

---

# Toward Optimal Active Learning through Monte Carlo Estimation of Error Reduction

---

**Nicholas Roy**

Robotics Institute, Carnegie Mellon University, Pittsburgh, PA 15213 USA

NICHOLAS.ROY@RI.CMU.EDU

**Andrew McCallum**

WhizBang! Labs - Research, 4616 Henry Street, Pittsburgh, PA 15213 USA

MCCALLUM@WHIZBANG.COM

## Abstract

This paper presents an active learning method that directly optimizes expected future error. This is in contrast to many other popular techniques that instead aim to reduce version space size. These methods are popular because for many learning models, closed form calculation of the expected future error is intractable. Our approach is made feasible by taking a Monte Carlo approach to estimating the expected reduction in error due to the labeling of a query. In experimental results on three real-world data sets we reach high accuracy with four times fewer labelled examples than competing methods.

## 1. Introduction

Traditional supervised learning methods set their parameters using whatever training data is given to them. By contrast, *active learning* is a framework in which the learner has the freedom to select which data points are added to its training set. An active learner may begin with a very small number of labelled examples, carefully select a few additional examples for which it requests labels, learn from the result of that request, and then using its newly-gained knowledge carefully choose which examples to request next. In this way the learner aims to reach high performance using as few labelled examples as possible. Active learning can be invaluable when, while it is often the case that *unlabelled* examples are cheap and plentiful, obtaining the labels for those examples is often time-consuming and difficult. Thus, the goal is for the learner to have the lowest error rate with the fewest possible labelings.

Cohn et al. (1996) describe a statistically optimal solution to this problem. Their method selects the training example that, once labelled and added to the training data, is expected to result in the lowest error on future test examples. They develop their method for two simple regression problems in which this question can be answered in closed form. Unfortunately there are many tasks and models for

which the optimal selection cannot efficiently be found in closed form.

Other, more widely used active learning methods attain practicality by optimizing a different, non-optimal criteria. For example, *uncertainty sampling* (Lewis & Gale, 1994) selects the example on which the current learner has lowest certainty; *Query-by-Committee* (Seung et al., 1992; Freund et al., 1997) selects examples that reduce the size of the version space (Mitchell, 1982) (the size of the subset of parameter space that correctly classifies the labelled examples). Tong and Koller’s method (2000) is also based on reducing version space size. None of these methods directly optimize the metric by which the learner will ultimately be evaluated—the learner’s expected error on future test examples. Uncertainty sampling often fails by selecting examples that are outliers—they have high uncertainty, but getting their labels doesn’t help the learner on the bulk of the test distribution. Version-space reducing methods, such as Query-by-Committee often fail by spending effort eliminating areas of parameter space that have no effect on the error rate. Thus these methods also are not immune to selecting outliers; see McCallum and Nigam (1998b) for examples.

This paper presents an active learning method that combines the best of both worlds. Our method selects the next example according to the optimal criteria (reduced error rate on future test examples), but solves the practicality problem by using Monte Carlo estimation. We describe our method in the framework of document classification with pool-based sampling, but it would also apply to other forms of classification or regression, and to generative sampling. We describe an implementation in terms of naive Bayes, but the same technique could apply to any learning method in which incremental training is efficient—for example support vector machines (SVMs) (Cauwenberghs & Poggio, 2000).

Our method estimates future error rate either by log-loss, using the entropy of the posterior class distribution on a sample of the unlabelled examples, or by 0-1 loss, using

the posterior probabilities of the most probable class for the sampled unlabelled examples. At each round of active learning, we select an example for labeling by sampling from the unlabelled examples, adding it to the training set with a sample of its possible labels, and estimating the resulting future error rate as just described. This seemingly daunting sampling and re-training can be made efficient through a number of rearrangements of computation, careful sampling choices, and efficient incremental training procedures for the underlying learner.

We show experimental results on two real-world document classification tasks, where, in comparison with weighted Query-by-Committee we reach 85% of full performance in one-quarter the number of training examples.

## 2. Optimal Active Learning and Monte Carlo Estimation

The optimal active learner is one that asks for labels on the examples that, once incorporated into training, will result in the lowest expected error on the test set. Note that this process is not first-order Markov—the optimal query depends on how many more queries can be made before testing time; for example, if this is the last query before testing, the optimal current query may be different than it is if there are fifty more queries available. However, most work in active learning ignores this factor, and makes the assumption that the next query is the last one before evaluation. In fact several methods take the approximation one step further, and independently select several examples for labeling in a batch.

