

Active Learning Using Uncertainty Information

Yazhou Yang*[†], Marco Loog*[‡]

*Pattern Recognition Laboratory, Delft University of Technology, Delft, The Netherlands

[†]College of Information System and Management, National University of Defense Technology, Changsha, China

[‡]The Image Section, University of Copenhagen, Copenhagen, Denmark

Email: {y.yang-4, m.loog}@tudelft.nl

Abstract—Many active learning methods belong to the retraining-based approaches, which select one unlabeled instance, add it to the training set with its possible labels, retrain the classification model, and evaluate the criteria that we base our selection on. However, since the true label of the selected instance is unknown, these methods resort to calculating the average-case or worse-case performance with respect to the unknown label. In this paper, we propose a different method to solve this problem. In particular, our method aims to make use of the uncertainty information to enhance the performance of retraining-based models. We apply our method to two state-of-the-art algorithms and carry out extensive experiments on a wide variety of real-world datasets. The results clearly demonstrate the effectiveness of the proposed method and indicate it can reduce human labeling efforts in many real-life applications.

I. INTRODUCTION

Over the past decade, a primary foundation of much progress in machine learning is the rapid growth of the number and size of data sets available, such as ImageNet [1] containing over 14 million labeled images for object recognition. In a practical scenario, we frequently encounter the situation where few labeled instances along with abundant unlabeled samples are available. Labeling a large amount of data is, however, very difficult due to the huge amount of time required or expensive because of the need of human experts [2]. Thus, it is very attractive to propose a proper labeling scheme to reduce the number of labels required in order to train a classifier.

Active learning has been put forward to overcome the above labeling problem. The main assumption behind active learning is that if an active learner can freely select any samples it wants, it can outperform random sampling with less labeling [2]. Thus, the main task of active learning is querying as little data as possible to minimize the annotation cost while maximizing the learning performance. Active learning tries to achieve this by selecting the most valuable samples. However, it is difficult to define or measure the value of one instance to the learning problem. We can view it as the amount of information carried which potentially promotes the learning performance, once its true label is known [3]. As a result of the fact that we do not have an exact measure of the value, there are a great number of selection criteria proposed from different perspectives on how to estimate the usefulness of each sample.

Most commonly used criteria in active learning include query-by-committee [4], uncertainty sampling [5]–[7], expected error reduction [8]–[11], expected model change [12]–[15], variance reduction [16]–[19] and “Min-max” view active learning [20], [21]. Query-by-committee put forward multiple mod-

els as the committees and selected the samples which receive highest level of disagreement from the committees [4]. Uncertainty sampling approach preferred the instances with maximum uncertainty. Based on the measurement of uncertainty, uncertainty sampling can be roughly divided two categories: maximum entropy of the estimated label [5] and minimum distance from the decision boundary [6], [7]. For example, Tong and Koller [6] proposed to query the instance which is closed to the current learning boundary using the classifier of support vector machines. Campbell *et al.* [7] shared the same idea with Tong and Koller [6].

Roy and McCallum [8] proposed the expected error reduction (EER), which is a popular active learning method. EER aimed to reduce the generalization error when labeling a new instance. Since we do not have access to the test data, Roy and McCallum suggested to compute the “future error” on the unlabeled pool under the assumption that the unlabeled data set is representative of the test distribution. In other words, the unlabeled pool can be viewed as a validation set. Also, we have no knowledge about the true labels of unlabeled samples. EER estimated the average-case criterion of potential loss instead. Expected model change followed the idea of EER, but turned to select the instance which leads to maximum change of the current model. The variance reduction methods tried to minimize the output variances [2]. Schein and Ungar [19] extended this approach to expected variance reduction method on logistic regression by following the idea of EER. “Min-max” view active learning was originally proposed by Hoi *et al.* [20], where “Min-max” indicates the worst-case criterion is adopted. The key idea behind is to select the sample which minimizes the gain of objective function no matter what its assigned label is. Huang *et al.* [21] extended this framework by taking into account all the unlabeled data when calculating the objective function.

Current active learning methods can be split in two classes: retraining-based and retraining-free active learning. Retraining-based active learning represents methods which measure the information of unlabeled sample by labeling it (any possible label) and adding it to the training set to retrain the classification model. Then, some appropriate criteria can be evaluated and used for the sample selection. The second class, retraining-free active learning, contains the remaining methods which not need repeatedly train the model for each unlabeled instance during one single selection. For example, uncertainty sampling and query-by-committee belong to this category.

