

Multi-Label Inference for Crowdsourcing

Jing Zhang

School of Computer Science and Engineering
 Nanjing University of Science and Technology
 200 Xiaolingwei Street
 Nanjing, Jiangsu 210094, China
 jzhang@njust.edu.cn

Xindong Wu

School of Computing and Informatics
 University of Louisiana at Lafayette
 222 James R. Oliver Hall
 Lafayette, LA 70504, U.S.
 xwu@louisiana.edu

ABSTRACT

When acquiring labels from crowdsourcing platforms, a task may be designed to include multiple labels and the values of each label may belong to a set of various distinct options, which is the so-called multi-class multi-label annotation. To improve the quality of labels, one task is independently completed by a group of heterogeneous crowdsourced workers. Then, the true values of the multiple labels of each task are inferred from these repeated noisy labels. In this paper, we propose a novel probabilistic method, which includes a *multi-class multi-label dependency* (MCMLD) model, to address this problem. The proposed method assumes that the label-correlation exists in both unknown true labels and noisy crowdsourced labels. Thus, it introduces a mixture of multiple independently multinoulli distributions to capture the correlation among the labels. Finally, the unknown true values of the multiple labels of each task, together with a set of confusion matrices modeling the reliability of the workers, can be jointly inferred through an EM algorithm. Experiments with three simulated typical crowdsourcing scenarios and a real-world dataset consistently show that our proposed MCMLD method significantly outperforms several competitive alternatives. Furthermore, if the labels are strongly correlated, the advantage of MCMLD will be more remarkable.

CCS CONCEPTS

- Computing methodologies → Maximum likelihood modeling; Mixture models; Latent variable models;
- Information systems → Crowdsourcing; Information integration;

KEYWORDS

Crowdsourcing, label aggregation, maximum likelihood estimation, mixture models, probabilistic graphical models

ACM Reference Format:

Jing Zhang and Xindong Wu. 2018. Multi-Label Inference for Crowdsourcing. In *KDD '18: The 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, August 19–23, 2018, London, United Kingdom. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3219819.3219958>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

KDD '18, August 19–23, 2018, London, United Kingdom

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 978-1-4503-5552-0/18/08...\$15.00
<https://doi.org/10.1145/3219819.3219958>

1 INTRODUCTION

Current intelligent systems usually require plenty of human-labeled datasets to build internal models. Human-annotation tasks are often too tedious and time-consuming for individuals to complete. The emergence of crowdsourcing systems provides a convenient and low-cost solution to obtain annotations from the Internet. Crowd-sourced annotation acquisition has been used in some important fields (such as data mining [35], machine learning [37], computer vision [14], information retrieval [16], and natural language processing [32]) in recent years. However, a significant challenge of using crowdsourcing to acquire labels is that their quality cannot be guaranteed due to the non-expert nature of crowdsourcing. To improve the quality of labels, we usually let each instance be labeled by multiple different heterogeneous workers, which is so-called *repeated labeling*. After sufficient crowdsourced labels are collected, an inference algorithm will be introduced to estimate the true labels of each instance from these noisy labels. This process is also known as *label aggregation* since the estimated true labels are the integration of multiple noisy labels. In addition to the integrated labels, this process may provide more information, such as the estimation of the workers' reliability. In crowdsourcing, the term *truth inference* is often used as a synonym of *label aggregation* as in this paper, but we should know that using *inference* implies that the process involves probabilistic modeling. A general inference algorithm in crowdsourcing should be agnostic, assuming that no prior knowledge can be exploited, because one of its objectives is to provide information for further usages, such as spammer detection [24] and incentive optimization [9].

Truth inference for crowdsourcing started from the study of majority voting [26]—a simple but efficient method. Subsequently, a large number of studies addressed this issue using different techniques [37]. Zheng et al. [39] reviewed ten-year efforts in this field by an empirical study, and observed that dozens of inference algorithms only focus on the *single-label* scenario, including *binary-class* and *multi-class* cases. However, in certain areas, it is quite common that each task associates with a set of appropriate labels. Figure 1 illustrates two examples of such applications. In the first application (a), workers are required to provide their judgments to the presence of seven emotions conveyed through the passages. The workers only need to take a binary option—yes or no—as with the problem settings in traditional *multi-label* classification [8]. A more general situation is that each label has more than two options (values) as the second application (b) illustrates. In this example, a description of an image contains three placeholders (labels) representing an animal (*dog*), an object (*soccer*), and a place (*lawn*), respectively. All of them have multiple options. We refer this application as *multi-class*

Instruction: read the passage and make judgments to the expression of each emotion.

Passage 1:

Martin said that 23 people were found dead inside the church, two were found dead outside the church and that another person who had been transported from the scene had died. Founded in 1926, the church has 267 members and has an average attendance of around 180 people, according to a listing on the Southern Baptist Convention website. "Everybody has a connection to that church," Wilson County Commissioner Albert Gamez Jr. said on "TODAY." "It's a family town. Anywhere you go, everybody recognizes everybody." Sherri Pomeroy, the wife of church pastor Frank Pomeroy, told NBC News in a statement that their daughter was among the victims. She said she and her husband were out of town in different states at the time of the shooting.

Emotions expressed by the content:

Anger	Calm	Disgust	Fear	Happy	Sadness	Surprise
<input checked="" type="radio"/> Yes	<input type="radio"/> Yes	<input type="radio"/> Yes	<input checked="" type="radio"/> Yes	<input type="radio"/> Yes	<input checked="" type="radio"/> Yes	<input type="radio"/> Yes
<input type="radio"/> No	<input checked="" type="radio"/> No	<input type="radio"/> No	<input type="radio"/> No	<input type="radio"/> No	<input type="radio"/> No	<input type="radio"/> No

Application (a): Sentiment Annotation

Instruction: choose a proper description of the image.



cat
dog
dolphin
hamster
monkey

basketball
plastic ball
soccer
tennis
volleyball

beach
creek
lawn
hill
sea

There is a playing a on (by, in, ...) the .

Application (b): Image Annotation

Figure 1: Two examples involving multi-label annotation in crowdsourcing. (a) In the sentiment judgment task, each label exactly has two choices. (b) In the image annotation task, each label has multiple choices.

multi-label annotation in this paper. The traditional multi-class learning usually transforms a multi-class dataset to various binary ones. However, in crowdsourcing, it will be more user-friendly and efficient if multi-class and binary-class multi-label inferences are handled uniformly.

A critical difference between the multi-label inference and traditional single-label inference in crowdsourcing lies in the exploration and exploitation of the correlation among the labels. For example, in the above sentiment annotation application, negative emotions such as *fear*, *disgust*, and *sadness* are usually strongly correlated. Leveraging the potential correlation among labels may improve correctness of the integrated labels on instances. This paper proposes a novel method to solve the multi-class multi-label inference problem in crowdsourcing. Within an EM procedure, it jointly infers the unknown true values for the multiple labels of instances as well as the reliability of each crowdsourced worker parameterized by a set of confusion matrices on multiple labels. During the truth inference, the correlation of the labels is captured and exploited using a mixture of independent multinoulli distributions. Both simulations on six synthetic datasets and evaluations on a real-world dataset confirm that our new method not only outperforms other single-label alternatives which successively perform inference on each label, but also outperforms our proposed baseline method that also jointly infers multiple labels.

2 RELATED WORK

Inferring true values of labels from multiple noisy labels is a fundamental research topic in crowdsourcing. Many previous studies modeled the complexity of crowdsourced labeling systems from different aspects. The earliest Dawid & Skene's (DS) algorithm [4] used a confusion matrix to model the reliability of uncertain sources. Some recent approaches are directly derived from the DS model by simplifying the parameters [25], adding constraints to workers [31] or optimizing initial settings [38]. Other approaches model the reliability of workers using fewer parameters [1, 5, 12, 34] but add additional parameters or variables to characterize workers' intention [1, 15] and bias [11, 33]. All above approaches are based on the probabilistic models and can be solved by an EM algorithm, gradient

descending, or Gibbs sampling. Some of them [1, 4, 5, 11, 15, 38] can be used in multi-class scenarios, while the others [12, 25, 31, 33, 34] are only suitable for binary cases, but none of them were proposed for multi-label annotation. In addition to the approaches based on the probabilistic models, some other studies introduced a variety of techniques to improve the accuracy of inference, including optimization [17, 18, 30, 40], clustering [36], and the injection of expert validation [20].

