

Document Classification Through Interactive Supervision of Document and Term Labels

Shantanu Godbole, Abhay Harpale, Sunita Sarawagi, and Soumen Chakrabarti

IIT Bombay
Powai, Mumbai, 400076, India
shantanu@it.iitb.ac.in

Abstract. Effective incorporation of human expertise, while exerting a low cognitive load, is a critical aspect of real-life text classification applications that is not adequately addressed by batch-supervised high-accuracy learners. Standard text classifiers are supervised in only one way: assigning labels to whole documents. They are thus deprived of the enormous wisdom that humans carry about the significance of words and phrases in context. We present HIClass, an interactive and exploratory labeling package that actively collects user opinion on feature representations and choices, as well as whole-document labels, while minimizing redundancy in the input sought. Preliminary experience suggests that, starting with essentially an unlabeled corpus, very little cognitive labor suffices to set up a labeled collection on which standard classifiers perform well.

1 Introduction

Motivated by applications like spam filtering, e-mail routing, Web directory maintenance, and news filtering, text classification has been researched extensively in recent years [1–3]. State-of-the-art classifiers now achieve up to 90% accuracy on well-known benchmarks. Almost all machine learning text classification research assumes some fixed, simple class of feature representation (such as bag-of-words), and at least a partially labeled corpus. Statistical learners also depend on the deployment scenario to be reasonably related to the training population.

Many of these assumptions do not hold in real-life applications. Discrimination between labels can be difficult unless features are engineered and selected with extensive human knowledge. Often, there is no labeled collection to start with. In fact, even the label set may not be specified up front, and must evolve with the user's understanding of the application. Several projects reported at the annual Operational Text Classification workshops [4] describe applications spanning law, journalism, libraries and scholarly publications in which automated, batch-mode techniques were not satisfactory; substantial human involvement was required before a suitable feature set, label system, labeled corpus, rule base, and resulting system accuracy were attained. However, not all the techniques used in commercial systems are publicly known, and few general principles can be derived from these systems.

There is much scope for building machine learning tools which engage the user in an active dialog to acquire human knowledge about features and document labels. When such supervision is available only as label assignments, *active learning* has provided some clear principles [5–7] and strategies for maximum payoffs from the dialog. We wish to extend the active learning paradigm significantly to include both feature engineering and document labeling conversations, exploiting rapidly increasing computing power to give the user immediate feedback on her choices.

Our contributions: In this paper we present the design of a system HIClass (Hyper Interactive text Classification) for providing this tight interaction loop. We extend SVMs to naturally absorb human inputs in the form of feature engineering, term inclusion/exclusion and term and document labels. In the past, such actions were performed through ad hoc means and as a distinct processing step before classification construction. We make these more effective by (1) providing the user easy access to a rich variety of summaries about the learnt model, the input data and aggregate performance measures, (2) drawing the user’s attention to terms, classes or documents in greatest need of inspection, and (3) helping the user assess the effect of every choice on the performance of the system on test data.

Outline: We describe the HIClass workbench in Section 2 and review the design choices and various modes of user interaction. Section 3 describes our method of active learning on documents with modifications to handle multi-labeled data and methods to reduce the cognitive load on the user. Section 4 introduces the idea of active learning on terms treating them as first class entities. We report our experiences with the workbench and experimental results in Section 5. We review related work in Section 6 and conclude in Section 7.

2 The HIClass Workbench for Text Classification

We present an overview of HIClass in Fig. 1. The lower layer shows the main data entities and the main processing units. There is a small pool of labeled documents (partitioned by sampling into a training and test set) and a large unlabeled pool. The feature extractor turns documents into feature vectors. Features are usually words, but the user can interactively refine features to be more complex; this is described next. The system can store and access by name multiple classifiers with their fitted parameters at any given time, assisting comparative analysis of parameters, performance on held-out data, and drill-down error diagnostics. The upper layer shows the prominent menus/modes in which a user can interact with the system. Next we describe the important building blocks shown in Fig. 1.

