

Multi-Label Active Learning from Crowds

Shao-Yuan Li, Yuan Jiang, Zhi-Hua Zhou*

*National Key Laboratory for Novel Software Technology
Nanjing University, Nanjing 210093, China*

Abstract

Multi-label active learning is a hot topic in reducing the label cost by optimally choosing the most valuable instance to query its label from an oracle. In this paper, we consider the pool-based multi-label active learning under the crowdsourcing setting, where during the active query process, instead of resorting to a high cost oracle for the ground-truth, multiple low cost imperfect annotators with various expertise are available for labeling. To deal with this problem, we propose the MAC (Multi-label Active learning from Crowds) approach which incorporate the local influence of label correlations to build a probabilistic model over the multi-label classifier and annotators. Based on this model, we can estimate the labels for instances as well as the expertise of each annotator. Then we propose the instance selection and annotator selection criteria that consider the uncertainty/diversity of instances and the reliability of annotators, such that the most reliable annotator will be queried for the most valuable instances. Experimental results demonstrate the effectiveness of the proposed approach.

Key words: Multi-Label learning, active learning, weak supervision, crowdsourcing

1. Introduction

Multi-label learning has received significant attention to deal with examples that are associated with multiple classes simultaneously [ZZ14]. It is worth noting that it is usually quite expensive to obtain the labels for training, as every possible label has to be checked whether it is proper for the instance. Thus multi-label active learning, which reduces the labeling cost by actively selecting the most valuable instances to query becomes a hot topic [LG13, HZ13].

*Corresponding author. Email: zhouzh@nju.edu.cn

Most multi-label active learning algorithms rely on the domain expert, or an oracle, to provide the ground-truth label for each query. While the labeling cost of domain experts is high, by employing multiple low cost non-expert labelers, crowdsourcing provides a low cost way to get annotations. A number of studies exploiting the wisdom of crowds have arisen recently [SOJN08, RYZ⁺10, WBBP10, ZLPM14], mainly focusing on single-label tasks where each instance is associated with only one label. Exploiting the wisdom of crowds for multi-label data has barely been touched to the best of our knowledge.

In this paper, we consider the problem of exploiting the wisdom of crowds for multi-label tasks. While multiple annotators may be available, in real-world applications, the labeling cost would still be high for multi-label objects considering that every label needs to be checked on each instance by several annotators, whereas the annotation budget is often limited. Thus we define our problem from the active learning perspective, where instances are actively selected to query the supervised information, but rather than a high cost oracle providing the ground-truth, multiple imperfect annotators are available for labeling. The goal of our work is to learn an effective multi-label classifier with as less annotation cost as possible.

By decomposing the multi-label task into a series of independent binary classification problems, Y. Yan's work [YRFD11] can be applied to each label independently, neglecting the fact that the information of one label may be helpful for learning another related label, especially when there are insufficient training data for every label. To handle the problem of multi-label active learning from crowds, we propose the MAC (Multi-label Active learning from Crowds) approach, which incorporates the local influence of label correlations to build a probabilistic model for multi-label classifier and annotators. Based on this model, we can estimate not only the labels of instances but also the expertise of annotators. Then we propose our instance selection criterion which considers both uncertainty and diversity of instances, and the most reliable annotator on the most valuable instance is queried.

In the following we start with a brief review of some related work. Then, we propose our MAC approach and report the experimental results. Finally, we conclude the paper .

2. Related Work

Multi-label Active Learning Multi-label learning deals with examples that are associated with multiple labels simultaneously; it has received significant attention during the past decades [ZZ14] and achieved successful application in various tasks such as image annotation [BLSB04], gene function classification [EW01] and text categorization [McC99]. To collect labels for multi-label training, each of the multiple labels should be checked whether it is proper for an instance; the labeling cost is quite expensive. Thus multi-label active learning, which reduces the labeling cost by actively selecting the most valuable instance to query from an oracle, has attracted great attention. Existing multi-label active learning research mainly focus on designing the criterion for instance selection, and can be roughly categorized into three categories: 1) uncertainty sampling [SCC08] which selects the instance that the classifier is most uncertain about ; 2) expected loss reduction [HL11] which queries instance that minimizes the expected loss of the classifier; and 3) combining multiple criteria [LG13, HZ13] to select instance. Given the selected instance, most work [HL11, LG13] query all the labels for the instance, which may lead to information redundancy and wasting of oracle’s effort, since labels are often correlated in multi-label learning.

Crowdsourcing With the advent of crowdsourcing platforms such as Amazon Mechanical Turk (AMT), crowdsourcing which makes use of crowdsourced annotations from multiple imperfect labelers is widely used for tasks including sentiment classification [SOJN08], medical diagnosis [RYZ⁺10], image tagging [WBBP10] and webpage categorization [ZLPM14]. As annotators may have different expertise on different tasks, the common wisdom is to distribute tasks to multiple annotators and then estimate the correct labels via some aggregation schemes. Simply taking the majority voting without considering the annotators ability variance may lead to poor impact on subsequent learning. Hence, in order to solve the above problem of annotators, a number of methods assessing abilities of annotators have been proposed. The annotator’s expertise on specific task are usually modeled by either measures with explicit explanation like accuracy [WRW⁺09, KOS11] and confusion matrix [RYZ⁺10, ZLPM14], or complex high multi-dimensional vectors [WBBP10]. A few work trying to get good learning performance with low crowds cost are proposed, for example, [YRFD11, FYT14] consider the cost by adaptively assigning the best annotator to specific task under the active learning framework and [ERH14] propose an approach to approximate the performance of majority voting with less annotations.

