

---

# Active Learning Using Pre-clustering

---

Hieu T. Nguyen  
Arnold Smeulders

TAT@SCIENCE.UVA.NL  
SMEULDERS@SCIENCE.UVA.NL

Intelligent Sensory Information Systems, University of Amsterdam, Faculty of Science, Kruislaan 403, NL-1098 SJ, Amsterdam, The Netherlands

## Abstract

The paper is concerned with two-class active learning. While the common approach for collecting data in active learning is to select samples close to the classification boundary, better performance can be achieved by taking into account the prior data distribution. The main contribution of the paper is a formal framework that incorporates clustering into active learning. The algorithm first constructs a classifier on the set of the cluster representatives, and then propagates the classification decision to the other samples via a local noise model. The proposed model allows to select the most representative samples as well as to avoid repeatedly labeling samples in the same cluster. During the active learning process, the clustering is adjusted using the coarse-to-fine strategy in order to balance between the advantage of large clusters and the accuracy of the data representation. The results of experiments in image databases show a better performance of our algorithm compared to the current methods.

## 1. Introduction

In recent years, research interest has been attracted to semi-supervised learning or learning in the condition that only a small initial amount of data is labeled while the majority of the data remain unlabeled. While many methods focus to improve the supervised learning by using the information from unlabeled data (Seeger, 2001), another important topic is a good strategy in selecting the data to label, considering that labeling data is a time-consuming job. The topic is known as *active learning* (Lewis & Gale, 1994).

Consider the problem of learning a binary classifier on a partially labeled database  $\mathcal{D} = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$ . Let

---

Appearing in *Proceedings of the 21<sup>st</sup> International Conference on Machine Learning*, Banff, Canada, 2004. Copyright 2004 by the authors.

$\mathcal{D}_\ell$  be the labeled set in which every sample is given a label  $y \in \{1, -1\}$ , and  $\mathcal{D}_u = \mathcal{D} \setminus \mathcal{D}_\ell$ . The active learning system comprises two parts: a learning engine and a selection engine. At every iteration the learning engine uses a supervised learning algorithm to train a classifier on  $\mathcal{D}_\ell$ . The selection engine then selects a sample from  $\mathcal{D}_u$  and requests a human expert to label the sample before passing it to the learning engine. The major goal is to achieve a good classifier as best as possible within a reasonable number of calls for labeling by human help.

Current methods on active learning can be characterized by their base learning algorithms which include probabilistic naive Bayes (Nigam et al., 2000; Roy & McCallum, 2001), combination of naive Bayes and logistic regression (Lewis & Gale, 1994), and the Support Vector Machine (SVM) (Campbell et al., 2000; Tong & Koller, 2001; Schohn & Cohn, 2000). The naive Bayes classifier suffers from two problems. First, the classifier assumes the independence between the component features of  $x$ . This assumption is often violated. The second problem is that naive Bayes is a generative model for which training relies on the estimation of the likelihood  $p(x|y)$ . This estimation is inaccurate in the case of active learning since the training data are not randomly collected. The paper focuses on discriminative models including logistic regression and SVM. These models aim to estimate the posterior probability  $p(y|x)$ . They are less sensitive to the way the training data is collected, and hence, are more suitable for active learning. A more theoretical consideration is given in (Zhang & Oles, 2000).

It is crucial to choose the most “valuable” training samples. Many methods choose the most uncertain samples which are closest to the current classification boundary. We name this approach the closest-to-boundary criterion. This simple and intuitive criterion performs well in some applications (Lewis & Gale, 1994; Tong & Chang, 2001; Schohn & Cohn, 2000; Campbell et al., 2000). Some other criteria have been proposed specifically for SVM. In (Campbell et al., 2000), it is proposed to select the sample that yields the largest decrease of the margin between the two classes. The method of (Tong & Koller, 2001) selects the

sample that halves the permitted region of the SVM parameters in the parameter space. Both (Campbell et al., 2000) and (Tong & Koller, 2001) need to predict the values of the SVM parameters for every possible case where a candidate sample might be added to the training set. Since it is hard to do this efficiently, the references finally resort to the closest-to-boundary criterion.

The closest-to-boundary methods ignore the prior data distribution which can be useful for active learning. In (Cohn et al., 1996), it is suggested to select samples that minimize the expected future classification error:

$$\int_{\mathbf{x}} E[(\hat{y}(\mathbf{x}) - y(\mathbf{x}))^2 | \mathbf{x}] p(\mathbf{x}) d\mathbf{x} \quad (1)$$