Facing the diversity of the inference algorithms, we carefully choose among them according to some previous comprehensive empirical studies. Sheshadri et al. [27] developed a benchmark tool SQUARE to evaluate six inference algorithms on ten real-world crowdsourcing datasets under both supervised and unsupervised paradigms. Muhammadi et al. [22] studied the performance of those algorithms on a set of synthetic datasets in a unified statistical framework. A most recent study [39] evaluated the performance of fourteen inference algorithms in terms of accuracy, F1-score, mean absolute error, and root mean square error on many real-world datasets. All these studies have shown some consistent conclusion that the approaches derived from the Dawid & Skene's model [4] have good probabilistic interpretability and exhibit strong robustness in most cases. Thus, our study also adopts the mainstream probabilistic models and uses a set of confusion matrices to model the reliability of crowdsourced workers. Compared with the single-label cases, the parameters of our multi-label model will greatly increase.

Several studies have touched the topic of multi-label annotation in crowdsourcing. Nowak et al. [23] studied the reliability of workers when they performed multi-label image annotation tasks, compared with expert annotation. Their study focused on the differences in annotation quality among multiple experts and non-experts workers when performing multi-label annotation tasks. As multi-label tasks consume more budgets under the repeated labeling scenario, Deng et al. [6] and Bragg et al. [3] designed some optimal annotation acquisition workflows that leverage the correlation among different labels. Hung et al. [10] proposed a Bayesian non-parametric model to aggregate the labels, which takes advantage of Clustering-based Bayesian Combination of Multi-label Classifiers

(cBCMC) [21] to exploit the correlation among labels. However, their model is not entirely unsupervised, requiring a particular portion of ground truths. Li et al. [19] also addressed the inference problem for multi-label annotation using probabilistic models under an active learning paradigm but ignoring the correlation among labels. In this paper, our proposed unsupervised method takes into account the correlation among labels and jointly infer the true values for all labels of each instance.

3 PROBLEM STATEMENT

In this section, we first formalize multi-label annotation in crowdsourcing and then discuss three critical issues to be addressed in this paper.

3.1 Multi-Label Annotation

The dataset posted on a crowdsourcing system to collect repeated noisy labels is denoted by $\mathcal{D} = \{\langle \mathbf{x}_i, \mathbf{y}_i \rangle\}_{i=1}^I$, where each instance \mathbf{x}_i associates with an unknown true class label set \mathbf{y}_i . Each \mathbf{y}_i consists of M labels, which is denoted by $\mathbf{y}_i = [y_i^{(1)}, \dots, y_i^{(M)}]$. All true label sets of the entire dataset are denoted by $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_I\}$. In a real-world application, because all instances are usually created from the same template, the present order of the labels on each instance is the same. Thus, we can simply say " $y_i^{(m)}$ is the m th label of instance \mathbf{x}_i ". For different labels, they usually have different numbers of values because each label reflects a specific aspect of the instance. For example, in Figure 1.(b), the first label (*animal*) may include dozens of values, but the third label (*place*) may include only several values. To simplify the formalization, we let each label associate with exactly K values, i.e., $y_i^{(m)} \in \{1, \dots, K\}$, where K is the maximum number of values that the labels have. This does not have any side effect in practice. For example, given $K = 5$, if the m th label only has three values, we can treat the $y_i^{(m)} = 4$ and $y_i^{(m)} = 5$ as two unobserved values. A probabilistic model will guarantee that all probabilities with respect to the unobserved values are always zero (i.e., $P(y_i^{(m)} = 4) = P(y_i^{(m)} = 5) = 0$), which does not affect the results of inference. Note that, in this simplification, the $k \in \{1, \dots, K\}$ is just an index of that value because values of different labels represent different things. However, we can easily make a discrimination of them by adding superscripts, i.e., $k^{(m)} \neq k^{(n)}$ if $m \neq n$.

There are totally J crowdsourced workers labeling the dataset. All collected noisy labels for each instance \mathbf{x}_i form a matrix $L_i^{J \times M} = [\mathbf{l}_{i1}, \dots, \mathbf{l}_{iJ}]^T$, where each element (noisy label) $l_{ij}^{(m)} = k \in \{0\} \cup \{1, \dots, K\}$ in a row vector \mathbf{l}_{ij} represents that worker j provides value k to the m th label of instance \mathbf{x}_i . Particularly, $l_{ij}^{(m)} = 0$ means that worker j does not provide any value to the m th label of instance \mathbf{x}_i . All crowdsourced labels of the entire dataset are denoted by $\mathbf{L} = \{L_1, \dots, L_I\}$.

3.2 Three Critical Issues

Our study aims to solve three critical issues in multi-class multi-label inference for crowdsourcing.

Issue 1: Jointly inferring truth labels. The primary objective of an inference algorithm for crowdsourcing is to infer the true

labels $\mathbf{Y}^{I \times M}$ for all instances from the collected noisy labels $\mathbf{L}^{I \times J \times M}$ and simultaneously minimize the overall inference error rate as follows:

$$\epsilon = \min \left\{ \frac{1}{I \cdot M} \sum_{i=1}^I \sum_{m=1}^M \mathbb{I}(\hat{y}_i^{(m)} \neq y_i^{(m)}) \right\}, \quad (1)$$

where \mathbb{I} is an indicator function. It outputs 1 when the test condition satisfies. Otherwise, it outputs 0.

A simple approach to solve this issue is successively running a *single-label* inference algorithm on each label of all instances. That is, the *single-label* inference algorithm totally runs M rounds. In the m th round, the inference error rate on the m th label is minimized. However, there exists a question that has never been answered. If we can jointly infer all M labels of each instance at once, will the total error rate be less than that obtained in the above manner?

Issue 2: Exploring and exploiting the correlation among the labels. In multi-label annotation, the correlation among the labels exhibits great complexity. On the one hand, there exists correlation among the true labels. On the other hand, this correlation, in turn, has some impact on the distribution of the crowdsourced labels. Workers can mutually check their own answers on different labels to improve the label quality if they are aware of the correlation among the true labels.

Considering instance \mathbf{x}_i and the J workers labeling it, if we use a joint probability of $P(\mathbf{l}_i, \mathbf{y}_i)$ to model the correlation among both crowdsourced and true labels, and use the chain rule of conditional probability to solve it, it will result in a complexity of $O(c \cdot I \cdot J \cdot K^M \cdot (K^J)^M)$ in total, which is computationally infeasible. Therefore, the second question is "Is there a simple method to model the correlation among the labels and does it improve the inference accuracy?"

Issue 3: Modeling the reliability of the workers. Estimates of the reliability of crowdsourced workers is an important function of an inference algorithm. A multi-class multi-label inference algorithm must model the reliability of the workers on M labels with $M \cdot K$ values. Furthermore, the upper-layer applications require the inference algorithm to provide fine-grained information as far as possible.

4 THE PROPOSED METHOD

In this section, we first present a novel probabilistic model for multi-label inference. Then, we solve the model with an EM algorithm. Finally, we present our inference algorithm.

4.1 Label-Dependent Model

We propose a novel probabilistic generative model, namely *multi-class multi-label dependent* (MCMLD) model, illustrated in Figure 2. In this model, the correlation among the (true) labels is modeled by *soft clusters*. The instances belonging to the same cluster suggest that their labels exhibit some correlation. We assign each instance \mathbf{x}_i a latent variable z_i to indicate its cluster membership. Given the total number R of clusters, the probability of instance \mathbf{x}_i being in cluster r ($1 \leq r \leq R$) is $P(z_i = r) = \omega_r$. Thus, for each instance \mathbf{x}_i , z_i is drawn from a multinoulli distribution, i.e.,

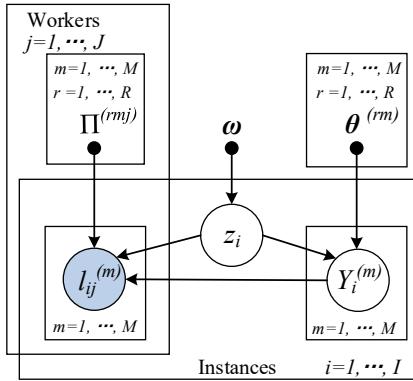


Figure 2: Probabilistic graphical model representation of the proposed MCMLD model.

