

Effective Multi-Label Active Learning for Text Classification

Bishan Yang^{†,*}, Jian-Tao Sun[‡], Tengjiao Wang[‡], Zheng Chen[‡]

[†]Key Laboratory of High Confidence Software Technologies (Peking Univ.), Ministry of Education, China
School of Electronics Engineering and Computer Science, Peking Univ., Beijing, China

[‡]Microsoft Research Asia, No. 49, Zhichun Road, Beijing, China

{bishan_yang, tjwang}@pku.edu.cn, {jtsun, zhengc@microsoft.com}

ABSTRACT

Labeling text data is quite time-consuming but essential for automatic text classification. Especially, manually creating multiple labels for each document may become impractical when a very large amount of data is needed for training multi-label text classifiers. To minimize the human-labeling efforts, we propose a novel multi-label active learning approach which can reduce the required labeled data without sacrificing the classification accuracy. Traditional active learning algorithms can only handle single-label problems, that is, each data is restricted to have one label. Our approach takes into account the multi-label information, and select the unlabeled data which can lead to the largest reduction of the expected model loss. Specifically, the model loss is approximated by the size of version space, and the reduction rate of the size of version space is optimized with Support Vector Machines (SVM). An effective label prediction method is designed to predict possible labels for each unlabeled data point, and the expected loss for multi-label data is approximated by summing up losses on all labels according to the most confident result of label prediction. Experiments on several real-world data sets (all are publicly available) demonstrate that our approach can obtain promising classification result with much fewer labeled data than state-of-the-art methods.

Categories and Subject Descriptors

H.3.3 [Information Systems]: Information Search and Retrieval; I.5.2 [Design Methodology]: Classifier Design and Evaluation

General Terms

Algorithms, Performance, Experimentation

*This work was finished when the first author conducted her internship at Microsoft Research Asia.

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. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

KDD'09, June 28–July 1, 2009, Paris, France.

Copyright 2009 ACM 978-1-60558-495-9/09/06 ...\$5.00.

Keywords

Active Learning, Text Classification, Multi-label Classification, Support Vector Machines

1. INTRODUCTION

As text data becomes a major information source in our daily life, many research efforts have been conducted in text classification to better organize text data, in applications like document filtering, email classification, Web search, etc. In particular, multi-label text classification problems have received considerable attention, since many text classification tasks are multi-labeled, i.e., each document can belong to more than one category. Take news classification as an example, one news article talking about the effect of Olympic games on tourism industry might belong to the following topic categories: *sports*, *economy* and *travel*.

In the literature, supervised learning algorithms are widely used in text classification. It requires a sufficient amount of labeled data for training a high quality model. However, labeling is usually a time-consuming and expensive process done by domain experts. Active learning is an approach to reduce the labeling cost. The active learner iteratively selects a sample of data to label based on some selection strategies suggesting that the data most deserves to be labeled. Thus it can achieve comparable performance with supervised learners while using much less labeled data. Active learning is particularly important for the multi-label text classification task. The reason is that, in the single-label case, a human judge can stop labeling an instance once its category is identified. But in the multi-label case, human judges need to decide all possible categories for each instance. Thus the effort of assigning labels for multi-label data is much larger than for the single-label data.

Despite the value and significance of this problem, there is very limited research on multi-label active learning. Most of the active learning research focuses on the single-label classification problem [10, 21, 14, 22]. The sample selection strategy strictly follows the assumption that each instance has only one label, and thus cannot directly applied in multi-label active learning. The reason can be explained by the following example. Suppose there are three categories c_1 , c_2 , c_3 in the multi-label classification task. The popular one-versus-all technique [3] is used and the classification probabilities on all possible classes are given. Assume the probabilities on instance x_1 are $[c_1:0.8, c_2:0.5, c_3:0.1]$ and on x_2 is $[c_1:0.7, c_2:0.1, c_3:0.1]$. x_1 actually has two labels c_1 and c_2 , and x_2 has one label c_1 . It can be found that correctly predicting labels for x_1 is harder than x_2 . However, if we

assume each instance only has one label and take the most uncertainty strategy, x_2 would be considered to be harder to classify, since the probability score on the predicted label of x_2 is 0.7, which is lower than that of x_1 0.8. Thus considering multi-label information in the sample selection strategy is very important.

In this paper, we propose a novel multi-label active **learning** approach for text classification. The sample selection strategy aims to label data which can help maximize the reduction rate of the expected model loss. To measure the loss reduction, we use Support Vector Machines (SVM) in terms of version space [21] due to the effectiveness of SVM active **learning** on text classification. In the original work, the loss is modeled for single-label case, and here we extend it to multi-label case. We also propose an effective method to predict labels for multi-label data. The expected loss is approximated with the loss associated with the most confident result of label prediction. We will show that a proper label prediction method is critical in measuring loss for multi-label data.

We empirically evaluate the effectiveness of the proposed approach using several real-world data sets that are publicly available. The results demonstrate that our method is superior to the state-of-the-art active **learning** algorithms for multi-label text classification, and can significantly reduce the demand of labeled data while maintaining promising classification results.

The remainder of this paper is organized as follows. Section 2 discusses related work. Section 3 presents the definition of multi-label active **learning** problem. Section 4 introduces our SVM-based active learner, including the loss optimization framework and the sample selection strategy. Section 5 shows experimental results of our algorithm on several real-world data sets compared with other baseline methods. Section 6 presents conclusions and future work.

2. RELATED WORK

Active **learning** on text classification has been well researched. Based on the adopted sample selection strategy, they can be grouped into three types: 1) Uncertainty sampling [10, 14]. The active learner iteratively labels the unlabeled data on which the current hypothesis is most uncertain. 2) Expected-error reduction [2, 18, 22]. The strategy aims to label data to minimize the expected error on the unlabeled data. Usually it requires expensive computational effort on estimating the expected error, since each of the unlabeled data associated with each possible labeling needs to be evaluated. 3) Committee-based active learner. It has the similar idea with uncertainty sampling strategy. The active learner selects data to label that have the largest disagreement among several committee members (classifiers) from the version space. The work of query by committee [19] is the first algorithm of this kind. In [21], the idea is extended to Support Vector Machine active **learning**, and it models the reduction of version space size with SVM.