Let  $P(y|x)$  be an unknown conditional distribution over inputs,  $x$ , and output classes,  $y \in \{y_1, y_2, \dots, y_n\}$ , and let  $P(x)$  be the marginal “input” distribution. The learner is given a labelled training set,  $\mathcal{D}$ , consisting of IID input/output pairs drawn from  $P(x)P(y|x)$ , and estimates a classification function that, given an input  $x$ , produces an estimated output distribution  $\hat{P}_{\mathcal{D}}(y|x)$ . We can then write the expected error of the learner as follows:

$$E_{\hat{P}_{\mathcal{D}}} = \int_x L(P(y|x), \hat{P}(y|x)) P(x) \quad (1)$$

where  $L$  is some loss function that measures the degree of our disappointment in any differences between the true distribution,  $P(y|x)$  and the learners prediction,  $\hat{P}_{\mathcal{D}}(y|x)$ . Two common loss functions are:

*log loss*  $L = \sum_{y \in \mathcal{Y}} P(y|x) \log(\hat{P}_{\mathcal{D}}(y|x))$

and *0/1 loss*  $L = 1 - \max_{y \in \mathcal{Y}} \hat{P}_{\mathcal{D}}(y|x)$ .

First-order Markov active learning thus aims to select a query,  $x^*$ , such that when the query is given label  $y^*$  and added to the training set, the learner trained on the resulting set ( $\mathcal{D} + (x^*, y^*)$ ) has lower error than any other  $x$ .

$$\forall x, E_{\hat{P}_{\mathcal{D}+(x^*, y^*)}} < E_{\hat{P}_{\mathcal{D}+(x, y)}} \quad (2)$$

We concern ourselves here with *pool-based active learning*, in which the learner has available a large pool,  $\mathcal{P}$ , of unlabelled examples sampled from  $P(x)$ , and the queries may be chosen only from this pool. The pool thus not only provides us with a finite set of queries, but also an estimate of  $P(x)$ .

This paper takes a Monte Carlo sampling approach to error estimation and the choice of query. Rather than estimating expected error over the full distribution,  $P(x)$ , we measure it over the sample in the pool. Furthermore, the true output distribution  $P(y|x)$  is unknown for each sample  $x$ , so we estimate it using the current learner. (For log loss this results in estimating the error by the entropy of the learner’s posterior distribution). Writing the labelled documents  $\mathcal{D} + (x^*, y^*)$  as  $\mathcal{D}^*$ , for log loss we have

$$\tilde{E}_{\hat{P}_{\mathcal{D}^*}} = \frac{1}{|\mathcal{P}|} \sum_{x \in \mathcal{P}} \sum_{y \in \mathcal{Y}} \hat{P}_{\mathcal{D}^*}(y|x) \log(\hat{P}_{\mathcal{D}^*}(y|x)), \quad (3)$$

and for 0/1 loss

$$\tilde{E}_{\hat{P}_{\mathcal{D}^*}} = \frac{1}{|\mathcal{P}|} \sum_{x \in \mathcal{P}} \left( 1 - \max_{y \in \mathcal{Y}} \hat{P}_{\mathcal{D}^*}(y|x) \right). \quad (4)$$

Of course, before we make the query, the true label for  $x^*$  is also unknown. Again, the current learned classifier gives us an estimate of the distribution from which the  $x^*$ ’s true label would be chosen,  $\hat{P}_{\mathcal{D}}(y|x)$ , and we can use this in an expectation calculation by calculating the estimated error for each possible label,  $y \in \{y_1, y_2, \dots, y_n\}$ , and taking an average weighted by the current classifier’s posterior,  $\hat{P}_{\mathcal{D}}(y|x)$ .

In the above formulation, we are using the current learner to estimate the true label probabilities, which may seem counter-intuitive. Using these loss functions will cause to the learner to select those examples which maximises the sharpness of learner’s posterior belief about the unlabelled examples. An example will be selected if it dramatically reinforces the learner’s existing belief over unlabelled examples for which it is currently unsure. In practice, selecting these instances for labelling is reasonable because the most useful (or informative) labellings are usually consistent with the learner’s prior belief over the majority (but not all) of unlabelled examples.

Our algorithm thus consists of the following steps:

1. train a classifier using the current labelled examples
2. consider each unlabelled example,  $x$ , in the pool as a candidate for the next labeling request
3. consider each possible label,  $y$ , for  $x$ , and add the pair  $(x, y)$  to the training set
  - (a) re-train the classifier with the enlarged training set,  $\mathcal{D} + (x, y)$
  - (b) estimate the resulting expected loss as in equation (3) or equation (4).

- (c) Assign to  $x$  the average expected losses for each possible labeling,  $y$ , weighted according to the current classifiers posterior,  $\hat{P}_{\mathcal{D}}(y|x)$
- 4. Select for labelling the unlabelled example  $x$  that generated the lowest expected error on all other examples.

If implemented naively, the above algorithm would be hopelessly inefficient. However, with some thought and a some rearrangements of computation, there are a number of optimizations and approximations that make this algorithm much more efficient and very tractable:

- Most importantly, many learning algorithms have algorithms for very efficient *incremental* training. That is, the cost of re-training after adding one more example to the training set is far less than re-training as if the entire set were new. For example, in a naive Bayes classifier, only a few event counts need be incremented. SVMs also have efficient re-training procedures (Cauwenberghs & Poggio, 2000).
- Furthermore, many learners have efficient algorithms for incremental *re-classification* of the examples in the pool. In incremental re-classification, the only parts of computation that need to be redone are those that would have changed as a result of the additional training instance. Again, naive Bayes and SVMs are two examples of algorithms that permit this.
- After adding a candidate query to the training set, the error associated with other examples in the pool does not need to be re-estimated—only those likely to be effected by the inclusion of the candidate in the training set. In many cases this means simply skipping examples not in the “neighborhood” of the candidate or skipping examples without any features that overlap with the features of the candidate. Inverted indices, in which all the examples containing a particular features are listed in an efficiently-accessed list, can make this extremely efficient.
- The pool of candidate queries can be reduced by random sub-sampling, or pre-filtering to remove outliers according to some efficient criteria. In fact, any of the suboptimal active learning methods might make good pre-filters.
- The expected error can be estimated using only a sub-sample of the pool. Especially when the pool is large, there is no need to use all examples—a good estimate may be formed with only few hundred examples.