While current work on crowdsourcing mainly focus on single-label tasks, to our best knowledge, using crowds in a cost economic way has not been studied in multi-label learning. We are presenting possibly the first approach to exploiting the wisdom of crowds efficiently in multi-label learning.

3. MAC: Multi-Label Active learning from Crowds

We consider the pool-based multi-label active learning from crowds, where a set of labeled annotated data of N_l examples $D_l = \{(\mathbf{x}_1, \mathbf{z}_1, \{\mathbf{y}_{1j}\}), \dots, (\mathbf{x}_{N_l}, \mathbf{z}_{N_l}, \{\mathbf{y}_{N_lj}\})\}$ and an unlabeled set of N_u examples $D_u = \{\mathbf{x}_{N_l+1}, \dots, \mathbf{x}_{N_l+N_u}\}$ are available for training. Each instance \mathbf{x}_i is a d -dimensional feature vector. Given an instance $\mathbf{x}_i \in D_l$, $\mathbf{z}_i \in \{+1, -1\}^{L \times 1}$ is its true label assignment, where L is the number of labels and $\mathbf{z}_i^l = +1(-1)$ indicates that the l -th label is tagged as positive(negative) for \mathbf{x}_i ; and $\mathbf{y}_{ij} \in \{+1, -1, 0\}^{L \times 1}$ is the label assignment given by annotator j , where $\mathbf{y}_{ij}^l = 0$ means that the annotator j gives no annotation for \mathbf{x}_i on the l -th label, and $\mathbf{y}_{ij}^l = +1(-1)$ means annotator j tags the l -th label as positive(negative) for \mathbf{x}_i . We assume a set of M annotators $\{\mathbf{w}_j\}_j$ where $j \in \{1, 2, \dots, M\}$ are available. We do not require each instance to be annotated by all annotators on every label, and the annotator set labeling instance \mathbf{x}_i on label l is denoted as M_i^l , i.e., $M_i^l = \{j | \mathbf{y}_{ij}^l = 1(-1)\}$. During the active learning process, rather than an oracle being available to provide the groundtruth labels \mathbf{z} , multiple imperfect annotators with various expertise are available for annotating. The target of our work is to actively select the most reliable annotator for the most valuable instance such that we can learn an effective multi-label classifier with as less labeling cost as possible.

Compared with traditional multi-label active learning, we can see that the key challenges in the active learning from crowds lie in that: 1) the annotators reliability can be various for specific instance, which requires the careful selection of the annotator with the best expertise for a query; 2) the labels provided by selected annotator for specific queried instance can be noisy, which requires a further aggregation step to get a better estimation of the groundtruth label. To solve the above challenges, we use a probabilistic model which can simultaneously inference the annotators expertise and groundtruth label at the same time. Furthermore, the local influence of neighborhoods' label correlations are incorporated in this probabilistic model to help the multi-label classifier learning. After we obtain the estimation of annotators expertise and instances

labels, we can design some criterion for instance and annotator selection to get the most reliable annotation for the most valuable instance. In the following subsections, we first propose our probabilistic multi-label crowdsourcing model and then the active selection strategy.

3.1. Probabilistic Multi-Label Crowdsourcing Model

We first introduce our probabilistic model on one single label to describe the classifier and annotators, and then encode the local influence of neighborhoods' label correlations in this model to deal with multi-label crowdsourcing.

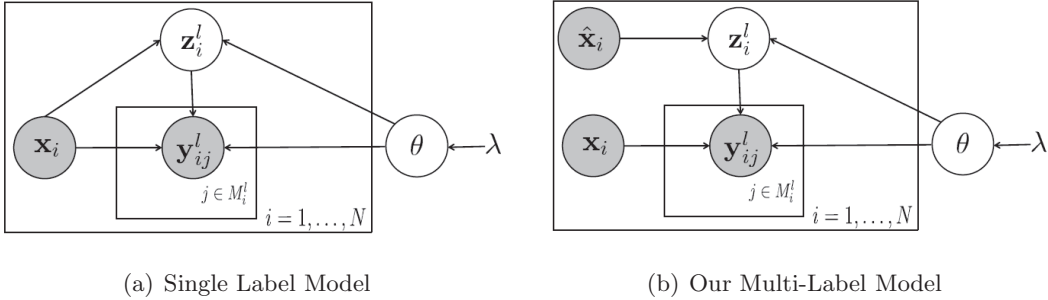


Figure 1: (a) The single label graphical model for instances $\{\mathbf{x}_i\}$, annotations $\{\mathbf{y}_{ij}^l\}$ and unknown groundtruth $\{\mathbf{z}_i^l\}$; (b) Our multi-label graphical model for instances $\{\mathbf{x}_i\}$, $\{\hat{\mathbf{x}}_i\}$, annotations $\{\mathbf{y}_{ij}^l\}$ and unknown groundtruth $\{\mathbf{z}_i^l\}$ on label l .