$P(z_i|\omega) = \prod_{r=1}^R (\omega_r)^{\mathbb{I}(z_i=r)}$, where $\omega = [\omega_1, \dots, \omega_R]$ are its parameters. Obviously, a soft cluster means that each instance belongs to cluster r with a probability ω_r and $\sum_{r=1}^R \omega_r = 1$.

Generation of true labels. For multi-class classification, each true label (supposed to be the m th label in multi-label annotation) is independently drawn from a multinoulli distribution with parameters $\theta^{(m)} = [\theta_1^{(m)}, \dots, \theta_K^{(m)}]$, where $\sum_{k=1}^K \theta_k^{(m)} = 1$. That is, for the m th label of an instance x_i , we have $P(y_i^{(m)}|\theta^{(m)}) = \prod_{k=1}^K \left(\theta_k^{(m)}\right)^{\mathbb{I}(y_i^{(m)}=k)}$. The probability of the M -dimensional true label vector y_i of the instance is

$$P(y_i|\Theta) = \prod_{m=1}^M \prod_{k=1}^K \left(\theta_k^{(m)}\right)^{\mathbb{I}(y_i^{(m)}=k)},$$

where $\Theta = [\theta^{(1)}, \dots, \theta^{(M)}]$ is a set of parameters for all M multinoulli distributions. In MCMLD, each true label vector y_i is generated from R independent clusters with the probabilities $[\omega_1, \dots, \omega_R]$, which is a *mixture of multiple independent multinoulli* distributions. Thus, the probability of y_i of instance x_i has the form

$$P(y_i|\Theta, \omega) = \sum_{r=1}^R \omega_r \prod_{m=1}^M \prod_{k=1}^K \left(\theta_{rk}^{(m)}\right)^{\mathbb{I}(y_i^{(m)}=k)}, \quad (2)$$

where $\omega = [\omega_1, \dots, \omega_R]$ are also called the mixing coefficients of multiple multinoulli distributions with the parameter set $\Theta = [\Theta_1, \dots, \Theta_R]$.

Generation of crowdsourced labels. In multi-class classification, a confusion matrix is a powerful tool that can comprehensively depict the distribution of a classifier's capability over all pairs of classes, providing fine-grained information to upper-layer applications. We use a set of confusion matrices $\Pi^{(j)} = [\Pi^{(1j)}, \dots, \Pi^{(Mj)}]$ to model the reliability of worker j with respect to M labels. We denote all sets of confusion matrices of totally J workers by a parameter set $\tilde{\Pi} = \{\Pi^{(1)}, \dots, \Pi^{(J)}\}$. In matrix $\Pi^{(mj)}$, each element $\pi_{kd}^{(mj)}$ ($1 \leq k, d \leq K$) represents the probability of worker j labeling (true) class k as class d on the m th label, which derives $P(l_{ij}^{(m)} = d|y_i^{(m)} = k) = \pi_{kd}^{(mj)}$. That is, $l_{ij}^{(m)}$ conditioning

on $y_i^{(m)} = k$ obeys a multinoulli distribution with parameters $[\pi_{kd}^{(mj)}]_{d=1}^K$ and $\sum_{d=1}^K \pi_{kd}^{(mj)} = 1$.

Consider each instance x_i is independently labeled by J workers, the likelihood of all observed noisy labels on the instance can be calculated as follows:

$$\begin{aligned} P(L_i|\Theta, \tilde{\Pi}, \omega) &= \sum_{k^{(1)}=1}^{k^{(1)}=K} \cdots \sum_{k^{(M)}=1}^{k^{(M)}=K} \left[\right. \\ &\quad P\left(y_i^{(1)} = k^{(1)}, \dots, y_i^{(M)} = k^{(M)} | \Theta, \omega\right) \\ &\quad \cdot P\left(L_i | y_i^{(1)} = k^{(1)}, \dots, y_i^{(M)} = k^{(M)}, \tilde{\Pi}\right) \left. \right]. \end{aligned} \quad (3)$$

In this equation, because MCMLD is a mixture model, we sum up all the probability of latent variable z_i , so that z_i disappears. According to the *conditional independency* principle in probabilistic graphical models, adding a latent variable z_i among the unobserved variables y_i will make them conditional mutually independent, which derives

$$P\left(y_i^{(1)} = k^{(1)}, \dots, y_i^{(M)} = k^{(M)} | \Theta, \omega\right) = \sum_{r=1}^R \omega_r \prod_{m=1}^M \theta_{rk}^{(m)}, \quad (4)$$

$$\begin{aligned} P\left(L_i | y_i^{(1)} = k^{(1)}, \dots, y_i^{(M)} = k^{(M)}, \tilde{\Pi}\right) \\ = \prod_{j=1}^J \left(\prod_{m=1}^M \prod_{d(m)=1}^K \left(\pi_{k^{(m)} d^{(m)}}^{(mj)}\right)^{\mathbb{I}(l_{ij}^{(m)}=d^{(m)})} \right). \end{aligned} \quad (5)$$

Here, because we need to represent the prior of y_i under the conditions that each of its elements is assigned with a different value, we use the superscript (m) to distinguish these ks on different labels. (Refer to Section 3.1 for the physical meaning of the superscript.) Similarly, the variable d is also decorated by the superscript (m) .

Finally, the log-likelihood of all crowdsourced labels of the entire dataset is

$$\begin{aligned} \ln P(L|\Theta, \omega, \tilde{\Pi}) &= \sum_{i=1}^I \ln \left\{ \sum_{k^{(1)}=1}^{k^{(1)}=K} \cdots \sum_{k^{(M)}=1}^{k^{(M)}=K} \left[\sum_{r=1}^R \omega_r \prod_{m=1}^M \theta_{rk}^{(m)} \right. \right. \\ &\quad \cdot \left. \prod_{j=1}^J \left(\prod_{m=1}^M \prod_{d(m)=1}^K \left(\pi_{k^{(m)} d^{(m)}}^{(mj)}\right)^{\mathbb{I}(l_{ij}^{(m)}=d^{(m)})} \right) \right] \left. \right\}. \end{aligned} \quad (6)$$

4.2 Inference with EM

To maximize the likelihood of the observed labels L defined by Eq.(6), we use an expectation-maximization (EM) algorithm. The difficulty of applying EM to MCMLD lies in that there are two types of latent variables Y and z . (Traditionally, a mixture model only has one type of latent variables.)

E-step. We first construct an expectation function of the log-likelihood of the observed data L conditioning on latent variables Y and z as follows:

$$Q(\Psi, \Psi^{old}) = \mathbb{E}_{Y, z|L, \Psi^{old}} [\ln P(L, Y, z|\Psi)], \quad (7)$$

where $\Psi = \{\Theta, \tilde{\Pi}, \omega\}$.

For the instance x_i with two types of latent variables y_i and z_i , when the crowdsourced labels L are observed, we must calculate

the joint posterior probability of \mathbf{y}_i and z_i by Bayes' theorem [2] as follows:

$$\begin{aligned} & P\left(y_i^{(1)} = k^{(1)}, \dots, y_i^{(M)} = k^{(M)}, z_i = r \mid \mathbf{L}, \Psi\right) \\ & \propto P\left(L_i \mid y_i^{(1)} = k^{(1)}, \dots, y_i^{(M)} = k^{(M)}, z_i = r, \Psi\right) \\ & \quad \cdot P\left(y_i^{(1)} = k^{(1)}, \dots, y_i^{(M)} = k^{(M)}, z_i = r \mid \Theta, \omega\right) \\ & = \omega_r \prod_{m=1}^M \theta_{r k^{(m)}}^{(m)} \prod_{j=1}^J \prod_{d^{(m)}=1}^K \left(\pi_{k^{(m)} d^{(m)}}^{(m)}\right)^{\mathbb{I}(l_{ij}^{(m)}=d^{(m)})}. \end{aligned} \quad (8)$$