However, since the true label of the selected unlabeled

instance is unknown, these methods resort to calculating the average-case or worse-case criteria with respect to the unknown label. In this paper, we propose a different criterion for retraining-based methods. We incorporate the uncertainty information (measured by the posterior probabilities within the min-max framework) for the selection. The proposed criterion can be seen as a trade-off of the exploration and the exploitation. The uncertainty information plays the role of the exploitation while the retraining-based models act as the exploration part. We concentrate on the pool-based active learning setting which assumes a large pool of unlabeled data along with a small set of labeled data already available [2]. We consider the myopic active learning which sequentially and iteratively selects unlabeled instance.

A. Outline

The rest of this paper is organized as follows. Section II firstly reviews the framework of retraining-based active learning. Then two state-of-the-art methods under the retraining framework are briefly described. Section III demonstrates the primary motivation of the proposed method and derives a general algorithm for retraining-based active learning in detail. It also illustrate how to extend the proposed criterion to current methods. Experimental design and results are reported in IV ; Section V concludes this work followed by some future issues.

II. RETRAINING-BASED ACTIVE LEARNING

In this section, we summarize a general framework of retraining-based active learning. Then we demonstrate two examples under this framework: Expected error reduction and Minimum Loss Increase.

A. Retraining-based Active Learning

Firstly, let us introduce some preliminaries and notation. Let $\mathcal{L} = \{(x_i, y_i)\}_{i=1}^m$ represent the training data set that consists of m labeled instances and \mathcal{U} be the pool of unlabeled instances $\{x_i\}_{i=m+1}^n$. Each $x_i \in \mathbb{R}^d$ is a d dimensional feature vector, and $y_i \in C = \{+1, -1\}$ is the class label of x_i . In this paper, let us focus on binary classification problem firstly, and it is easy to extend this work to multi-class problem by extending C to multi-labels set. We denote $P_{\mathcal{L}}(y|x)$ be the conditional probability of y given x according to a classifier trained on \mathcal{L} .

For the retraining-based active learning, its framework can be summarized in Algorithm 1, where $V(x_i, y_i)$ represents any selection criterion associated with (x_i, y_i) . The main procedure contains the loops which checks all the points in unlabeled pool \mathcal{U} over all the possible labels. For example, we firstly select one instance from the unlabeled pool and assign it any possible label. Then we update the labeled set (since we acquire a new labeled sample) and retrain the classifier we use. Based on the new trained classifier, we can measure some kind of selection criteria (*e.g.*, generalization error in EER [8]). However, since the true label information of last selected sample is unknown, we need calculate some kind of performance, *e.g.*, the average-case in [8], [13], [19], worst-case in [21], or even the best-case criteria in [9]. Finally, we will query the instance which leads

Algorithm 1 General Retraining-based Active Learning Procedure

- 1: **Input:** Labeled data \mathcal{L} , unlabeled data \mathcal{U}
 - 2: **repeat**
 - 3: Train the classifier on \mathcal{L} and calculate $P_{\mathcal{L}}(y_i|x_i)$ for each $x_i \in \mathcal{U}$, each $y_i \in C$;
 - 4: **for** each $x_i \in \mathcal{U}$ **do**
 - 5: **for** each $y_i \in C$ **do**
 - 6: Re-train the model on $\mathcal{L} \cup \{x_i, y_i\}$;
 - 7: Calculate some criterion $V(x_i, y_i)$, (*e.g.*, error or variance);
 - 8: **end for**
 - 9: **end for**
 - 10: Compute some kind of performance based on $P_{\mathcal{L}}(y_i|x_i)$ and $V(x_i, y_i)$;
 - 11: Query the instance x^* which leads to the best performance and label it y^* , update $\mathcal{L} \leftarrow \mathcal{L} \cup \{x^*, y^*\}, \mathcal{U} \leftarrow \mathcal{U} \setminus \{x^*\}$;
 - 12: **until** Stopping criterion is satisfied
-

to maximum or minimum value in terms of the criterion we are interested in.

EER is one example of retraining-based active learning, which uses the generalization error as $V(x_i, y_i)$. We get expected model change [12]–[15] by adopting model change as the criterion. By adopting variance and logistic regression as the classifier, we get expected variance reduction [19]. Similarly, if we want to minimize the value of objective function after labeling a new instance and use the worst-case performance (corresponding to min-max framework), then we can get [20], [21]. Clearly, the retraining-based approaches may suffer from high computational cost due to the fact that they need go over all the unlabeled data and all the possible labels.