2.1 Document and Classification Models

The first step of the design of HIClass is to choose a flexible classification model that (1) suits state-of-the-art automated learners and (2) can be easily interpreted and tuned by the user.

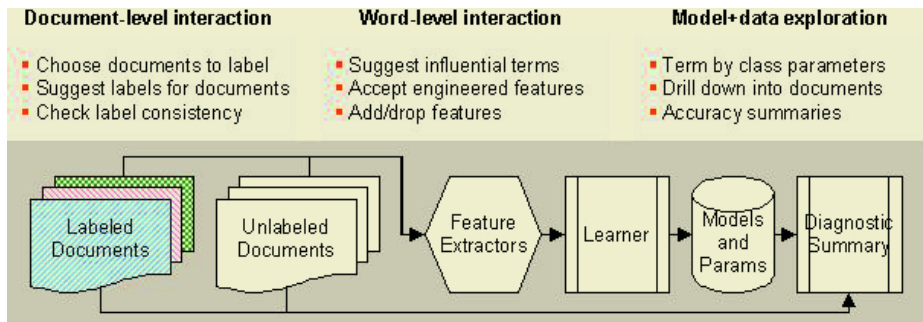


Fig. 1. The architecture of HIClass.

A document is a bag of features. Usually, features are words after minor processing like stemming and case-normalization. But the user can also (dynamically) define features to reflect domain knowledge. E.g., month names or currency names may be conflated into synthetic features. On the other hand, the user may notice harmful conflation between “blood bank” and “bank”, and define “blood bank” as a single compound feature. We will continue to use *term*, *word* and *feature* interchangeably where no confusion can result. Documents are represented as unit vectors.

Labeled documents can be associated with more than one class in general. HIClass supports **linear additive** classifier models, where each class c is associated with a set of weights w_1^c, \dots, w_T^c corresponding to the T terms in a vocabulary. Each document is represented by a vector of non-negative weights $\mathbf{x} = (x_1, \dots, x_T)$, each component corresponding to a feature. The classifier assigns a document all class labels c for which $\mathbf{w}^c \cdot \mathbf{x} + b_c \geq 0$ where b_c is a scalar per-class bias parameter. As documents vectors have only non-negative components, both magnitude and sign of components of w^c give natural interpretations of salience of terms.

The linear additive model generalizes a number of widely-used classifiers, including naive Bayes (NB), maximum entropy, logistic regression, and support vector machines (SVMs). Here we focus on SVMs. Given documents d_i with labels $y_i \in \{-1, +1\}$, a two-class linear SVM finds a vector \mathbf{w} and a scalar constant b , such that for all documents $y_i(\mathbf{w}_c \cdot d_i + b) \geq 1$, and $\|\mathbf{w}_c\|$ is minimized.

When the application demands more than two classes, one can (1) rewrite the above optimization slightly, with one \mathbf{w} vector per class, so that the discriminant $\mathbf{w}_{c_j} \cdot d_i + b_j$ is largest for the correct class c_j ; or (2) build an ensemble of SVMs, each playing off one class against another (“one-vs-one”), and assigning the document to the class that wins the largest number of matches; or (3) build an ensemble of SVMs, as many as there are classes, each predicting an yes/no label for its corresponding class (“one-vs-rest” or “one-vs-others”). In practice, all these approaches are comparable in accuracy [8]. We use one-vs-others as it is easily extended to make multi-labeled prediction and is efficient.

2.2 Exploration of Data, Model, Performance Summaries

HIClass provides support for viewing the trained classifier scores, aggregate as well as drill-down statistics about terms, documents and classes, and standard accuracy measures. Aggregate statistics like per-class population, similarity distribution, and uncertainty distribution can be viewed.

After building an initial classifier using the starting labeled set L , the user can view the learnt model as a matrix of class-by-term scores. Simple inspection of term-class scores in an OLAP-like tool enables the expert to propose changes to per-class classification models like including and excluding certain terms. The user-interface allows easy movement from a term-centric to a class-centric analysis where the user can see documents with terms indicative of belongingness to classes.