Figure 1(a) illustrate the probabilistic graphical model on label l over the instances $\{\mathbf{x}_i\}$, the observed annotations given by annotators $\{\mathbf{y}_{ij}^l\}$, and the unobserved groundtruth labels $\{\mathbf{z}_i^l\}$. We assume that the annotations given by the annotators depend both on what input they observe and what task they are expected to finish, i.e., the input instance \mathbf{x} and the target groundtruth label \mathbf{z} . The groundtruth label of the instance is believed to reflect specific property of the instances, for example, $\mathbf{z}_i = 1/ - 1$ may denote the presence/absence of the tree in a picture. Using θ to denote the parameters involved in this model and $p_r(\theta)$ its prior probability distribution, the probabilistic graphical model can be represented using the following joint probability distribution:

$$P(\{\mathbf{y}_{ij}^l\}_{ij}, \{\mathbf{z}_i^l\}_i | \{\mathbf{x}_i\}_i, \theta) = \prod_i p(\mathbf{z}_i^l | \mathbf{x}_i, \theta) \prod_{j \in M_i^l} p(\mathbf{y}_{ij}^l | \mathbf{z}_i^l, \mathbf{x}_i, \theta) p_r(\theta). \quad (1)$$

For each annotator j , we define an expertise level variable \mathbf{e}_{ij}^l to build the relationship between its annotation \mathbf{y}_{ij}^l and the input $(\mathbf{x}_i, \mathbf{z}_i^l)$, which is mathematically define as the probability that

annotator j provides the groundtruth as its annotation for instance i . Here we exploit a *logistic sigmoid* acting on some function $f_{j,\theta}^l$ over the instance \mathbf{x}_i , that is,

$$\mathbf{e}_{ij}^l = p(\mathbf{y}_{ij}^l = \mathbf{z}_i^l | \mathbf{x}_i, \theta) = \sigma(f_{j,\theta}^l(\mathbf{x}_i)), \quad (2)$$

where the logistic sigmoid function is defined as $\sigma(z) = 1/(1 + \exp(-z))$. Using $I(\cdot)$ to denote the indicator function, then the distribution $p(\mathbf{y}_{ij}^l | \mathbf{z}_i^l, \mathbf{x}_i, \theta)$ can be formulated as the following Bernoulli distribution,

$$p(\mathbf{y}_{ij}^l | \mathbf{z}_i^l, \mathbf{x}_i, \theta) = (1 - \mathbf{e}_{ij}^l)^{[1 - I(\mathbf{y}_{ij}^l = \mathbf{z}_i^l)]} (\mathbf{e}_{ij}^l)^{[I(\mathbf{y}_{ij}^l = \mathbf{z}_i^l)]}. \quad (3)$$

Since we are dealing with classification, we use a *logistic sigmoid* acting on some function $f_{0,\theta}^l$ over instance \mathbf{x}_i to model the Bernoulli distribution $p(\mathbf{z}_i^l | \mathbf{x}_i, \theta)$ for \mathbf{z} , that is,

$$p(\mathbf{z}_i^l = 1 | \mathbf{x}_i, \theta) = \sigma(f_{0,\theta}^l(\mathbf{x}_i)) \quad \text{and} \quad p(\mathbf{z}_i^l = -1 | \mathbf{x}_i, \theta) = \sigma(-f_{0,\theta}^l(\mathbf{x}_i)). \quad (4)$$

While any function can be used to implement $f_{0,\theta}^l$ and $f_{j,\theta}^l$, for ease of exposition, we consider the linear discriminating functions, i.e., $f_{0,\theta}^l(\mathbf{x}_i) = (\mathbf{w}_0^l)' \mathbf{x}_i$ and $f_{j,\theta}^l(\mathbf{x}_i) = (\mathbf{w}_j^l)' \mathbf{x}_i$. Given this, the parameters becomes the classifier/annotator parameters $\theta = \{\mathbf{w}_0^l, \{\mathbf{w}_j^l\}\}$. To overcome overfitting, we introduce a zero-mean λ -variance Gaussian prior for $\{\mathbf{w}_0^l\}$ and $\{\mathbf{w}_j^l\}$, i.e., :

$$p_r(\theta) = p(\mathbf{w}_0^l | \lambda) \prod_{j \in M_i^l} p(\mathbf{w}_j^l | \lambda) \quad \text{where} \quad p(\mathbf{w}_0^l | \lambda) = p(\mathbf{w}_j^l | \lambda) = \mathcal{N}(\mathbf{0}, \lambda^{-1} \mathbf{I}) \quad (5)$$

We can estimate the parameters by exploiting the maximum likelihood criterion, to solve which a standard Expectation-Maximization (EM) approach [DLR77] can be used with missing variables $\{\mathbf{z}_i^l\}$ and observed variables $\{\mathbf{y}_{ij}^l\}_{ij}$.