In the remainder of the paper we describe a naive Bayes implementation of our method, discuss related work, and present experimental results on three real-world data sets showing that our method significantly outperforms methods optimize indirect criteria, such as version-space reduction and query uncertainty. We also outline some future work.

### 3. Naive Bayes Text Classification

Text classification is not only a task of tremendous practical significance, but is also an interesting problem in machine learning because it involves an unusually large number of features, and thus requires estimating an unusually large number of parameters. It is also a domain in which obtaining labelled data is expensive, since the human effort of reading and assigning documents to categories is almost always required. Hence, the large number of parameters often must be estimated from a small amount of labelled data.

When little training data is being used to estimate the parameters for a large number of features, it is often best to use a simple learning model. There isn’t enough data to support estimations of feature correlations and other complex interactions. One such classification method that performs surprisingly well given its simplicity is *naive Bayes*. Naive Bayes is not always the best performing classification algorithm for text (McCallum et al., 2000; Joachims, 1998), but it continues to be widely used for the purpose because it is efficient and simple to implement, and even against significantly more complex methods, it rarely trails far behind in accuracy.

This paper’s Monte Carlo approach to active learning could be applied to several different learners. We apply it here to naive Bayes for the sake of simplicity of explanation and implementation. Experiments with other learners is an item of future work.

Naive Bayes is a Bayesian classifier based on a generative model in which data are produced by first selecting a class,  $y \in \mathcal{Y}$ , and then generating features of the instance,  $x \in \mathcal{X}$ , independently given the class. For text classification, the common variant of naive Bayes has unordered word counts for features, and uses a per-class multinomial to generate the words (McCallum & Nigam, 1998a). Let  $w_i$  be the  $i$ th word in the dictionary of words  $D$ , and  $\theta = (\theta_{y_j}, \theta_{w_t|y_j})_{y_j \in \mathcal{Y}, w_t \in D}$  be the parameters of the model, where  $\theta_{y_j}$  is the prior probability of class  $y_j$ ,  $P(y_j)$ , and where  $\theta_{w_t|y_j}$  is the probability of generating word  $w_t$  from the multinomial associated with class  $y_j$ , and  $\sum_t \theta_{w_t|y_j} = 1$  for all  $y_j$ .

Thus the probability of generating the  $i$ th instance,  $x_i$  is

$$P(x_i|\theta) \propto \sum_{j=1}^{|\mathcal{Y}|} P(y_j|\theta) \prod_{k=1}^{|x_i|} P(w_{x_{ik}}|c_j; \theta) \quad (5)$$

where  $x_{ik}$  is the  $k$ th word in document  $x_i$ . Then, by Bayes rule, the probability that document  $x_i$  was generated by class  $y_j$  is

$$P(y_j|x_i; \theta) = \frac{P(y_j|\theta) \prod_{k=1}^{|x_i|} P(w_{x_{ik}}|y_j; \theta)}{\sum_{r=1}^{|\mathcal{Y}|} P(y_r|\theta) \prod_{k=1}^{|x_i|} P(w_{x_{ik}}|y_r; \theta)} \quad (6)$$

Maximum *a posteriori* parameter estimation is performed

by taking ratios of counts:

$$\hat{\theta}_{w_t|y_j} = \frac{1 + \sum_{i=1}^{|\mathcal{D}|} N(w_t, x_i)P(y_j|x_i)}{|V| + \sum_{s=1}^{|\mathcal{V}|} \sum_{i=1}^{|\mathcal{D}|} N(w_s, x_i)P(y_j|x_i)} \quad (7)$$

$$\hat{\theta}_{y_j} = \frac{\sum_{i=1}^{|\mathcal{D}|} P(y_j|x_i)}{|\mathcal{D}|} \quad (8)$$

where  $N((w_t, x_i))$  is the number of times word  $w_t$  occurs in document  $x_i$ .

### 3.1 Fast Naive Bayes Updates

Equations (3) and (4) show how the pool of unlabelled documents can be used to estimate the change in classifier error if we label one document. However, in order to choose the best candidate from the pool of unlabelled documents  $\mathcal{P}$ , we have to train the classifier  $|\mathcal{P}|$  times, and classify  $|\mathcal{P}| - 1$  documents. Performing  $\mathcal{O}(|\mathcal{P}|^2)$  classifications for every query by the learner can be computationally infeasible for large document pools.