B. Expected Error Reduction

Expected error reduction has demonstrated its effectiveness on text classification domain [8]. There are also some follow-up work of EER contributed by other researchers [9] [10] [11]. EER aims to select the sample which will reduce the future generalization error. Since we can not see the test data, the unlabeled pool can be used as the validation set to predict the future test error. We encounter a new problem since we do not know the true labels of the pool. Roy and McCallum [8] suggested, in practice, we can approximately estimate the error using the expected log-loss or 0/1 loss over the pool. For example, if we adopt the log loss, EER can be written as follows:

$$\arg \min_{x \in \mathcal{U}} \sum_{y \in C} P_{\mathcal{L}}(y|x) \left(- \sum_{x_i \in \mathcal{U}} \sum_{y_i \in C} P_{\mathcal{L}^+}(y_i|x_i) \log P_{\mathcal{L}^+}(y_i|x_i) \right)$$

where $\mathcal{L}^+ = \mathcal{L} \cup (x, y)$ means that the selected instance x is labeled y and added to \mathcal{L} . Note that the first term $P_{\mathcal{L}}(y|x)$ contains the pre-trained label information. The second term is the sum of potential entropy over the unlabeled data set \mathcal{U} .

C. Minimum Loss Increase

We can find that EER attempts to reduce the future generalization error, however, it is not easy due to the missing of

test data and true label information of unlabeled data. There are some researchers which try to solve this problem from a different perspective. Hoi *et al.* [20] presented a so called “min-max” view active learning. It prefers the instance which results in a small value of an objective function in spite of its assigned label. This is because the smaller the value of an objective function, the better the learning model, at least in high probability. Assume $G_{\mathcal{L}}$ is the value of an objective function on current labeled data \mathcal{L} . When we label a new instance and update the training data $\mathcal{L}^+ = \mathcal{L} \cup \{x_i, y_i\}$, we get a new value of objective function $G_{\mathcal{L}^+}$. What we want is the minimum increase of objective function, *i.e.*, $G_{\mathcal{L}^+} - G_{\mathcal{L}}$, when adding one more labeled sample. However, because the second term $G_{\mathcal{L}}$ is independent of the next queried instance, so we can ignore it and focus on minimizing $G_{\mathcal{L}^+}$. Since we expect a minimum value of $G_{\mathcal{L}^+}$ regardless of the assigned label of x_i , we adopt the worst-case performance as follows, instead of the average-case version.

$$\arg \min_{x_i \in \mathcal{U}} \max_{y_i \in \mathcal{C}} G_{\mathcal{L}^+}$$

Note that we can view $G_{\mathcal{L}^+}$ as one choice of $V(x_i, y_i)$ mentioned in Algorithm 1.

Let us consider an unconstrained optimization problem using L_2 -loss regularized classifier with arbitrary loss $l(w; x_i, y_i)$: $g(w) = \frac{1}{2\lambda} \|w\|^2 + \sum_{x_i \in \mathcal{L}} l(w; x_i, y_i)$, where w is the parameter of learning classifier. If we adopt the Hinge loss $l(w; x_i, y_i) = \max(0, 1 - y_i w^T x_i)$, we can derive the same model with “min-max” view active learning described in [20], but without extend it to batch model setting. If we use square loss $l(w; x_i, y_i) = (y_i - w^T x_i)^2$, we can get the same model with [21]. Note that, as is stated in [22], though [21] includes all the unlabeled data when calculating the objective function, the unlabeled examples play no role since [21] relaxes the constraint of the labels of unlabeled pool in the end. This operation can guarantee *zero* contribution of unlabeled data to the objection function. Thus, [21] is also one special case using the square loss. Moreover, we can conclude that the main idea of min-max view active learning is to minimize the increase of the value of an objective function.

In our paper, we consider the logistic loss $l(w; x_i, y_i) = \log(1 + \exp^{-y_i w^T x_i})$, which results in:

$$\arg \min_{x_i \in \mathcal{U}} \max_{y_i \in \mathcal{C}} \frac{1}{2\lambda} \|\hat{w}\|^2 + \sum_{x_i \in \mathcal{L}^+} -\log P_{\mathcal{L}^+}(y_i|x_i) \quad (1)$$

where \hat{w} is estimated parameter of L_2 -regularized logistic regression model. Logistic regression is chosen as the base classifier since it is generally widely used in many fields and can output the conditional probability straightly, which can be used in active learning [22]. We call this method Minimum Loss Increase (MLI) in this paper. EER tries to minimize the error on unlabeled data while MLI aims to minimize the loss on data already labeled.

III. A NEW RETRAINING-BASED ACTIVE LEARNER

In this section, we motivate our proposed method and, subsequently, describe a general adaptation for retaining-based active learning models.