E-step: Given current estimation of the parameters $\theta = \{\mathbf{w}_0^l, \{\mathbf{w}_j^l\}\}$ from last M step, the posterior probability of ground truth label $\{\mathbf{z}_i^l\}$ is computed:

$$\begin{aligned} p(\mathbf{z}_i^l) &= p(\mathbf{z}_i^l | \mathbf{x}_i, \{\mathbf{y}_{ij}^l\}_{ij}, \mathbf{w}_0^l, \{\mathbf{w}_j^l\}) \\ &\propto p(\mathbf{z}_i^l | \mathbf{x}_i, \mathbf{w}_0^l) \prod_{j \in M_i^l} p(\mathbf{y}_{ij}^l | \mathbf{z}_i^l, \mathbf{x}_i, \mathbf{w}_j^l). \end{aligned} \quad (6)$$

Substituting Eq. 3 and Eq. 4 into Eq. 6, E-step reduces to:

$$p(\mathbf{z}_i^l = 1) \propto \sigma(\mathbf{w}_0^l' \mathbf{x}_i) \prod_{\mathbf{y}_{ij}^l = 1} \sigma(\mathbf{w}_j^l' \mathbf{x}_i) \prod_{\mathbf{y}_{ij}^l = -1} \sigma(-\mathbf{w}_j^l' \mathbf{x}_i) \quad (7)$$

$$p(\mathbf{z}_i^l = -1) \propto \sigma(-\mathbf{w}_0^l' \mathbf{x}_i) \prod_{\mathbf{y}_{ij}^l = -1} \sigma(\mathbf{w}_j^l' \mathbf{x}_i) \prod_{\mathbf{y}_{ij}^l = 1} \sigma(-\mathbf{w}_j^l' \mathbf{x}_i) \quad (8)$$

M-step: To estimate the parameters $\theta = \{\mathbf{w}_0^l, \{\mathbf{w}_j^l\}\}$, we maximize the expectation of the joint log-likelihood of $(\{\mathbf{y}_{ij}^l\}_{ij}, \{\mathbf{z}_i^l\}_i)$ over parameters $\theta = \{\mathbf{w}_0^l, \{\mathbf{w}_j^l\}\}$, with respect to the posterior probabilities of $\{\mathbf{z}_i^l\}$ computed by last E step:

$$\theta = \arg \max_{\theta} Q(\theta),$$

where $Q(\theta)$ is

$$\begin{aligned} Q(\theta) &= E_{\mathbf{z}}[\ln P(\{\mathbf{y}_{ij}^l\}_{ij}, \{\mathbf{z}_i^l\}_i | \{\mathbf{x}_i\}_i, \mathbf{w}_0^l, \{\mathbf{w}_j^l\}) p_r(\theta)] \\ &= E_{\mathbf{z}}[\ln \prod_i p(\mathbf{z}_i^l | \mathbf{x}_i, \mathbf{w}_0^l) p(\mathbf{w}_0^l | \lambda) \prod_j p(\mathbf{y}_{ij}^l | \mathbf{z}_i^l, \mathbf{x}_i, \{\mathbf{w}_j^l\}) p(\mathbf{w}_j^l | \lambda)] \\ &= \sum_i E_{\mathbf{z}}[\ln p(\mathbf{z}_i^l | \mathbf{x}_i, \mathbf{w}_0^l)] + \sum_{ij} E_{\mathbf{z}}[\ln p(\mathbf{y}_{ij}^l | \mathbf{z}_i^l, \mathbf{x}_i, \mathbf{w}_j^l)] \\ &\quad + \frac{\lambda}{2} \|\mathbf{w}_0^l\|^2 + \frac{\lambda}{2} \sum_j \|\mathbf{w}_j^l\|^2. \end{aligned} \quad (9)$$

Substituting Eq. 3 and Eq. 4 into Eq. 9, $Q(\theta)$ reduces to:

$$Q(\mathbf{w}_0^l) = \sum_i (\mathbf{w}_0^{l'} \mathbf{x}_i) p(\mathbf{z}_i^l = 1) - \sum_i [\ln(1 + \exp(\mathbf{w}_0^{l'} \mathbf{x}_i))]. \quad (10)$$

$$Q(\mathbf{w}_j^l) = \sum_i (\mathbf{w}_j^{l'} \mathbf{x}_i) p(\mathbf{z}_i^l = 1) \cdot \mathbf{y}_{ij}^l - \sum_i \ln[1 + \exp(\mathbf{w}_j^{l'} \mathbf{x}_i \cdot \mathbf{y}_{ij}^l)]. \quad (11)$$

The gradient of Q with respect to parameters $\mathbf{w}_0^l, \mathbf{w}_j^l$ is computed as:

$$\frac{\partial Q}{\partial \mathbf{w}_0^l} = \sum_i \mathbf{x}_i [p(\mathbf{z}_i^l = 1) - \sigma(\mathbf{w}_0^{l'} \mathbf{x}_i)]. \quad (12)$$

$$\frac{\partial Q}{\partial \mathbf{w}_j^l} = \sum_i \mathbf{x}_i [p(\mathbf{z}_i^l = 1) - \sigma(\mathbf{w}_j^{l'} \mathbf{x}_i \cdot \mathbf{y}_{ij}^l)] \cdot \mathbf{y}_{ij}^l. \quad (13)$$

To deal with multi-label data, one straightforward approach is to treat each label independently and applying the above model on each label separately. But by this, the label correlations are ignored, which however is widely verified to be helpful for multi-label learning. In this paper, we derive a local code to enhance the feature representation of instances, which reflects the local influence of its neighborhoods' label correlations. We use the simple idea that instances similar in the feature space should also be similar in the label space, and utilize the information from the label space of labeled instances to construct extra features. In detail, for each instance \mathbf{x} , its code vector \mathbf{c} is constructed as the label average mean of its k nearest neighbors in the initial labeled training set, i.e.,

$$\mathbf{c} = \frac{1}{k} \sum_{\mathbf{x}_i \in N(\mathbf{x})} \mathbf{z}_i. \quad (14)$$