Eq.(8) provides the expected values of two indicator functions (calculated as two marginal probabilities) that will be used in the M-step as follows:

$$\mathbb{E}\left[\mathbb{I}\left(y_i^{(m)} = k\right)\right] = \sum_{y_i^{(m')} (\forall m' \in \{1, \dots, M\} \setminus m), z_i} P\left(y_i^{(m)} = k, \dots \mid \mathbf{L}, \Psi\right), \quad (9)$$

$$\mathbb{E}\left[\mathbb{I}\left(z_i = r\right)\right] = \sum_{y_i^{(m)} (\forall m \in \{1, \dots, M\})} P\left(z_i = r, \dots \mid \mathbf{L}, \Psi\right). \quad (10)$$

M-step. We determine the revised parameter estimate Ψ by maximizing the objective function Eq.(7), i.e., $\Psi^{new} = \text{argmax}_{\Psi}(\Psi, \Psi^{old})$. Because we have obtained the posterior probabilities of the latent variables \mathbf{Y} and \mathbf{z} , we omit the condition of \mathbf{L} and Ψ^{old} and eliminate the constant factors, i.e.,

$$\begin{aligned} \mathbb{E}_{\mathbf{Y}, \mathbf{z}} [\ln P(\mathbf{L}, \mathbf{Y}, \mathbf{z} \mid \Psi)] &= \sum_{i=1}^I \mathbb{E}_{\mathbf{Y}, \mathbf{z}} [\ln P(L_i \mid \mathbf{y}_i, z_i, \Pi_i)] \\ &= \sum_{i=1}^I \mathbb{E}_{\mathbf{Y}, \mathbf{z}} \left[\ln P(L_i \mid \mathbf{y}_i, z_i, \Pi_i) \underbrace{P(\mathbf{y}_i \mid z_i, \Theta) P(z_i \mid \omega)}_{\text{constants w.r.t. } \Psi \text{ in derivatives}} \right]. \end{aligned}$$

Thus, we only need to maximize $\sum_{i=1}^I \mathbb{E}_{\mathbf{Y}, \mathbf{z}} [\ln P(L_i \mid \mathbf{y}_i, z_i, \Pi_i)]$, which derives

$$\begin{aligned} & \sum_{i=1}^I \mathbb{E}_{\mathbf{Y}, \mathbf{z}} \left\{ \ln \prod_{r=1}^R \left[\omega_r \prod_{m=1}^M \prod_{k=1}^K \left(\theta_{r k}^{(m)} \right. \right. \right. \\ & \quad \cdot \left. \left. \left. \prod_{j=1}^J \prod_{d=1}^K \left(\pi_{k d}^{(m)} \right)^{\mathbb{I}(l_{ij}^{(m)}=d)} \right)^{\mathbb{I}(y_i^{(m)}=k)} \right]^{\mathbb{I}(z_i=r)} \right\} \\ &= \sum_{i=1}^I \sum_{r=1}^R \mathbb{E}[\mathbb{I}(z_i = r)] \left\{ \ln \omega_r + \sum_{m=1}^M \sum_{k=1}^K \mathbb{E}\left[\mathbb{I}\left(y_i^{(m)} = k\right)\right] \right. \\ & \quad \cdot \left. \left[\ln \theta_{r k}^{(m)} + \sum_{j=1}^J \sum_{d=1}^K \mathbb{I}(l_{ij}^{(m)}=d) \ln \pi_{k d}^{(m)} \right] \right\}. \end{aligned} \quad (11)$$

(1) Using a Lagrange multiplier to optimize Eq.(11) with respect to $\theta_{r k}^{(m)}$, we construct a function

$$\mathcal{F}_1 = \sum_{i=1}^I \mathbb{E}_{\mathbf{Y}, \mathbf{z}} [\ln P(L_i \mid \mathbf{y}_i, z_i, \Pi_i)] + \lambda \left(\sum_{k=1}^K \theta_{r k}^{(m)} - 1 \right). \quad (12)$$

We let the partial derivative of Eq.(12) with respect to $\theta_{r k}^{(m)}$ be zero. Applying sum-up-to-one condition $\left(\sum_{k=1}^K \theta_{r k}^{(m)} = 1\right)$, we have

$$\frac{\partial \mathcal{F}_1}{\partial \theta_{r k}^{(m)}} = \frac{\sum_{i=1}^I \left\{ \mathbb{E}[\mathbb{I}(z_i = r)] \mathbb{E}\left[\mathbb{I}\left(y_i^{(m)} = k\right)\right] \right\}}{\theta_{r k}^{(m)}} + \lambda = 0 \quad (13)$$

$$\begin{aligned} & \Rightarrow \sum_{k=1}^K \sum_{i=1}^I \left\{ \mathbb{E}[\mathbb{I}(z_i = r)] \mathbb{E}\left[\mathbb{I}\left(y_i^{(m)} = k\right)\right] \right\} = -\lambda \sum_{k=1}^K \theta_{r k}^{(m)} \\ & \Rightarrow \sum_{i=1}^I \mathbb{E}[\mathbb{I}(z_i = r)] = -\lambda. \end{aligned} \quad (14)$$

Plugging Eq.(14) into Eq.(13), we have

$$\hat{\theta}_{r k}^{(m)} = \frac{\sum_{i=1}^I \left\{ \mathbb{E}[\mathbb{I}(z_i = r)] \mathbb{E}\left[\mathbb{I}\left(y_i^{(m)} = k\right)\right] \right\}}{\sum_{i=1}^I \mathbb{E}[\mathbb{I}(z_i = r)]}. \quad (15)$$

(2) Using a Lagrange multiplier to optimize Eq.(11) with respect to $\pi_{k d}^{(m)}$, we construct a function

$$\mathcal{F}_2 = \sum_{i=1}^I \mathbb{E}_{\mathbf{Y}, \mathbf{z}} [\ln P(L_i \mid \mathbf{y}_i, z_i, \Pi_i)] + \lambda \left(\sum_{k=1}^K \pi_{k d}^{(m)} - 1 \right). \quad (16)$$

We let the partial derivative of Eq.(16) with respect to $\pi_{k d}^{(m)}$ be zero. Applying sum-up-to-one condition $\left(\sum_{d=1}^K \pi_{k d}^{(m)} = 1\right)$, we have

$$\frac{\partial \mathcal{F}_2}{\partial \pi_{k d}^{(m)}} = \frac{\sum_{i=1}^I \left\{ \mathbb{E}\left[\mathbb{I}\left(y_i^{(m)} = k\right)\right] \mathbb{E}\left[\mathbb{I}\left(l_{ij}^{(m)} = d\right)\right] \right\}}{\pi_{k d}^{(m)}} + \lambda = 0 \quad (17)$$

$$\begin{aligned} & \Rightarrow \sum_{d=1}^K \sum_{i=1}^I \left\{ \mathbb{E}\left[\mathbb{I}\left(y_i^{(m)} = k\right)\right] \mathbb{E}\left[\mathbb{I}\left(l_{ij}^{(m)} = d\right)\right] \right\} = -\lambda \sum_{d=1}^K \pi_{k d}^{(m)} \\ & \Rightarrow \sum_{i=1}^I \mathbb{E}\left[\mathbb{I}\left(y_i^{(m)} = k\right)\right] = -\lambda. \end{aligned} \quad (18)$$

Plugging Eq.(18) into Eq.(17), we have

$$\hat{\pi}_{k d}^{(m)} = \frac{\sum_{i=1}^I \left\{ \mathbb{E}\left[\mathbb{I}\left(y_i^{(m)} = k\right)\right] \mathbb{E}\left[\mathbb{I}\left(l_{ij}^{(m)} = d\right)\right] \right\}}{\sum_{i=1}^I \mathbb{E}\left[\mathbb{I}\left(y_i^{(m)} = k\right)\right]}. \quad (19)$$

(3) Using a Lagrange multiplier to optimize Eq.(11) with respect to ω_r , we construct a function

$$\mathcal{F}_3 = \sum_{i=1}^I \mathbb{E}_{\mathbf{Y}, \mathbf{z}} [\ln P(L_i \mid \mathbf{y}_i, z_i, \Pi_i)] + \lambda \left(\sum_{r=1}^R \omega_r - 1 \right). \quad (20)$$