With every proposed change, the user can study the impact of the change by observing its performance on the test data. The user can inspect graphs for the accuracy and $F1$ of the whole system or for individual classes across iterations. The user can identify classes which hamper the overall performance of the system and can concentrate on them further by adding more labeled documents or fine tuning important relevant terms. The user can also inspect a confusion matrix between any two classes that reveals their overlap, allowing the user to inspect and tune discriminating terms by inspecting the results of a binary SVM on the classes.

2.3 Feature Engineering

Fast evaluation over a variety of test data enables a user to easily identify limitations of a trained model and perhaps the associated feature set. Most users, on inspection of this set of scores, will be able to propose a number of modifications to the classifiers. Some of these modifications may not impact performance on the available test set but could be beneficial in improving the robustness and performance of the classifier in the long term. For the Reuters dataset, close inspection of some of the terms shown to have a high positive weight for the class *crude* reveals:

- “Reagan” is found to be a positive indicator of the class *crude* though proper names should be identified and treated differently.
- “Ecuador” and “Ecuadorean” reveal insufficient stemming.
- “World bank” and “Buenos Aires” should always occur together as a bigram; “Union”, a high weight term for *crude*, should be associated with “Pacific Union” in *crude*, but as “Soviet Union” in other classes.
- Month names, currencies, date formats, proper nouns should be recognized and grouped into appropriate aggregate features.

2.4 Document Labeling Assistant

When unlabeled data is abundant and labeled data is limited, a user can choose to add labels to some of the unlabeled documents. Active learning has proven

to be highly effective in interactively choosing documents for labeling so that the total number of documents to be labeled is minimized. HIClass provides a number of mechanisms to lighten the user’s cognitive load in the document labeling process. Details of appear in Section 3.

2.5 Term-Level Active Learning

The high accuracy of linear SVMs at text classification [1] suggests that class membership decision depends on the combination of “soft” evidence from a class-conditional subset of the corpus vocabulary. E.g., high rate of occurrence of one or more of the words *wicket*, *run*, *stump*, and *ball* leads us to believe a document is about (the game) cricket. Given enough training documents, good classifiers can learn the importance of these terms. However, in the initial stages of bootstrapping a labeled corpus, it is far more natural for the user to directly specify these important features as being positively associated with the class “cricket”, rather than scan and label long documents containing these words.

HIClass allows users to label terms with classes just like documents. We expect the cognitive load of labeling terms to be lower because the user does not have to waste time reading long documents. We help the user in spotting such terms by doing active learning on terms. This is elaborated in Section 4.

3 Active Labeling of Documents

The system starts with a small training pool of labeled documents L and a large pool of unlabeled documents U . Assume that the number of class labels is k and each document can be assigned multiple labels. We train k one-vs-others SVMs on L . Our goal during active learning is to pick some unlabeled documents about whose predictions the classifier is *most uncertain*. Various measures are used for calculating uncertainty with SVMs [6]. However, these assume binary, single-labeled documents. We extend these to the multi-class, multi-labeled setting as described next.

3.1 Uncertainty

Each unlabeled document gets k discriminant values, one from each SVM in the one-vs-others ensemble. We arrange these values on the number line, and find the largest gap between adjacent values. A reasonable policy for multilabel classification using one-vs-others SVMs is that discriminant values to the right of the gap (larger values) correspond to SVMs that should be assigned a positive label to the document and the rest should be negative.

We need this policy because, in our experience with one-vs-others ensembles, as many as 30% of documents may be labeled negative by all members of the ensemble. For single label classification, it is common to pick the maximum discriminant even if it is negative. Our policy may be regarded as an extension of this heuristic to predict multiple labels.

With this policy, we declare that document to be most uncertain whose this largest gap is the smallest among all documents. When documents are restricted to have one label, this reduces to defining certainty (confidence) in terms of the gap between the highest scoring and the second highest scoring class.

3.2 Bulk-Labeling

The user could label these uncertain documents one by one. But experience suggests that we can do better: often, many of these document are quite similar, and if we could present tight clusters that the user can label all at once, we can reduce the cognitive load on the user and speed up the interaction.