However, most of the previous research targets single-label classification problems. The sample selection strategy evaluates each unlabeled data by assuming it has only one label. For instance, the uncertainty sampling strategy will focus on measuring the confidence of the most likely class, and the error reduction strategy will estimate the expected error by just considering single-label cases. Thus these strategies can not be directly applied in multi-label text classification.

There is very limited research on multi-label active **learning**. The research work of [9] is the one most related to our paper with respect to the studied problem. It decomposes the multi-label classification problem to several binary ones using one-versus-all approach. The selection strategy minimizes the smallest SVM margin among all binary classification problems. The approach does not consider the multi-label information, and treats all classes equally. In [12], an SVM active **learning** method was proposed for multi-label image classification. It selects unlabeled data which has the maximum mean loss value over the predicted classes. The multi-label classification problem is also viewed as several binary classification tasks. A threshold of loss value is estimated for each binary classifier, and then used to decide the predicted classes for unlabeled data. According to our experiments, this threshold cutting method is not effective on the text classification data sets we used. Also for image classification, [16, 17] developed a two-dimensional active **learning** algorithm, which selects sample-label pairs to minimize the Bayesian classification error bound. It is reasonable to label picture-category pairs since judging a picture's label is very efficient. However, this method is not suitable for text classification task. Because it will introduce much additional cost if a document is read several times. Obviously, the cost of reading a document and judging its label is much bigger than that of a picture. Recently, [4] proposed several active **learning** strategies for multi-label text classification. Each selection strategy consists of one rule to combine the output of individual binary classifiers, including three orthogonal dimensions: "evidence", "class" and "weight". Also, they do not take account of the label prediction result for each instance in the selection strategy.

3. PROBLEM DEFINITION

Multi-label text classification is the task of automatically classifying text documents into a subset of predefined classes. Denote training examples as $\mathbf{x}_1, \dots, \mathbf{x}_n$ and the k classes as $1, \dots, k$. We represent the label set of \mathbf{x}_i by a binary vector $\mathbf{y}_i = [y_i^1, \dots, y_i^k]$, $y_i^j \in \{-1, +1\}$, where $y_i^j = 1$ if \mathbf{x}_i belongs to class j , otherwise $y_i^j = -1$. Denote the set of all possible class combinations as \mathcal{Y} ($|\mathcal{Y}| = 2^k$). The multi-label classifier can be expressed as a decision function $f: \mathbf{X} \rightarrow \mathcal{Y}$.

In our active **learning** study, we consider SVM as the basic multi-label classifier, since SVM has demonstrated significant success on text classification tasks [7, 23]. Usually, multi-label SVM adopts the one-versus-all approach, which trains a separate binary classifier for each possible class against the rest of classes, and combines the output of all the binary classifiers to determine the final labels of the given data. In binary classification, SVM tries to find the hyperplane that can separate the training data by a maximal margin. Denote f^i as the binary classifier associated with target class i . Given a test instance \mathbf{x}' , if $f^i(\mathbf{x}') > 0$, then \mathbf{x}' belongs to class i , otherwise, the labels of \mathbf{x}' will not include class i .

In this paper, we adopt the pool-based active **learning** approach which is usually used in the literature. Assume we are given a pool of partially labeled data. Denote the data with labels by D_l , which is typically small in size, and the remaining data without labels by D_u . At the beginning, a classifier is trained using the initial labeled set D_l . Based on this classifier, the learner selects a sample from D_u and queries for its true labels according to some criterion. Then

the newly labeled data is incorporated into D_l . The training and labeling process runs iteratively after a certain number of iterations or when the classifier reaches a sufficient accuracy.

The key issue of active learning is how to select the most informative data examples to label, which is also called sample selection strategy. So, the research problem studied in this work can be described as follows: in order to train an effective multi-label active learner, how to design the sample selection strategy to reduce the human labeling cost as much as possible?

4. SVM-BASED ACTIVE LEARNING FOR MULTI-LABEL TEXT CLASSIFICATION

In this section, we will first introduce the optimization framework for multi-label active learning. Next we will describe our sample selection strategy with SVM.

4.1 Optimization Framework for Multi-label Active Learning

The optimization goal of our multi-label active learner is to label data which can contribute the largest reduction of the expected model loss.

Let $P(\mathbf{x})$ be the input distribution. Denote the multi-label prediction function given training set D_l as f_{D_l} . The predicted label set of \mathbf{x} is $f_{D_l}(\mathbf{x})$. Suppose the true label set of \mathbf{x} is \mathbf{y} , then the estimated loss on \mathbf{x} can be written as $L(f_{D_l}(\mathbf{x}), \mathbf{y})$ (we will simplify $L(f_{D_l}(\mathbf{x}), \mathbf{y})$ by writing $L(f_{D_l})$ in the following discussion), and the expected loss of the learner can be expressed as follows:

$$\widehat{\sigma}_{D_l} = \int_{\mathbf{x}} \left(\sum_{\mathbf{y} \in \mathcal{Y}} L(f_{D_l}) P(\mathbf{y}|\mathbf{x}) \right) P(\mathbf{x}) d\mathbf{x} \quad (1)$$

As it is rather difficult to estimate $P(\mathbf{x})$ directly, a practical way to estimate $\widehat{\sigma}_{D_l}$ is to measure it over all the examples in D_u . Therefore we have

$$\widehat{\sigma}_{D_l} = \frac{1}{|D_u|} \sum_{\mathbf{x} \in D_u} \sum_{\mathbf{y} \in \mathcal{Y}} L(f_{D_l}) P(\mathbf{y}|\mathbf{x}) \quad (2)$$

The active learner will evaluate each possible set of unlabeled data D_s to find the optimal query set D_s^* . When D_s obtains its labels, it can be incorporated to the training set. Denote the new training set as $D'_l = D_l + D_s$, and the expected loss for the classifier trained on D'_l as $\widehat{\sigma}_{D'_l}$. The optimization problem is to find the optimal query set D_s^* , which once added, will generate the largest reduction on expected loss.