While we cannot reduce the total number of classifications for every query, we can take advantage of the structure in a naive Bayes classifier to allow more efficient retraining of the classifier and relabeling of each unlabelled document, for each putative labelling.

Recall from equation (5) that the each label probability for each unlabelled document is a product of multinomials; given a particular label, there is one multinomial for each word in the document. Retraining a naive Bayes classifier is a process of updating each multinomial for the new label that corresponds to a word in the newly-labelled document.

Notice that when we compute the class probabilities for each unlabelled document using the new classifier, we modifying *some* multinomials in the product of equation (5), but not all. By propagating only changes to the multinomials through the label distributions, we can gain substantial computational savings compared to naively retraining and reclassifying each document from scratch. We update each label probability (which again, is a product of multinomials) by dividing out the changed multinomials, and factoring in the new multinomial values.

More specifically, given a classifier learned from labelled documents  $\mathcal{D}$ , we can label a document  $x_p$  with label  $y_p$  and update the class probabilities (for that class  $y_p$ ) of each unlabelled document  $x_i$  in  $\mathcal{P}$  using equation (9):

$$P(y_p|x_i; \hat{\theta}') = \frac{P(y_p|x_i; \hat{\theta}) \prod_{k=1}^{|x_p|} P'(w_{x_{pk}}|y_p)P(w_{x_{pk}}|x_i)}{\prod_{k=1}^{|x_p|} P(w_{x_{pk}}|y_p)P(w_{x_{pk}}|x_i)} \quad (9)$$

where  $w_{pk}$  is the  $k^{\text{th}}$  word in document  $x_p$ ,  $P'(w_{x_{pk}}|y_p)$  are the new word probabilities given  $D + (x_p, y_p)$ , and  $P(w_{x_{pk}}|y_p)$  are the old word probabilities given only  $D$ .

The product on the right hand side of the numerator is the product of the new multinomials  $P(W|Y)$  that result from

labelling document  $x_p$ . The denominator is the step of dividing out the old multinomials from the previous classifier. The term  $P(w_{x_{pk}}|x_i)$  in top and bottom is simply an indicator variable with value  $\{0, 1\}$  that determines whether  $k$ -th word in document  $x_p$  is also in document  $x_i$ .

The old multinomials that are divided out are the same as in equation (6):

$$P(w_{x_{pk}}|y_p) = \frac{1 + \sum_{l=1}^{|\mathcal{L}|} N(w_{x_{pk}}, x_l)P(y_p|x_l)}{|V| + \sum_{s=1}^{|\mathcal{V}|} \sum_{l=1}^{|\mathcal{L}|} N(w_s, x_l)P(y_p|x_l)} \quad (10)$$

We can update each multinomial rapidly by updating the word counts:

$$P'(w_{x_{pk}}|y_p) = \frac{1 + N(w_{x_{pk}}, x_p) + \sum_{l=1}^{|\mathcal{L}|} N(w_{x_{pk}}, x_l)P(y_p|x_l)}{|V| + |x_p| + \sum_{s=1}^{|\mathcal{V}|} \sum_{l=1}^{|\mathcal{L}|} N(w_s, x_l)P(y_p|x_l)} \quad (11)$$

where  $N(w_{x_{pk}})$  is the word count for the  $k$ -th word in putatively labelled document  $x_p$ . Again, we only do this for the label probabilities of the putative label  $y_p$ ; all other label probabilities remain unchanged.

### 3.2 Obtaining Smoother Posteriors for Naive Bayes

Our active learning method relies on obtaining reasonably accurate class posteriors from the classification procedure. It is well-known that naive Bayes, with its violated independence assumption, gives overly sharp posteriors—the probability of the winning class tends to be very close to 1.0, and the losing classes have probabilities close to 0.

We address this problem with a sampling-based approach to variance reduction, otherwise known as *bagging* (Breiman, 1996). From our original labelled training set of size  $s$ , a different training set is created by sampling  $s$  times with replacement from the original. The learner then creates a new classifier from this sample, this procedure is repeated  $m$  times, and the final class posterior for an instance is taken to be the unweighted average of the class posteriors for each of the classifiers. In regions of uncertain classification is it often the case that the classifiers from different samples give different answers. Thus if the posteriors from any individual classifier are completely extreme, the bagged posterior is more smooth and reflective of the true uncertainty. This approach has been shown not necessarily to reduce overfitting (Domingos, 2000), but it does certainly give better posterior probabilities.

One interesting aspect of this approach is that it can be applied to absolutely any classifier—even ones that don't give class posterior probabilities at all, or for which the distribution over classifier parameters is unclear. This “bagging approach” to sampling from the distribution over classifiers has been used in previous work related to QBC (Abe &

Mamitsuka, 1998); see the related work section for more details.

## 4. Related Work

Cohn et al. (1996) proposed one of the first statistical analyses of active learning, demonstrating how to construct queries that maximize the error reduction by minimizing the learner's variance. They take advantage of the fact that an unbiased learner that minimizes the expected error given as the expected sum of squared error is equivalent to an unbiased learner that minimizes its variance. Such a learner can then use the estimated distribution of  $\hat{y}$  to estimate  $\langle \hat{\sigma}_y^2 \rangle$ , the expected variance of the learner after querying at  $\hat{x}$ . However, a closed-form solution for the expected variance of the text classifier is exceedingly difficult to compute. Furthermore, they construct exactly the query that maximizes this reduction, rather than choosing from a pool of possible queries.