We pick the u most uncertain documents and compute pairwise vector-space similarity between documents in the uncertain set, and prepare for the user a cluster/subset of fixed size (set by the parameter s) that has the largest sum of pairwise similarities.

When showing these uncertain clusters to the user, we also provide an ordered list of suggested labels. The ordering is created by taking the centroid of each uncertain cluster and finding its similarity to the k centroids of positive training data of the k classes. Fig. 2 summarizes the active bulk-labeling process for documents.

```

Start with a labeled pool  $L$  and an unlabeled pool  $U$ .
while user wants to continue with active labeling do
  Train a  $A$ -vs-not $A$  SVM ensemble on  $T$ 
  Calculate uncertainty on all documents in  $U$ :
  for all documents  $d \in U$  do
    Get  $k$  scores by applying the  $k$  SVMs to  $d$ . Find the largest gap in score
    values.
  end for
  Sort the  $|U|$  gaps in ascending order and add top  $u$  to the uncertain set.
  Select the  $s$  most similar documents from top  $u$ 
  Suggesting ranked list of labels for the group  $s$ :
  for all  $k$  classes do
    Find similarity between centroid of  $s$  and centroid of positive training data of
    class  $k$ 
  end for
  Sort these distances in a suggested list of classes
  Present  $s$  and the ranked list of  $k$  suggestions to the user for active labeling
  Accept multi-labeled suggestions for all documents in  $s$ . Check for conflicts
  Add these  $s$  documents to  $L$  with user provided labels and remove from  $U$ 
end while

```

Fig. 2. The algorithm for active learning on documents.

(An alternative is to use the existing classifier itself to propose suggestions based on the confidence with which the documents in the uncertain cluster are classified into various classes. However, we feel keeping the same suggestion list

for all documents in each uncertain cluster reduces the cognitive load on the user. Also, empirically we found in the initial stages this provides better suggestion than the SVMs.)

The user provides feedback to the system by labeling all documents in an uncertain cluster in one shot. The labeled documents are inspected by a conflict check module for consistency. We defer discussion of this topic due to lack of space. Once the user confirms the labels, the newly labeled documents are removed from U and added to L . The system then iterates back to re-training the SVM ensemble.

4 Active Learning Involving Terms

As mentioned in Section 2.5, users generally find it easier to bootstrap the labeled set using trigger terms (that they already know) rather than tediously scrutinize lengthy documents for known triggers. We demonstrate this with an example from the Reuters dataset. We trained two SVMs using the *interest* class in Reuters; the first trained with a single document per class and the second trained with 50 documents per class. For each SVM, we report some terms corresponding to the maximum positive weights in the table. The SVM using more data contains terms like “rate” “fe” (foreign exchange), “pct” (percent), and “interest”: that a user can readily recognize as being positively associated with the label *interest* that are missing from the first SVM.

Num labeled=1		Num labeled=50	
Term	w	Term	w
forecast	0.40	rate	2.08
bank	0.29	fe	1.97
noon	0.20	pct	1.65
account	0.20	market	1.26
oper	0.14	custom	1.01
market	0.14	interest	0.92
england	0.09	forecast	0.92
		stg	0.87
		bank	0.83

We allow a direct process of proposing trigger terms within the additive linear framework. We believe such manual addition of terms will be most useful in the initial phases to bootstrap a starting classifier which is subsequently strengthened using document-level active learning. We propose a mechanism analogous to active learning on documents to help a user spot such terms. SVMs treat labeled terms as mini-documents whose vector representation has a 1 at the term’s position and 0 everywhere else, resulting in standard unit length document vectors.

We develop a criterion for term active learning that is based on the theoretically optimum criterion of minimizing uncertainty on the unlabeled set but avoids the exhaustive approach required to implement it [5–7] by exploiting the special nature of single-term documents.

Consider adding a term t whose current weight is w_t in the trained SVM. For terms not in any of the labeled documents $w_t = 0$. Suppose we add t as a “mini-document” with the user-assigned label y_t . Let the new SVM weight vector be w' . Since the term t is a mini-document whose vector has $x_t = 1$ and $\forall t' \neq t, x_{t'} = 0$, we can assume that in the new w' only w_t is changed to a new w'_t and no other $w_{t'}$ is affected significantly. This is particularly true for terms that do not already appear in the labeled set. From the formulation of SVMs, $y_t(w'_t + b) \geq 1$.