$$\begin{aligned} D_s^* &= \arg \max_{D_s} (\widehat{\sigma}_{D_l} - \widehat{\sigma}_{D'_l}) \\ &= \arg \max_{D_s} \left(\sum_{\mathbf{x} \in D_u} \sum_{\mathbf{y} \in \mathcal{Y}} (L(f_{D_l}) - L(f_{D'_l})) P(\mathbf{y}|\mathbf{x}) \right) \end{aligned} \quad (3)$$

As in [1], we assume that any \mathbf{x} in $D_u - D_s$ has equal impact on the learner trained from D_l and D'_l . Then we will have

$$D_s^* = \arg \max_{D_s} \left(\sum_{\mathbf{x} \in D_s} \sum_{\mathbf{y} \in \mathcal{Y}} (L(f_{D_l}) - L(f_{D'_l})) P(\mathbf{y}|\mathbf{x}) \right) \quad (4)$$

4.2 Sample Selection Strategy with SVM

According to Equation 4, the optimization problem can be divided into two parts: how to measure the loss reduction of

the multi-label classifier and how to provide a good probability estimation for the conditional probability $p(\mathbf{y}|\mathbf{x})$. We will address these two issues respectively in the following subsections.

4.2.1 Estimate Loss Reduction

As discussed in Section 3, we decompose the multi-label problem to one-versus-all subproblems and use SVM as the base binary classifier in active learning. By decomposing the classifier into several binary ones, the overall model loss can be measured by gathering the model loss of all binary classifiers.

$$L(f) = \sum_{i=1}^k l(f^i), \quad (5)$$

where $l(f^i)$ is the model loss of binary classifier f^i . So the problem becomes how to estimate the model loss of each binary classifier. As suggested by S. Tong et al. [21], we measure the model loss by the size of version space of a binary SVM. According to [21], the version space of SVM can be defined as follows:

$$V = \{\mathbf{w} \in W \mid \|\mathbf{w}\| = 1, y_i(\mathbf{w} \cdot \mathbf{x}_i) > 0, i = 1, \dots, n\} \quad (6)$$

where W denotes the parameter space. The size of a version space is defined as the surface area of the hypersphere $\|\mathbf{w}\| = 1$ in W .

Based on the work in [21], we can use SVM margin as the measure of the version space size. When a new labeled example is added, we can approximate the new version space size by computing the SVM margin of the updated classifier. However, it is too expensive in computation when each data in the unlabeled pool associated with each possible label set needs to be evaluated. To make it more practical, we apply the heuristics idea in [20] to simplify the approximation by mapping the SVM margin of the current classifier to the size of the new version space.

In multi-label settings, denote $V_{D_l}^i$ as the size of version space of the binary classifier $f_{D_l}^i$ associated with target class i and learnt from labeled data D_l . After adding new data point (\mathbf{x}, y^i) , where $y^i \in \{-1, +1\}$ is the true label for data \mathbf{x} on class i , the new model loss versus the old one on the binary classifier $f_{D_l}^i$, can be approximated by:

$$\frac{l(f_{D_l+\mathbf{x}, y^i}^i)}{l(f_{D_l}^i)} \approx \frac{V_{D_l+\mathbf{x}, y^i}^i}{V_{D_l}^i} \approx \frac{1 + y^i f_{D_l}^i(\mathbf{x})}{2} \quad (7)$$

Then the loss reduction part in Equation 4 can be re-written by:

$$\begin{aligned} L(f_{D_l}) - L(f_{D'_l}) &= \sum_{i=1}^k (l(f_{D_l}^i) - l(f_{D'_l}^i)) \\ &= \sum_{i=1}^k (l(f_{D_l}^i) \cdot (1 - \frac{l(f_{D'_l}^i)}{l(f_{D_l}^i)})) \end{aligned} \quad (8)$$

Note that $l(f_{D_l}^i)$ has nothing to do with the selected unlabeled example \mathbf{x} , so we can focus on optimizing the reduction rate, which can be approximated as

$$\sum_{i=1}^k \left(\frac{1 - y^i f_{D_l}^i(\mathbf{x})}{2} \right) \quad (9)$$

Intuitively, the idea of the above estimation can be explained as follows. Consider an unlabeled data example \mathbf{x} ,

if \mathbf{x} can be correctly predicted by the binary classifier f^i , then the smaller the value of $|f^i(\mathbf{x})|$ is, the more uncertain the classifier is on \mathbf{x} , and \mathbf{x} deserves more to label. This is consistent with the result of the above measure, since \mathbf{x} will contribute more in reducing the size of the version space. On the other hand, if the classifier provides wrong prediction result for \mathbf{x} , then the larger $|f^i(\mathbf{x})|$ is, the more mistake the classifier will make, and in another view, adding \mathbf{x} will greatly help reduce the size of the version space.

4.2.2 Label Prediction

Now we come to the issue of estimating the conditional probability $p(\mathbf{y}|\mathbf{x})$, $\mathbf{y} \in \mathcal{Y}$. Note that for k labels, there are 2^k possible label combinations. It is intractable for active learner to provide estimation on all these possibilities. Particularly, it will become harder when the training data is quite limited, which is common in active **learning**. To simplify the estimation, we approximate the expected loss with the loss on the most possible label combination, since the predicted labels with the largest confidence will be most likely to be correct. Thus the problem becomes how to produce better label prediction on the unlabeled data. We propose a novel prediction approach to address this problem. Instead of directly estimating the possible labels for each data, we first try to decide the possible label number each data may have, and then determine the final labels based on the probability on each label obtained by the corresponding binary classifier.

Suppose there are k classes. Using the one-versus-all approach, we can have k binary classifiers. Given data \mathbf{x} , denote $p(y^i = 1|\mathbf{x})$ as the probability of \mathbf{x} belonging to class i . We can obtain k classification probabilities on \mathbf{x} produced by the k binary classifiers. Sort these k probabilities in decreasing order. If \mathbf{x} actually has m labels, the first m probabilities are expected to be large while the other $k - m$ probabilities are expected to be small. Based on this assumption, we want to predict the number of labels for each data based on the probabilities output by the binary classifiers.