After the code vector is computed, the *enhanced representation* of instance \mathbf{x} becomes $\hat{\mathbf{x}} = [\mathbf{x}, \mathbf{c}]$. Then we learn the above probabilistic model on each label separately, but using different instance representations for classifier and annotators. For label l , the multi-label classifier parameter \mathbf{w}_0 is learned on the *enhanced representation*, but for the annotator model parameters $\{\mathbf{w}_j^l\}$, the *original representations* \mathbf{x} are used, i.e., \mathbf{w}_0 is of dimension $(d + L)$ and $\{\mathbf{w}_j^l\}$ are of dimension d . Utilizing the local influence of label correlations, our multi-label probabilistic crowdsourcing model can be described by Figure 1(b). The reason we use different instance representations for the classifier and annotators is that, the *enhanced representation* of instances are deduced from the initial labeled data, which are generated by the true multi-label data distribution, so it is expected that only for cases where the training data are consistent with the multi-label data distribution, the *enhanced representation* are to be helpful for learning. While the groundtruth estimation aggregated from the multiple annotators should be trustable and consistent with the multi-label data distribution, the annotations provided by each annotator can be of low quality and far away from the multi-label data distribution, thus the *enhanced representation* generated from the labeled training data would be inconsistent with its labeling behavior. Though here the way we compute the *enhanced representation* is simple, we believe the idea of constructing new representations for the classifier using information from the label space would be a good direction for multi-label crowdsourcing.

3.2. Active Selection

Based on the above model, we can learn an estimation of the multi-label classifier, the annotators expertise and the labels for instances from the crowds annotations. In traditional multi-label active learning, most work select the instance and query its supervised information on all labels, which may lead to information redundancy and wasting of annotation effort since labels are often correlated in multi-label learning. To avoid information redundancy, we exploit the idea of first choosing the most valuable instance according to its label estimation, and then selecting its most uncertain label to query the supervised information. Considering that the supervised information in our problem are provided by crowds whose reliability on different instances may be different, therefore, the selection of the most reliable annotator for the specific instance-label pair should also be carefully considered.

Instance-Label Pair Selection We design our instance selection criterion by combining the uncertainty criterion with the diversity of queried annotations. LCI (instance Label Cardinality Inconsistency) is a commonly used instance uncertain measure in multi-label active learning and its combination with other measures to conduct instance selection has shown promising results in several works [LG13, HZ13]. It is defined as the inconsistency between the number of predicted positive labels of instances and the average label cardinality on the labeled data,

$$LCI(\mathbf{x}_i) = \left(\sum_{l=1}^L I(\hat{\mathbf{z}}_i^l = 1) - \frac{1}{N_l} \sum_{j=1}^{N_l} \sum_{l=1}^L I(\mathbf{z}_j^l = 1) \right)^2,$$

where $\hat{\mathbf{z}}_i^l$ is the label estimation of instance \mathbf{x}_i on l . We extend LCI to incorporate the query diversity during the active process and define a new criterion

$$i^* = \arg \max_i CI(\mathbf{x}_i) = \frac{\left| \sum_{l=1}^L I(\hat{\mathbf{z}}_i^l = 1) - \frac{1}{N_l} \sum_{j=1}^{N_l} \sum_{l=1}^L I(\mathbf{z}_j^l = 1) \right|}{\max\{\xi, \text{anno}(\mathbf{x}_i)\}}, \quad (15)$$

where we set $\hat{\mathbf{z}}_i^l = 1(-1)$ if its probability predicted by Eq. 6 is larger (no more) than 0.5, $\text{anno}(\mathbf{x}_i)$ is the number of queried annotations of instance \mathbf{x}_i and $\xi \in (0, 1)$ is a small constant to avoid zero divisor. We use the number of queried annotation to tradeoff LCI to avoid querying the same instance, and giving less queried instance more chance to be queried which may contain more unknown information. $\xi \in (0, 1)$ is set as 0.5 in this paper.

Given the selected instance \mathbf{x}_* by Eq. 15, the label whose estimation is most uncertain is selected:

$$l^* = \arg \min_l CL(\mathbf{x}_{i^*}, l) = |p(\mathbf{z}_{i^*}^l = 1 | \mathbf{x}_{i^*}) - 0.5|. \quad (16)$$

Annotator Selection Given the selected instance and label (\mathbf{x}_*, l^*) by Eq. 15 and Eq. 16, we need to select the most reliable annotator on it to collect high quality annotations. Eq. 2 provides the information about how reliable each annotator would be on each instance and label, which is inferred from observed annotations by our crowdsourcing model. So we compute the expertise level of all annotators on the selected instance-label pair and select the most reliable one:

$$j^* = \arg \max_j CA(\mathbf{x}_{i^*}, l^*, \mathbf{w}_j^{l^*}) = \mathbf{e}_{i^*j}^{l^*} = \sigma((\mathbf{w}_j^{l^*})' \mathbf{x}_{i^*}). \quad (17)$$

After the instance-label pair and the best annotator on it are selected, the annotation is queried and added to the training data to update the multi-label crowdsourcing model. The MAC (Multi-label Active learning from Crowds) algorithm is summarized in Algorithm 1.