If the current w_t is such that $|w_t + b| \geq 1$ then adding t will probably not have any affect. So we consider only those ts where $|w_t + b| < 1$. Adding t with a label $+1$ will enforce $w'_t + b = 1$ i.e., $w'_t = 1 - b$ and with a label of -1 will make it $w'_t = -1 - b$. For each possible value of $y_t = c$, we get a new value of $w'_t(c)$. Thus we can directly compute the new uncertainty of each unlabeled document x by computing the *change* in the distance from separator value as $(w'_t(c) - w_t)x_t$, since uncertainty is inversely proportional to distance from the separator. Let $Pr(c, t)$ be the probability that the term t will be assigned to class c , as our weighing factor. We estimate $Pr(c, t)$ by the fraction of documents containing term t which have been predicted to belong to class c . We then compute the weighted uncertainty $WU(t)$ for a term t as $WU(t) = \sum_c U(c, t)Pr(c, t)$ and then select the term with the smallest $WU(t)$ for labeling. Other details and approximate variants can be found in [9]. This gives us a way to compute the total uncertainty over the unlabeled set without retraining a SVM for each candidate term.

5 Experimental Study

We have experimented with several text classification tasks ranging from well-established benchmarks like Reuters and 20-newsgroups to more noisy classification tasks, like the Outdoors dataset, chosen from Web directories [11]. It is difficult to quantify the many ways in which HIClass is useful. Therefore we pick a few measures like the benefits of active learning with terms and document to report as performance numbers. We also present some results which quantify the cognitive load on the user and try to show how HIClass eases the user’s interaction and labeling process.

HIClass consists of roughly 5000 lines of C++ code for the backend and 1000 lines of PHP scripts to manage frontend user interactions. The frontend is a web browser, readily available on any user’s desktop. XML is used to pass messages between the frontend and the server backend. LibSVM [12] is used as the underlying SVM classifier.

All our development and experiments were done on a dual-processor P3 server running Debian Linux and with 2GB RAM. Due to space limitations we report numbers for fixed settings of some of our system parameters. Further experiments can be found in [9]. Unless otherwise stated, the number of initial documents per class is set to 1, the number of documents selected for bulk labeling is 5 and the number of uncertain documents over which we pick similar clusters (the parameter u of Section 3) is set to 75.

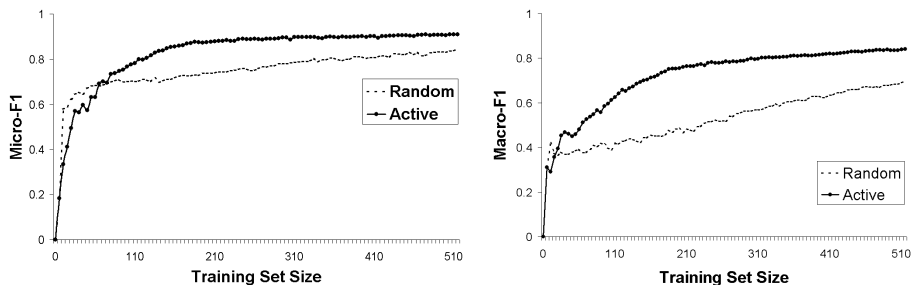


Fig. 3. Reuters - Micro and Macro-averaged F1 on held-out test data while increasing training set size, randomly versus using document level active learning.

5.1 Document-Level Active Learning

We now show how active learning on documents can reduce the number of documents for which the user needs to provide labels in a multiclass, multi-labeled setting. We started with one document in each class and added 5 documents in each round. All graphs are averaged over 30 random runs. Fig. 3 compares the micro and macro averaged $F1$, of selecting 5 documents per round using active learning and using random selection for Reuters (similar results with other datasets omitted due to lack of space). We see that active learning outperforms randomly adding documents to L and reaches its peak $F1$ levels faster.