This "Constructive Query Generation" approach is contrasted with "Query-Filtering" (or Seung et al. (1992)'s "Selective Sampling"), in which unlabelled data is presented to the learner from some distribution, and the learner chooses whether or not to query the true label. From this data-oriented perspective, Lewis and Gale (1994) presented the *uncertainty sampling* algorithm for choosing the example with the greatest uncertainty in predicted label, although label distribution uncertainty is a poor way of estimating classifier uncertainty. Freund et al. (1997) showed that the *uncertainty sampling* algorithm does not converge to the optimal classifier as quickly as the "Query-By-Committee" algorithm (Seung et al., 1992).

In the "Query By Committee" (QBC) approach, the approach is to reduce the error of the learner by choosing the instance to be labelled that will minimize the size of the version space (Mitchell, 1982) that is consistent with the labelled examples. Instead of explicitly determining the size of the version space, predicted labels for each unlabelled example (a "committee" of labels) are generated by drawing hypotheses probabilistically from the version space, according to a distribution over the concepts in the version space, and using these hypotheses to predict the example label. Examples arrive from a stream, and are labelled whenever the 'committee' of hypotheses disagree on the predicted label. This approach chooses examples that "split the version space into two parts of comparable size" with a degree of probability that guarantees data efficiency that is logarithmic in the desired probability of error.

A number of others have made use of QBC-style algorithms; in particular, Liere and Tadepalli (1997) use committees of Winnow learners for text classification, and Argamon-Engelson and Dagan (1999) use a Monte Carlo style QBC for natural language processing. Our algorithm differs from theirs in that we are estimating the error re-

duction, whereas Argamon et al. are simply estimating the example disagreement.

Abe and Mamitsuka (1998) use a bagging and boosting approach for maximising the classifier accuracy on the test data. This approach suggests that by maximising the margin on training data, accuracy on test data is improved, an approach that is not always successful (Grove & Schuurmans, 1998). Furthermore, like the QBC algorithms before it, the QBC-by-boosting approach fails to maximize the margin on *all* unlabelled data, instead choosing to query the single instance with the smallest margin.

McCallum and Nigam (1998b) extended the earlier QBC approach by not only using pool-based QBC, but also using a novel disagreement metric. Whereas the stream-based approaches classify whenever a level of disagreement (possibly any) occurs, in pool-based QBC, the *best* unlabelled example should be chosen. Argamon-Engelson and Dagan (1999) suggest using a probabilistic measure based on vote-entropy of the committee, whereas McCallum & Nigam explicitly measure disagreement using the Jensen-Shannon divergence (Lin, 1991; Pereira et al., 1993). However, they recognize that this error metric does not measure the impact that a labelled document had on classifier uncertainty on *other* unlabelled documents. They therefore factored document density into their error metric, to decrease the likelihood of uncertain documents that are outliers. Nevertheless, document density is a heuristic that is specific to text classification, and does not directly measure the impact of a document's *label* on other predicted labelings.

Lindenbaum et al. (1999) examined active learning by minimizing the expected error, however, they used nearest neighbour classifiers which requires an implicit distance metric.

Most recently, Tong and Koller (2000) have used active learning with Support Vector Machines for text classification. In a return to the original QBC formulation, their approach is to estimate the reduction in classifier uncertainty by estimating the reduction in version space size as a function of querying instances. Thus, like QBC, it explicitly reduces version space size, implicitly reducing future expected error. The active learning technique they propose also makes strong assumptions about the linear separability of the data.

## 5. Experimental Results

### NEWSGROUP DOMAIN

The first set of experiments used Ken Lang's News-groups, containing 20,000 articles evenly divided among 20 UseNet discussion groups (McCallum & Nigam, 1998b). We performed two experiments to perform binary classification. The first experiment used the 2 classes `comp.graphics` and `comp.windows.x` as the data

set. The data was pre-processed to remove UseNet headers, UU-encoded binary data, and punctuation. Additionally, words are removed from a stoplist of common words, and words that appeared in less than 3 documents. As in McCallum and Nigam (1998b), no feature selection or stemming was performed. The resulting vocabulary was 10,205 words long. The data set of 2000 documents was split into a training set of 1000 documents, and 1000 test documents. All results reported are the average of 10 trials.

We tested 4 different active learning algorithms:

- *Random* – choosing the query document at random.
- *Most Uncertain Example* – choosing the document with the largest label uncertainty, as in (Lewis & Gale, 1994).
- *Most Disagreed example* – choosing the document with the greatest committee disagreement in the predicted label, as measured by Jensen-Shannon divergence, weighted by document density, as in (McCallum & Nigam, 1998b).
- *Entropy-reducing example* – the method introduced in this paper – choosing the document that maximizes the reduction in the total predicted label entropy, as in equation (1), with error as given in equation (3).