We let the partial derivative of Eq.(20) with respect to ω_r be zero. Applying sum-up-to-one condition $\left(\sum_{r=1}^R \omega_r = 1\right)$, we have

$$\frac{\partial \mathcal{F}_3}{\partial \omega_r} = \frac{\sum_{i=1}^I \mathbb{E}[\mathbb{I}(z_i = r)]}{\omega_r} + \lambda = 0, \quad (21)$$

$$\Rightarrow \sum_{r=1}^R \sum_{i=1}^I \mathbb{E}[\mathbb{I}(z_i = r)] = -\lambda \sum_{r=1}^R \omega_r \Rightarrow I = -\lambda. \quad (22)$$

Plugging Eq.(22) into Eq.(21), we have

$$\hat{\omega}_r = \sum_{i=1}^I \mathbb{E}[\mathbb{I}(z_i = r)] / I. \quad (23)$$

After convergence. The integrated labels of instance \mathbf{x}_i should be the sequence of class values with the maximum marginal posterior probability. That is,

$$\begin{aligned} \hat{y}_i^{(1)}, \dots, \hat{y}_i^{(M)} &= \operatorname{argmax}_{k^{(1)}, \dots, k^{(M)}} \left\{ \sum_{r=1}^R P(y_i^{(1)} = k^{(1)}, \right. \\ &\quad \left. \dots, y_i^{(M)} = k^{(M)}, z_i = r | \mathbf{L}, \Psi \right\}. \end{aligned} \quad (24)$$

The *hard* value of z_i also can be determined like Eq.(24).

Algorithm 1 The MCMLD Inference for Crowdsourcing

Input: Multi-Class Multi-Label Crowdsourced Dataset \mathcal{D}, ρ

Output: Integrated labels $\hat{\mathbf{y}}_i$ for each instance \mathbf{x}_i

- 1: Create a two-dimensional matrix $U^{(I,J) \times M}$ from \mathcal{D} .
 - 2: Apply PCA on U , obtaining normalized eigenvalues of the covariance matrix $\boldsymbol{\eta} = \{\eta_1, \dots, \eta_N\}$ in descending order.
 - 3: Set $R = \operatorname{argmin}_{1 \leq n \leq N} |\rho - \sum_{b=1}^n \eta_b|$, where $\rho \in (0, 1]$.
 - 4: **repeat**
 - 5: E-step: calculate expectation by Eqs.(9) and (10).
 - 6: M-step: update parameters by Eqs.(15), (19) and (23).
 - 7: **until** convergence
 - 8: For each instance \mathbf{x}_i , calculate $\hat{\mathbf{y}}_i$ by Eq.(24).
 - 9: **return** all $\hat{\mathbf{y}}_i$ s.
-

4.3 Algorithm

To apply the proposed mixture model into a practical usage, we propose a novel algorithm, where a critical issue of setting the value of parameter R under an unsupervised scenario is addressed by exploring the correlation among the crowdsourced labels. First, it combines all crowdsourced labels into a two-dimensional matrix $U^{(I,J) \times M}$. Each row in U can be viewed as a data point, whose features consist of M labels provided by a worker. Then, we apply the principal components analysis (PCA) algorithm on matrix U . The dimension of the features will be reduced from M to N ($N \leq M$), which stands for the number of irrelevant labels. In practice, we use the normalized eigenvalues ($\boldsymbol{\eta}$) of the covariance matrix in PCA, which present the contributions of the components. If $\boldsymbol{\eta}$ is in the descending order, step 3 in the algorithm will determine the value of R , which means that R components contribute approximately ρ portion of the variance. Note that setting ρ is much easier than setting R . For example, in our experiments, we simply deem that if the R components contribute 70% of the variance, they are irrelevant.

5 EXPERIMENTS

In this section, we evaluate the performance of the proposed MCMLD method on both simulated and real-world datasets.

5.1 Metric and algorithms in comparisons

Evaluation metric. Because our primary objective is to accurately infer the true values of all labels on each instance, we merely use

the label-based accuracy (i.e., the inverse of Eq.(1)) as the evaluation metric.

Algorithms in comparisons. To the best of our knowledge, there is no entirely unsupervised multi-class multi-label annotation inference algorithm for crowdsourcing. Therefore, we choose four alternatives and another baseline proposed by us to conduct comparisons.

- **MV** is usually served as a baseline for label inference.
- **DS** [4] uses the confusion matrix to model workers and makes inference using maximum likelihood estimation. Previous large-scale empirical studies [22, 27, 39] have shown its robustness and efficiency.
- **SpectralDS** [38] is a variant of DS, where the initial values of the confusion matrices are carefully set.
- **iBCC** [13] is a non-EM based inference method, where a Bayesian generative model of the crowdsourced annotation is solved by Gibbs sampling.
- **MCMLI**¹ is a simple baseline model for multi-class multi-label inference, which can jointly infer the M true labels and reliability of workers. If we remove node z_i (and its links) from the probabilistic graphic model of MCMLD (in Figure 2), we obtain the MCMLI model. It addresses issues 1 and 3 in Section 3.2.

The four alternatives are chosen because they are not only efficient and robust but also suitable for both binary and multi-class cases [39]. For each of these *single-label* algorithms, from the 1st labels of all instances to the M th labels, we repeatedly run the algorithm to infer their true values.

5.2 Evaluation by simulations

The advantage of simulation lies in its ability of revealing the differences among compared algorithms under various predefined conditions, especially, when there are no sufficient real-world datasets to cover different situations.

Datasets for simulation. We selected four *real-world* multi-label classification datasets from the multi-label learning library MULAN², because their labels naturally exhibit some degree of correlation. However, these labels, collected from experts and served as the ground truth, only have binary values, since they were created for the traditional multi-label classification. To make a comprehensive investigation, we created another two datasets *pentacorrel* and *decasynth*, in which each label has three and four values, respectively. Dataset *pentacorrel* was created with strong label-correlation, where the values of the first, third and fifth labels always appear three patterns (1, 2, 3), (2, 3, 1) and (3, 1, 2). For dataset *decasynth*, we relaxed the correlation. The values of labels are approximately uniformed distributed and generated regardless of others. The detailed information of these datasets for simulation is listed in Table 1.

Scenario 1: Annotation errors are uniformly distributed.

In our simulation, the overall accuracy of each worker is independently drawn from a uniform distribution $U(0.60, 0.80)$. Previous

¹The multi-class multi-label independent (MCMLI) model is also proposed by us. Since it does not directly model the correlation, we only use it as a baseline method in this paper. In fact, MCMLI has made an obvious progress as our experiments show.

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

Table 1: Six multi-label datasets for simulations

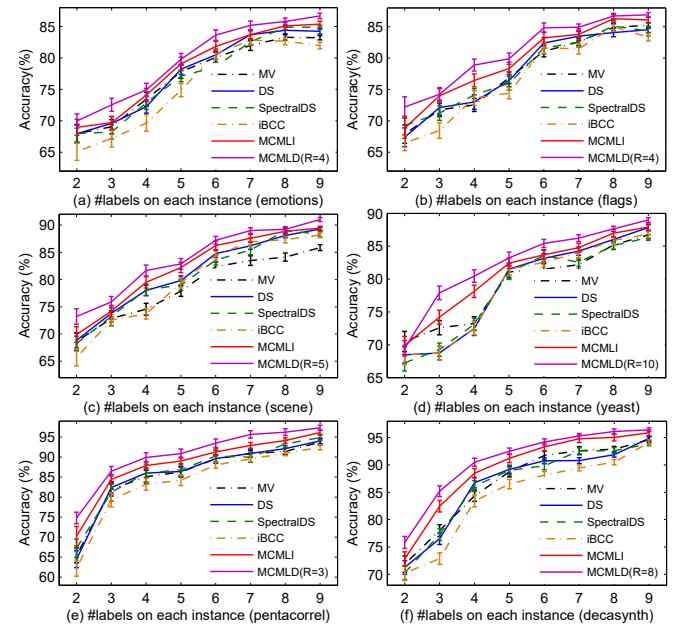
Dataset	#inst.	#lb.	#cls.	card.	dens.
<i>emotions</i>	593	6	2	1.869	0.311
<i>flags</i>	194	7	2	3.392	0.485
<i>scene</i>	2407	6	2	1.074	0.179
<i>yeast</i>	2417	14	2	4.237	0.303
<i>pentacorrel</i>	600	5	3	5.00	1.000
<i>decasynth</i>	1000	10	4	10.00	1.000

Note: #inst.: number of instances; #lb.: number of labels; #cls.: number of classes on each label; #card. and #dens.: cardinality and density [8].

empirical studies [27, 39] have shown that the accuracies of workers are normally within this range. Under Scenario 1, we simulate that workers are never aware of the label-correlation. That is, their annotation errors are uniformly distributed. When workers provide values to the m th labels of all instances, they randomly select a fixed portion of the instances to provide true values. Furthermore, on a multi-class label, the errors are also uniformly distributed across the incorrect values. For MCMLD, R is set using the method in Section 4.3. Specifically, we use `sklearn.decomposition.pca` in `scikit-learn` to perform PCA. Then, R is set to the value that the first R components of the returned result `explained_variance_ratio_` contribute 70% variance. For all confusion matrices, their diagonal elements are randomly initialized with the values within the range of [0.70, 0.90] and their non-diagonal elements are set to the remainder (i.e., one minus the diagonal value) averaged by $K - 1$. The number of workers increases from 2 to 9. For each setting, we first generate a synthetic dataset and then run the six algorithms. The procedure is repeated 20 times and the means of the accuracies and their standard deviations are reported.