A. Motivation

Obviously, not knowing the true labels of the unlabeled data complicates calculating the final score of each instance in step 10 in Algorithm 1. One simple possibility is computing the average-case [8] or worst-case performance [21], or even the best-case criterion [9]. These choices, however, may fail to take into account some potentially valuable information: Firstly, although the average-case criterion makes use of the label distribution information $P_{\mathcal{L}}(y_i|x_i)$ already known, the expectation calculation can hide or underestimate some outstanding samples due to the re-weighting by $P_{\mathcal{L}}(y_i|x_i)$. For example, the true label of instance x_i is +1 but the estimated $P_{\mathcal{L}}(+1|x_i) = 0.1$, and the $V(x_i, +1)$ has a maximum value compared with other instances. Then the average-case criterion of x_i , namely $\sum_{y_i} P_{\mathcal{L}}(y_i|x_i)V(x_i, y_i)$, is highly likely to be surpassed by other instances. Secondly, as to the worst-case criterion, it suffers from not taking advantage of label distribution information at all. Worst-case analysis is a safe analysis since it is never underestimated. However, making no use of the available label information $P_{\mathcal{L}}(y_i|x_i)$ can lose sight of some valuable information.

Thus, to overcome the shortcomings mentioned, a new criterion for retraining-based active learning is proposed. The main motivation is that we want to incorporate the uncertainty information (*e.g.*, known label distribution information) within min-max framework for retraining-based models. The proposed criterion is therefore as follows:

$$\min_{x_i \in \mathcal{U}} \max_{y_i \in \mathcal{C}} P_{\mathcal{L}}(y_i|x_i)V(x_i, y_i) \quad (2)$$

where $P_{\mathcal{L}}(y_i|x_i)$ contains the pre-trained label information and $V(x_i, y_i)$ represents any criteria we are interested. Note that for some classifiers like logistic regression, we can use the estimated posterior probability as $P_{\mathcal{L}}(y_i|x_i)$. For classifiers which do not produce a probabilistic output, *e.g.*, SVMs, we can transform their output to some probability using Platt’s [23] or Duin & Tax’s method [24]. And for $V(x_i, y_i)$, various choices are possible, such as the test error on the unlabeled pool in EER, the output variance as in [19], or the value of an objective function [21].

The proposed method can be interpreted as follows: it utilizes the pre-trained label information, although this kind of information might be inaccurate due to limited labeled data we have, it still shows some underlying or potential useful clues which may promote active learning. Firstly, it improves upon the average-case criterion since it does not compute the expected value. The calculation of expectation tends to ruin the discriminative information contained in the data due to its averaging manner. Secondly, it outperforms the worst-case criterion because it takes advantage of the knowledge of the potential label distribution while worst-case analysis does not use this at all. Thus, it avoids the disadvantages of average-case and worst-case criteria. It can be seen as a trade-off between the average-case and the worst-case criteria. Lastly, it can be considered as incorporating uncertainty sampling (encoded by the posterior probabilities) for retraining-based model. If all $V(x_i, y_i)$ become one constant term like 1 or $P_{\mathcal{L}}(y_i|x_i)$ itself,

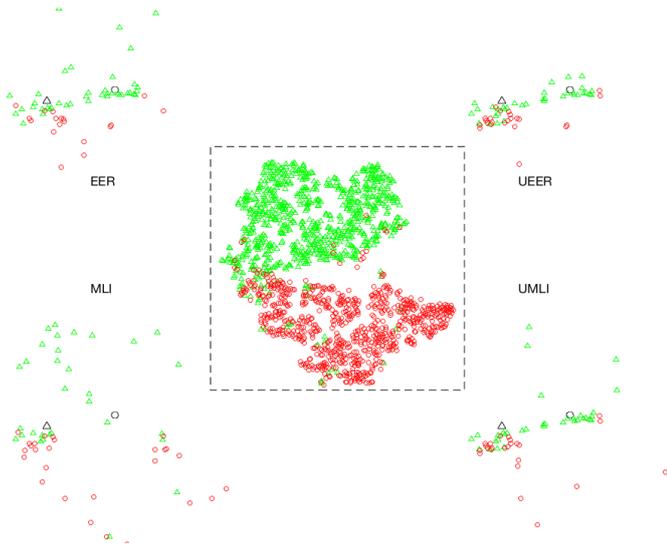


Fig. 1. Illustration of the inherent characteristics of the proposed method. The middle is the distribution of a synthetic binary data set. Four corners represent the performance of four active learning methods, EER, UEER, MLI and UMLI, respectively. One black triangle and circle represent the initial labeled set.

then the proposed method will turn into exactly the uncertainty sampling. More specifically, $\min_{x_i \in \mathcal{U}} \max_{y_i \in \mathcal{C}} P_{\mathcal{L}}(y_i|x_i)$ or $\min_{x_i \in \mathcal{U}} \max_{y_i \in \mathcal{C}} [P_{\mathcal{L}}(y_i|x_i)]^2$ will act as totally same as uncertainty sampling since they will select the instance whose posterior probability comes closest to 0.5 on the binary problem. This shows that our proposed method actually fuses uncertainty sampling with retraining-based models.

