Talagrand's contraction lemma [47], we have

$$\begin{aligned} \mathcal{R}_n(\tilde{\mathcal{L}}_h \circ \mathcal{F}) &\leq q \mathcal{R}_n(\tilde{\phi} \circ \mathcal{F}) \\ &\leq q K_\rho \mathcal{R}_n(\mathcal{F}), \end{aligned}$$

which completes the proof. \square

Without loss of generality, assume that $\forall j \in [q], \mu = \max_j \frac{1}{1 - \rho_{-1}^j - \rho_{+1}^j}$.

Lemma 4. For any $\delta > 0$, with probability at least $1 - \delta$,

$$\max_{\mathbf{f} \in \mathcal{F}} |\hat{R}_{\tilde{\mathcal{L}}_h}(\mathbf{f}) - R_{\tilde{\mathcal{L}}_h}(\mathbf{f})| \leq 2q K_\rho \mathcal{R}_n(\mathcal{F}) + q\mu\Theta \sqrt{\frac{\ln \frac{2}{\delta}}{2n}}.$$

Proof. Since both two directions can be proved in the same way, we consider one single direction $\sup_{\mathbf{f}_1, \dots, \mathbf{f}_q \in \mathcal{F}} (\hat{R}_{\tilde{\mathcal{L}}}(\mathbf{f}) - R_{\tilde{\mathcal{L}}}(\mathbf{f}))$. Note that the change in \mathbf{x}_i leads to a perturbation of at most $\frac{q\mu\Theta}{n}$ by replacing a single point $(\mathbf{x}_i, \tilde{\mathbf{y}}_i)$ with $(\mathbf{x}'_i, \tilde{\mathbf{y}}'_i)$, since the change in any $\tilde{y}_j^{(i)}$ leads to a perturbation as Eq.(10). By using McDiarmid's inequality [48] to the single-direction uniform deviation $\sup_{\mathbf{f}_1, \dots, \mathbf{f}_k \in \mathcal{F}} (\hat{R}_{\tilde{\mathcal{L}}}(\mathbf{f}) - R_{\tilde{\mathcal{L}}}(\mathbf{f}))$, we have

$$\begin{aligned} & p \left\{ \sup_{\mathbf{f} \in \mathcal{F}} (\hat{R}_{\tilde{\mathcal{L}}}(\mathbf{f}) - R_{\tilde{\mathcal{L}}}(\mathbf{f})) - \mathbb{E} \left[\sup_{\mathbf{f} \in \mathcal{F}} (\hat{R}_{\tilde{\mathcal{L}}}(\mathbf{f}) - R_{\tilde{\mathcal{L}}}(\mathbf{f})) \right] \geq \epsilon \right\} \\ & \leq \exp \left(-\frac{2\epsilon^2}{n(\frac{q\mu\Theta}{n})^2} \right), \end{aligned}$$

or equivalently, with probability at least $1 - \delta$,

$$\begin{aligned} & \sup_{\mathbf{f} \in \mathcal{F}} (\hat{R}_{\tilde{\mathcal{L}}}(\mathbf{f}) - R_{\tilde{\mathcal{L}}}(\mathbf{f})) \\ & \leq \mathbb{E} \left[\sup_{\mathbf{f} \in \mathcal{F}} (\hat{R}_{\tilde{\mathcal{L}}}(\mathbf{f}) - R_{\tilde{\mathcal{L}}}(\mathbf{f})) \right] + q\mu\Theta \sqrt{\frac{\ln \frac{1}{\delta}}{2n}}. \end{aligned}$$

According to [48], it is straightforward to show that

$$\mathbb{E} \left[\sup_{\mathbf{f} \in \mathcal{F}} (\hat{R}_{\tilde{\mathcal{L}}}(\mathbf{f}) - R_{\tilde{\mathcal{L}}}(\mathbf{f})) \right] \leq 2\mathcal{R}_n(\tilde{\mathcal{L}} \circ \mathcal{F}).$$

With the lemma 3, we complete the proof. \square

Based on the Lemma 4, with $\mathbf{f}^* = \operatorname{argmin}_{\mathbf{f} \in \mathcal{F}} R_{\mathcal{L}, \mathcal{D}}(\mathbf{f})$, it is obvious to prove the generalization error bound as follows:

$$\begin{aligned} & R_{\mathcal{L}, \mathcal{D}}(\hat{\mathbf{f}}) - R_{\mathcal{L}, \mathcal{D}}(\mathbf{f}^*) \\ &= R_{\tilde{\mathcal{L}}, \mathcal{D}_\rho}(\hat{\mathbf{f}}) - R_{\tilde{\mathcal{L}}, \mathcal{D}_\rho}(\mathbf{f}^*) \\ &= (\hat{R}_{\tilde{\mathcal{L}}}(\hat{\mathbf{f}}) - \hat{R}_{\tilde{\mathcal{L}}}(\mathbf{f}^*)) + (R_{\tilde{\mathcal{L}}, \mathcal{D}_\rho}(\hat{\mathbf{f}}) - \hat{R}_{\tilde{\mathcal{L}}}(\hat{\mathbf{f}})) \\ &\quad + (\hat{R}_{\tilde{\mathcal{L}}}(\mathbf{f}^*) - R_{\tilde{\mathcal{L}}, \mathcal{D}_\rho}(\mathbf{f}^*)) \\ &\leq 0 + 2 \max_{\mathbf{f} \in \mathcal{F}} |\hat{R}_{\tilde{\mathcal{L}}}(\mathbf{f}) - R_{\tilde{\mathcal{L}}, \mathcal{D}_\rho}(\mathbf{f})|. \end{aligned}$$

The first equation holds due to the unbiasedness of the estimator and for the last line, we used the fact $\hat{R}_{\tilde{\mathcal{L}}}(\hat{\mathbf{f}}) \leq \hat{R}_{\tilde{\mathcal{L}}}(\mathbf{f}^*)$ by the definition of $\hat{\mathbf{f}}$. \square

7.3 Proof for Theorem 3

Before providing the proof, the definition of multi-label consistency can be formulated as follows.

Definition 1. [49] Given a below-bounded surrogate loss \mathcal{L} , where $\mathcal{L}(\cdot, \mathbf{y})$ is continuous for every $\mathbf{y} \in \mathcal{Y}$, \mathcal{L} is said to be multi-label consistent w.r.t. the loss L if it holds, for every \mathbf{p} , that

$$W^*(\mathbf{p}) < \inf_{\mathbf{f}} \{W(\mathbf{p}, \mathbf{f}) : \mathbf{f} \notin \mathcal{A}(\mathbf{p})\},$$

where $\mathcal{A}(\mathbf{p}) = \{\mathbf{f} : l(\mathbf{p}, \mathbf{f}) = \inf_{\mathbf{f}'} l(\mathbf{p}, \mathbf{f}')\}$ is the set of Bayes decision functions.

Based on the definition, the following theorem can be further established.