Figure 3 shows the comparison results under Scenario 1. We first focus on the performance of jointly inferring the true labels. Both MCMLD and MCMLI perform joint inference. Clearly, even the simple model MCMLI outperforms DS. The average excess is around 2% (absolute value), and the maximum is over 5% (on *yeast* when #labels on each instance is 4). Especially, on two multi-class multi-label datasets, both MCMLD and MCMLI exhibit obvious advantages. Furthermore, MCMLD outperforms MCMLI in most cases with the average excess around 2%. Although errors are uniformly distributed, the crowdsourced labels are generated according to their true values, so that their values may inevitably inherit a certain degree of correlation from their true values. This results in two consequences: (1) jointly inferring multiple labels can achieve a better local optimum, comparing to running a *single-label* inference multiple times on a dataset, which makes both MCMLD and MCMLI outperform the others. (2) MCMLD exploits the correlation and outperforms MCMLI.

Scenario 2: Annotation with strong label-correlation. Under this scenario, workers are aware of the label-correlation. When labeling, the workers can mutually check their own answers on different labels with the help of the correlation. We strengthen the correlation as follows. For the four real-world datasets and *decasynth*, we run PCA on the true label matrix $Y^{I \times M}$ to obtain

**Figure 3: Comparison results in accuracy under Scenario 1: annotation errors are uniformly distributed.**

a correlation matrix, then we select $t = (30\% \times \#lb.)$ pairs of labels with the top- t large correlation coefficients from the matrix (i.e., for *emotions*, *flags*, and *scene*, $t = 2$; for *yeast*, $t = 4$; and for *decasynth*, $t = 3$). For dataset *pentacorrel*, its strong-correlated labels are the first, third and fifth ones. Suppose (l_a, l_b) is a pair of strong-correlated labels on an instance. If the value of label l_a has not been generated, we generate its value with a random correct probability within the range of [0.75, 0.90]. The strong-correlated labels can help workers verify their answers, which improves their accuracies. Thus, when label l_a is set to its true value, label l_b has 95% probability of being its true value. Although the workers have higher accuracies when labeling the strong-correlated labels, we still let their overall accuracies within the range of [0.60, 0.80] by tuning the accuracies on the non-correlated labels.

Figure 4 shows the comparison results under Scenario 2. When workers are aware of label-correlation, the correlation among crowdsourced labels is strengthened, which directly results in the decline of the calculated values of the parameter R . On datasets *flags*, *yeast* and *decasynth*, the values of R s individually reduce one compared with those values under Scenario 1. When comparing Figures 4.(b), 4.(c), 4.(d), and 4.(f) with Figures 3.(b), 3.(c), 3.(d), and 3.(f), respectively, we find that under Scenario 1, MCMLD does not significantly outperforms MCMLI at some settings of the numbers of labels on each instance (represented as #labels in the following), for examples, on *flags* when #labels = 3 or 8, on *scene* when #labels = 5 or 8, on *yeast* when #labels = 2, 5, or 8, and on *decasynth* when #labels = 7 or 9. Comparatively, under Scenario 2, MCMLD significantly outperforms MCMLI at all points above. The superiority of MCMLD is obviously enlarged on all dataset. Especially, on the strong-correlated dataset *pentacorrel*, the performance of MCMLD exceeds that of MCMLI by more than 4%

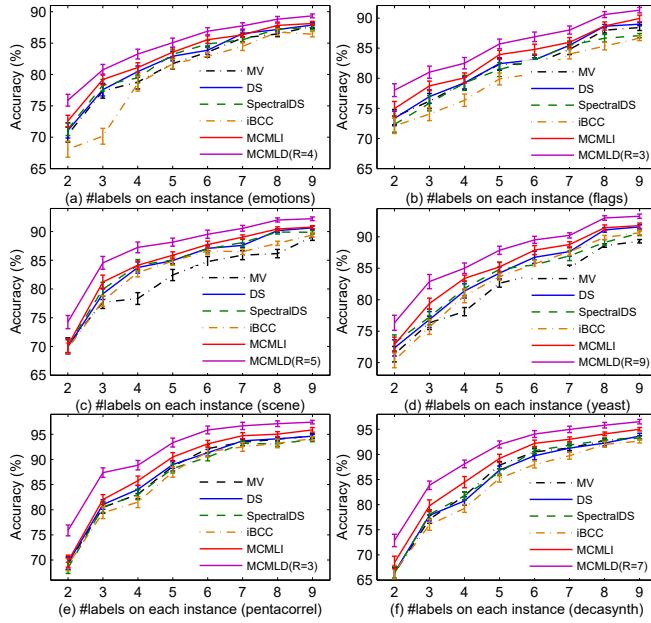


Figure 4: Comparison results in accuracy under Scenario 2: annotation with strong label-correlation.

on average. Thus, MCMLD can benefit a lot from the growth of the crowdsourced label-correlation. In contrast, the advantage of MCMLI is not enlarged under this scenario.

Scenario 3: Biased annotation. The crowdsourced annotation often suffers from the biases of workers. This phenomenon has been extensively studied in recent years [7, 11, 36]. The causes for the biased annotation may be due to workers’ subjective social reality, lack of expertise, common sense decisions, and even the spam submissions. In our study, we simply explain the biased annotation as a phenomenon that the workers tend to provide some classes overwhelming the others. For the binary annotation, we assume that workers tend to provide the negative class, which results in the increase of their accuracies on the true negative labels and the decrease of their accuracies on the true positive ones. For datasets *emotions*, *flags*, *scene*, and *yeast*, whose $K = 2$ (binary annotation), we let the workers’ accuracies on the negative class be within the range of [0.80, 0.90] and those on the positive class be within the range of [0.40, 0.50]. For the multi-class datasets *pentacorrel* (*decaysynth*), we let the workers’ accuracies on the first class (the first two classes for *decaysynth*) be within the range of [0.80, 0.90] and the others be within the range of [0.40, 0.50]. Again, we simulate that there are at least two and at most nine biased workers.