B. Two Examples of the Proposed Method

To provide valuable insights on the underlying characteristic of the proposed method, we apply it to two state-of-the-art retraining-based models EER and MLI. We also demonstrate its advantage on a synthetic data set in Figure 1.

Since our method tries to make use of the uncertainty information, the following adapted methods are termed uncertainty retraining-based active learners. It is easy to extend EER to uncertainty-based error reduction by adopting our method in Equation 2 as follows:

$$\arg \min_{x \in \mathcal{U}} \max_{y \in \mathcal{C}} P_{\mathcal{L}}(y|x) \left(- \sum_{x_i \in \mathcal{U}} \sum_{y_i \in \mathcal{C}} P_{\mathcal{L}^+}(y_i|x_i) \log P_{\mathcal{L}^+}(y_i|x_i) \right)$$

This method is called UEER for short. We can also apply our proposed criterion on MLI. The new approach is called UMLI in this paper. Note that the regularization parameter $\frac{1}{2\lambda}$ in Equation 1 is usually quite small, so we ignore it in our adapted criterion:

$$\arg \min_{x \in \mathcal{U}} \max_{y \in \mathcal{C}} P_{\mathcal{L}}(y|x) \sum_{x_i \in \mathcal{L}^+} -\log P_{\mathcal{L}^+}(y_i|x_i)$$

As is shown in Figure 1, we construct a synthetic binary data set and two colours represent different classes. We demonstrate the performance of four retraining-based active learners EER, UEER, MLI and UMLI on four corners, respectively. One black triangle and circle in each corner represent two initial labeled

points. When we compare UEER with EER, it is obvious that UEER selects a number of instances near the decision boundary while EER explores points in a wider range. This is because our method helps UEER make use of the uncertainty information and uncertainty information makes UEER focus on the region which is least certain about. Similar results can also be found between UMLI and MLI. MLI explores over the data space and queries the points around the border while UMLI balances the exploration and the exploitation. UMLI concentrates on the central part (exploitation) and also searches around the edge. Therefore, we can see that our method enhances retraining-based model by balancing the exploration and the exploitation.

IV. EXPERIMENTS

In this section, we investigate the performance of our proposed methods to examine the effectiveness and robustness of our new criterion. The following experiments are limited to binary classification problems. Firstly, we show the experimental setting, then present the extensive experiment results, followed by further discussion and analysis.

A. Experimental setting

We compare the our proposed methods UEER and UMLI against their original version EER and MLI, respectively. Random sampling is also included in this comparison. In all the experiments, we use L_2 -regularized logistic regression included in LIBLINEAR package [25] as default classifier with the same regularization parameter, $\lambda = 100$, for all methods.

The classification accuracy is used as the comparison criterion in our experiment. However, since active learning is a iteratively labeling procedure, we care about the performance during the whole learning process. Thus, it is not reasonable to merely compare the accuracy at some single points. Instead, we generate the learning curve of classification accuracy versus the number of labeled instances. Then, we calculate the area under the learning curve (ALC) as a measure of evaluation.

We test on totally 49 real-world data sets from various real-life applications, including many UCI data sets [26], MNIST handwritten digit dataset [27] and 20 Newsgroups dataset [28]. There are 39 datasets from UCI benchmark datasets, such as breast, vehicle, heart and so on. These datasets are pre-processed according to [29]. For wine data set, we conduct class 2 against class 1 and 3 as binary problem. For glass data set, we also split it into two groups (class 1-3 vs. class 5-7) to build binary case. We randomly sub-sample 1000 instances from mushroom for computing efficiency. We select six pairs of letters from Letter Recognition Data Set [26], *i.e.*, D vs. P, E vs. F, I vs. J, M vs. N, V vs. Y and U vs. V since these pairs look similar to each other and distinguishing them is a little challenging. 3 vs. 5, 5 vs. 8 and 7 vs. 9 are three difficult pairs taken from MNIST data set ¹ and used as the binary classification data set. We randomly sub-sample 1500 instances from the three data sets for computing efficiency. We also test the performance on 20 Newsgroups dataset which is a common benchmark used for text classification ². Following the work of

¹<http://yann.lecun.com/exdb/mnist/>

²<http://qwone.com/~jason/20Newsgroups/>

TABLE I
DATA SETS INFORMATION: IT SHOWS THE NUMBER OF INSTANCES (# INS)
AND THE FEATURE DIMENSIONALITY (# FEA)