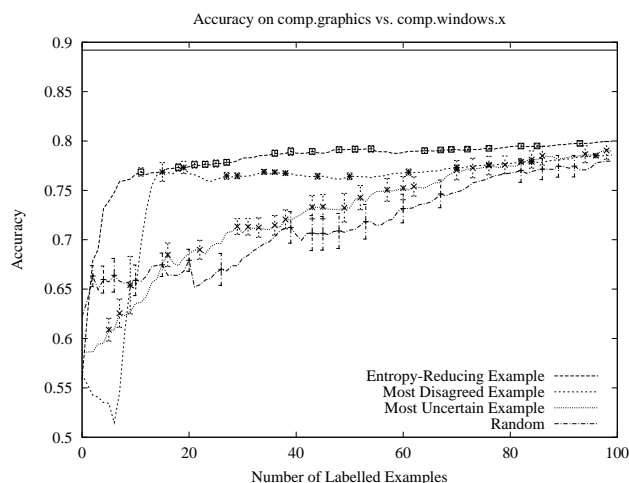
Two of the four algorithms are QBC algorithms, and are given a committee of 3. The algorithms are given 6 labelled examples, 3 from each class, to allow 1 labelled example per class per committee member.

At each iteration, 250 documents (25% of the unlabelled documents) are randomly sampled from the larger pool of unlabelled documents, as candidates for labelling<sup>1</sup>. The error metric was then computed for each putative labelling against all remaining unlabelled documents (not just the sampled pool.) Figure 1 shows the active learning process. The results are reported in terms of accuracy of classification on the held-out test set, up to 100 queries. All results reported are the average of 10 trials.

The solid grey line shows the maximum possible accuracy of 89.2% accuracy after all the unlabelled data has been labelled. After 16 queries, the *Entropy-reducing algorithm* has reached 77.2%, or 85% of the maximum possible accuracy. The *Most Disagreed Example* takes 68 queries to reach the same point (four times more slowly), and maintains a lower accuracy for the remainder of the queries. The *Entropy-reducing algorithm* reaches 80.0% accuracy after 100 queries.

It is also interesting to compare the documents chosen by the two algorithms for initial labelling. Looking at the documents chosen in the first 10 queries, over the 10 trials, the

<sup>1</sup>The subsampling was performed in the interests of these experimental results. In a real active learning setting, all algorithms would be run over as much unlabelled data as was computationally feasible in that setting.



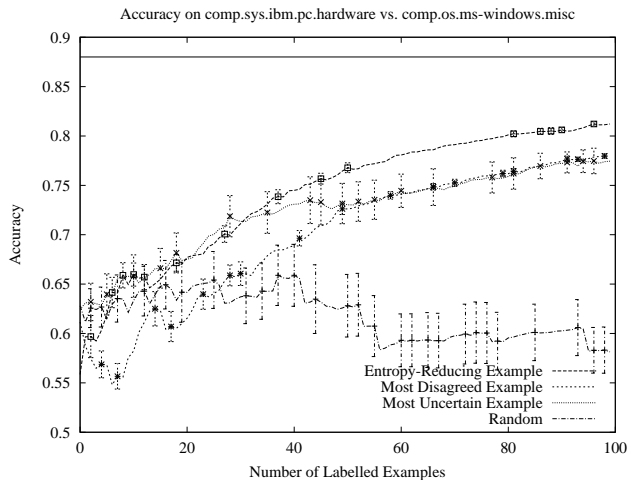
**Figure 1.** Average test set accuracy for `comp.graphics` vs. `comp.windows.x`. The *Entropy-reducing* algorithm reaches 85% of maximum in 16 documents, compared to 68 documents for the *Most Disagreed* algorithm. The error bars are placed at local maximum to reduce clutter.

first 10 documents chosen by the *Entropy-reducing* algorithm are an FAQ, tutorial or How-To 9.8 times out of ten. By comparison, the first 10 documents chosen by the *Most Disagreed* algorithm are an FAQ only 5.8 times out of 10. While the high incidence of FAQs in the initial phases is not quantitatively meaningful, it does suggest that the learner’s behaviour is somewhat intuitive.

The particular newsgroups in the preceding experiment were chosen because they are relatively easy to distinguish. A more difficult text-categorization problem is classifying the newsgroups `comp.sys.ibm.pc.hardware` and `comp.os.ms-windows.misc`. The documents were pre-processed as before, resulting in a vocabulary of 9,895 words. The data set of 2000 documents was split into a training set of 1000 documents, and 1000 test documents. Also as before, the unlabelled data was sampled randomly down to 250 documents for candidate labellings at each iteration, although the error was measured against all unlabelled documents.

The solid grey line shows the maximum possible accuracy of 88% accuracy after all the unlabelled data has been labelled. After 42 queries, the *Entropy-reducing algorithm* has reached 75%, or 85% of the maximum possible accuracy, reaching 80% accuracy after 100 queries. The *Most Disagreed* algorithm takes 70 queries to reach the level of 75% accuracy.