where  $y(\mathbf{x})$  is the true label of  $\mathbf{x}$  and  $\hat{y}(\mathbf{x})$  is the classifier output.  $E[\cdot | \mathbf{x}]$  denotes the expectation over  $p(y | \mathbf{x})$ . Due to the complexity of the integral, the direct implementation of eq. (1) is usually difficult. However, it shows that the data uncertainty should be weighted with the prior density  $p(\mathbf{x})$ . If  $p(\mathbf{x})$  is uniform or unknown, the expectation under the integral is the contribution by a sample into the classification error. The expectation can then be used to measure the value of the sample in the condition that the computation of the integral is complex. Under the assumption that the current classification boundary is good, it is easy to show that the error expectation is maximal for the samples lying on the classification boundary, see section 3.3. When  $p(\mathbf{x})$  is known and non-uniform, the information about the distribution can be used to select better data. In this paper  $p(\mathbf{x})$  is obtained via clustering which can be done offline without the interaction with human. The clustering information is then useful for active learning in two ways. First, the representative samples located in center of clusters are more important than the other, and should be selected first in labeling. Secondly, samples in the same cluster are likely to have the same label, (Seeger, 2001; Chapelle et al., 2002). This assumption should be used to accelerate active learning by reducing the number of labeling samples from the same cluster.

The idea to combine clustering and active learning has appeared in previous work. In (McCallum & Nigam, 1998), a naive Bayes classifier is trained over both labeled and unlabeled data using an EM algorithm. Under the condition that the overwhelming majority of the data is unlabeled, that training algorithm amounts to clustering the data set, and the role of the labeled data is for initialization only. Clustering information also contributes to the selection where an uncertainty measure is weighted with the density of the sample. The referenced approach does not match, however, the objective of this paper to combine clustering with a discriminative model. Several other active learning schemes also weigh the uncertainty with the data density (Zhang & Chen, 2002; Tang et al., 2002). Some

methods put more emphasis on the sample representativeness by selecting cluster centers from a set of most interesting samples. In the representative sampling by (Xu et al., 2003), the algorithm uses the k-means algorithm to cluster the samples lying within the margin of a SVM classifier trained on the current labeled set. The samples at cluster centers are then selected for human labeling. The method of (Shen & Zhai, 2003) has a similar idea, but applies the k-medoid algorithm for the top relevant samples. In general, heuristic methods have been proposed to balance between the uncertainty and the representativeness of the selected sample. They encourage the selection of cluster centers. However, no measure has been taken to avoid repeatedly labeling samples in same cluster. In addition, there are important questions that remain open, namely, how to adapt the classification model for a training set that contains only cluster centers? and, how to classify samples that are disputed by several clusters? This paper presents a solution for these issues using a mathematical model that explicitly takes clustering into account.

The organization of the paper is as follows. Section 2 describes the incorporation of the clustering information into the data model, and provides the theoretical framework for the data classification. Section 3 presents our active learning algorithm. Section 4 shows the results of the algorithm for the classification of images in test databases.

## 2. Probabilistic framework

### 2.1. Data model

In the standard classification, data generation is described by the joint distribution  $p(\mathbf{x}, y)$  of the data  $\mathbf{x}$  and the class label  $y \in \{-1, +1\}$ . The clustering information is explicitly incorporated by introducing the hidden cluster label  $k \in \{1, \dots, K\}$ , where  $K$  is the number of clusters in the data.  $k$  indicates that the sample belongs to the  $k$ -th cluster. Assume that all information about the class label  $y$  is already encoded in the cluster label  $k$ . This implies that once  $k$  is known,  $y$  and  $\mathbf{x}$  are independent. The joint distribution is written as:

$$p(\mathbf{x}, y, k) = p(y | \mathbf{x}, k) p(\mathbf{x} | k) p(k) = p(y | k) p(\mathbf{x} | k) p(k) \quad (2)$$

The simple Bayesian belief net representing the model is

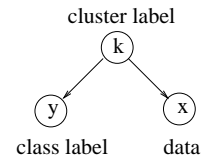


Figure 1. The Bayesian net for the data model.

depicted in Figure 1.

Before giving the specific form for the three distributions in eq. (2) we remark that a similar scheme has been proposed for **the passive semi-supervised learning** (Miller & Uyar, 1996; Seeger, 2001). The conceptual difference, however, between their approach and ours is in the definition of  $p(y|k)$ . In the references,  $p(y|k)$  is defined within individual clusters. As a consequence, the estimation of the parameters of  $p(y|k)$  can be unreliable due to insufficient labeled data in a cluster. In our model,  $p(y|k)$  is defined for all clusters with the same parameters.

We use logistic regression for  $p(y|k)$ :

$$p(y|k) = \frac{1}{1 + \exp\{-y(c_k \cdot a + b)\}} \quad (3)$$

Here,  $c_k$  is a representative of the  $k$ -th cluster which is determined via clustering.  $a \in \mathbb{R}^d$  and  $b \in \mathbb{R}$  are the logistic regression parameters. In essence,  $p(y|k)$  is the label model for a representative subset of the database.