Figure 5 shows the comparison results under Scenario 3. We notice that although the overall accuracies of the workers are within the same range as the those of workers under Scenarios 1 and 2, the accuracies of inference on all datasets consistently decrease. The decrease of the accuracies is more than 15% (absolute value), which suggests that the biased annotation has a negative impact on the performance of inference algorithms. Under this scenario, we find that on the most datasets, except *scene* whose underlying true labels are not correlated, the performance of MCMLI obviously

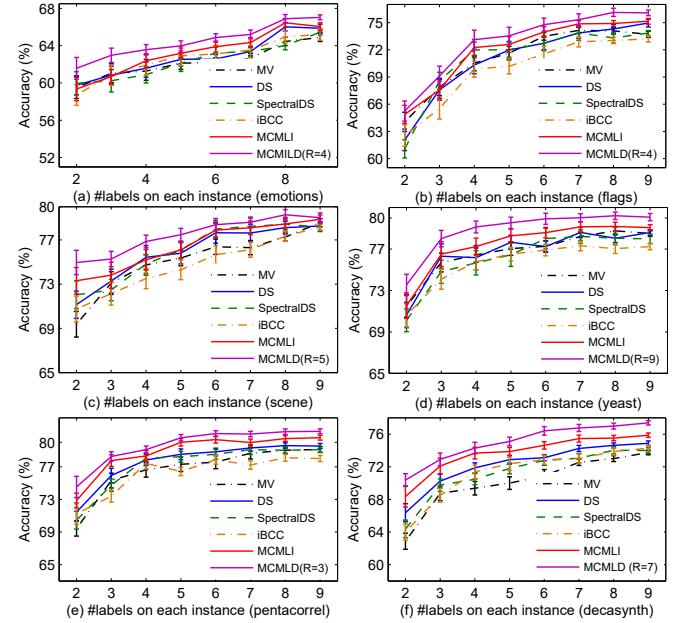


Figure 5: Comparison results in accuracy under Scenario 3: biased annotation.

outperforms that of the compared *single-label* inference algorithms. Furthermore, the performance of MCMLD obviously outperforms that of MCMLI in most cases. It seems that the biased annotation actually increases the degree of the crowdsourced label-correlation, which results in MCMLD’s significantly outperforming MCMLI. This fact is also indicated by the calculation of the parameter R on *yeast* (also on *decaysynth*), whose values are 10 and 9 (8 and 7 for *decaysynth*) under Scenarios 1 and 3, respectively. Both MCMLD and MCMLI are consistently more advantageous in multi-class annotation. In addition, when the number of labels (i.e. M) is greater, e.g., on *yeast*, MCMLD is more advantageous.

Other observations. SpectralDS does not exhibit obvious advantages over DS, indicating that it is not easy to find optimal initial settings for an EM algorithm suitable for different circumstances. iBCC, which uses Gibbs sampling for inference, performs worst. This indicates that as a brutal approximate algorithm, Gibbs sampling usually performs worse than EM in our simulations.

Summary. Under the simulated scenarios, both proposed models MCMLI and MCMLD consistently outperform the existing compared alternatives. Compared with our baseline MCMLI, MCMLD performs better in most cases because it can explore and exploit the label-correlation.

5.3 Evaluation on real-world datasets

To evaluate the performance of the proposed method in a real-world usage, we utilize the crowdsourcing dataset *affectee* [28]. The original *affectee* dataset (SemEval) was constructed as the annotation task [29], asking experts to provide numeric judgments rating the headline of a piece of news for six emotions: *anger*, *disgust*, *fear*, *joy*, *sadness*, and *surprise*. These expert judgments serve as the ground truth. Snow et al. [28] selected a 100-headline sample from

Table 2: Accuracy on the real-world datasets with different label values (in percent)

Dataset	MV	iBCC	DS	SpectralDS	MCMLI	MCMLD($R = 4$)
$D2 (K = 2)$	69.71 ± 0.29	70.87 ± 0.98	72.83 ± 0.11	71.71 ± 1.05	<u>74.60 ± 0.21</u>	76.28 ± 0.26
$D3 (K = 3)$	63.33 ± 0.13	63.83 ± 0.46	64.21 ± 0.12	63.41 ± 0.59	<u>66.14 ± 0.34</u>	67.33 ± 0.31
$D4 (K = 4)$	53.85 ± 0.10	53.50 ± 0.36	52.96 ± 0.24	54.06 ± 0.32	<u>56.32 ± 0.11</u>	58.68 ± 0.21
$D5 (K = 5)$	43.33 ± 0.12	44.16 ± 0.14	44.33 ± 0.09	43.88 ± 0.22	44.49 ± 0.24	<u>45.52 ± 0.09</u>

the original dataset and collected 1000 affect labels from the Amazon Mechanical Turk for each of six label types, where each instance was labeled by ten unique workers from six aspects (emotions). For each emotion, the workers were asked to provide scores within the range of $[0, 100]$.

Conversion to multi-class multi-label datasets. In above *affective* dataset, each instance obtained six labels from a worker, which forms a typical multi-label dataset. The value of each label in this dataset is a real number. This provides an opportunity to change the dataset into several classification datasets by dividing the real label value into intervals. The analysis of the scores in this dataset shows that the scores provided by experts are rather strict. Therefore, for the true values (provided by the experts), we let those equaling to zero be the negative class (no emotion) and others be the positive class when creating a binary dataset (i.e., $K = 2$). From its positive class, we extract out the instances with the original true label values greater than 50 to form a new class "strong" emotion. Then, we have a second dataset, where each label has three choices (negative, moderate, and strong), i.e., $K = 3$. We continue to divide the instances into two and three categories by averaging the interval of the class "moderate" emotion, obtaining two datasets with $K = 4$ and $K = 5$, respectively. For the scores provided by the crowdsourced workers, we simply divide the values that are greater than zero into two, three, and four intervals, forming three, four, and five classes, respectively. Finally, we have four multi-class multi-label datasets whose the numbers of labels are the same (i.e., $M = 6$) but the values of each label are different (i.e., $K = 2, 3, 4$, and 5, respectively).

Results. We repeatedly run the compared algorithms 20 times on these real-world datasets. The parameter R in MCMLD is set to 4, which is calculated by our PCA-method algorithm. The experimental results are listed in Table 2. The best results with statistical significance (i.e., $p - value < 0.05$) are in bold and the runner-ups are underlined. Obviously, on all datasets, MCMLD significantly outperforms the other methods, and MCMLI is the runner-up. Among the existing *single-label* alternatives (i.e., MV, DS, SpectralDS, and iBCC), DS performs best on $D2, D3$, and $D5$, while SpectralDS performs best on $D4$. The differences in the average accuracies between our MCMLD (MCMLI) method and the best one of *single-label* alternatives are 3.45 (1.77), 3.12 (1.93), 4.62 (2.26) and 1.19 (0.16) on the four datasets, respectively. Therefore, the observations on the real-world datasets are consistent with those in our simulation. Furthermore, we find that if we divide the values of labels into more intervals, the quality of the crowdsourced labels will decrease. Because on the one hand there will be fewer instances with the crowdsourced labels whose values lies in the interval of true values, on the other hand, each interval of true values will contain fewer

instances. That is why on $D5$ all methods have low accuracies, and MCMLI cannot significantly outperform the existing alternatives.

6 CONCLUSION AND FUTURE WORK

A real-world crowdsourcing annotation application may involve a natural multi-class multi-label scenario. This study takes the first step toward the multi-label crowdsourcing learning. Specifically, we proposed a novel probabilistic inference method MCMLD based on the mixture of multiple independently multinoulli distributions, which can explore and exploit the label-correlation to jointly infer the true values of multiple labels of each instance in an EM procedure as well as models the reliability of crowdsourced workers. We conducted a set of simulations on six synthetic datasets under three typical crowdsourcing scenarios with different numbers of labels and classes on each label as well as a set of real-world dataset evaluations. Both the simulations and real-world evaluations have shown that our MCMLD method significantly outperforms the existing *single-label* alternatives and our baseline multi-class multi-label method MCMLI.

Several promising topics are worth further investigation. As for MCMLD, its efficiency may deteriorate when the numbers of labels and classes are large. It becomes especially important to reduce its complexity or parallelize it. As a cost-effective tool, MCMLD may help develop novel incentive mechanisms that motivate workers to provide higher-quality answers. Incorporating multi-label model learning and MCMLD is also interesting and full of challenges.

ACKNOWLEDGMENTS

This research has been supported by the National Natural Science Foundation of China under grant 61603186, the Natural Science Foundation of Jiangsu Province, China, under grant BK20160843, the China Postdoctoral Science Foundation under grants 2017T100370 and 2016M590457, the Science Foundation of the Science and Technology Commission of the Central Military Commission (CMC Youth Program), China, and the US National Science Foundation (NSF) under grant IIS-1613950.

REFERENCES

- [1] Wei Bi, Liwei Wang, James T Kwok, and Zhuowen Tu. 2014. Learning to predict from crowdsourced data. In *UAI*. 82–91.
- [2] Christopher M Bishop. 2006. *Pattern Recognition and Machine Learning*. Springer, Chapter 9, 440–441.
- [3] Jonathan Bragg, Mausam, and Daniel S Weld. 2013. Crowdsourcing multi-label classification for taxonomy creation. In *Proceedings of the First AAAI Conference on Human Computation and Crowdsourcing (HCOMP)*. 25–33.
- [4] Alexander P Dawid and Allan M Skene. 1979. Maximum likelihood estimation of observer error-rates using the EM algorithm. *Applied statistics* 28, 1 (1979), 20–28.