We can again examine the documents chosen by the different algorithms during the initial phases for the incidence of FAQs and tutorials. The *Entropy-reducing* algorithm had an average incidence of 7.3 FAQs in the first 10 documents, compared with 2.6 for the *Most Disagreed* algorithm. In this experiment, however, we see that the apparently intuitive behaviour is not sufficient, and the learners require



**Figure 2.** Average test set accuracy for the `comp.sys.ibm.pc.hardware` vs. `comp.os.ms-windows.misc`. The *Entropy-reducing* algorithm reaches 85% of maximum in 45 documents, compared to 72 documents for the *Most Disagreed* algorithm. The error bars again are placed at local maximum.

several more documents to begin to achieve a reasonable accuracy.

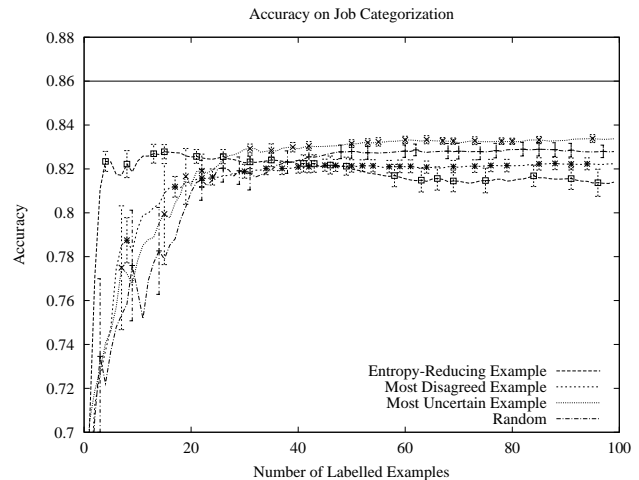
#### JOB CATEGORY DOMAIN

The third set of experiments uses a data set collected at WhizBang! Labs. The *Job Category* data set contains 112,643 documents containing job descriptions for 16 different categories such as *Clerical*, *Educational* or *Engineer*. The data set was collected by automatic spidering of the web, and labelled by hand. The 16 different categories are then broken down into as many as 9 subcategories. We selected the *Engineer* job category, and took 500 articles from each of the 6 *Engineer* categories: *Chemical*, *Civil*, *Industrial*, *Electrical*, *Mechanical*, *Operations* and *Other*. The documents were pre-processed to remove the job title, as well as rare or stoplist words.

This experiment trained the naive Bayes classifier to distinguish one job category from the remaining five. Each data point is an average of 10 trials per job category, averaged over all 6 job categories. In this example, the *Entropy-reducing* algorithm reaches 82% accuracy (94% of the maximum accuracy at 86%) in 5 documents. The *Job Category* data set is easily distinguishable, however, since similar accuracy is achieved after choosing 36 documents at random. The region of interest for evaluating this domain is the initial stages as shown by figure 4. Although the other algorithms do catch up, the *Entropy-reducing* algorithm reaches very high accuracy very quickly.

## 6. Summary

Unlike existing algorithms which typically consider each unlabelled example in isolation, we use the entire pool of



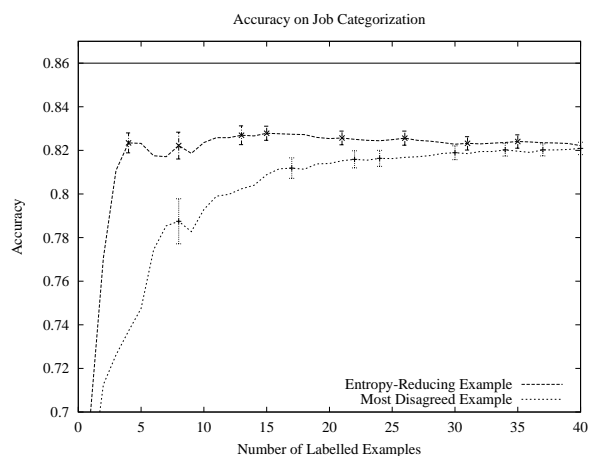
**Figure 3.** Average test set accuracy for the *Job Category* domain, distinguishing one job category from 5 others. The *Entropy-reducing* algorithm reaches 82% accuracy in 5 documents, compared to 36 documents for both the *Most Disagreed* and *Random* algorithms.

unlabelled data to estimate the expected error of the current learner, and we determine the impact on the expected error of each potential labelling request. This approach can be compared to existing statistical techniques (Cohn et al., 1996) that compute the reduction in error (or some equivalent quantity) in closed form; however, we approximate the reduction in error by repeated Monte Carlo sampling of the labelled data and then computing error terms with respect to the unlabelled data on the resultant committee of learners. In this respect, we have attempted to bridge the gap between closed-form statistical active learning and more recent work in *Query-By-Committee* (Freund et al., 1997; McCallum & Nigam, 1998b).

We presented results on two domains, the *Newsgroups* domain and also the *Job Category* domain. Our results show that *Entropy-reducing* algorithm outperforms the best existing algorithm substantially, achieving a high level of accuracy in up to 25% of the labelling queries required by the *Most Disagreed* algorithm.