Theorem 6. [49] The surrogate loss \mathcal{L} is multi-label consistent w.r.t. the loss L if and only if it holds for any sequence $\{\mathbf{f}_n\}_{n \geq 1}$ that

$$\text{if } R_{\mathcal{L}}(\mathbf{f}_n) \rightarrow R_{\mathcal{L}}^* \text{ then } R(\mathbf{f}_n) \rightarrow R^*.$$

$$\frac{1}{n} \left| \frac{(1 - \rho_{-1}^j)\phi(f_j) - \rho_{+1}^j\phi(-f_j) - [(1 - \rho_{+1}^j)\phi(-f_j) - \rho_{-1}^j\phi(f_j)]}{1 - \rho_{-1}^j - \rho_{+1}^j} \right| = \frac{1}{n} \left| \frac{\phi(f_j) - \phi(-f_j)}{1 - \rho_{-1}^j - \rho_{+1}^j} \right| \leq \frac{\mu\Theta}{n}. \quad (10)$$

$$\begin{bmatrix} \alpha_{+1}^j \alpha_{-1}^k & \rho_{+1}^j \rho_{-1}^k & \rho_{+1}^j \alpha_{-1}^k & \rho_{-1}^k \alpha_{+1}^j \\ \rho_{-1}^j \rho_{+1}^k & \alpha_{-1}^j \alpha_{+1}^k & \rho_{-1}^j \alpha_{+1}^k & \rho_{+1}^k \alpha_{-1}^j \\ \alpha_{+1}^j \rho_{+1}^k & \rho_{+1}^j \alpha_{+1}^k & \rho_{+1}^j \rho_{+1}^k & \alpha_{+1}^j \alpha_{+1}^k \\ \rho_{-1}^j \alpha_{-1}^k & \alpha_{-1}^j \rho_{-1}^k & \alpha_{-1}^j \alpha_{-1}^k & \rho_{-1}^j \rho_{-1}^k \end{bmatrix} \begin{bmatrix} \tilde{\phi}((f_j, f_k), (+1, -1)) \\ \tilde{\phi}((f_j, f_k), (-1, +1)) \\ \tilde{\phi}((f_j, f_k), (+1, +1)) \\ \tilde{\phi}((f_j, f_k), (-1, -1)) \end{bmatrix} = \begin{bmatrix} \phi(f_j - f_k) \\ \phi(f_k - f_j) \\ \phi(0) \\ \phi(0) \end{bmatrix}. \quad (11)$$

where $\alpha_{+1}^j = 1 - \rho_{+1}^j$ and $\alpha_{-1}^j = 1 - \rho_{-1}^j$.

$$\begin{aligned} \tilde{\phi}((f_j, f_k), (+1, -1)) &= \kappa_{jk} \left[(1 - \rho_{-1}^j)(1 - \rho_{+1}^k)\phi(f_{jk}) + \rho_{+1}^j \rho_{-1}^k \phi(-f_{jk}) \right] \\ \tilde{\phi}((f_j, f_k), (-1, +1)) &= \kappa_{jk} \left[(1 - \rho_{+1}^j)(1 - \rho_{-1}^k)\phi(-f_{jk}) + \rho_{-1}^j \rho_{+1}^k \phi(f_{jk}) \right] \\ \tilde{\phi}((f_j, f_k), (+1, +1)) &= -\kappa_{jk} \left[\rho_{+1}^j (1 - \rho_{-1}^k)\phi(-f_{jk}) + \rho_{+1}^k (1 - \rho_{-1}^j)\phi(f_{jk}) \right] \\ \tilde{\phi}((f_j, f_k), (-1, -1)) &= -\kappa_{jk} \left[\rho_{-1}^j (1 - \rho_{+1}^k)\phi(f_{jk}) + \rho_{-1}^k (1 - \rho_{+1}^j)\phi(-f_{jk}) \right]. \end{aligned} \quad (12)$$

This theorem tells us that the multi-label consistency is a necessary and sufficient condition for the convergence of \mathcal{L} -risk to the Bayes \mathcal{L} -risk, implying $R(\mathbf{f}) \rightarrow R^*$. The proof of Theorem 3 is presented as follows.

With respect to $\tilde{\mathcal{L}}_h$, the conditional surrogate loss can be defined as

$$\begin{aligned} \widetilde{W}(\mathbf{p}, \mathbf{f}) &= \mathbb{E}_{\tilde{\mathbf{y}}}[\tilde{\mathcal{L}}(\mathbf{f}(\mathbf{x}), \tilde{\mathbf{y}})] = \sum_{\mathbf{y} \in \mathcal{Y}} p_{\mathbf{y}} \mathbb{E}_{\tilde{\mathbf{y}}|\mathbf{y}}[\tilde{\mathcal{L}}(\mathbf{f}(\mathbf{x}), \tilde{\mathbf{y}})] \\ &= \sum_{\mathbf{y} \in \mathcal{Y}} p_{\mathbf{y}} \mathcal{L}(\mathbf{f}(\mathbf{x}), \mathbf{y}) = \sum_{j=1}^q \sum_{\mathbf{y} \in \mathcal{Y}} p_{\mathbf{y}} \phi(y_j f_j(\mathbf{x})) \\ &= W(\mathbf{p}, \mathbf{f}) = \sum_{j=1}^q p_j^+ \phi(f_j(\mathbf{x})) + p_j^- \phi(-f_j(\mathbf{x})). \end{aligned}$$

where the second equality is based on the law of total probability. Here, $p_j^+ = \sum_{\mathbf{y}: y_j = +1} p_{\mathbf{y}}$ and $p_j^- = \sum_{\mathbf{y}: y_j = -1} p_{\mathbf{y}}$. Accordingly, we have

$$\widetilde{W}^*(\mathbf{p}) = W^*(\mathbf{p}) = \inf_{\mathbf{f}} W(\mathbf{p}, \mathbf{f}) = \sum_{j=1}^q \inf_{f_j} W_j(p_j^+, f_j),$$

where $W_j(p_j^+, f_j) = p_j^+ \phi(f_j) + p_j^- \phi(-f_j)$. This yields that minimizing $W(\mathbf{p}, \mathbf{f})$ is equivalent to minimizing $W_j(p_j^+, f_j)$ for every $1 \leq j \leq q$. The consistency for binary classification has been well-studied [50], [51]. Based on their results, it is easy to prove that learning from CCMN data with $\tilde{\mathcal{L}}_h$ is consistent with respect to hamming loss. Therefore, based on Corollary 25 in [52], we have

$$R_{L_h}(\hat{\mathbf{f}}) - R_{L_h}^* \leq \xi(R_{L_h}(\hat{\mathbf{f}}) - R_{L_h}^*),$$

where ξ is a non-negative concave function with $\xi(0) = 0$. \square

7.4 Proof for Lemma 2

For each instance, regarding each pair of noisy labels $(\tilde{y}_j, \tilde{y}_k)$, there may exist four cases as follows: $(\tilde{y}_j = +1, \tilde{y}_k = -1)$, $(\tilde{y}_j = -1, \tilde{y}_k = +1)$, $(\tilde{y}_j = +1, \tilde{y}_k = +1)$ and $(\tilde{y}_j = -1, \tilde{y}_k = -1)$. By considering these four cases separately, we

have four linear equations as presented in Eq.(11). Solving these four equations for $\tilde{\phi}(j, k)$ gives the solution as shown in Eq.(12). With simple computation for these four equations, we obtain the unbiased estimator defined in Eq.(6), which completes the proof. \square

7.5 Proof for Theorem 4

We first propose the following two lemmas, which are useful for proving Theorem 4.

Lemma 5. Let $\mathcal{R}_n(\tilde{\mathcal{L}}_r \circ \mathcal{F})$ be the Rademacher complexity of $\tilde{\mathcal{L}}_r$ and \mathcal{F} over S_ρ with n training data drawn from \mathcal{D}_ρ , which can be defined as

$$\mathcal{R}_n(\tilde{\mathcal{L}}_r \circ \mathcal{F}) = \mathbb{E}_{S_\rho} \mathbb{E}_\sigma \left[\sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \tilde{\mathcal{L}}_r(\mathbf{f}(\mathbf{x}_i), \tilde{\mathbf{y}}_i) \right].$$

Then,

$$\mathcal{R}_n(\tilde{\mathcal{L}}_r \circ \mathcal{F}) \leq 2q(q-1)K_\rho \mathcal{R}_n(\mathcal{F}),$$

where K_ρ is the Lipschitz constant of $\tilde{\mathcal{L}}_r$.