Algorithm 1 The MAC Algorithm

- 1: **Input:**
 - 2: training set $D = \{D_l, D_u\}$, parameters λ, k
 - 3: **Train:**
 - 4: **initialization**
 - 5: get the enhanced representation $\hat{\mathbf{x}} = [\mathbf{x}, \mathbf{c}]$ for instance $\mathbf{x} \in D_l$ by Eq. 14
 - 6: get the initial estimation of $\{\mathbf{w}_0^l, \mathbf{w}_j^l\}$ using LIBLINEAR on D_l
 - 7: **repeat:**
 - 8: get the label estimations for instances \mathbf{x} in D_u by Eq. 6 with $\{\mathbf{w}_0^l\}$ and $\{\mathbf{w}_j^l\}$
 - 9: select instance \mathbf{x}_{i^*} by Eq. 15
 - 10: select label l^* for instance \mathbf{x}_{i^*} by Eq. 16
 - 11: select annotator j^* for instance \mathbf{x}_{i^*} on label l^* by Eq. 17
 - 12: get the annotation $\mathbf{y}_{i^*j^*}^{l^*}$ for \mathbf{x}_{i^*} on label l^* from annotator j^*
 - 13: remove annotator j^* for instance \mathbf{x}_{i^*} on label l^*
 - 14: remove label l^* from instance \mathbf{x}_{i^*} if all annotators on l^* are removed
 - 15: remove instance \mathbf{x}_{i^*} from D_u if all labels of \mathbf{x}_{i^*} are removed
 - 16: add \mathbf{x}_{i^*} to D_l if \mathbf{x}_{i^*} is not in D_l , add $\mathbf{y}_{i^*j^*}^{l^*}$ to D_l for \mathbf{x}_{i^*} on label l^*
 - 17: update $\{\mathbf{w}_0^l, \mathbf{w}_j^l\}$ on D_l
 - 18: **until** the maximum number of queries is reached or D_u is empty.
 - 19: **Test:**
 - 20: for instance \mathbf{x}_t , get its enhanced representation $\hat{\mathbf{x}}_t = [\mathbf{x}_t, \mathbf{c}_t]$ by Eq. 14
 - 21: get its prediction on label l \mathbf{z}_t^l by Eq. 4
-

3.3. Computational Complexity

The computational complexity of MAC composes of three parts: 1) the initialization step, we use LIBLINEAR [FCH⁺08] to get the initial estimation of classifier and annotator parameters which is fast and accurate; 2) the active selection part by Eq. 15-17, which is linear in the number of instance, the number of labels and the number of annotators; 3) model updating, after the queried annotation is added to the training data, the crowdsourcing model is updated using EM. The parameters are initialized by results from last step which makes the EM converges fast, usually in less than 100 iterations. The computational complexity of E-step in Eq. 6 is linear in the number

of involved instances and the number of involved annotations; the computation complexity of M-step in Eq. 9 depends on the employed optimization approach. We exploit gradient ascent to solve the M-step, which needs to iteratively compute the gradient of parameters until convergence or maximum iteration number is reached. At each iteration, the computation complexity is linear in the number of involved annotators, the number of involved instance and the number of involved annotation.

4. Experiments

In this section, we compare our MAC approach with a number baseline methods on two natural scene classification datasets in which each picture can be categorized into one or more classes. The Image dataset [BLSB04]¹ contains 2000 natural scene images and 5 possible labels $\{desert, mountains, sea, sunset, trees\}$, each image having on average 1.24 ± 0.44 labels. The Scene dataset [ZZ07]² contains 2407 images and 6 possible labels $\{beach, field, foliage, mountain, sunset, urban\}$, each image having on average 1.07 ± 0.26 labels. The images belonging to more than one class in these two datasets comprise about 22% of the whole data.

The parameter λ is set as $1e^{-3}$ and k is set as 10 in the experiments. As there exists no other method for the problem of multi-label active learning from crowds, we compare our MAC approach with the following 2 groups of 7 baselines.

Group1 exploits the label correlations of multi-label data, i.e., the multi-label classifier and annotators are learned respectively on enhanced and original instance representations: 1) MCR+RD: learn the classifier and annotators using our multi-label crowdsourcing model in Figure 1(b); randomly select (instance, label, annotator) for labeling; 2) MV+ACT: instead of learning by our multi-label crowdsourcing model, use majority voting on the collected annotations and train a logistic regression classifier; select the (instance, label) by our active selection component in section 3.2, randomly select one annotator for labeling; 3) MV+RD: instead of learning by our multi-label crowdsourcing model, use majority voting on the available annotations and train a logistic regression classifier; randomly select (instance, label, annotator) for labeling;

¹ <http://mulan.sourceforge.net/datasets.html>

² <http://cse.seu.edu.cn/people/zhangml/Resources.htm>

Group2 does not exploit the label correlations of multi-label data, i.e., both the multi-label classifier and annotators are learned on original instance representations: 4) SCR+ACT: learn the classifier and annotators using the probabilistic model in Figure 1(a) on each label independently; actively select (instance, label, annotator) by our active selection component in section 3.2; 5) SCR+RD: learn the classifier and annotators using the probabilistic model in Figure 1(a) on each label independently; randomly select (instance, label, annotator) for labeling; 6) SMV+ACT: use majority voting on the available annotations and train a logistic regression classifier; select the (instance, label) by our active selection component in section 3.2, randomly select one annotator for labeling; 7) SMV+RD: use majority voting on the available annotations and train a logistic regression classifier; randomly select (instance, label, annotator) for labeling.