In the ideal case where data is well clustered, once all the parameters of  $p(y|k)$  are determined, one could use this probability to determine the label of the cluster representatives, and then assign the same label to the remaining samples in the cluster. In practice, however, clustering can be inaccurate and we will have problems with classification of samples at border between the clusters. To achieve better classification for those samples, we use **a soft cluster membership** which allows a sample to be connected to more than one clusters (representatives) with a probability. The noise distribution  $p(x|k)$  is then used to propagate information of label  $y$  from the representatives into the remaining majority of the data, see Figure 2. We use the isotropic Gaussian model:

$$p(x|k) = (2\pi)^{-d/2} \sigma^{-d} \exp\left\{-\frac{1}{2\sigma^2} \|x - c_k\|^2\right\} \quad (4)$$

where  $\sigma^2$  is the variance assumed to be the same for all clusters.

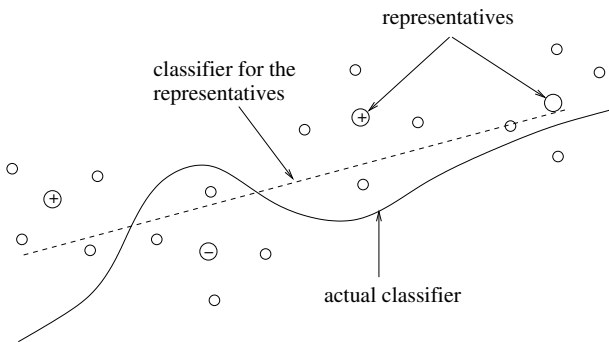


Figure 2. The classification model

Let  $p(k) = \alpha_k$ . Then,  $p(x)$  is a mixture of  $K$  Gaussians with the weights  $\alpha_k$ .

In the presented method, the parameters  $c_k$ ,  $\alpha_k$ ,  $a$  and  $b$  are estimated from the data. The scale parameter  $\sigma^2$  is given initially. It can be changed during active learning when a different clustering setting is needed.

## 2.2. Data classification

Given the above model, one calculates  $p(y|x)$ , the posterior probability of label of a sample as follows:

$$p(y|x) = \sum_{k=1}^K p(y, k|x) = \sum_{k=1}^K p(y|k)p(k|x) \quad (5)$$

where  $p(k|x) = p(x|k)p(k)/p(x)$ .

Data are then classified using the Bayes decision rule:

$$\hat{y}(x) = \begin{cases} 1 & \text{if } p(y=1|x; \hat{a}, \hat{b}) > p(y=-1|x; \hat{a}, \hat{b}) \\ -1 & \text{otherwise} \end{cases} \quad (6)$$

where  $\hat{a}, \hat{b}$  denote the current estimates of the parameters.

Observe from eq. (5) that the classification decision is a weighted combination of the classification decision for the representatives. Well clustered samples will be assigned the same label as the nearest representative. Samples disputed by several clusters, on the other hand, will be assigned the label of the cluster, which has the highest confidence. Note that the weights  $p(k|x)$  are fixed unless the data are re-clustered whereas  $p(y|k)$  is updated upon the arrival of new training data.

## 3. Description of algorithm

The parameters of the model proposed in section 2.1 are estimated via likelihood maximization. The data likelihood comprises two parts: the likelihood of the labeled data and the likelihood of the unlabeled data:

$$L = \sum_{i \in \mathcal{I}_\ell} \ln p(x_i, y_i) + \sum_{i \in \mathcal{I}_u} \ln p(x_i) \quad (7)$$

where  $\mathcal{I}_\ell$  and  $\mathcal{I}_u$  denote the set of indices of labeled and unlabeled samples respectively. Expanding  $\ln p(x_i, y_i)$  as the sum of  $\ln p(x_i)$  and  $\ln p(y_i|x_i)$ , the likelihood (7) can be written with explicit dependence on the parameters as follows:

$$\begin{aligned} L(c_1, \dots, c_K, \alpha_1, \dots, \alpha_K, a, b) = & \sum_{i \in \mathcal{I}_\ell \cup \mathcal{I}_u} \ln p(x_i; c_1, \dots, c_K, \alpha_1, \dots, \alpha_K) + \\ & \sum_{i \in \mathcal{I}_\ell} \ln p(y_i|x_i; c_1, \dots, c_K, a, b) \end{aligned} \quad (8)$$

As the amount of the unlabeled data is overwhelming over the labeled data, the parameters  $c_1, \dots, c_K$  and

$\alpha_1, \dots, \alpha_K$  are determined mainly by maximizing the first term in eq. (8). The maximization of each term can therefore be done separately. The clustering algorithm maximizes the likelihood of the data samples to obtain the cluster representatives and the cluster prior. The maximization of the label likelihood follows to estimate the parameters  $a$  and  $b$ .