Proof. Recalling $\tilde{\mathcal{L}}_r(\mathbf{f}(\mathbf{x}), \tilde{\mathbf{y}}) = \sum_{j,k} \tilde{\phi}((f_j, f_k), (y_j, y_k))$, we have

$$\begin{aligned} \mathcal{R}_n(\tilde{\mathcal{L}}_r \circ \mathcal{F}) &= \mathbb{E}_{S_\rho} \mathbb{E}_\sigma \left[\sup_{f_1, \dots, f_q \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \sum_{1 \leq j < k \leq q} \tilde{\phi}((f_j, f_k), (y_j, y_k)) \right] \\ &= \mathbb{E}_{\mathcal{X}} \mathbb{E}_\sigma \left[\sup_{f_1, \dots, f_q \in \mathcal{F}} \frac{1}{n} \sum_{\mathbf{x}_i \in \mathcal{X}} \sigma_i \sum_{1 \leq j < k \leq q} \tilde{\phi}((f_j, f_k), (y_j, y_k)) \right] \\ &\leq \sum_{1 \leq j < k \leq q} \mathbb{E}_{\mathcal{X}} \mathbb{E}_\sigma \left[\sup_{f_1, \dots, f_q \in \mathcal{F}} \frac{1}{n} \sum_{\mathbf{x}_i \in \mathcal{X}} \sigma_i \tilde{\phi}((f_j, f_k), (y_j, y_k)) \right]. \end{aligned} \quad (13)$$

Sequentially, let (y, y') be the current label pair to be cumulated, then, according to Talagrand's contraction lemma [47], we have

$$\begin{aligned} & \mathbb{E}_{\mathcal{X}} \mathbb{E}_{\sigma} \left[\sup_{f_y, f_{y'} \in \mathcal{F}} \frac{1}{n} \sum_{\mathbf{x}_i \in \mathcal{X}} \sigma_i \tilde{\phi}((f_y, f_{y'}), (y, y')) \right] \\ & \leq K_{\rho} \mathbb{E}_{\mathcal{X}} \mathbb{E}_{\sigma} \left[\sup_{f_y, f_{y'} \in \mathcal{F}} \frac{1}{n} \sum_{\mathbf{x}_i \in \mathcal{X}} \sigma_i (f_y(\mathbf{x}_i) - f_{y'}(\mathbf{x}_i)) \right] \\ & \leq K_{\rho} \mathbb{E}_{\mathcal{X}} \mathbb{E}_{\sigma} \left[\sup_{f_y \in \mathcal{F}} \frac{1}{n} \sum_{\mathbf{x}_i \in \mathcal{X}} \sigma_i f_y(\mathbf{x}_i) \right] \\ & \quad + K_{\rho} \mathbb{E}_{\mathcal{X}} \mathbb{E}_{\sigma} \left[\sup_{f_{y'} \in \mathcal{F}} \frac{1}{n} \sum_{\mathbf{x}_i \in \mathcal{X}} \sigma_i f_{y'}(\mathbf{x}_i) \right] \\ & = 2K_{\rho} \mathcal{R}_n(\mathcal{F}), \end{aligned}$$

where f_y represent the classifier corresponding class label y . Then, by combining with Eq.(13), it is easy to prove that $\mathcal{R}_n(\tilde{\mathcal{L}}_r \circ \mathcal{F}) \leq 2q(q-1)K_{\rho} \mathcal{R}_n(\mathcal{F})$. \square

Without loss of generality, assume that $\forall j \in [q], \mu = \max_j \frac{1 - \min(\rho_{-1}^j - \rho_{+1}^j, \rho_{+1}^j - \rho_{-1}^j)}{(1 - \rho_{-1}^j - \rho_{+1}^j)^2}$.

Lemma 6. For any $\delta > 0$, with probability at least $1 - \delta$,

$$\begin{aligned} & \max_{\mathbf{f} \in \mathcal{F}} |\hat{R}_{\tilde{\mathcal{L}}}(\mathbf{f}) - \mathcal{R}_{\tilde{\mathcal{L}}}(\mathbf{f})| \\ & \leq 2q(q-1)K_{\rho} \mathcal{R}_n(\mathcal{F}) + q(q-1)\mu\Theta \sqrt{\frac{\ln \frac{2}{\delta}}{2n}}. \end{aligned}$$

We omit the proof, since it can be proved similarly to Lemma 4. \square

Finally, similar to theorem 2, based on lemma 5 and 6, it is easy to obtain theorem 4. \square

7.6 Proof for Theorem 5

For notational simplicity, we introduce the following notation:

$$\Delta_{j,k} = \sum_{\mathbf{y}: y_j = -1, y_k = +1} p_{\mathbf{y}}.$$

With respect to $\tilde{\mathcal{L}}$, the conditional surrogate loss can be defined as

$$\begin{aligned} \tilde{W}(\mathbf{p}, \mathbf{f}) &= \mathbb{E}_{\tilde{\mathbf{y}}}[\tilde{\mathcal{L}}(\mathbf{f}(\mathbf{x}), \tilde{\mathbf{y}})] = \sum_{\mathbf{y} \in \mathcal{Y}} p_{\mathbf{y}} \mathbb{E}_{\tilde{\mathbf{y}}|\mathbf{y}}[\tilde{\mathcal{L}}(\mathbf{f}(\mathbf{x}), \tilde{\mathbf{y}})] \quad (14) \\ &= \sum_{\mathbf{y} \in \mathcal{Y}} p_{\mathbf{y}} \mathcal{L}(\mathbf{f}(\mathbf{x}), \mathbf{y}) = W(\mathbf{p}, \mathbf{f}) \\ &= \sum_{1 \leq j < k \leq q} \Delta_{j,k} \phi(f_k - f_j) + \Delta_{k,j} \phi(f_j - f_k), \end{aligned}$$

where the second equality is based on the law of total probability.

Following Theorem 10 in [49], if ϕ is a differential and non-increasing function with $\phi'(0) < 0$ and $\phi(t) + \phi(-t) = 2\phi(0)$, it suffices to prove that $f_j > f_k$ if $\Delta_{j,k} < \Delta_{k,j}$ for every \mathbf{f} such that $W^*(\mathbf{p}) = W(\mathbf{p}, \mathbf{f}) = \tilde{W}(\mathbf{p}, \mathbf{f})$, where $W^*(\mathbf{p}) = \inf_{\mathbf{f}} W(\mathbf{p}, \mathbf{f})$ is the conditional Bayes \mathcal{L} -risk. Therefore, by minimizing $\tilde{W}(\mathbf{p}, \mathbf{f})$, we obtain the Bayes decision function \mathbf{f}^* , which proves $\tilde{\mathcal{L}}_r$ is consistent to

ranking loss. Accordingly, based on Corollary 25 in [52], we have

$$R_{L_r}(\hat{\mathbf{f}}) - R_{L_r}^* \leq \xi(R_{\mathcal{L}_r}(\hat{\mathbf{f}}) - R_{\mathcal{L}_r}^*),$$

where ξ is a non-negative concave function with $\xi(0) = 0$. \square

8 EXPERIMENT

The experiments for CCMN data are first reported, followed by the experiments for PML data.

8.1 Experimental Settings

Datasets We evaluate our method on eight benchmark multi-label datasets: music_emotion, music_style, mirflickr, tmc2007², Multi-MNIST³ [53], Multi-Kuzushiji-MNIST (Multi-KMNIST for short), Multi-Fashion-MNIST⁴ (Multi-FMNIST for short), and VOC2007⁵ [54]. It is noteworthy that the first three datasets, i.e., music_emotion, music_style and mirflickr are real-world PML datasets [32]. For these datasets, candidate are first collected from web users and then these candidate labels are further examined by human labelers to specify the ground-truth labels. For tmc2007, we use its shorter version which contains 28,596 instances and 500 features for each instance. For three Multi-MNIST-style datasets, each dataset contains 10,000 images. VOC2007 (VOC for short) contains 9,963 images for 20 object categories, which are divided into train, val and test sets. Following [55], [56], we use the *trainval* set to train the models, and evaluate the performance on the test set. For each dataset except for VOC, we randomly sample 60% examples for training and 40% examples for testing.

Metrics We evaluate the performance of the proposed method based on multiple standard multi-label criterion: hamming loss, ranking loss, coverage and average precision. For hamming loss, ranking loss, coverage, the smaller value, the better performance; for average precision, the larger value, the better performance. The detail of these criterion can be found in [40].