3 labelers with different expertise are simulated to generate the crowdsourcing annotations. For each dataset, on one specific label l , we proceed as follows: we first train a logistic regression model for this label on the whole dataset by LIBLINEAR [FCH⁺08], then according to the probability output of the logistic regression model on the dataset, we cluster the data into three subsets by k-means. After that, each of the 3 simulated annotators i , $i = 1, 2, 3$ performs as an expert on the i -th cluster and gives the groundtruth label for annotation; for the remaining data belonging to the other clusters, its probability of correctly annotating the groundtruth label is 75% (the labels for these data are randomly switched with probability 25%). By such an annotation generation process, the annotations provided by crowds for specific label is dependent on the semantic of this label. Furthermore, on each label, the expertise of annotators are different on different instances.

For each dataset, we randomly partition the data into three parts which compose 2.5%, 47.5% and 50% of the whole dataset to respectively construct the initial annotated labeled training data, the unlabeled training data and the test data. The random partition is repeated for five times on each dataset and the average performance is reported. At each active query, the (instance, label, annotator) triple is selected to get the annotation and then added into the annotated data. After every 200 annotation queries the performance of the multi-label classifier on the test data is reported. The query process terminates after the number of queried annotations reaches 8,000. The commonly used multi-label performance measure micro-F1 [ZZ14] is used in our experiments.

In Figure 2(a)/3(a) and Figure 2(b)/3(b), we test our crowdsourcing learning component and the active selection component between methods within Goup1 and Group2. In Figure 2(c)/3(c)

we test the label relationship role by comparing our MAC with the best baseline in Group1 and Group2.

From Figure 2(a) and Figure 2(b), we can see that: 1) compared to random sampling, our active selection component clearly helps improve learning for both our crowdsourcing learning model and majority voting (MAC vs MCR+RD; MV+ACT vs MV+RD; SCR+ACT vs SCR+RD; SMV+ACT vs SMV+RD); 2) compared to majority voting, our crowdsourcing model aggregates crowds better, especially when the annotations quality can not be guaranteed (MCR+RD vs MV+RD; SCR+RD vs SMV+RD). From Figure 2(c) in which MAC is compared with the best method in Group1 (MV+ACT) and Group2 (SCR+ACT, SMV+ACT), it can be seen that utilizing the label relationship for classifier learning greatly boost learning (MAC vs SCR+ACT; MV+ACT vs SMV+ACT), especially at early stage when the number of training data is small. Similar results for our MAC and baselines on Scene dataset are also obtained in Figure 3. It's not surprising that by utilizing the label relationship and active selection, our MAC approach achieves the best learning performance.

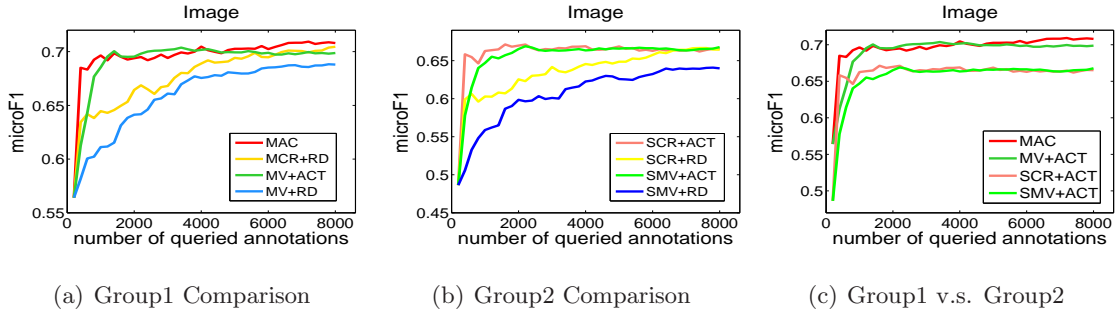


Figure 2: The average results over 5 runs in terms of micro-F1(the larger the better) on Image.

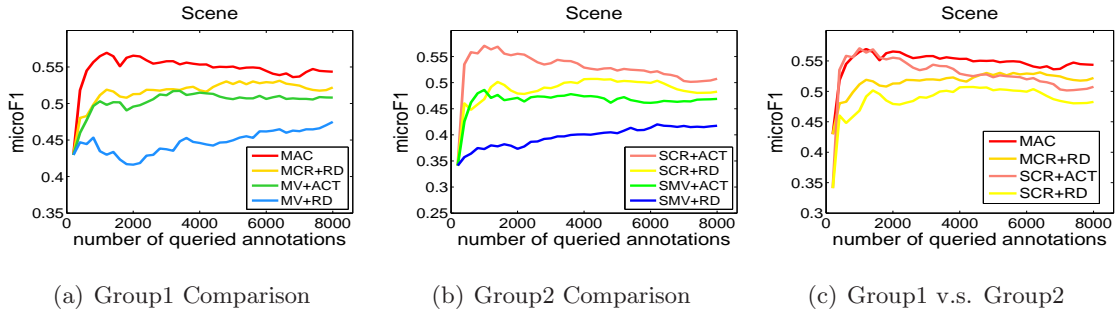


Figure 3: The average results over 5 runs in terms of micro-F1(the larger the better) on Scene.