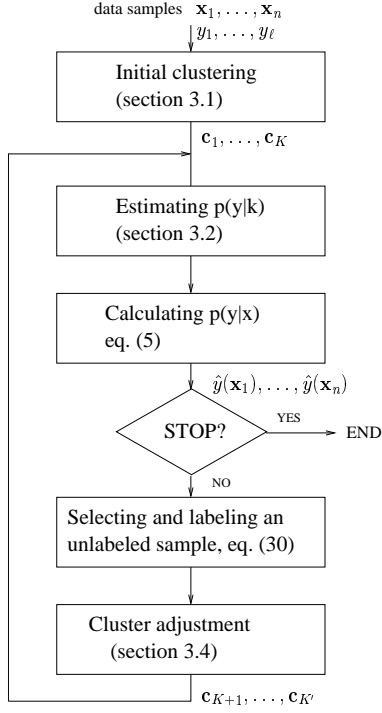


Figure 3. The proposed active learning algorithm.

The block scheme of the algorithm is illustrated in Figure 3.

### 3.1. Initial clustering

In the presented algorithm, the goal of clustering is data representation rather than data classification. We therefore use the  $K$ -medoid algorithm of (Kaufman & Rousseeuw, 1990). The algorithm finds  $K$  representatives  $c_1, \dots, c_K$  of the data set  $x_1, \dots, x_n$  so as to minimize the sum of the distance from the data samples to the nearest representative. See (Struyf et al., 1997) for the detailed implementation.

The  $K$ -medoid algorithm is computationally expensive when either  $n$  or  $K$  is large. In practical applications both numbers are very large indeed. The following simplifications are employed to reduce computations. First, the data set is split into smaller subsets. The  $K$ -medoid algorithm is then applied to cluster every subset into a limited number  $K_0$  of clusters. Clustering is continued by subsequently breaking the cluster with the largest radius  $r_k$  into

two smaller ones with:

$$r_k = \max_{i \in \mathcal{I}_k} \|x_i - c_k\| \quad (9)$$

where  $\mathcal{I}_k$  denotes the set of indices of the samples in the  $k$ -th cluster. The process of cluster fission is completed when:

$$\max_k r_k < \kappa \sigma \quad (10)$$

where  $\kappa$  is a predefined constant. We have used  $\kappa = 8$ . Thus, the cluster size and the final number of clusters  $K$  is controlled by the scale parameter  $\sigma$ .

Once the cluster representatives  $c_1, \dots, c_K$  have been determined, the cluster prior  $\alpha_k$  is obtained by iterating the following two equations until stability:

$$p(k|x_i) = \frac{\alpha_k \exp\{-\frac{1}{2\sigma^2}\|x_i - c_k\|^2\}}{\sum_{k'=1}^K \alpha_{k'} \exp\{-\frac{1}{2\sigma^2}\|x_i - c_{k'}\|^2\}} \quad (11)$$

$$\alpha_k = \frac{1}{n} \sum_{i=1}^n p(k|x_i) \quad (12)$$

### 3.2. The estimation of the class label model

This section presents the estimation of the distribution  $p(y|k)$  based on the maximization of the second likelihood in eq. (8). Fixing the cluster representatives  $c_k$ , the likelihood depends only on the parameters  $a$  and  $b$ :

$$\max_{a,b} \sum_{i \in \mathcal{I}_\ell} \ln p(y_i|x_i; a, b) \quad (13)$$

From eq.(5),  $p(y_i|x_i; a, b)$  can be written as a mixture of  $K$  logistic distributions  $p(y_i|k; a, b)$  with the weights  $p(k|x_i)$ .

In case the dimensionality  $d$  is higher than the number of the labeled samples, the optimization in (13) is numerically unstable. The conventional approach to overcome the problem is to add a regularization term  $\frac{1}{2}\lambda\|a\|^2$  where  $\lambda$  is a predefined parameter. This leads to the minimization of the following objective function:

$$\mathcal{L}(a, b) = \frac{1}{2}\lambda\|a\|^2 - \sum_{i \in \mathcal{I}_\ell} \ln \left\{ \sum_{k=1}^K p(k|x_i) p(y_i|k; a, b) \right\} \quad (14)$$

Remark that eq. (14) is the extension of the regularized logistic regression (Zhang & Oles, 2001; Zhu & Hastie, 2001) for the mixture of the logistic distributions.

The minimization of  $\mathcal{L}$  is implemented using Newton's algorithm which guarantees to find a local minimum. Furthermore, since  $\mathcal{L}$  is convex, it has only one local minimum which is also the global minimum.