Implementation For experiments on all datasets except for VOC, we train a linear model by using Adam [57] optimizer with learning rate chosen from $\{0.01, 0.001\}$. We added an ℓ_2 -regularization term, with the regularization parameter of 0.0001. For experiments on VOC, we use an Alexnet [58] pre-trained with the ILSVRC2012 dataset on Pytorch platform [59]. The Alexnet is trained by using stochastic gradient descent (SGD) with learning rate chosen from $\{0.01, 0.001\}$. An ℓ_2 -regularization term is added with the regularization parameter of 0.0001. The batch size for all datasets is set as 200 except for tmc and mirflickr, where the batch size is set as 400. All the experiments are conducted on GeForce RTX 2080 GPUs

2. See <http://mulan.sourceforge.net/datasets-mlc.html> for tmc2007.

3. See <https://github.com/shaohua0116/MultiDigitMNIST> for Multi-MNIST.

4. Similar to Multi-MNIST, we construct Multi-Kuzushiji-MNIST and Multi-Fashion-MNIST for two commonly used datasets Kuzushiji-MNIST and Fashion-MNIST, respectively.

5. See <http://host.robots.ox.ac.uk/pascal/VOC/voc2007/> for VOC2007.

TABLE 1

Comparison results between Ub-HL, Ub-RL and their baselines with diverse loss functions using linear models, where \bullet/\circ indicates whether the proposed method is significantly superior/inferior to the comparing methods via paired t -test (at 0.05 significance level).

Data	music_emotion	music_style	mirflickr	tmc2007	Multi-MNIST	Multi-KMNIST	Multi-FMNIST	VOC
Hamming loss (the smaller, the better)								
Ub-HL/Hinge	.207 \pm .003 \bullet	.125 \pm .003 \bullet	.103 \pm .007 \bullet	.059 \pm .002 \bullet	.258 \pm .008	.256 \pm .007 \bullet	.199 \pm .005	.063 \pm .002 \bullet
Ub-HL/Square	.201 \pm .002	.122 \pm .002 \bullet	.104 \pm .009	.060 \pm .001 \bullet	.258 \pm .006	.256 \pm .007 \bullet	.205 \pm .005 \bullet	.092 \pm .006 \bullet
Ub-RL/Sigmoid	.199 \pm .004 \bullet	.123 \pm .001	.111 \pm .017	.062 \pm .001 \bullet	.256 \pm .007 \bullet	.249 \pm .006 \bullet	.198 \pm .007 \bullet	.077 \pm .000 \bullet
Ub-RL/Hinge	.203 \pm .001	.123 \pm .002 \bullet	.107 \pm .010 \bullet	.061 \pm .002 \bullet	.257 \pm .007	.251 \pm .006 \bullet	.203 \pm .005	.068 \pm .001 \bullet
Ub-RL/Square	.203 \pm .002	.124 \pm .003 \bullet	.115 \pm .013 \bullet	.063 \pm .002 \bullet	.259 \pm .005	.259 \pm .002 \bullet	.218 \pm .005 \bullet	.092 \pm .015 \bullet
B-HL/Hinge	.277 \pm .028 \bullet	.189 \pm .043 \bullet	.151 \pm .021 \bullet	.124 \pm .020 \bullet	.278 \pm .011 \bullet	.282 \pm .009 \bullet	.235 \pm .009 \bullet	.172 \pm .017 \bullet
B-HL/Square	.261 \pm .022 \bullet	.171 \pm .030 \bullet	.130 \pm .015 \bullet	.112 \pm .017 \bullet	.274 \pm .008 \bullet	.280 \pm .010 \bullet	.232 \pm .010 \bullet	.235 \pm .021 \bullet
B-RL/Sigmoid	.280 \pm .008 \bullet	.188 \pm .045 \bullet	.167 \pm .025 \bullet	.130 \pm .021 \bullet	.275 \pm .010 \bullet	.281 \pm .012 \bullet	.229 \pm .006 \bullet	.098 \pm .007 \bullet
B-RL/Hinge	.263 \pm .023 \bullet	.166 \pm .026 \bullet	.133 \pm .015 \bullet	.112 \pm .017 \bullet	.275 \pm .009 \bullet	.280 \pm .011 \bullet	.231 \pm .009 \bullet	.183 \pm .019 \bullet
B-RL/Square	.255 \pm .018 \bullet	.164 \pm .028 \bullet	.131 \pm .014 \bullet	.112 \pm .015 \bullet	.276 \pm .008 \bullet	.280 \pm .010 \bullet	.236 \pm .007 \bullet	.215 \pm .019 \bullet
Ranking loss (the smaller, the better)								
Ub-HL/Hinge	.245 \pm .003 \bullet	.168 \pm .006 \bullet	.100 \pm .012 \bullet	.062 \pm .004 \bullet	.340 \pm .010 \bullet	.330 \pm .014 \bullet	.234 \pm .008	.130 \pm .016 \bullet
Ub-HL/Square	.237 \pm .002 \bullet	.160 \pm .005	.100 \pm .013	.067 \pm .003 \bullet	.343 \pm .009 \bullet	.333 \pm .017 \bullet	.245 \pm .009 \bullet	.180 \pm .015 \bullet
Ub-RL/Sigmoid	.239 \pm .003	.153 \pm .004 \bullet	.104 \pm .019	.068 \pm .003 \bullet	.333 \pm .012 \bullet	.313 \pm .014 \bullet	.223 \pm .012 \bullet	.114 \pm .013 \bullet
Ub-RL/Hinge	.238 \pm .003	.160 \pm .004	.105 \pm .016	.064 \pm .004 \bullet	.338 \pm .011 \bullet	.323 \pm .016 \bullet	.241 \pm .012 \bullet	.133 \pm .015 \bullet
Ub-RL/Square	.240 \pm .005	.163 \pm .006	.111 \pm .020	.070 \pm .003 \bullet	.348 \pm .014 \bullet	.336 \pm .012 \bullet	.270 \pm .011 \bullet	.227 \pm .018 \bullet
B-HL/Hinge	.350 \pm .036 \bullet	.309 \pm .114	.153 \pm .031 \bullet	.193 \pm .040 \bullet	.388 \pm .018 \bullet	.389 \pm .028 \bullet	.294 \pm .021 \bullet	.248 \pm .038 \bullet
B-HL/Square	.332 \pm .034 \bullet	.284 \pm .080 \bullet	.139 \pm .044	.185 \pm .035 \bullet	.380 \pm .014 \bullet	.386 \pm .029 \bullet	.290 \pm .021 \bullet	.270 \pm .035 \bullet
B-RL/Sigmoid	.392 \pm .024 \bullet	.298 \pm .053 \bullet	.181 \pm .031 \bullet	.240 \pm .037 \bullet	.390 \pm .014 \bullet	.384 \pm .038 \bullet	.279 \pm .020 \bullet	.241 \pm .034 \bullet
B-RL/Hinge	.342 \pm .035 \bullet	.288 \pm .086 \bullet	.139 \pm .039	.185 \pm .033 \bullet	.380 \pm .014 \bullet	.384 \pm .033 \bullet	.286 \pm .019 \bullet	.250 \pm .034 \bullet
B-RL/Square	.330 \pm .030 \bullet	.259 \pm .053 \bullet	.143 \pm .040	.186 \pm .031 \bullet	.381 \pm .013 \bullet	.385 \pm .029 \bullet	.296 \pm .017 \bullet	.285 \pm .029 \bullet
Coverage (the smaller, the better)								
Ub-HL/Hinge	.415 \pm .007 \bullet	.230 \pm .008 \bullet	.107 \pm .010	.154 \pm .005 \bullet	.491 \pm .009 \bullet	.480 \pm .014 \bullet	.393 \pm .012	.189 \pm .019 \bullet
Ub-HL/Square	.403 \pm .004 \bullet	.222 \pm .006	.107 \pm .011	.164 \pm .004 \bullet	.493 \pm .010 \bullet	.481 \pm .017 \bullet	.401 \pm .009 \bullet	.250 \pm .015 \bullet
Ub-RL/Sigmoid	.403 \pm .004	.211 \pm .004 \bullet	.109 \pm .016	.161 \pm .005 \bullet	.483 \pm .010 \bullet	.466 \pm .016 \bullet	.375 \pm .011 \bullet	.171 \pm .016 \bullet
Ub-RL/Hinge	.407 \pm .006	.221 \pm .006	.111 \pm .014	.156 \pm .005	.488 \pm .012 \bullet	.477 \pm .019 \bullet	.400 \pm .018	.191 \pm .017 \bullet
Ub-RL/Square	.410 \pm .009	.224 \pm .008	.116 \pm .017 \bullet	.168 \pm .004 \bullet	.499 \pm .013 \bullet	.488 \pm .018 \bullet	.429 \pm .014 \bullet	.293 \pm .018 \bullet
B-HL/Hinge	.502 \pm .029 \bullet	.354 \pm .098 \bullet	.151 \pm .026 \bullet	.324 \pm .040 \bullet	.539 \pm .016 \bullet	.532 \pm .029 \bullet	.449 \pm .023 \bullet	.317 \pm .037 \bullet
B-HL/Square	.497 \pm .031 \bullet	.336 \pm .069 \bullet	.141 \pm .037	.323 \pm .036 \bullet	.531 \pm .013 \bullet	.530 \pm .030 \bullet	.446 \pm .022 \bullet	.339 \pm .032 \bullet
B-RL/Sigmoid	.578 \pm .016 \bullet	.349 \pm .040 \bullet	.174 \pm .027 \bullet	.391 \pm .033 \bullet	.546 \pm .010 \bullet	.527 \pm .039 \bullet	.436 \pm .027 \bullet	.310 \pm .032 \bullet
B-RL/Hinge	.506 \pm .032 \bullet	.343 \pm .072 \bullet	.140 \pm .033	.325 \pm .033 \bullet	.530 \pm .013 \bullet	.527 \pm .033 \bullet	.440 \pm .024 \bullet	.319 \pm .031 \bullet
B-RL/Square	.497 \pm .028 \bullet	.317 \pm .046 \bullet	.144 \pm .034	.327 \pm .033 \bullet	.530 \pm .013 \bullet	.527 \pm .028 \bullet	.452 \pm .021 \bullet	.353 \pm .028 \bullet
Average precision (the greater, the better)								
Ub-HL/Hinge	.632 \pm .004 \bullet	.700 \pm .007 \bullet	.834 \pm .014 \bullet	.790 \pm .009 \bullet	.510 \pm .016	.517 \pm .014 \bullet	.644 \pm .012	.650 \pm .034 \bullet
Ub-HL/Square	.640 \pm .002 \bullet	.709 \pm .008 \bullet	.832 \pm .017	.786 \pm .007 \bullet	.511 \pm .011	.514 \pm .016 \bullet	.629 \pm .014 \bullet	.573 \pm .030 \bullet
Ub-RL/Sigmoid	.636 \pm .008	.703 \pm .002	.820 \pm .029	.779 \pm .008 \bullet	.516 \pm .015 \bullet	.533 \pm .014 \bullet	.650 \pm .017 \bullet	.668 \pm .028 \bullet
Ub-RL/Hinge	.638 \pm .002	.706 \pm .007	.825 \pm .020 \bullet	.782 \pm .008 \bullet	.515 \pm .014	.526 \pm .016 \bullet	.636 \pm .013 \bullet	.628 \pm .038 \bullet
Ub-RL/Square	.634 \pm .003 \bullet	.702 \pm .009 \bullet	.813 \pm .026 \bullet	.771 \pm .009 \bullet	.506 \pm .013 \bullet	.508 \pm .007 \bullet	.597 \pm .016 \bullet	.458 \pm .055 \bullet
B-HL/Hinge	.469 \pm .047 \bullet	.498 \pm .140 \bullet	.750 \pm .036 \bullet	.505 \pm .091 \bullet	.463 \pm .022 \bullet	.458 \pm .024 \bullet	.566 \pm .021 \bullet	.444 \pm .070 \bullet
B-HL/Square	.510 \pm .048 \bullet	.554 \pm .102 \bullet	.784 \pm .035	.560 \pm .080 \bullet	.474 \pm .018 \bullet	.463 \pm .024 \bullet	.572 \pm .023 \bullet	.423 \pm .061 \bullet
B-RL/Sigmoid	.482 \pm .043 \bullet	.506 \pm .114 \bullet	.716 \pm .045 \bullet	.472 \pm .092 \bullet	.468 \pm .019 \bullet	.461 \pm .030 \bullet	.583 \pm .016 \bullet	.484 \pm .074 \bullet
B-RL/Hinge	.502 \pm .051 \bullet	.564 \pm .102 \bullet	.780 \pm .032 \bullet	.559 \pm .076 \bullet	.472 \pm .018 \bullet	.463 \pm .027 \bullet	.577 \pm .020 \bullet	.457 \pm .064 \bullet
B-RL/Square	.520 \pm .042 \bullet	.584 \pm .083 \bullet	.781 \pm .031 \bullet	.556 \pm .072 \bullet	.469 \pm .015 \bullet	.464 \pm .024 \bullet	.564 \pm .018 \bullet	.393 \pm .045 \bullet