- [5] Gianluca Demartini, Djellel Eddine Difallah, and Philippe Cudré-Mauroux. 2012. ZenCrowd: Leveraging probabilistic reasoning and crowdsourcing techniques for large-scale entity linking. In *WWW*. 469–478.
- [6] Jia Deng, Olga Russakovsky, Jonathan Krause, Michael S Bernstein, Alex Berg, and Li Fei-Fei. 2014. Scalable multi-label annotation. In *SIGCHI*. 3099–3102.
- [7] Carsten Eickhoff. 2018. Cognitive biases in crowdsourcing. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining (WSDM)*. 162–170.
- [8] Eva Gibaja and Sebastián Ventura. 2015. A tutorial on multilabel learning. *Comput. Surveys* 47, 3 (2015), 52.
- [9] Chien-Ju Ho, Aleksandrs Slivkins, Siddharth Suri, and Jennifer Wortman Vaughan. 2015. Incentivizing high quality crowdwork. In *WWW*. 419–429.
- [10] Nguyen Quoc Viet Hung, Huynh Hui Viet, Nguyen Thanh Tam, Matthias Weidlich, Hongzhi Yin, and Xiaofang Zhou. 2018. Computing crowd consensus with partial agreement. *IEEE Trans. Knowledge and Data Engineering* 30, 1 (2018), 1–14.
- [11] Ece Kamar, Ashish Kapoor, and Eric Horvitz. 2015. Identifying and accounting for task-dependent bias in crowdsourcing. In *Proceedings of the Third AAAI Conference on Human Computation and Crowdsourcing (HCOMP)*. 92–101.
- [12] David R Karger, Sewoong Oh, and Devavrat Shah. 2011. Iterative learning for reliable crowdsourcing systems. In *NIPS*. 1953–1961.
- [13] Hyun-Chul Kim and Zoubin Ghahramani. 2012. Bayesian classifier combination. In *AISTATS*. 619–627.
- [14] Adriana Kovashka, Olga Russakovsky, Li Fei-Fei, and Kristen Grauman. 2016. Crowdsourceing in computer vision. *Foundations and Trends® in Computer Graphics and Vision* 10, 3 (2016), 177–243.
- [15] Aditya Kurve, David J Miller, and George Kesidis. 2015. Multicategory crowdsourcing accounting for variable task difficulty, worker skill, and worker intention. *IEEE Trans. Knowledge and Data Engineering* 27, 3 (2015), 794–809.
- [16] Matthew Lease and Emine Yilmaz. 2012. Crowdsourcing for information retrieval. In *ACM SIGIR Forum*, Vol. 45. 66–75.
- [17] Qi Li, Yaliang Li, Jing Gao, Lu Su, Bo Zhao, Murat Demirbas, Wei Fan, and Jiawei Han. 2014. A confidence-aware approach for truth discovery on long-tail data. *Proceedings of the VLDB Endowment* 8, 4 (2014), 425–436.
- [18] Qi Li, Yaliang Li, Jing Gao, Bo Zhao, Wei Fan, and Jiawei Han. 2014. Resolving conflicts in heterogeneous data by truth discovery and source reliability estimation. In *SIGMOD*. 1187–1198.
- [19] Shao-Yuan Li, Yuan Jiang, and Zhi-Hua Zhou. 2015. Multi-label active learning from crowds. *arXiv preprint arXiv:1508.00722* (2015).
- [20] Mengchen Liu, Liu Jiang, Junlin Liu, Xiting Wang, Jun Zhu, and Shixia Liu. 2017. Improving learning-from-crowds through expert validation. In *IJCAI*. 2329–2336.
- [21] Pablo G Moreno, Antonio Artés-Rodríguez, Yee Whye Teh, and Fernando Perez-Cruz. 2015. Bayesian nonparametric crowdsourcing. *Journal of Machine Learning Research* 16 (2015), 1607–16271.
- [22] Jafar Muhammadi, Hamid R Rabiee, and Abbas Hosseini. 2015. A unified statistical framework for crowd labeling. *Knowledge and Information Systems* 45, 2 (2015), 271–294.
- [23] Stefanie Nowak and Stefan Rüger. 2010. How reliable are annotations via crowdsourcing: A study about inter-annotator agreement for multi-label image annotation. In *ICMR*. 557–566.
- [24] Vikas C Raykar and Shipeng Yu. 2012. Eliminating spammers and ranking annotators for crowdsourced labeling tasks. *Journal of Machine Learning Research* 13, Feb (2012), 491–518.
- [25] Vikas C Raykar, Shipeng Yu, Linda H Zhao, Gerardo Hermosillo Valadez, Charles Florin, Luca Bogoni, and Linda Moy. 2010. Learning from crowds. *Journal of Machine Learning Research* 11, Apr (2010), 1297–1322.
- [26] Victor S Sheng, Foster Provost, and Panagiotis G Ipeirotis. 2008. Get another label? improving data quality and data mining using multiple, noisy labelers. In *SIGKDD*. 614–622.
- [27] Aashish Sheshadri and Matthew Lease. 2013. Square: A benchmark for research on computing crowd consensus. In *Proceedings of the First AAAI Conference on Human Computation and Crowdsourcing (HCOMP)*. 156–164.
- [28] Rion Snow, Brendan O’Connor, Daniel Jurafsky, and Andrew Y Ng. 2008. Cheap and fast—but is it good?: Evaluating non-expert annotations for natural language tasks. In *EMNLP*. 254–263.
- [29] Carlo Strapparava and Rada Mihalcea. 2007. SemEval-2007 task 14: Affective text. In *Proceedings of the 4th International Workshop on Semantic Evaluations*. Association for Computational Linguistics, 70–74.
- [30] Tian Tian and Jun Zhu. 2015. Max-margin majority voting for learning from crowds. In *NIPS*. 1621–1629.
- [31] Matteo Venanzi, John Guiver, Gabriella Kazai, Pushmeet Kohli, and Milad Shokouhi. 2014. Community-based Bayesian aggregation models for crowdsourcing. In *WWW*. 155–164.
- [32] Aobo Wang, Cong Duy Vu Hoang, and Min-Yen Kan. 2013. Perspectives on crowdsourcing annotations for natural language processing. *Language Resources and Evaluation* 47, 1 (2013), 9–31.
- [33] Peter Welinder, Steve Branson, Pietro Perona, and Serge J Belongie. 2010. The multidimensional wisdom of crowds. In *NIPS*. 2424–2432.
- [34] Jacob Whitehill, Ting-fan Wu, Jacob Bergsma, Javier R Movellan, and Paul L Ruvolo. 2009. Whose vote should count more: Optimal integration of labels from labelers of unknown expertise. In *NIPS*. 2035–2043.
- [35] Guo Xintong, Wang Hongzhi, Yangqiu Song, and Gao Hong. 2014. Brief survey of crowdsourcing for data mining. *Expert Systems with Applications* 41, 17 (2014), 7987–7994.
- [36] Jing Zhang, Victor S Sheng, Jian Wu, and Xindong Wu. 2016. Multi-class ground truth inference in crowdsourcing with clustering. *IEEE Trans. Knowledge and Data Engineering* 28, 4 (2016), 1080–1085.
- [37] Jing Zhang, Xindong Wu, and Victor S Sheng. 2016. Learning from crowdsourced labeled data: A survey. *Artificial Intelligence Review* 46, 4 (2016), 543–576.
- [38] Yuchen Zhang, Xi Chen, Denny Zhou, and Michael I Jordan. 2014. Spectral methods meet EM: A provably optimal algorithm for crowdsourcing. In *NIPS*. 1260–1268.
- [39] Yudian Zheng, Guoliang Li, Yuanbing Li, Caihua Shan, and Reynold Cheng. 2017. Truth inference in crowdsourcing: Is the problem solved? *Proceedings of the VLDB Endowment* 10, 5 (2017), 541–552.
- [40] Denny Zhou, Sumit Basu, Yi Mao, and John C Platt. 2012. Learning from the wisdom of crowds by minimax entropy. In *NIPS*. 2195–2203.