Starting with an initial guess  $a_0$  and  $b_0$ , the parameters  $a$  and  $b$  are updated iteratively. At each iteration, the param-

eter increment in the steepest direction is:

$$\begin{bmatrix} \Delta \mathbf{a} \\ \Delta b \end{bmatrix} = -\mathbf{H}^{-1} \nabla \mathcal{L} \quad (15)$$

where  $\nabla \mathcal{L}$  is the Jacobian of  $\mathcal{L}$ , and  $\mathbf{H}$  is a positive definite approximation of the hessian matrix of  $\mathcal{L}$ . Using eq.(3), it can be shown that:

$$\nabla \mathcal{L} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{a}} \\ \frac{\partial \mathcal{L}}{\partial b} \end{bmatrix} = \lambda \begin{bmatrix} \mathbf{a} \\ 0 \end{bmatrix} - \sum_{k=1}^K \gamma_k \begin{bmatrix} \mathbf{c}_k \\ 1 \end{bmatrix} \quad (16)$$

where

$$\gamma_k = \sum_{i \in \mathcal{I}_\ell} y_i p(y_i | k; \mathbf{a}, b) [1 - p(y_i | k; \mathbf{a}, b)] \frac{p(k | \mathbf{x}_i)}{p(y_i | \mathbf{x}_i; \mathbf{a}, b)} \quad (17)$$

For the hessian matrix, the following approximation can be used:

$$\mathbf{H} = \lambda \begin{bmatrix} \mathbf{I} & 0 \\ 0 & 0 \end{bmatrix} + \sum_{k=1}^K \tau_k \begin{bmatrix} \mathbf{c}_k \mathbf{c}_k^T & \mathbf{c}_k \\ \mathbf{c}_k^T & 1 \end{bmatrix} \quad (18)$$

where  $\mathbf{I}$  is the identity matrix and:

$$\tau_k = \sum_{i \in \mathcal{I}_\ell} p(y_i | k; \mathbf{a}, b)^2 [1 - p(y_i | k; \mathbf{a}, b)] \frac{p(k | \mathbf{x}_i)}{p(y_i | \mathbf{x}_i; \mathbf{a}, b)} \quad (19)$$

To get more insight into eq.(15), let:

$$\mathbf{z} = \sum_{k=1}^K \tau_k \mathbf{c}_k / \sum_{k=1}^K \tau_k \quad (20)$$

$$\mathbf{D} = [\sqrt{\tau_1} \mathbf{c}_1, \dots, \sqrt{\tau_K} \mathbf{c}_K] \quad (21)$$

$$\mathbf{G} = \mathbf{D} \mathbf{D}^T + \lambda \mathbf{I} \quad (22)$$

Here,  $\mathbf{D}$  is the  $d \times K$  matrix whose columns are the vectors  $\sqrt{\tau_k} \mathbf{c}_k$ . One can show that:

$$\Delta \mathbf{a} = -\mathbf{G}^{-1} \left[ \frac{\partial \mathcal{L}}{\partial \mathbf{a}} - \mathbf{z} \frac{\partial \mathcal{L}}{\partial b} \right] \quad (23)$$

$$\Delta b = \mathbf{z}^T \mathbf{G}^{-1} \left[ \frac{\partial \mathcal{L}}{\partial \mathbf{a}} - \mathbf{z} \frac{\partial \mathcal{L}}{\partial b} \right] - \frac{1}{\sum_{k=1}^K \tau_k} \frac{\partial \mathcal{L}}{\partial b} \quad (24)$$

If  $d$  is high, it is efficient to invert  $\mathbf{G}$  using the Woodbury formula:

$$\mathbf{G}^{-1} = \frac{1}{\lambda} [\mathbf{I} - \mathbf{D}(\lambda \mathbf{I} + \mathbf{D}^T \mathbf{D})^{-1} \mathbf{D}^T] \quad (25)$$

If all  $\tau_k$  are non-zero, the size of  $\mathbf{D}^T \mathbf{D}$  is  $K \times K$ . Since  $K$  is large, the inverting of  $\mathbf{G}$  as in eq. (25) would be computationally expensive still. However, remark that for a sample  $\mathbf{x}_i$  there are only few values of  $k$  such that  $p(k | \mathbf{x}_i)$  is different from zero, especially if  $\mathbf{x}_i$  is a cluster representative. In the latter case, the sample typically belongs to one cluster

only. As will be seen in the next subsection, the presented algorithm tends to select the training data from the cluster representatives. The number of non-zero  $\tau_k$  is then small, approximately the same as the number of labeled samples. The computation of  $\mathbf{G}^{-1}$  in eq. (25) can then be done efficiently by suppressing the columns in  $\mathbf{D}$  which correspond to  $\tau_k$  that equal to zero.

### 3.3. Criterion for data selection

The selection criterion gives priority to two types of samples: samples close to the classification boundary and samples which are cluster representatives. Furthermore, within the set of cluster representatives, one should start with the highest density clusters first.

We have noted that the computation of the future classification error in eq.(1) is complicated. So, instead of choosing the sample that produces the smallest future error, we select the sample that has the largest contribution to the current error. Although such approach does not guarantee the smallest future error, there is a good chance for a large decrease of the error. The selection criterion is:

$$s = \arg \max_{i \in \mathcal{I}_u} E[(\hat{y}_i - y_i)^2 | \mathbf{x}_i] p(\mathbf{x}_i) \quad (26)$$

where  $s$  denotes the index of the selected sample.