8.2 Study on CCMM data

To validate the effectiveness of the proposed unbiased estimators, we perform experiments on multi-label data with class-conditional multiple noisy labels. To inject label noise into training examples, each class label is flipped according to noise rate ρ_{-1}^j and ρ_{+1}^j randomly sampled from $\{0.1, 0.2, 0.3, 0.4, 0.5\}$. We repeat experiments 5 times with different noise rates and report their averaging results.

We compare the performance of the proposed method with their baselines. Specifically, two variants of the proposed method, i.e., **Unbiased-HL** ($\tilde{\mathcal{L}}_h$ defined in Eq.(4), Ub-HL for short) and **Unbiased-RL** ($\tilde{\mathcal{L}}_r$ defined in Eq.(6), Ub-RL for short) are compared with their baselines, i.e., **Biased-HL** (\mathcal{L}_h defined in Eq.(3), B-HL for short) and **Biased-RL** (\mathcal{L}_r defined in Eq.(5), B-RL for short). For $\tilde{\mathcal{L}}_h$ and \mathcal{L}_h , the surrogate loss function ϕ is defined as the hinge loss (Hinge for short) and least square loss (Square for short), respectively; For $\tilde{\mathcal{L}}_r$ and \mathcal{L}_r , besides the hinge loss and least square loss, we also adopt sigmoid loss (Sigmoid for short), since it is a consistent surrogate loss, which satisfies the conditions in Theorem 4.

Table 1 reports comparison results of the proposed methods against baseline methods by using linear models. For each data set, paired t -test based on 5 repeats (at

0.05 significance level) is conducted to show whether the proposed unbiased estimator is significantly different from the comparing methods. From the results, it is obvious that the proposed method outperforms their corresponding baselines with significant superiority in almost all cases, which validates the effectiveness of the proposed method. The performance of Ub-RL is generally superior to Ub-HL, which validates the label correlation is beneficial for learning with class-conditional multi-label noise. From the results, regarding Ub-RL, it can be observed that Sigmoid surrogate loss generally achieves better performance than the other losses, which provides an empirical validation of Theorem 4, since Sigmoid loss satisfies the conditions in Theorem 4 and is thus consistent for Ub-RL with respect to ranking loss while others are not [49].

8.3 Study on PML data

To validate the practical usefulness of the proposed unbiased estimators, we perform experiments on PML tasks. Regarding the multi-label datasets, for training instances, each negative label y_j can be flipped into a candidate label according to noise rate ρ^j randomly sampled from $\{0.2, 0.3, 0.4, 0.5, 0.6\}$. The noise rates are unknown in experiments and to be

TABLE 2

Comparison results between the proposed method (using linear model) and PML methods, where \bullet/\circ indicates whether the proposed method is significantly superior/inferior to the comparing method via paired t -test (at 0.05 significance level).