Data set (# Ins, # Fea)	Data set (# Ins, # Fea)	Data set (# Ins, # Fea)
ac-inflam (120, 6)	acute (120, 6)	australian (690, 14)
blood (748, 4)	breast (683, 10)	credit (690, 15)
cylinder (512, 35)	diabetes (768, 8)	fertility (100, 9)
german (1000, 24)	glass (214, 9)	haberman (306, 3)
heart (270, 13)	hepatitis (255, 19)	hill (606, 100)
ionosphere (351, 34)	liver (345, 6)	mushrooms (1000, 112)
mammographic (961, 5)	musk1 (476, 166)	ooctris2f (912, 25)
ozone (1000, 72)	parkinsons (195, 22)	pima (768, 8)
planning (182, 12)	sonar (208, 60)	splice (1000, 60)
tictactoe (958, 9)	ve2 (310, 6)	vehicle (435, 18)
wine (178, 13)	wisc (699, 9)	wdbc (569, 31)
d vs p (1608, 16)	e vs f (1543, 16)	i vs j (1502, 16)
m vs n (1575, 16)	v vs y (1577, 16)	u vs v (1550, 16)
3 vs 5 (1500, 784)	5 vs 8 (1500, 784)	7 vs 9 (1500, 784)
base-hockey (1993, 500)	pc-mac (1945, 500)	misc-atheism (1427, 500)
autos (3970, 8014)	motorcycles (3970, 8014)	baseball (3970, 8014)
hockey (3970, 8014)		

TABLE II
PERFORMANCE COMPARISON ON THE AREAS UNDER THE LEARNING
CURVE (ALC)

Dataset	Random	EER	UEER	MLI	UMLI
hill	0.581	0.616	0.599	0.626	0.612
planning	0.586	0.58	0.578	0.614	0.586
cylinder	0.586	0.61	0.597	0.608	0.617
liver	0.627	0.635	0.626	0.615	0.607
splice	0.659	0.679	0.682	0.65	0.666
german	0.664	0.673	0.679	0.691	0.703
ooctris2f	0.679	0.678	0.673	0.686	0.663
musk1	0.682	0.699	0.71	0.702	0.688
fertility	0.693	0.706	0.712	0.727	0.711
haberman	0.711	0.712	0.715	0.694	0.7
sonar	0.713	0.715	0.707	0.708	0.712
pima	0.716	0.706	0.714	0.711	0.722
pcmac	0.717	0.715	0.716	0.747	0.751
diabetes	0.719	0.723	0.723	0.726	0.728
religionatheism	0.72	0.708	0.718	0.691	0.739
hepatitis	0.731	0.753	0.75	0.73	0.738
blood	0.743	0.74	0.718	0.73	0.732
baseball	0.753	0.765	0.872	0.832	0.847
motorcycles	0.763	0.78	0.883	0.854	0.859
autos	0.768	0.768	0.872	0.838	0.835
heart	0.774	0.791	0.795	0.797	0.799
hockey	0.775	0.787	0.901	0.875	0.882
ionosphere	0.779	0.818	0.806	0.674	0.766
credit	0.779	0.793	0.814	0.797	0.809
mammographic	0.78	0.774	0.795	0.766	0.779
basehockey	0.793	0.785	0.801	0.817	0.847
vc2	0.807	0.815	0.812	0.825	0.82
parkinsons	0.811	0.824	0.821	0.83	0.826
australian	0.823	0.832	0.84	0.842	0.83
letterIJ	0.853	0.879	0.853	0.865	0.874
letterVY	0.855	0.878	0.884	0.861	0.867
3vs5	0.856	0.903	0.897	0.859	0.872
vehicle	0.859	0.878	0.888	0.883	0.89
5vs8	0.864	0.907	0.901	0.85	0.87
7vs9	0.876	0.914	0.921	0.841	0.874
ozone	0.882	0.86	0.899	0.892	0.882
tictactoe	0.894	0.912	0.899	0.853	0.88
glass	0.904	0.914	0.914	0.917	0.912
wine	0.906	0.936	0.943	0.94	0.939
letterMN	0.916	0.944	0.941	0.927	0.932
mushrooms	0.931	0.969	0.974	0.971	0.972
letterEF	0.933	0.954	0.961	0.956	0.957
wdbc	0.938	0.953	0.956	0.958	0.957
letterDP	0.939	0.963	0.969	0.967	0.966
letterUV	0.945	0.972	0.979	0.974	0.974
wisc	0.949	0.951	0.954	0.956	0.956
breast	0.95	0.956	0.959	0.962	0.962
ac-inflam	0.955	0.981	0.984	0.98	0.983
acute	0.978	0.971	0.984	0.992	0.992
Mean	0.798	0.812	0.822	0.812	0.818
Average Rank	4.143	3.102	2.388	2.857	2.510
Win/tie/loss counts	-	29/7/13		27/11/11	

[30], we also evaluate three binary tasks from 20 Newsgroups dataset: baseball vs. hockey, pc vs. mac, and religion.misc vs. alt.atheism. And the three pairs represent easy, moderate and difficult classification problems, respectively. We apply PCA to reduce the dimensionality of the above three datasets to 500 for computation efficiency. We also use the pre-processed data autos, motorcycles, baseball, hockey used in [18].