The error expectation for an unlabeled  $\mathbf{x}_i$  is calculated over the distribution  $p(y_i | \mathbf{x}_i)$ :

$$\begin{aligned} E[(\hat{y}_i - y_i)^2 | \mathbf{x}_i] &= p(y_i = 1 | \mathbf{x}_i) (\hat{y}_i - 1)^2 + p(y_i = -1 | \mathbf{x}_i) (\hat{y}_i + 1)^2 \\ &= 2 - 2\hat{y}_i (p(y_i = 1 | \mathbf{x}_i) - p(y_i = -1 | \mathbf{x}_i)) \end{aligned} \quad (27)$$

It should be noted that the probability  $p(y_i = 1 | \mathbf{x}_i)$  is unknown and needs to be approximated. An obvious choice is to use the current estimation  $p(y_i = 1 | \mathbf{x}_i; \hat{\mathbf{a}}, \hat{b})$ , assuming  $\hat{\mathbf{a}}, \hat{b}$  are good enough. Letting

$$f(\mathbf{x}_i) = p(y_i = 1 | \mathbf{x}_i; \hat{\mathbf{a}}, \hat{b}) - p(y_i = -1 | \mathbf{x}_i; \hat{\mathbf{a}}, \hat{b}) \quad (28)$$

it follows from eq. (5) that:

$$E[(\hat{y}_i - y_i)^2 | \mathbf{x}_i] = 2(1 - |f(\mathbf{x}_i)|) \quad (29)$$

Observe that if  $\mathbf{x}_i$  lies on the current classification boundary, the quantity  $|f(\mathbf{x}_i)|$  is minimal, and hence the expected error is maximal.

Eq. (26) becomes:

$$s = \arg \max_{i \in \mathcal{I}_u} (1 - |f(\mathbf{x}_i)|) p(\mathbf{x}_i) \quad (30)$$

where

$$p(\mathbf{x}_i) = \sum_{k=1}^K \alpha_k \exp\left\{-\frac{1}{2\sigma^2} \|\mathbf{x}_i - \mathbf{c}_k\|^2\right\} \quad (31)$$



Figure 4. Example view of images in the first database.

The resulting criterion indeed satisfies the demands put in the beginning of the subsection. The term  $(1 - |f(x_i)|)$  gives priority to the samples at the boundary. Meantime,  $p(x_i)$  gives priority to the representatives of dense clusters.

### 3.4. Coarse-to-fine adjustment of clustering

The labeling of high density clusters promises a substantial move of the classification boundary. It is therefore advantageous to group the data into large clusters in the initial clustering. This is achieved by setting a high value for the initial scale parameter  $\sigma_0$ . When the classification boundary reaches the border between the global clusters, a finer clustering with smaller cluster size is better to obtain a more accurate classification boundary. The maximum in eq. (30) can be used as the indication for the need to adjust the clustering. If this quantity drops below a threshold  $t_e$ :

$$\max_{i \in \mathcal{I}_u} (1 - |f(x_i)|) p(x_i) < t_e \quad (32)$$

the scale parameter is decreased:

$$\sigma_{new} = \eta \sigma_{old} \quad (33)$$

where  $0 < \eta < 1$ . The data set is then re-clustered. The parameters  $t_e$  and  $\eta$  are predefined. We have used  $t_e = 0.2$  and  $\eta = 0.9$ . Note that clustering the data set with different scales can be done offline. Furthermore, change of the scale takes place not in every iteration, but only few times during the learning process.

## 4. Experiments

We have performed two experiments to test the performance of the proposed algorithm. In the first experiment, the algorithm is applied to find human face images in a database containing 2500 images of size  $20 \times 20$ . See (Pham et al., 2002) for details on how the images were created. Example views of some images are shown in Figure 4. In the second experiment, a test database was made of images of handwritten digits taken from the MNIST database

(<http://yann.lecun.com/exdb/mnist/>). The size of images is  $28 \times 28$ . The objective is to separate the images of a given digit against the other nine.

In the experiments the following setting was used. The images are considered as the vectors composed of the pixel grey values which range from 0 to 255. The initial training set contains equal numbers of object and non-object images. The initial size of this set was  $0.4\%n$  and was increased to  $3.5\%n$  during active learning, where  $n$  is the number of samples in the database. For clustering, the databases were split into subsets per 1250 samples. The  $K$ -medoid algorithm was applied for each subset with the initial number of clusters  $K_0 = 20$ . The initial value of  $\sigma$  was  $\sigma_0 = 2.2d$  where  $d$  is the number of pixels in one image. For the estimation of the class label distribution, we use the regularization coefficient  $\lambda = 80$ .