Data	music_emotion	music_style	mirflickr	tmc2007	Multi-MNIST	Multi-KMNIST	Multi-FMNIST	VOC
Hamming loss (the smaller, the better)								
uPML-HL	.200 \pm .002	.114 \pm .002	.156 \pm .008	.053 \pm .000	.234 \pm .000	.231 \pm .001	.171 \pm .000 \bullet	.052 \pm .000
uPML-RL	.199 \pm .002	.115 \pm .001	.157 \pm .003	.054 \pm .000 \bullet	.234 \pm .000	.230 \pm .000	.166 \pm .000	.066 \pm .011
PMLNI	.228 \pm .001 \bullet	.115 \pm .002	.171 \pm .002 \bullet	.081 \pm .003 \bullet	.258 \pm .004 \bullet	.263 \pm .004 \bullet	.198 \pm .005 \bullet	.096 \pm .001 \bullet
PMLLS	.231 \pm .000 \bullet	.121 \pm .001 \bullet	.175 \pm .008	.078 \pm .003 \bullet	.260 \pm .004 \bullet	.261 \pm .005 \bullet	.197 \pm .007 \bullet	.094 \pm .001 \bullet
PARVLS	.256 \pm .001 \bullet	.126 \pm .001 \bullet	.174 \pm .002 \bullet	.095 \pm .001 \bullet	.320 \pm .005 \bullet	.325 \pm .007 \bullet	.289 \pm .004 \bullet	.128 \pm .000 \bullet
PARMAP	.251 \pm .004 \bullet	.125 \pm .004 \bullet	.167 \pm .003	.097 \pm .002 \bullet	.283 \pm .004 \bullet	.288 \pm .004 \bullet	.251 \pm .003 \bullet	.120 \pm .002 \bullet
fPML	.232 \pm .001 \bullet	.122 \pm .002 \bullet	.176 \pm .009 \bullet	.080 \pm .003 \bullet	.258 \pm .003 \bullet	.261 \pm .007 \bullet	.198 \pm .006 \bullet	.104 \pm .003 \bullet
Ranking loss (the smaller, the better)								
uPML-HL	.236 \pm .005 \bullet	.144 \pm .005	.124 \pm .007	.051 \pm .000	.303 \pm .000 \bullet	.287 \pm .001 \bullet	.180 \pm .000 \bullet	.089 \pm .004 \bullet
uPML-RL	.226 \pm .003	.142 \pm .002	.118 \pm .002	.055 \pm .001 \bullet	.298 \pm .000	.272 \pm .000	.172 \pm .000	.083 \pm .003
PMLNI	.250 \pm .003 \bullet	.145 \pm .005	.124 \pm .002 \bullet	.109 \pm .010 \bullet	.334 \pm .007 \bullet	.339 \pm .013 \bullet	.223 \pm .009 \bullet	.160 \pm .006 \bullet
PMLLS	.260 \pm .001 \bullet	.151 \pm .000 \bullet	.127 \pm .005 \bullet	.106 \pm .010 \bullet	.340 \pm .007 \bullet	.340 \pm .016 \bullet	.221 \pm .011 \bullet	.155 \pm .004 \bullet
PARVLS	.356 \pm .003 \bullet	.243 \pm .008 \bullet	.170 \pm .004 \bullet	.186 \pm .004 \bullet	.462 \pm .004 \bullet	.467 \pm .010 \bullet	.401 \pm .008 \bullet	.291 \pm .002 \bullet
PARMAP	.317 \pm .006 \bullet	.197 \pm .006 \bullet	.139 \pm .003 \bullet	.197 \pm .013 \bullet	.405 \pm .010 \bullet	.410 \pm .017 \bullet	.339 \pm .010 \bullet	.283 \pm .003 \bullet
fPML	.264 \pm .003 \bullet	.157 \pm .004 \bullet	.127 \pm .007 \bullet	.118 \pm .011 \bullet	.339 \pm .009 \bullet	.338 \pm .014 \bullet	.223 \pm .012 \bullet	.200 \pm .011 \bullet
Coverage (the smaller, the better)								
uPML-HL	.411 \pm .002 \bullet	.205 \pm .005	.125 \pm .006	.138 \pm .001	.456 \pm .000 \bullet	.435 \pm .000 \bullet	.327 \pm .000 \bullet	.140 \pm .004 \bullet
uPML-RL	.392 \pm .003	.199 \pm .003	.120 \pm .002	.142 \pm .001 \bullet	.450 \pm .000	.419 \pm .000	.319 \pm .000	.133 \pm .003
PMLNI	.413 \pm .002 \bullet	.204 \pm .006	.124 \pm .002 \bullet	.231 \pm .014 \bullet	.483 \pm .006 \bullet	.486 \pm .016 \bullet	.373 \pm .010 \bullet	.230 \pm .008 \bullet
PMLLS	.419 \pm .002 \bullet	.209 \pm .001 \bullet	.128 \pm .004	.229 \pm .015 \bullet	.488 \pm .007 \bullet	.485 \pm .019 \bullet	.371 \pm .012 \bullet	.224 \pm .006 \bullet
PARVLS	.492 \pm .002 \bullet	.276 \pm .008 \bullet	.160 \pm .003 \bullet	.332 \pm .007 \bullet	.551 \pm .002 \bullet	.546 \pm .004 \bullet	.490 \pm .008 \bullet	.355 \pm .003 \bullet
PARMAP	.479 \pm .003 \bullet	.257 \pm .006 \bullet	.140 \pm .003 \bullet	.356 \pm .015 \bullet	.551 \pm .006 \bullet	.558 \pm .020 \bullet	.498 \pm .013 \bullet	.362 \pm .001 \bullet
fPML	.423 \pm .002 \bullet	.215 \pm .004 \bullet	.127 \pm .006 \bullet	.245 \pm .012 \bullet	.491 \pm .007 \bullet	.487 \pm .017 \bullet	.376 \pm .016 \bullet	.272 \pm .013 \bullet
Average precision (the greater, the better)								
uPML-HL	.646 \pm .005	.730 \pm .005	.753 \pm .014	.818 \pm .001	.561 \pm .000 \bullet	.570 \pm .001 \bullet	.707 \pm .000 \bullet	.746 \pm .004
uPML-RL	.649 \pm .005	.728 \pm .005	.755 \pm .004	.810 \pm .002 \bullet	.562 \pm .000	.575 \pm .000	.717 \pm .000	.744 \pm .007
PMLNI	.601 \pm .003 \bullet	.728 \pm .006	.752 \pm .004	.723 \pm .017 \bullet	.516 \pm .010 \bullet	.505 \pm .010 \bullet	.651 \pm .011 \bullet	.621 \pm .011 \bullet
PMLLS	.586 \pm .003 \bullet	.708 \pm .002 \bullet	.746 \pm .013	.732 \pm .018 \bullet	.513 \pm .009 \bullet	.507 \pm .012 \bullet	.654 \pm .012 \bullet	.632 \pm .007 \bullet
PARVLS	.524 \pm .004 \bullet	.658 \pm .004 \bullet	.731 \pm .004 \bullet	.631 \pm .005 \bullet	.458 \pm .008 \bullet	.451 \pm .008 \bullet	.530 \pm .006 \bullet	.470 \pm .002 \bullet
PARMAP	.533 \pm .007 \bullet	.671 \pm .008 \bullet	.752 \pm .004	.612 \pm .015 \bullet	.454 \pm .004 \bullet	.445 \pm .012 \bullet	.533 \pm .006 \bullet	.456 \pm .010 \bullet
fPML	.583 \pm .005 \bullet	.703 \pm .008 \bullet	.744 \pm .014	.721 \pm .018 \bullet	.516 \pm .007 \bullet	.508 \pm .014 \bullet	.653 \pm .010 \bullet	.561 \pm .014 \bullet

TABLE 3

Friedman statistics F_F in terms of each evaluation metric and the critical value at 0.05 significance level (# comparing algorithms $k = 7$, # data sets $N = 8$).

Evaluation metric	F_F	critical value
Hamming Loss	29.6784	
Ranking loss	55.7200	
Coverage	53.8932	2.3240
Average Precision	34.5364	

estimated by using the estimator proposed in Section 6. We repeat experiments 5 times with different noise rate and report the averaging results.