The algorithm presented here is more computationally complex than most existing QBC algorithms; however, we have shown several ways to make the active learning steps more competitive. Ultimately, the trade-off between computational complexity and query complexity should always be decided in favour of reducing query complexity, which requires humans in the loop. A human labeller typically requires 30 seconds or more to label a document, during which time a computer active learner can select an example in a very large pool of documents. The results presented here typically required less than a second of computation time per query.

Furthermore, our algorithm uses subsampling of the unlabelled data to generate a pool of candidates at each itera-



**Figure 4.** A magnified view of the average test set accuracy for the Job Category domain, distinguishing one job category from 5 others. The Entropy-reducing algorithm clearly reaches a high level of accuracy much before the Most Disagreed algorithm.

tion. By initially using a fairly restrictive pool of candidates for labelling, and increasing the pool as time permits, our algorithm can be considered an anytime algorithm.

Our technique should perform even more strongly with models that are complex and have complex parameter spaces—not all of which are relevant to performance on a particular data set. In these situations active learning methods based on version space reduction would not be expected to work as well, since they will expend effort excluding portions of version space that have no impact on expected error.

We plan to extend this active learning technique to other classifiers, such as Support Vector Machines. The recent work by Cauwenberghs and Poggio (2000) describes techniques for efficient incremental updates to SVMs and will apply to our approach. In the absence of sufficient efficiency in their technique, an intermediate heuristic might be to perform data selection via the naive Bayes approach, and measure the test set accuracy using a SVM trained on the same data.

## Acknowledgements

Nicholas Roy was supported by Whizbang! Labs for this work. The authors would like to thank Andrew Ng for his advice and suggestions. Kamal Nigam, Simon Tong and Pedro Domingos provided much valuable feedback and suggestions of earlier versions of this paper, as did the anonymous reviewers.

## References

Abe, N., & Mamitsuka, H. (1998). Query learning strategies using boosting and bagging. *International Conference on Machine Learning (ICML)*.

Argamon-Engelson, S., & Dagan, I. (1999). Committee-based sample selection for probabilistic classifiers. *Journal of Artificial Intelligence Research*, 11, 335–460.

Breiman, L. (1996). Bagging predictors. *Machine Learning*, 24, 123–140.

Cauwenberghs, G., & Poggio, T. (2000). Incremental and decremental support vector machine learning. *Advances in Neural Information Processing 13 (NIPS)*. Denver, CO.

Cohn, D. A., Ghahramani, Z., & Jordan, M. I. (1996). Active learning with statistical models. *Journal of Artificial Intelligence Research*, 4, 129–145.

Domingos, P. (2000). Bayesian averaging of classifiers and the overfitting problem. *Proceedings of the International Conference on Machine Learning (ICML)* (pp. 223–230).

Freund, Y., Seung, H., Shamir, E., & Tishby, N. (1997). Selective sampling using the Query By Committee algorithm. *Machine Learning*, 28, 133–168.

Grove, A., & Schuurmans, D. (1998). Boosting in the limit: Maximizing the margin of learned ensembles. *Proc. of the Fifteenth National Conference on Artificial Intelligence (AAAI-98)*.

Joachims, T. (1998). Text categorization with support vector machines: Learning with many relevant features. *Proceedings of the European Conference on Machine Learning*.

Lewis, D., & Gale, W. (1994). A sequential algorithm for training text classifiers. *Proceedings of the International ACM-SIGIR Conference on Research and Development in Information Retrieval* (pp. 3–12).

Liere, R., & Tadepalli, P. (1997). Active learning with committees for text categorization. *Proceedings of the National Conference in Artificial Intelligence (AAAI-97)*. Providence, RI.

Lin, K. (1991). Divergence measures based on the Shannon entropy. *IEEE Transactions on Information Theory*, 37, 145–151.

Lindenbaum, M., Markovitch, S., & Rusakov, D. (1999). Selective sampling for nearest neighbor classifiers. *Proceedings of the Sixteenth National Conference on Artificial Intelligence (AAAI-99)* (pp. 366–371).

McCallum, A., Freitag, D., & Pereira, F. (2000). Maximum entropy markov models for information extraction and segmentation. *Proceedings of International Conference on Machine Learning (ICML)*.

McCallum, A., & Nigam, K. (1998a). A comparison of event models for Naive Bayes text classification. *AAAI-98 Workshop on "Learning for Text Categorization"*.

McCallum, A., & Nigam, K. (1998b). Employing EM and pool-based active learning for text classification. *Proc. of the Fifteenth Inter. Conference on Machine Learning* (pp. 359–367).

Mitchell, T. M. (1982). Generalization as search. *Artificial Intelligence*, 18.

Pereira, F., Tishby, N., & Lee, L. (1993). Distributional clustering of English words. *Proceedings of the 31st ACL*.

Seung, H. S., Opper, M., & Sompolinsky, H. (1992). Query by committee. *Proceedings of the Fifth Annual ACM Workshop on Computational Learning Theory* (pp. 287–294).

Tong, S., & Koller, D. (2000). Support vector machine active learning with applications to text classification. *Proc. of the Seventeenth International Conference on Machine Learning*.