Specifically, we predict the number of labels by tackling a multi-class classification problem. Logistic regression (LR) algorithm is used to train a multi-class model and predict the probabilities of having different number of labels for each data. For k classes, there are k possible number of labels: $1, \dots, k$. So we have k classes in the multi-class classification problem. Before LR is used, we transform the decision output on the training data to classification probabilities. Here, we use the sigmoid function [13] to transform the SVM output to probability values. For a data example \mathbf{x} , we have

$$p(y^i = 1|\mathbf{x}) = \frac{1}{1 + \exp(Af^i(\mathbf{x}) + B)}$$

where f^i is the binary SVM classifier associated with class i , A and B are scalar values fit by maximum likelihood estimation.

The process of predicting number of labels can be described as follows:

1. Use the SVM classifier to assign classification probabilities for all data examples.
2. For each instance \mathbf{x} , sort the classification probabilities in decreasing order, $p(y^{i_1} = 1|\mathbf{x}) \geq p(y^{i_2} = 1|\mathbf{x}) \geq \dots \geq p(y^{i_k} = 1|\mathbf{x})$. Normalize the classification probabilities and obtain $q_1(\mathbf{x}), \dots, q_k(\mathbf{x})$, where

$$q_p(\mathbf{x}) = \frac{p(y^{i_p} = 1|\mathbf{x})}{\sum_{t=1}^k p(y^{i_t} = 1|\mathbf{x})}.$$

3. Train logistic regression classifier. For each training data \mathbf{x} , present $[1 : q_1(\mathbf{x}), 2 : q_2(\mathbf{x}), \dots, k : q_k(\mathbf{x})]$ as the training features for LR model. The number of labels of \mathbf{x} is used as the category to train a multi-class classifier.
4. For each data in the unlabeled pool, apply the LR classifier to predict the probabilities of having different number of labels, and output the label with the largest probability to be the predicted number of labels for the data.

Suppose the most possible number of labels for data \mathbf{x} is m , and i_1, \dots, i_m are the m classes associated with the largest probabilities produced by the m corresponding binary SVM classifiers. Then the predicted label vector $\hat{\mathbf{y}}$ can be represented by the binary vector $[\hat{y}^{i_1} = 1, \dots, \hat{y}^{i_j} = 1, \hat{y}^{i_{j+1}} = -1, \dots, \hat{y}^{i_k} = -1]$. We call this approach *LR-based* label prediction.

By incorporating the predicted label vector into the expected loss estimation, we obtain our data selection strategy, Maximum loss reduction with Maximal Confidence(MMC). It can be written as

$$D_s^* = \arg \max_{D_s} \left(\sum_{\mathbf{x} \in D_s} \sum_{i=1}^k \left(\frac{1 - \hat{y}^i f^i(\mathbf{x})}{2} \right) \right), \quad (10)$$

Based on the above discussion, the proposed active **learning** algorithm is described in Algorithm 1.

Algorithm 1 Multi-label Active **Learning**

Input: Labeled set D_l

Unlabeled set D_u

Number of classes k

Number of iterations T

Number of selected examples per iteration S

- 1: **for** $t = 1$ to T **do**
 - 2: Train k binary SVM classifiers f^1, \dots, f^k based on training data D_l
 - 3: **for** each instance \mathbf{x} in D_u **do**
 - 4: Predict its label vector using the LR-based prediction method described in Section 4.2.2.
 - 5: Calculate the expected loss reduction with the most confident label vector $\hat{\mathbf{y}}$, $score(\mathbf{x}) = \sum_{i=1}^k \left(\frac{1 - \hat{y}^i f^i(\mathbf{x})}{2} \right)$
 - 6: Sort $score(\mathbf{x})$ in decreasing order for all \mathbf{x} in D_u
 - 7: Select a set of S examples D_s^* with the largest scores, and update the training set $D_l \leftarrow D_l + D_s^*$
-

5. EXPERIMENTS

In this section, we will evaluate our proposed multi-label active **learning** approach for multi-label text classification task on seven real-world data sets, comparing with the state-of-the-art active **learning** approaches.

Table 1: Statistics on RCV1-V2 and Yahoo! Data Sets

Data sets	#Samples	#Features	#Label
RCV1-V2	3,000	47,236	101
Arts&Humanities	7,484	23,146	26
Business&Economy	11,214	21,924	30
Computers&Internet	12,444	34,096	33
Education	12,030	27,534	33
Entertainment	12,730	32,001	21
Health	9,205	30,605	32

5.1 Data Sets and Experiment Settings

The first data set¹ we used is the RCV1-V2 [11] text data set, which has been widely used as a benchmark data set to evaluate text classification algorithms. It contains Reuters newswire stories which are organized by three different category sets: Topics, Industries, and Regions. Each document is assigned with at least one label in the related category set. A sample of 3,000 documents in the Topics category set is chosen for our experiments, including 101 labels.

The other 6 data sets² are web pages collected through the hyperlinks from Yahoo!’s top directory (www.yahoo.com). They are used in [15, 8] to evaluate multi-label text classification algorithms. Each data set is associated with one of Yahoo!’s top categories, and each page is labeled with one or more second level sub-categories. We choose 6 data sets for our experiments, which are Arts&Humanities, Business&Economy, Computers&Internet, Education, Entertainment, and Health.

The details of all the 7 data sets are given in Table 1. “#Samples” is the number of samples in each data set. “#Features” is the feature dimension of each data set. “#Label” is the number of labels in each data set.

On all data sets, the documents are transformed to vectors with TF-IDF format, and each vector has unit modulus with L-2 length normalization. One-versus-all classification is conducted for each category and the multi-label classification problem is treated as several binary classification problems, where the documents from the target category are given positive label (i.e. $y = 1$), and the rest of the documents are given negative label (i.e. $y = -1$). *SVM^{Light}* package [7] is downloaded and used to train the binary classifier. Linear kernel is used due to its good performance in text classification task [6]. The penalty parameter C is set to 1.0 by default.