We compare with five state-of-the-art PML algorithms: PML-NI [37], PARTICLE (including two implementations: PAR-VLS and PAR-MAP) [32], PML-LRS [34] and fPML [60]. It is noteworthy that to make a fair comparison, we use a linear classifier as the base model for uPML in all datasets except for VOC, where we only fine-tune parameters of the last layer while the other parameters are frozen. For other comparing methods, parameters are determined by the performance on validation set if no default value given in their literature.

Table 2 reports the comparison result of two variants of the proposed uPML method, i.e., uPML-HL ($\hat{\mathcal{L}}_h$ defined in Eq.(8)) with hinge loss and uPML-RL ($\hat{\mathcal{L}}_r$ defined in Eq.(9)) with Sigmoid loss, against comparing PML methods. For each data set, paired t -test based on five repeats (at 0.05 significance level) is conducted to show whether the proposed unbiased estimator is significantly different from the comparing methods. The proposed uPML method signif-

icantly outperforms the comparing PML methods in almost all cases. In particular, on three real-world PML datasets, i.e., music_emotion, music_style and mirflickr, uPML achieves the best performances on all cases, which demonstrates the practical usefulness of the proposed method. In general, uPML-RL is better than uPML-HL due to the label correlation is considered for the former one.

Furthermore, the commonly used *Friedman test* [61] is employed as the statistical test to analyze the relative performance among the comparing approaches. Table 3 reports the Friedman statistics F_F and the corresponding critical value with respect to each evaluation metric (# comparing algorithms $k = 7$, # data sets $N = 8$). For each evaluation metric, the null hypothesis of indistinguishable performance among the comparing algorithm is rejected at 0.05 significance level.

Finally, the post-hoc *Bonferroni-Dunn test* [61] is utilized to illustrate the relative performance among comparing approaches. Here, uPML is regarded as the control method whose average rank difference against the comparing algorithm is calibrated with the *critical difference* (CD). Accordingly, uPML is deemed to have significantly different performance to one comparing algorithm if their average ranks differ by at least one CD (CD = 2.8494 in our experiment: # comparing algorithms $k = 7$, # data sets $N = 8$). Figure 1 shows the CD diagrams ([61]) on each evaluation metric, where the average rank of each comparing algorithm is marked along the axis (lower ranks to the right). In each subfigure, any comparing algorithms whose average rank is within one CD to that of uPML is interconnected to each other with a thick line. From the figure, it can be observed that: 1)

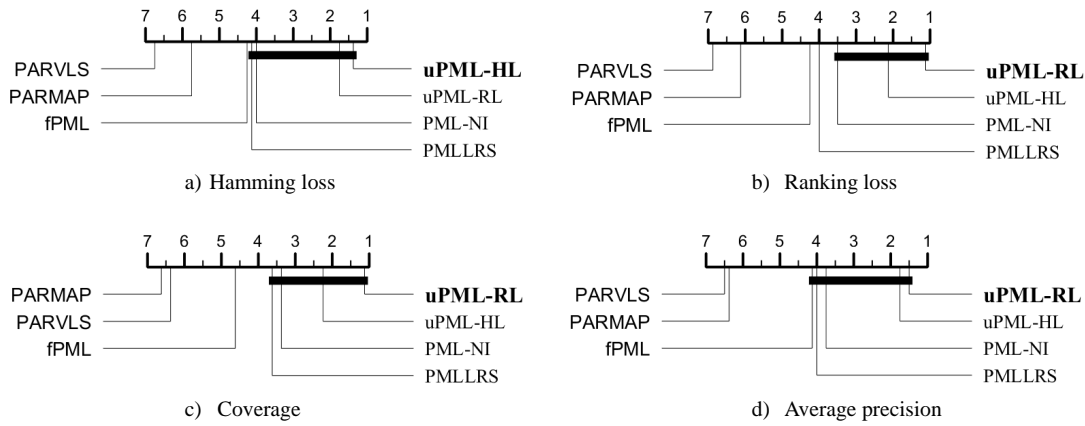


Fig. 1. Comparison of uPML (control algorithm) against five comparing algorithms with the Bonferroni-Dunn test. Algorithms not connected with uPML in the CD diagram are considered to have a significantly different performance from the control algorithm ($CD = 2.8494$ at 0.05 significance level).

uPML achieves the best (lowest) average rank in terms of all evaluation metrics; 2) uPML is significantly better than the PARVLS and PARMAP in terms of all evaluation metrics; 3) uPML significantly outperforms comparing methods other than PML-NI and PMLLRs in terms of *hamming loss*, *ranking loss* and *coverage*. These experimental results convincingly demonstrate the significance of the superiority for our uPML approach.

9 CONCLUSION

In this paper, we study the problem of multi-label classification with class-conditional multiple noisy labels, where multiple class labels assigned to each instance may be corrupted simultaneously with class-conditional probabilities. From the perspective of the unbiased estimator, we derive efficient methods for solving CCMN problems with theoretical guarantee. Generally, we prove that learning from class-conditional multiple noisy labels with the proposed unbiased estimators is consistent with respect to hamming loss and ranking loss. Furthermore, we propose a novel method called uPML for solving PML problems, which can be regarded as a special case of CCMN framework. Empirical studies on multiple data sets validate the effectiveness of the proposed method. In the future, we will study CCMN with more loss functions.

REFERENCES

- [1] B. Frénay and M. Verleysen, "Classification in the presence of label noise: A survey," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 25, no. 5, pp. 845–869, 2014.
- [2] T. Xiao, T. Xia, Y. Yang, C. Huang, and X. Wang, "Learning from massive noisy labeled data for image classification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 2691–2699.
- [3] I. Jindal, D. Pressel, B. Lester, and M. Noleby, "An effective label noise model for dnn text classification," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 2019, pp. 3246–3256.
- [4] B. Biggio, B. Nelson, and P. Laskov, "Support vector machines under adversarial label noise," in *Asian conference on machine learning*, 2011, pp. 97–112.
- [5] D. Angluin and P. Laird, "Learning from noisy examples," *Machine Learning*, vol. 2, no. 4, pp. 343–370, 1988.
- [6] G. Stempfel and L. Ralaivola, "Learning svms from sloppily labeled data," in *International conference on artificial neural networks*. Springer, 2009, pp. 884–893.
- [7] N. Natarajan, I. S. Dhillon, P. K. Ravikumar, and A. Tewari, "Learning with noisy labels," in *Advances in neural information processing systems*, 2013, pp. 1196–1204.
- [8] Z.-H. Zhou, "A brief introduction to weakly supervised learning," *National Science Review*, vol. 5, no. 1, pp. 44–53, 2018.
- [9] T. Cour, B. Sapp, and B. Taskar, "Learning from partial labels," *Journal of Machine Learning Research*, vol. 12, pp. 1501–1536, 2011.
- [10] M. Xie and S. Huang, "Partial multi-label learning," in *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI-18)*, 2018, pp. 4302–4309.
- [11] C. Gong, T. Liu, Y. Tang, J. Yang, J. Yang, and D. Tao, "A regularization approach for instance-based superset label learning," *IEEE transactions on cybernetics*, vol. 48, no. 3, pp. 967–978, 2017.
- [12] Y. Sun, Y. Zhang, and Z. Zhou, "Multi-label learning with weak label," in *Proceedings of the 24th AAAI Conference on Artificial Intelligence*, 2010.
- [13] Y. Freund, "A more robust boosting algorithm," *stat*, vol. 1050, p. 13, 2009.
- [14] P. M. Long and R. A. Servedio, "Random classification noise defeats all convex potential boosters," *Machine learning*, vol. 78, no. 3, pp. 287–304, 2010.
- [15] M. Dredze, K. Crammer, and F. Pereira, "Confidence-weighted linear classification," in *Proceedings of the 25th international conference on Machine learning*, 2008, pp. 264–271.
- [16] K. Crammer, A. Kulesza, and M. Dredze, "Adaptive regularization of weight vectors," in *Advances in neural information processing systems*, 2009, pp. 414–422.
- [17] K. Crammer and D. D. Lee, "Learning via gaussian herding," in *Advances in neural information processing systems*, 2010, pp. 451–459.
- [18] N. Manwani and P. Sastry, "Noise tolerance under risk minimization," *IEEE transactions on cybernetics*, vol. 43, no. 3, pp. 1146–1151, 2013.
- [19] C. Scott, G. Blanchard, and G. Handy, "Classification with asymmetric label noise: Consistency and maximal denoising," in *Conference On Learning Theory*, 2013, pp. 489–511.
- [20] G. Blanchard and C. Scott, "Decontamination of mutually contaminated models," in *Artificial Intelligence and Statistics*, 2014, pp. 1–9.
- [21] A. Vahdat, "Toward robustness against label noise in training deep discriminative neural networks," in *Advances in Neural Information Processing Systems*, 2017, pp. 5596–5605.
- [22] Y. Li, J. Yang, Y. Song, L. Cao, J. Luo, and L.-J. Li, "Learning from noisy labels with distillation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 1910–1918.
- [23] X. Xia, T. Liu, N. Wang, B. Han, C. Gong, G. Niu, and M. Sugiyama, "Are anchor points really indispensable in label-noise learning?"