5. Conclusion

In this paper, we consider the problem of reducing the labeling cost of multi-label learning by actively learning from crowdsourcing annotations, where during the multi-label active query process, rather than resorting to a high cost oracle for the ground-truth, the labels are provided by multiple low-cost imperfect annotators. To deal with this problem, we propose the MAC (Multi-label Active learning from Crowds) approach which encodes the local influence of label correlations to build a probabilistic multi-label crowdsourcing model and then propose an active selection criterion to query the best annotator for the most valuable instance. Experimental results show that our approach is more effective than baselines. Currently, the computational complexity of MAC is linearly dependent on the number of labels; in the future, we plan to develop more efficient methods for large number of labels.

References

- [BLSB04] M. R. Boutell, J. Luo, X. Shen, and C. M. Brown. Learning multi-label scene classification. *Pattern Recognition*, 37(9):1757–1771, 2004.
- [DLR77] A. Dempster, N. Laird, and D. Rubin. Maximum likelihood estimation from incomplete data. *Journal of the Royal Statistical Society - B*, 39(1):1–38, 1977.
- [ERH14] Seyda Ertekin, Cynthia Rudin, and Haym Hirsh. Approximating the crowd. *Data Mining and Knowledge Discovery*, 28(5-6):1189–1221, 2014.
- [EW01] A. Elisseff and J. Weston. A kernel method for multi-labelled classification. In *Advances in Neural Information Processing Systems 14*, pages 681–687, 2001.
- [FCH⁺08] Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, and Chih-Jen Lin. Liblinear: A library for large linear classification. *Journal of Machine Learning Research*, 9:1871–1874, 2008.
- [FYT14] Meng Fang, Jie Yin, and Dacheng Tao. Active learning for crowdsourcing using knowledge transfer. In *Proceedings of the 28th AAAI Conference on Artificial Intelligence*, pages 1809–1815, 2014.

- [HL11] C.-W. Hung and H.-T. Lin. Multi-label active learning with auxiliary learner. In *Proceedings of the 3th Asian Conference on Machine Learning*, pages 315–330, 2011.
- [HZ13] S.-J. Huang and Z.-H. Zhou. Active query driven by uncertainty and diversity for incremental multi-label learning. In *Proceedings of the 13th IEEE International Conference on Data Mining*, pages 1079–1084, 2013.
- [KOS11] D. R. Karger, S. Oh, , and D. Shah. Iterative learning for reliable crowdsourcing systems. In *Advances in Neural Information Processing Systems 24*, pages 1953–1961, 2011.
- [LG13] X. Li and Y. Guo. Active learning with multi-label svm classification. In *Proceedings of the 23rd International Joint Conference on Artificial Intelligence*, pages 1479–1485, 2013.
- [McC99] A. K. McCallum. Multi-label text classification with a mixture model trained by EM. In *Working Notes of the AAAI’99 Workshop on Text Learning*, pages 17–26, 1999.
- [RYZ⁺10] V.C. Raykar, S. Yu, L.H. Zhao, G.H. Valadez, C. Florin, L. Bogoni, and L. Moy. Learning from crowds. *Journal of Machine Learning Research*, 11:1297–1322, 2010.
- [SCC08] M. Singh, E. Curran, and P. Cunningham. Active learning for multi-label image annotation. In *Proceedings of the 19th Irish Conference on Artificial Intelligence and Cognitive Science*, 2008.
- [SOJN08] R. Snow, B. OConnor, D. Jurafsky, and A. Y. Ng. Cheap and fastbut is it good? evaluating non-expert annotations for natural language tasks. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 254–263, 2008.
- [WBBP10] P. Welinder, S. Branson, S. Belongie, and P. Perona. The multidimensional wisdom of crowds. In *Advances in Neural Information Processing Systems 23*, pages 2024–2432, 2010.
- [WRW⁺09] J. Whitehill, P. Ruvolo, T. Wu, J. Bergsma, and J. Movellan. Whose vote should

- count more: Optimal integration of labels from labelers of unknown expertise. In *Advances in Neural Information Processing Systems 22*, pages 2035–2043, 2009.
- [YRFD11] Y. Yan, R. Rosales, G. Fung, and J. Dy. Active learning from crowds. In *Proceedings of the 28th International Conference on Machine Learning*, pages 1161–1168, 2011.
- [ZLPM14] D. Zhou, Q. Liu, J. C. Platt, and C. Meek. Aggregating ordinal labels from crowds by minimax conditional entropy. In *Proceedings of the 31st International Conference on Machine Learning*, pages 262–270, 2014.
- [ZZ07] M.-L. Zhang and Z.-H. Zhou. Ml-knn: A lazy learning approach to multi-label learning. *Pattern Recognition*, 40(7):2038–2048, 2007.
- [ZZ14] M.-L. Zhang and Z.-H. Zhou. A review on multi-label learning algorithms. *IEEE Transactions on Knowledge and Data Engineering*, 26(8):1819–1837, 2014.