In our active **learning** experiments on each data set, we first randomly selected a small set of documents to form the initial labeled set, and left the remaining documents as the unlabeled pool. Then the active learner selects a given number of examples from the unlabeled pool in each iteration, and then add them to the labeled set with their labels. We performed several active **learning** iterations on each data set until the learner achieves sufficient accuracy. In every iteration, once the selected data being incorporated, the active learner retrained a new classifier on the expanded labeled set and its performance was evaluated on the remaining data examples. We used Micro-Average F1 score as the evaluation

measure, since it is a standard evaluation used in most previous text classification research. As defined in [23], micro-F1 score in multi-label case is given as follows

$$\frac{2 \sum_{j=1}^k \sum_{i=1}^n \hat{y}_i^j y_i^j}{\sum_{j=1}^k \sum_{i=1}^n \hat{y}_i^j + \sum_{j=1}^k \sum_{i=1}^n y_i^j}$$

where n is the number of test data, \mathbf{y}_i is the true label vector of the i -th data instance, $y_i^j = 1$ if the instance belongs to category j ; otherwise $y_i^j = -1$. $\hat{\mathbf{y}}_i$ is the predicted label vector. We computed the average of micro-F1 scores for each active **learning** iteration based on 10 randomized experiments.

In our experiments, we will evaluate and compare four active **learning** methods:

- MMC. The active **learning** method proposed in this paper.
- Random. The sample selection strategy is to randomly select data examples from the unlabeled pool to label.
- BinMin. This is a sample selection strategy proposed in [9], which is most related to our research work with respect to the problem studied. In this work, one-versus-all approach is used for multi-label classification, and SVM is used as the basic binary classifier. The optimal unlabeled example is selected according to

$$\arg \min_x \min_{i=1, \dots, k} |f^i(x)|$$

where f^i is the binary classifier on the binary problem associated with class i . That is, it selects unlabeled examples with respect to the most uncertain label. As stated in Section 2, this method does not take advantages of the multi-label information.

- Mean Max Loss(MML). This strategy is to select unlabeled data which has the maximum mean loss value over all the predicted labels [12]. For each predicted label j , the loss is measured as

$$\sum_{i=1}^k \max[(1 - m_{ij} f^i(x)), 0]$$

where $m_{ij} = 1$ if $i = j$, else $m_{ij} = -1$, and f^i is the binary SVM classifier on class i . The algorithm uses a threshold cutting method to decide the predicted labels. However, according to our experiments on the text data sets, this method is usually unable to pick out predicted labels correctly. Thus we replace the label prediction part with our LR-based prediction method in Section 4.2.2, and focus on evaluating the effectiveness of the loss optimization.

5.2 Results and Discussions

In this section, we will present and discuss the experiment results on the RCV1-V2 data set as well as the 6 Yahoo! data sets.

¹<http://trec.nist.gov/data/reuters/reuters.html>

²<http://www.kecl.ntt.co.jp/as/members/ueda/yahoo.tar.gz>

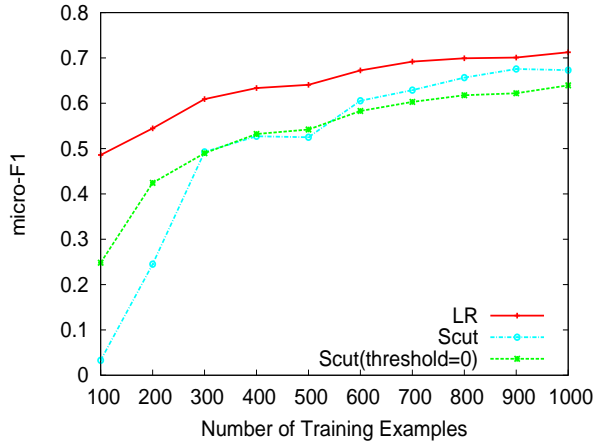


Figure 1: Comparison between label prediction methods on RCV1-V2 data set (no active learning)

Experimental Results with RCV1-V2 data set.

In the first experiment, we would like to verify whether our method of label prediction (presented in Section 4.2.2) is effective when only a small amount of training data is available, as this is very typical in active learning. Two popular prediction methods for multi-label classification are implemented for comparison purposes. In previous studies, the SCut method is widely used and proved very effective for predicting labels in multi-label classification tasks [11]. In [11], a binary classifier is first trained for each label. A threshold score is tuned for each binary classification task and then used to decide if an unlabeled data example belongs to the corresponding class or not. The second prediction method is simply setting the threshold score to be zero for each binary problem. If the classification score is positive, then the data belongs to this class, and vice versa. This simple method has its theoretical foundation, as when SVM is used, zero score corresponds with the classification hyperplane induced from statistical learning theory.

In order to verify the effectiveness of the LR-based method in predicting labels, we varied the number of training data from 100 to 1,000 (with 100 as step size). The corresponding micro-F1 curves for predicting labels are plot in Figure 1. We can observe that, as the number of training data varies, the LR-based method achieves substantially better performance than both baseline methods. When less training data is available, the advantage of LR is more obvious. This demonstrates that the LR-based method is more effective for label prediction in multi-label active learning framework.