To objectively evaluate the performance, each data set is randomly divided into training and test data set of equal size. At the very beginning of active learning, we assume that only two instances randomly picked up from the training data are labeled, and one of them is from the positive class and the other is from the negative class. We run each active learning algorithm 20 times on each real-world dataset. The average performance of each active learning method is reported in the following section.

B. Results

Table II shows the experimental results on 49 data sets. The datasets in Table II are sorted with respect to the performance of random sampling. We can find that the comparisons contain the datasets which vary from very difficult problems (*e.g.*, hill) to easy tasks (*e.g.*, acute). To clearly demonstrate the advantage of the proposed method, we do pairwise comparison between the original algorithm and its counterpart, *e.g.*, EER vs. UEER and MLI vs. UMLI, respectively. On each data set, a paired t-tests at 95% significance level is used to determine which method has the best performance or provides comparable outcome. These methods are highlighted in bold face. Over all the experiments, average performances are reported in Table II. “Average Rank” shows the average rank of all the methods with regard to their performances on all the experiments. The lower the value of average rank, the better the method. The “win/tie/loss counts” represents times of our proposed methods versus its counterparts over all the 49 datasets.

As is shown in Table II, our proposed methods UEER and UMLI evidently outperform their counterparts EER and MLI, respectively. UEER surpasses EER in terms of average accuracy, and improves its performance from 0.812 to 0.822. UEER also outperforms EER in terms of “average rank”, which demonstrates the effectiveness of our method. Similar results can be found between UMLI and MLI. UMLI is superior to MLI on the overall performance. Moreover, it is interesting to

observe that UEER attains the best overall performance among all the active learning methods. Over all the experimental data sets, the “win/tie/loss” counts of UEER versus EER is 29/7/13, meaning that UEER is the preferred active learner in over half the cases. With regard to UMLI and MLI, the “win/tie/loss” count is 27/11/11, which also shows the clear benefit of our scheme nonetheless. We also notice that even random sampling can surpass all the other methods, *e.g.*, on the blood data set, indicating that, generally, one might not want to use active learners in a blind way.

To investigate the robustness of our method, we also apply the worst-case criterion on EER and the average-case criterion on MLI, respectively. Due to the lack of space, we omit the results on each data set and only report the average performances. The average performance (ALC) of the worst-case on EER is 0.771 while that of the average-case on MLI is 0.710. To our surprise, they definitely show poorer performances in comparison with our method and even perform worse than random sampling. The possible reason may be that: EER computes the error on the unlabeled data and none of the true label are known, the average-case criterion is a safe choice for EER. Since MLI estimates the loss on the enlarged labeled set $\mathcal{L} \cup \{x_i, y_i\}$ and only the true label of x_i is unknown, the worst-case criterion is more appropriate for MLI than the average-case criterion. However, since the proposed method is a trade-off of the two criteria, it can adapt to both settings and show a robust performance for different retraining-based models.

V. CONCLUSIONS

In this paper, we propose a new general method for retraining-based active learning. The proposed method can balance a trade-off of the average-case and worst-case criteria by incorporating uncertainty information (carried by the pre-trained posterior probabilities) within min-max framework. It drives current retraining-based models to pay more attention to the exploitation. We employ the new idea on two state-of-the-art methods to investigate its effectiveness. The synthetic data demonstrates that our method prefers to select the instances which are near the decision boundary in comparison with the original retraining-based approaches. Moreover, extensive experiments on 49 real-world datasets also prove that the proposed method is a promising approach for promoting retraining-based active learners.