5.2 Reducing Labeling Effort

We next show the effectiveness of the two techniques that we proposed in Section 3 for reducing the effort spent for labeling a document. For lack of space we only show results with Reuters in this sub-section.

Quality of Suggestions. We quantify the quality of suggestions provided to the user by the average rank of the true labels in the suggested list. We see in Figure 4 that even in the initial stages of active learning the true classes on an average are within rank 4 whereas the total number of possible classes is 20 for this dataset. We also see that the suggestions with u fixed at 75 are better than at 10 as expected.

Bulk-Labeling. We quantify the benefit of bulk-labeling by measuring **inverse similarity**, defined as the number of true distinct labels in a batch of s documents as a fraction of the total number of document-label pairs in the batch. So, if $s = 5$ and each document in a batch has one label and all of them are the same, then the inverse similarity is $\frac{1}{5}$.

It is reasonable to assume that the cognitive load of labeling is proportional to the number of distinct labels that the user has to assign. Thus, Fig. 5 establishes that our chosen set of similar documents reduce cognitive load by a

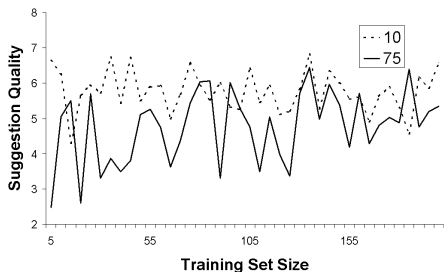


Fig. 4. Reuters – Quality of suggestion measured as the rank at which correct labels are found in the suggested labels.

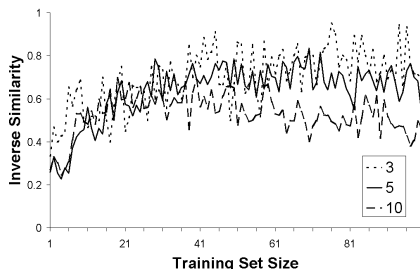


Fig. 5. Reuters – Benefits of bulk-labeling measured as inverse similarity defined in section 5.2.

factor of 2. The benefits are higher in the initial stages because then there are several documents with high uncertainty to choose from. With higher number of documents per batch, the benefits get larger.

We cannot set s to be very high because there is a tradeoff between *reducing effort per label* by bulk labeling similar documents and *increasing number of labels* by possibly including redundant documents per batch. If we calculated labeling cost in terms of *number of documents* to be labeled, the optimum strategy is to label the most uncertain single document per batch. But the effort the user has to spend in deciding on the right label for rapidly changing document contexts will be high. The right tradeoff can only be obtained through experience and will vary with different classification tasks and also the user’s experience and familiarity with the data.

5.3 Term-Level Active Learning

Our goal here is to evaluate the efficacy of training with labeled terms. We take all available labeled documents for a class and train a one-vs-rest SVM for that class. All single-term documents that are predicted as positive or negative with very large margins (above $b/3$ here) are labeled with the predicted class and the rest are not labeled. We then start with a SVM trained initially with a single labeled document on each side and keep adding these collected labeled terms in order of the magnitude of their weights (the AllData method). We also evaluate the performance of our term level active learning described in 4. However, we use an approximation algorithm which is less time-intensive and computationally efficient. We select terms with higher values of $f(t)$ where $f(t) = (\sum_{i \in pos} x_{it} - \sum_{i \in neg} x_{it}) * (N - (pos - neg)(b + w_t))$ where N is the total number of unlabeled documents, and pos and neg refer to number of positive and negative unlabeled documents.