In the following we will report the active learning experimental results. We randomly selected 500 examples as the initial labeled data. Active learning was iteratively performed for 50 iterations, selecting 20 examples from the unlabeled pool each time. Figure 2 and Table 2 show the experimental results of micro-F1 scores averaging over 10 random trials. The proposed MMC strategy outperforms other baseline methods by a large margin. Surprisingly, we can see that MML performs even worse than Random at the beginning, and worse than MMC and BinMin for all cases. Since MML adopts the same label prediction approach as MMC, the observation above indicates that the loss optimization

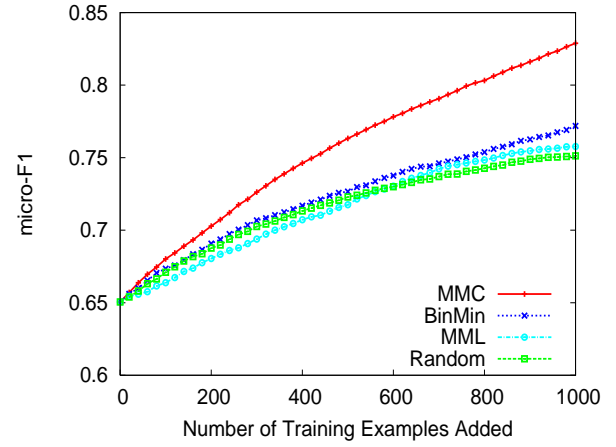


Figure 2: Micro-F1 score on RCV1-V2 data set

Table 2: Micro-F1 score at different iterations on RCV1-V2 data set(%)

K	MMC	BinMin	MML	Random
100	68.02	67.35	66.38	67.10
200	70.28	69.08	68.05	68.77
300	72.62	70.68	69.39	70.26
400	74.62	71.69	70.72	71.33
500	76.33	72.66	71.75	72.29
600	77.81	73.76	73.04	72.99
700	79.07	74.61	74.23	73.69
800	80.32	75.37	74.84	74.27
900	81.62	76.25	75.47	74.89
1000	82.88	77.19	75.77	75.12

approach used in MML is not effective on multi-label text data. Instead, our approach optimizes the loss reduction rate over all labels based on the most confident label vector, and it can successfully pick out useful data examples to label. We can also find that our proposed method outperforms all baseline methods more significantly than BinMin. An explanation is that the BinMin strategy does not take advantage of the multi-label information, while our approach effectively estimates possible labels for each instance and incorporates the multi-label information to optimize the expected loss reduction.

Table 2 shows the performance results with the number of training samples added. We can find that as the number of selected data increases, the improvement becomes more and more significant. For example, when 1,000 examples are added, the micro-F1 score of our method achieves 82.88%, while that of BinMin, MML and Random are 77.19%, 75.77% and 75.12% respectively. We can find that MMC achieves the similar performance with BinMin by using about 600 selected examples, while BinMin needs to select 1,000 examples. It indicates that MMC can save about 40% labeling effort compared with BinMin.

In order to investigate if our MMC algorithm is sensitive to the size of initial labeled data set, we varied the number of initial training data from 100 to 1,000, with 100 as step size. For each fixed initial labeled set, we applied active learning and selected 20 examples at each iteration. Then we com-

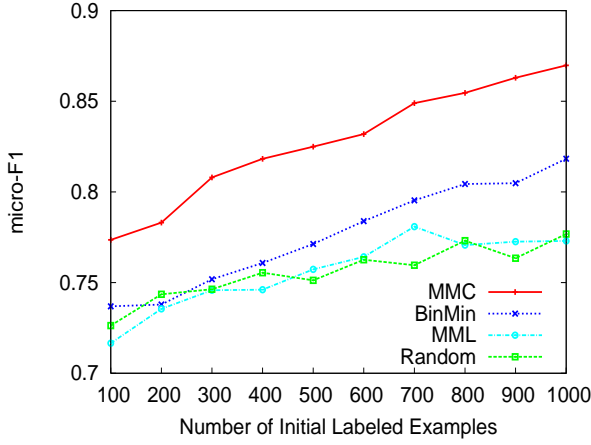


Figure 3: Micro-F1 score on RCV1-V2 data set after adding 1000 examples

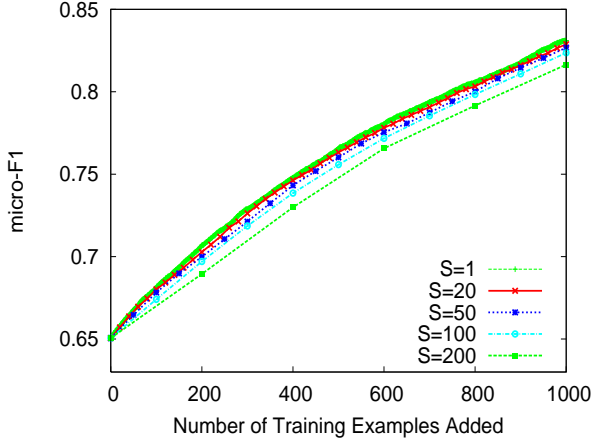


Figure 4: Micro-F1 score of MMC on RCV1-V2 data set with different sampling sizes per run

pared the performance of the final classifier after 50 active learning iterations. Figure 3 presents the micro-F1 scores of final classifiers with the size of initial training data set. We can see that our proposed MMC algorithm consistently outperforms all other methods when the initial training data set varies in size. The consistent improvement indicates that our MMC strategy is robust with different size of the initial labeled data set.

We also varied the sampling size per run and investigated its impact on the performance of the active learner. In this experiment, we started with 500 training examples and stopped after 1,000 examples are added. The sampling size S was set to 1, 20, 50, 100 and 200. The results of MMC with various sampling size are depicted in Figure 4. We can see that generally the performance improves as the sampling size decreases. A possible explanation is that having more chances to query labels enables the learner to make better evaluation on unlabeled examples, and to choose more informative examples to label.

Table 3: Micro-F1 score on the Yahoo! data sets with 2,500 training samples added (%)

Data sets	MMC	BinMin	MML	Random
Arts&Humanities	65.03	61.67	60.26	58.74
Business&Economy	80.54	78.37	77.08	75.90
Computers&Internet	77.13	75.05	74.37	73.94
Education	71.28	69.29	66.97	67.65
Entertainment	75.46	74.46	73.14	70.82
Health	81.05	79.56	74.74	74.60

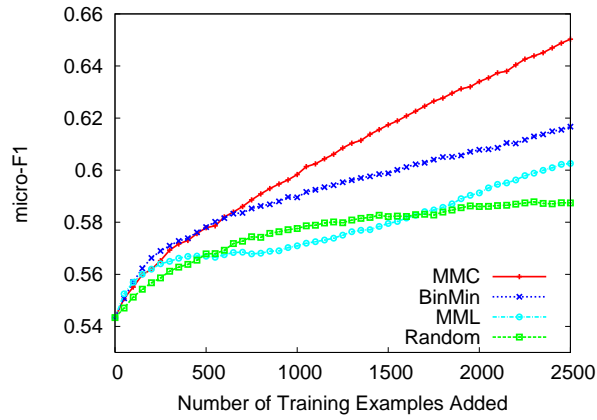
Experimental Results with Yahoo! data sets.