- Advances in Neural Information Processing Systems*, vol. 32, pp. 6838–6849, 2019.
- [24] G. Patrini, A. Rozza, A. Krishna Menon, R. Nock, and L. Qu, “Making deep neural networks robust to label noise: A loss correction approach,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1944–1952.
 - [25] B. Han, J. Yao, G. Niu, M. Zhou, I. Tsang, Y. Zhang, and M. Sugiyama, “Masking: A new perspective of noisy supervision,” in *Advances in Neural Information Processing Systems*, 2018, pp. 5836–5846.
 - [26] M. Ren, W. Zeng, B. Yang, and R. Urtasun, “Learning to reweight examples for robust deep learning,” in *International Conference on Machine Learning*, 2018, pp. 4334–4343.
 - [27] B. Han, Q. Yao, X. Yu, G. Niu, M. Xu, W. Hu, I. Tsang, and M. Sugiyama, “Co-teaching: Robust training of deep neural networks with extremely noisy labels,” *arXiv preprint arXiv:1804.06872*, 2018.
 - [28] J. Li, R. Socher, and S. C. Hoi, “Dividemix: Learning with noisy labels as semi-supervised learning,” *arXiv preprint arXiv:2002.07394*, 2020.
 - [29] A. Ghosh, H. Kumar, and P. Sastry, “Robust loss functions under label noise for deep neural networks,” in *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
 - [30] Z. Zhang and M. Sabuncu, “Generalized cross entropy loss for training deep neural networks with noisy labels,” in *Advances in neural information processing systems*, 2018, pp. 8778–8788.
 - [31] H. Wang, W. Liu, Y. Zhao, C. Zhang, T. Hu, and G. Chen, “Discriminative and correlative partial multi-label learning,” in *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, August 10-16, 2019*, 2019, pp. 3691–3697.
 - [32] M.-L. Zhang and J.-P. Fang, “Partial multi-label learning via credible label elicitation,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
 - [33] W. Liu, X. Shen, H. Wang, and I. W. Tsang, “The emerging trends of multi-label learning,” *arXiv preprint arXiv:2011.11197*, 2020.
 - [34] T. W. C. L. Lijuan Sun, Songhe Feng and Y. Jin, “Partial multi-label learning by low-rank and sparse decomposition,” in *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence (AAAI-19)*, 2019.
 - [35] Y. Yan and Y. Guo, “Adversarial partial multi-label learning,” *arXiv preprint arXiv:1909.06717*, 2019.
 - [36] T. Yu, G. Yu, J. Wang, C. Domeniconi, and X. Zhang, “Partial multi-label learning using label compression,” in *2020 IEEE International Conference on Data Mining (ICDM)*. IEEE, 2020, pp. 761–770.
 - [37] M.-K. Xie and S.-J. Huang, “Partial multi-label learning with noisy label identification,” in *AAAI*, 2020, pp. 6454–6461.
 - [38] —, “Semi-supervised partial multi-label learning,” in *2020 IEEE International Conference on Data Mining (ICDM)*. IEEE, 2020, pp. 691–700.
 - [39] Z.-S. Chen, X. Wu, Q.-G. Chen, Y. Hu, and M.-L. Zhang, “Multi-view partial multi-label learning with graph-based disambiguation,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 3553–3560.
 - [40] M. Zhang and Z. Zhou, “A review on multi-label learning algorithms,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 26, no. 8, pp. 1819–1837, 2014.
 - [41] M. R. Boutell, J. Luo, X. Shen, and C. M. Brown, “Learning multi-label scene classification,” *Pattern Recognition*, vol. 37, no. 9, pp. 1757–1771, 2004.
 - [42] A. Elisseeff and J. Weston, “A kernel method for multi-labelled classification,” in *Advances in neural information processing systems*, 2002, pp. 681–687.
 - [43] N. Lu, G. Niu, A. K. Menon, and M. Sugiyama, “On the minimal supervision for training any binary classifier from only unlabeled data,” in *International Conference on Learning Representations*, 2018.
 - [44] J. Read, B. Pfahringer, G. Holmes, and E. Frank, “Classifier chains for multi-label classification,” *Machine learning*, vol. 85, no. 3, p. 333, 2011.
 - [45] H. Yu, P. Jain, P. Kar, and I. S. Dhillon, “Large-scale multi-label learning with missing labels,” in *Proceedings of the 31th International Conference on Machine Learning*, 2014, pp. 593–601.
 - [46] T. Liu and D. Tao, “Classification with noisy labels by importance reweighting,” *IEEE Transactions on pattern analysis and machine intelligence*, vol. 38, no. 3, pp. 447–461, 2015.
 - [47] M. Ledoux and M. Talagrand, *Probability in Banach Spaces: isoperimetry and processes*. Springer Science & Business Media, 2013.
 - [48] M. Mohri, A. Rostamizadeh, and A. Talwalkar, *Foundations of machine learning*. MIT press, 2018.
 - [49] W. Gao and Z.-H. Zhou, “On the consistency of multi-label learning,” *Artificial Intelligence*, vol. 199, no. 1, pp. 22–44, 2013.
 - [50] T. Zhang, “Statistical behavior and consistency of classification methods based on convex risk minimization,” *Annals of Statistics*, pp. 56–85, 2004.
 - [51] P. L. Bartlett, M. I. Jordan, and J. D. McAuliffe, “Convexity, classification, and risk bounds,” *Journal of the American Statistical Association*, vol. 101, no. 473, pp. 138–156, 2006.
 - [52] T. Zhang, “Statistical analysis of some multi-category large margin classification methods,” *Journal of Machine Learning Research*, vol. 5, no. Oct, pp. 1225–1251, 2004.
 - [53] C. Finn, P. Abbeel, and S. Levine, “Model-agnostic meta-learning for fast adaptation of deep networks,” PMLR, pp. 1126–1135, 2017.
 - [54] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, “The pascal visual object classes (voc) challenge,” *International journal of computer vision*, vol. 88, no. 2, pp. 303–338, 2010.
 - [55] T. Chen, Z. Wang, G. Li, and L. Lin, “Recurrent attentional reinforcement learning for multi-label image recognition,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.
 - [56] Z.-M. Chen, X.-S. Wei, P. Wang, and Y. Guo, “Multi-label image recognition with graph convolutional networks,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5177–5186.
 - [57] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2015.
 - [58] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” *Advances in neural information processing systems*, vol. 25, pp. 1097–1105, 2012.
 - [59] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, “Pytorch: An imperative style, high-performance deep learning library,” *arXiv preprint arXiv:1912.01703*, 2019.
 - [60] G. Yu, X. Chen, C. Domeniconi, J. Wang, Z. Li, Z. Zhang, and X. Wu, “Feature-induced partial multi-label learning,” in *2018 IEEE International Conference on Data Mining (ICDM)*. IEEE, 2018, pp. 1398–1403.
 - [61] J. Demsar, “Statistical comparisons of classifiers over multiple data sets,” *J. Mach. Learn. Res.*, vol. 7, pp. 1–30, 2006.