Every time a new training sample is added, the classifier is re-trained and tested on the rest of the database. The classification error is calculated as the sum of the missed positives and false alarms relative to  $n$ . The performance evaluation is based on the decrease of the classification error as the function of the amount of training samples.

For comparison, we have also implemented three other active learning algorithms. They use the standard linear SVM for classification. The first algorithm selects training data randomly. In the second algorithm, data are selected according to the closest-to-boundary criterion. The third algorithm uses the representative sampling of (Xu et al., 2003) which selects the medoid centers in the SVM margin. As we select one sample per iteration, this leads to selection of the most representative sample among the samples in the margin.

Figure 5 shows the result of the first experiment for different proportions between the numbers of face and non-face images in the database. Figure 6 shows the result of the second experiment for the different sizes of the database. Both figures show the average of the classification error obtained by repeating the experiments with three different initial training sets that are picked up randomly. The results of Figure 6 are also an average over the ten digits.

The proposed algorithm outperforms all three other algorithms. The most significant improvement is observed in Figure 5a with equal numbers for object and non-object samples in the database. The improvement decreases when the amount of the object samples is small relative to the non-object samples, see Figure 5c. In this case, since there are no clusters of the object class, the proposed algorithm is not advantageous over the closest-to-boundary algorithm in finding object samples. Nevertheless, the proposed algorithm remains better as it still benefits from the clustering of non-object samples. Representative sampling turns out

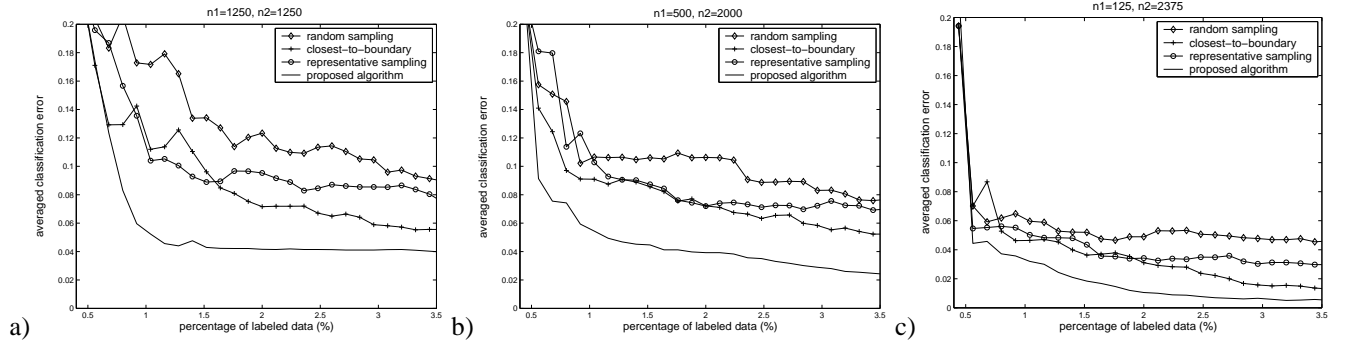


Figure 5. The results for the classification of face images.  $n_1$  and  $n_2$  are the number of face and non-face images in the databases respectively.

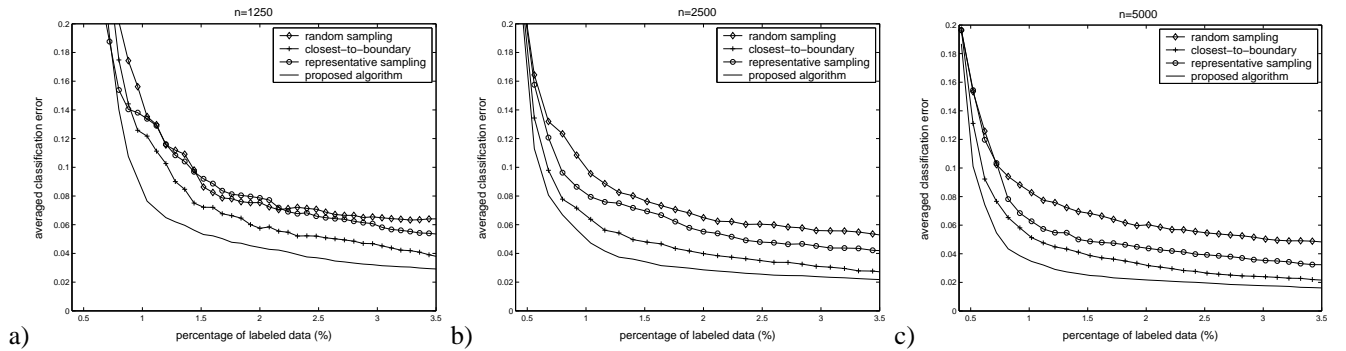


Figure 6. The classification results of handwritten digits from the MNIST database.  $n$  is the database size.

to perform better only than random sampling. A possible reason could be the undervaluation of the uncertainty and the lack of a proper classification model.