The following experiments are conducted with the 6 Yahoo! data sets. On each data set, we randomly selected 500 data instances as the initial training data, and set the sampling size in each active learning run to 50. The learning process was repeated for 50 rounds. The active learning results were averaged over 10 random trials. Fig 5 presents the performance of all active learners with the number of training data added. We can observe that our proposed method MMC outperforms other baseline methods on all six data sets. The most noticeable case is the Computers&Internet data set, where the BinMin method only provides slight improvement over the Random method. However, MMC achieves substantially better performance. It can be observed that MMC only requires labeling 900 examples to achieve the similar performance with BinMin and MML which require labeling about 1,600 and 2,100 examples respectively. We can also see that MML has worse performance compared with Random on most of the cases. This implies that the loss optimization framework of MML is worse than that of MMC on multi-label text data. Compared with MMC, BinMin is less effective to enhance the active learner as the training example grows. This underscores the importance of considering multi-label information when evaluating unlabeled examples. The promising results of MMC confirm that the proposed method can provide proper evaluation on the unlabeled data examples, and select the informative ones which can help enhance the learner more effectively. Table 3 summarizes the classification results measured by micro-F1 after 50 active learning iterations on the six Yahoo! data sets. It shows that the proposed MMC method provides more favorable performance than all other baseline methods for all six data sets, and the improvement of MMC over Random is more significant than that of BinMin and MML.

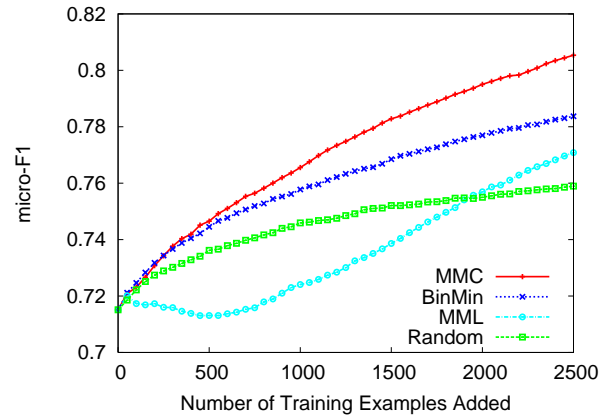
From the above experiments, we can observe that MMC provides promising performance on diverse data sets. This indicates that it is more effective and robust for training multi-label text classifier than the state-of-the-art active learning methods.

6. CONCLUSIONS

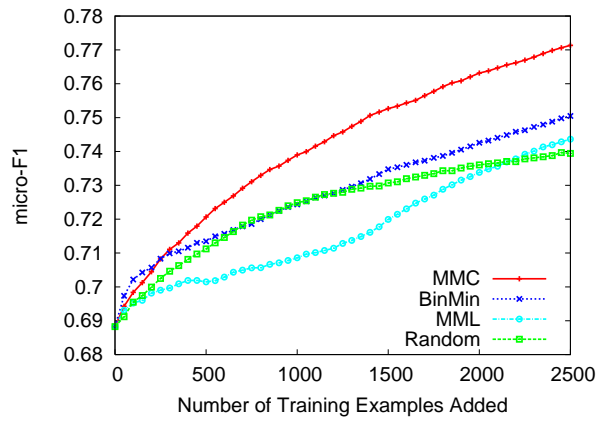
In this paper, we try to address the problem of multi-label active learning for text classification. The goal is to reduce the required size of labeled data in multi-label classification while maintaining favorable accuracy performance. We propose a novel multi-label active learning algorithm with Support Vector Machines (SVM). The optimization goal is to select data to label which can maximize the reduction in the expected model loss. Our approach provides proper ap-