REFERENCES

- [1] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “ImageNet: A Large-Scale Hierarchical Image Database,” in *CVPR09*, 2009.
- [2] B. Settles, “Active learning literature survey,” *University of Wisconsin, Madison*, vol. 52, no. 55-66, p. 11, 2010.
- [3] R. Chattopadhyay, Z. Wang, W. Fan, I. Davidson, S. Panchanathan, and J. Ye, “Batch mode active sampling based on marginal probability distribution matching,” *ACM Transactions on Knowledge Discovery from Data (TKDD)*, vol. 7, no. 3, p. 13, 2013.
- [4] H. S. Seung, M. Opper, and H. Sompolinsky, “Query by committee,” in *Proceedings of the fifth annual workshop on Computational learning theory*. ACM, 1992, pp. 287–294.
- [5] D. D. Lewis and J. Catlett, “Heterogeneous uncertainty sampling for supervised learning,” in *Proceedings of the eleventh international conference on machine learning*, 1994, pp. 148–156.
- [6] S. Tong and D. Koller, “Support vector machine active learning with applications to text classification,” *The Journal of Machine Learning Research*, vol. 2, pp. 45–66, 2002.
- [7] C. Campbell, N. Cristianini, A. Smola *et al.*, “Query learning with large margin classifiers,” in *ICML*, 2000, pp. 111–118.
- [8] N. Roy and A. McCallum, “Toward optimal active learning through monte carlo estimation of error reduction,” *ICML, Williamstown*, 2001.
- [9] Y. Guo and R. Greiner, “Optimistic active-learning using mutual information,” in *IJCAI*, vol. 7, 2007, pp. 823–829.
- [10] A. Holub, P. Perona, and M. C. Burl, “Entropy-based active learning for object recognition,” in *Computer Vision and Pattern Recognition Workshops, 2008. CVPRW’08. IEEE Computer Society Conference on*. IEEE, 2008, pp. 1–8.
- [11] Y. Guo and D. Schuurmans, “Discriminative batch mode active learning,” in *Advances in neural information processing systems*, 2008, pp. 593–600.
- [12] B. Settles, M. Craven, and S. Ray, “Multiple-instance active learning,” in *Advances in neural information processing systems*, 2008, pp. 1289–1296.
- [13] A. Freytag, E. Rodner, and J. Denzler, “Selecting influential examples: Active learning with expected model output changes,” in *Computer Vision—ECCV 2014*. Springer, 2014, pp. 562–577.
- [14] W. Cai, Y. Zhang, S. Zhou, W. Wang, C. Ding, and X. Gu, “Active learning for support vector machines with maximum model change,” in *Machine Learning and Knowledge Discovery in Databases*. Springer, 2014, pp. 211–226.
- [15] C. Kading, A. Freytag, E. Rodner, P. Bodesheim, and J. Denzler, “Active learning and discovery of object categories in the presence of unnameable instances,” in *Computer Vision and Pattern Recognition (CVPR), 2015 IEEE Conference on*. IEEE, 2015, pp. 4343–4352.
- [16] S. C. Hoi, R. Jin, J. Zhu, and M. R. Lyu, “Batch mode active learning and its application to medical image classification,” in *Proceedings of the 23rd international conference on Machine learning*. ACM, 2006, pp. 417–424.
- [17] T. Zhang and F. Oles, “The value of unlabeled data for classification problems,” in *Proceedings of the Seventeenth International Conference on Machine Learning (Langley, P., ed.)*. Citeseer, 2000, pp. 1191–1198.
- [18] K. Yu, J. Bi, and V. Tresp, “Active learning via transductive experimental design,” in *Proceedings of the 23rd international conference on Machine learning*. ACM, 2006, pp. 1081–1088.
- [19] A. I. Schein and L. H. Ungar, “Active learning for logistic regression: an evaluation,” *Machine Learning*, vol. 68, no. 3, pp. 235–265, 2007.
- [20] S. C. Hoi, R. Jin, J. Zhu, and M. R. Lyu, “Semi-supervised svm batch mode active learning for image retrieval,” in *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*. IEEE, 2008, pp. 1–7.
- [21] S.-J. Huang, R. Jin, and Z.-H. Zhou, “Active learning by querying informative and representative examples,” in *Advances in neural information processing systems*, 2010, pp. 892–900.
- [22] Y. Yang and M. Loog, “A benchmark and comparison of active learning methods for logistic regression,” *arXiv preprint*, 2016.
- [23] J. Platt *et al.*, “Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods,” 1999.
- [24] R. P. Duin and D. M. Tax, “Classifier conditional posterior probabilities,” in *Advances in pattern recognition*. Springer, 1998, pp. 611–619.
- [25] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin, “Lib-linear: A library for large linear classification,” *The Journal of Machine Learning Research*, vol. 9, pp. 1871–1874, 2008.
- [26] M. Lichman, “UCI machine learning repository,” 2013. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [27] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [28] K. Lang, “Newsweeder: Learning to filter netnews,” in *Proceedings of the Twelfth International Conference on Machine Learning*, 1995, pp. 331–339.
- [29] M. Fernández-Delgado, E. Cernadas, S. Barro, and D. Amorim, “Do we need hundreds of classifiers to solve real world classification problems?” *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 3133–3181, 2014.
- [30] X. Zhu, J. Lafferty, and Z. Ghahramani, “Combining active learning and semi-supervised learning using gaussian fields and harmonic functions,” in *ICML 2003 workshop on the continuum from labeled to unlabeled data in machine learning and data mining*, 2003, pp. 58–65.