In Fig. 6 we show the resulting accuracy on 8 classes of Reuters and 3 classes of the 20-newsgroups dataset. Active learning on terms clearly works as expected though the gains are small. This is to be expected since SVMs are trained with very few terms instead of entire documents. Random selection performs much

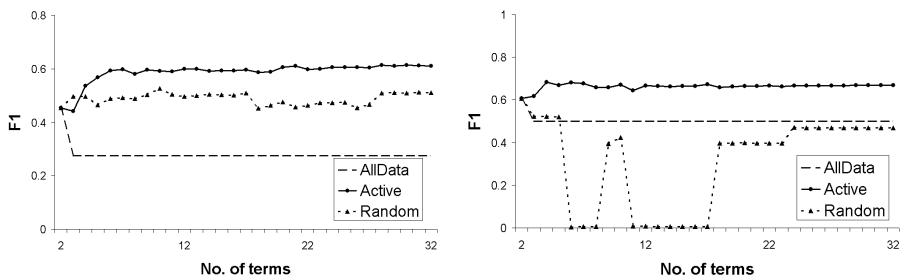


Fig. 6. Adding labeled terms in score order Reuters (left) and 20-newsgroups (right).

worse in both the datasets. This confirmed our intuition that term-level active learning is best viewed as a bootstrapping technique followed by document-level active learning.

6 Related Work

Most earlier work on applying active learning to text categorization [6, 10] assume a single binary SVM whereas our proposed scheme is for multiple one-vs-others SVMs and for multi-labeled classification. Active learning has also recently been applied to the problem of selecting missing attributes of labeled instances whose values should be filled in by the user [13]. This is different from our setting of term active learning because our goal is to add terms as additional labeled instances. The notion of labeling terms is used in [14] for building lexicons of terms related to a concept. So the goal there is not to assign documents to categories but to exploit the co-occurrence patterns of terms in documents to categorize terms.

7 Conclusion

We have described HIClass, an interactive workbench for text classification which combines the cognitive power of humans with the power of automated learners to make statistically sound decisions. The system is based on active learning, starting with a small pool of labeled documents and a large pool of unlabeled documents. We introduce the novel concept of active learning on terms for text classification. We describe our OLAP-like interface for browsing the term-class matrix of the classifier cast as a linear additive model. The user can tune weights of terms in classes leading to better, more understandable classifiers. HIClass provides user continuous feedback on the state of the system, drawing her attention to classes, documents, and terms which would benefit by manual tuning.

References

1. T. Joachims. Text categorization with support vector machines: learning with many relevant features. In *Proceedings of ECML-98*.
2. K. Nigam, J. Lafferty, and A. McCallum. Using maximum entropy for text classification. In *IJCAI'99 Workshop on Information Filtering*.

3. J. Zhang and Y. Yang. Robustness of regularized linear classification methods in text categorization. In *SIGIR*, 2003.
4. Third workshop on Operational Text Classification OTC 2003. In conjunction with *SIGKDD 2003*.
5. D. A. Cohn, Z. Ghahramani, and M. I. Jordan. Active learning with statistical models. In *Advances in Neural Information Processing Systems*, 1995.
6. S. Tong and D. Koller. Support vector machine active learning with applications to text classification. *Journal of Machine Learning Research*, 2:45–66, Nov. 2001.
7. Y. Freund, H. S. Seung, E. Shamir, and N. Tishby. Selective sampling using the query by committee algorithm. *Machine Learning*, 28(2-3):133–168, 1997.
8. C. Hsu and C. Lin. A comparison of methods for multi-class support vector machines. In *IEEE Transactions on Neural Networks*, 13(2002), 415–425, 2001.
9. A. Harpale. Practical alternatives for active learning with applications to text classification. Master’s Thesis, IIT Bombay, 2004.
10. A. K. McCallum and K. Nigam. Employing EM in pool-based active learning for text classification. In *Proceedings of ICML-98*.
11. S. Sarawagi, S. Chakrabarti, and S. Godbole. Cross-training: Learning probabilistic mappings between topics. In *SIGKDD*, 2003.
12. C.C. Chang, and C.J. Lin. LIBSVM: a library for support vector machines, 2001. <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>
13. D. Lizotte, O. Madani, and R. Greiner. Budgeted learning of naive-bayes classifiers. In *UAI*, 2003.
14. H. Avancini, A. Lavelli, B. Magnini, F. Sebastiani, and R. Zanolì. Expanding domain-specific lexicons by term categorization. In *SAC*, 2003.