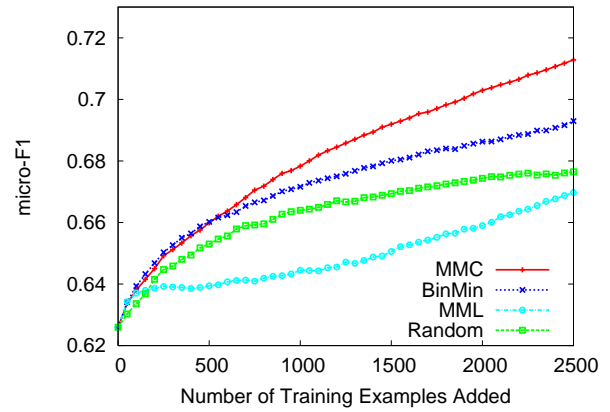
(a) Arts&Humanities



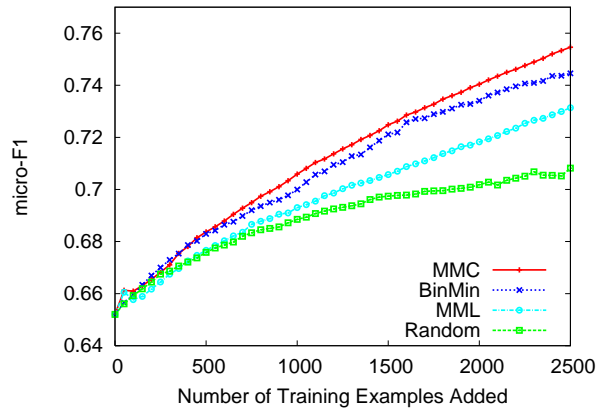
(b) Business&Economy



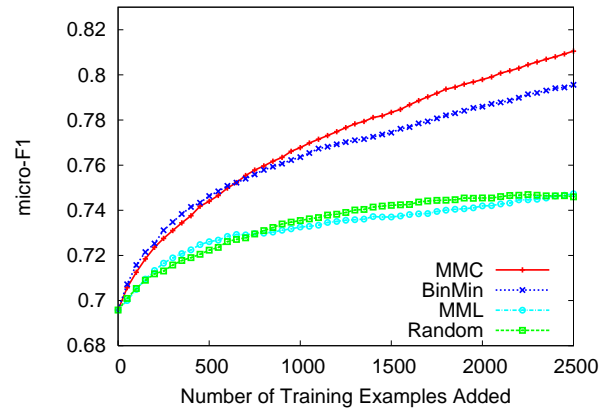
(c) Computers&Internet



(d) Education



(e) Entertainment



(f) Health

Figure 5: Micro-F1 score on Yahoo! data sets

proximation on the loss reduction and the expected loss in the optimization framework. Experiments on several real-world data sets show that our proposed method outperforms the state-of-the-art active **learning** techniques on multi-label text classification by a large margin and can significantly reduce the labeling cost.

Note that our active **learning** approach should evaluate each of the unlabeled data at every active **learning** iteration. The computation would be expensive when the size of unlabeled pool is very large and the number of categories is very big. So it would be interesting to study how to evaluate only a subset of the unlabeled pool and also be able to pick out informative data to label. We plan to explore this extension in the future. Also, we will apply our method on other multi-label classification tasks, e.g., image classification.

7. ACKNOWLEDGMENTS

We express our grateful thanks to Prof. Jian Pei from Simon Fraser University for his valuable suggestions on this work.

8. REFERENCES

- [1] C. Campbell, N. Cristianini, and A. J. Smola. Query **learning** with large margin classifiers. In *Proceedings of the 7th International Conference on Machine Learning (ICML'00)*, pages 111–118, 2000.
- [2] D. A. Cohn, Z. Ghahramani, and M. I. Jordan. Active **learning** with statistical models. In *Advances in Neural Information Processing Systems*, volume 7, pages 705–712. The MIT Press, 1995.
- [3] C. Cortes and V. Vapnik. Support vector networks. In *Machine Learning*, pages 273–297, 1995.
- [4] A. Esuli and F. Sebastiani. Active **learning** strategies for multi-label text classification. In *Proceedings of the 31th European Conference on Information Retrieval (ECIR'09)*, pages 102–113, 2009.
- [5] R.-E. Fan and C.-J. Lin. A study on threshold selection for multi-label classification. *Technical Report, National Taiwan University*, 2007.
- [6] T. Joachims. Text categorization with support vector machines: **Learning** with many relevant features. pages 137–142. Springer Verlag, 1998.
- [7] T. Joachims. *Learning to Classify Text Using Support Vector Machines: Methods, Theory and Algorithms*. Kluwer Academic Publishers, Norwell, MA, USA, 2002.
- [8] H. Kazawa, T. Izumitani, H. Taira, and E. Maeda. Maximal margin labeling for multi-topic text categorization. In *Advances in Neural Information Processing Systems (NIPS'05)*, pages 649–656, 2005.
- [9] K. Brinker. *On Active Learning in Multi-label Classification. "From Data and Information Analysis to Knowledge Engineering" of Book Series "Studies in Classification, Data Analysis, and Knowledge Organization"*, Springer, 2006. 1, 2.
- [10] D. D. Lewis and W. A. Gale. A sequential algorithm for training text classifiers. In *Proceedings of the 17th annual international ACM SIGIR conference on Research and development in information retrieval (SIGIR'94)*, pages 3–12, 1994.
- [11] D. D. Lewis, Y. Yang, T. G. Rose, G. Dietterich, F. Li, and F. Li. Rcv1: A new benchmark collection for text categorization research. *Journal of Machine Learning Research*, 5:361–397, 2004.
- [12] X. Li, L. Wang, and E. Sung. Multi-label svm active **learning** for image classification. In *International Conference on Image Processing*, pages 2207–2210, 2004.
- [13] H.-T. Lin, C.-J. Lin, and R. C. Weng. A note on Platt's probabilistic outputs for support vector machines. *Journal of Machine Learning Research*, 68(3):267–276, 2007.
- [14] T. Luo, K. Kramer, D. B. Goldgof, L. O. Hall, S. Samson, A. Remsen, and T. Hopkins. Active **learning** to recognize multiple types of plankton. *Journal of Machine Learning Research*, 6:589–613, 2005.
- [15] N. Ueda and K. Saito. Single-shot detection of multiple categories of text using parametric mixture models. In *Proceedings of the 8th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD'02)*, pages 626–631, 2002.
- [16] G.-J. Qi, X.-S. Hua, Y. Rui, J. Tang, and H.-J. Zhang. Two-dimensional active **learning** for image classification. *IEEE Conference on Computer Vision and Pattern Recognition*, 2008.
- [17] G.-J. Qi, X.-S. Hua, Y. Rui, J. Tang, and H.-J. Zhang. Two-dimensional multi-label active **learning** with an efficient online adaptation model for image classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 99(1), 2008.
- [18] N. Roy and A. McCallum. Toward optimal active **learning** through sampling estimation of error reduction. In *Proceedings of the 8th International Conference on Machine Learning (ICML'01)*, pages 441–448, 2001.
- [19] H. S. Seung, M. Oppor, and H. Sompolinsky. Query by committee. In *Proceedings of the 5th annual workshop on Computational learning theory (COLT'92)*, pages 287–294, 1992.
- [20] S. Tong. *Active Learning: Theory and Applications*. PhD thesis, Stanford University, CA, 2001.
- [21] S. Tong and D. Koller. Support vector machine active **learning** with applications to text classification. *Journal of Machine Learning Research*, 2:45–66, 2002.
- [22] R. Yan, J. Yang, and A. Hauptmann. Automatically labeling video data using multi-class active **learning**. In *Proceedings of the 9th IEEE International Conference on Computer Vision (ICCV'03)*, page 516, 2003.
- [23] Y. Yang. A study on thresholding strategies for text categorization. In *Proceedings of 24th International Conference on Research and Development in Information Retrieval (SIGIR'01)*, pages 137–145, 2001.