## 5. Conclusion

The paper has proposed a formal model for incorporation of clustering into active learning. The model allows to select most representative training examples as well as to avoid repeatedly labeling samples in same cluster, leading to better performance than the current methods. To take the advantage of the similarity between the class label of data in the same cluster, the method first constructs a classifier over the population of the cluster representatives. We use regularized logistic regression which is a discriminative model with state-of-the-art performance and which is naturally fitted into a probabilistic framework. The gaussian noise model is then used to infer the class label for non-representative samples. New training data are selected from the samples having the maximal contribution to the current expected error. In addition to closeness to the classification boundary, the selection criterion gives priority also to the representatives of the dense clusters, making the training set statistically stable.

The method was restricted to linear logistic regression as

the main purpose of the paper is to show the advantage of using clustering information. We have succeeded in that goals for the given datasets.

## References

- Campbell, C., Cristianini, N., & Smola, A. (2000). Query learning with large margin classifiers. *Proc. 17th International Conf. on Machine Learning* (pp. 111–118). Morgan Kaufmann, CA.
- Chapelle, O., Weston, J., & Scholkopf, B. (2002). Cluster kernels for semi-supervised learning. *Advances in Neural Information Processing Systems*.
- Cohn, D. A., Ghahramani, Z., & Jordan, M. I. (1996). Active learning with statistical models. *Journal of Artificial Intelligence research*, 4, 129–145.
- Kaufman, L., & Rousseeuw, P. (1990). *Finding groups in data: An introduction to cluster analysis*. John Wiley & Sons.
- Lewis, D. D., & Gale, W. A. (1994). A sequential algorithm for training text classifiers. *Proceedings of SIGIR-94, 17th ACM International Conference on Research*

- and Development in Information Retrieval (pp. 3–12). Springer Verlag.
- McCallum, A. K., & Nigam, K. (1998). Employing EM in pool-based active learning for text classification. *Proc. 15th International Conf. on Machine Learning* (pp. 350–358). Morgan Kaufmann, CA.
- Miller, D., & Uyar, H. (1996). A mixture of experts classifier with learning based on both labelled and unlabelled data. *Advances in Neural Information Processing Systems 9* (pp. 571–577).
- Nigam, K., McCallum, A., Thrun, S., & Mitchell, T. (2000). Text classification from labeled and unlabeled documents using EM. *Machine Learning*, 39, 103–134.
- Pham, T., Worring, M., & Smeulders, A. (2002). Face detection by aggregated bayesian network classifiers. *Pattern Recogn. Letters*, 23, 451–461.
- Roy, N., & McCallum, A. (2001). Toward optimal active learning through sampling estimation of error reduction. *Proc. 18th International Conf. on Machine Learning* (pp. 441–448). Morgan Kaufmann, CA.
- Schohn, G., & Cohn, D. (2000). Less is more: Active learning with support vector machines. *Proc. 17th International Conf. on Machine Learning* (pp. 839–846). Morgan Kaufmann, CA.
- Seeger, M. (2001). *Learning with labeled and unlabeled data* (Technical Report). Edinburgh University.
- Shen, X., & Zhai, C. (2003). Active feedback - UIUC TREC-2003 HARD experiments. *The 12th Text Retrieval Conference, TREC*.
- Struyf, A., Hubert, M., & Rousseeuw, P. (1997). Integrating robust clustering techniques in s-plus. *Computational Statistics and Data Analysis*, 26, 17–37.
- Tang, M., Luo, X., & Roukos, S. (2002). Active learning for statistical natural language parsing. *Proc. of the Association for Computational Linguistics 40th Anniversary Meeting*. Philadelphia, PA.
- Tong, S., & Chang, E. (2001). Support vector machine active learning for image retrieval. *Proceedings of the 9th ACM int. conf. on Multimedia* (pp. 107–118). Ottawa.
- Tong, S., & Koller, D. (2001). Support vector machine active learning with applications to text classification. *Journal of Machine Learning Research*, 2, 45–66.
- Xu, Z., Yu, K., Tresp, V., Xu, X., & Wang, J. (2003). Representative sampling for text classification using support vector machines. *25th European Conf. on Information Retrieval Research, ECIR 2003*. Springer.
- Zhang, C., & Chen, T. (2002). An active learning framework for content-based information retrieval. *IEEE trans on multimedia*, 4, 260–268.
- Zhang, T., & Oles, F. (2000). A probability analysis on the value of unlabeled data for classification problems. *Proc. Int. Conf. on Machine Learning*.
- Zhang, T., & Oles, F. J. (2001). Text categorization based on regularized linear classification methods. *Information Retrieval*, 4, 5–31.
- Zhu, J., & Hastie, T. (2001). Kernel logistic regression and the import vector machine. *Advances in Neural Information Processing Systems*.