

Active Learning for Mention Detection: A Comparison of Sentence Selection Strategies

Nitin Madnani
nmadnani@umiacs.umd.edu

Institute for Advanced Computer Studies
Dept. of Computer Science
University of Maryland
College Park, MD, 20770

Abstract. We propose and compare various sentence selection strategies for active learning for the task of detecting mentions of entities. The best strategy employs the sum of confidences of two statistical classifiers trained on different views of the data. Our experimental results show that, compared to the random selection strategy, this strategy reduces the amount of required labeled training data by over 50% while achieving the same performance. The effect is even more significant when only named mentions are considered: the system achieves the same performance by using only 42% of the training data required by the random selection strategy.

1 Introduction

Human annotation is expensive, yet it is needed in many tasks in order to create training data. Given the high cost, it is critical to improve the efficiency of such annotation. Active learning involves selecting samples “intelligently” rather than randomly, for human annotation. By doing so, it is possible for systems to attain better performance with the same amount of annotation or achieve the same level of performance with a lot less annotated data.

In this paper, we present several new active learning strategies for the task of Mention Detection (MD). Here, we employ the terminology used in the Automatic Content Extraction Conferences [1]. Mentions are references to real-world entities that can be *named* (e.g. “John”), *nominal* (e.g. “survivor”) or *pronominal* (e.g. “he”).

We propose and investigate a variety of sentence selection criteria for active learning, including various sentence scoring metrics that combine uncertainty-based and query-by-committee like measurements. Experimental results show that these sentence selection strategies are quite effective for mention detection: compared to the random selection strategy, the best strategy reduces the amount of required annotated training data by over 50% while achieving the same performance. The effect is even more significant when only named mentions are considered: the system achieves the same performance by using only 42% of the training data required by the random strategy.

In the next section, we discuss related work on active learning. Section 3 describes our framework in detail and presents our experimental setup. Section 4 presents the results of the different experiments. Finally, we discuss observations from our experiments and present some ideas for future work.

2 Related Work

Active learning has been utilized for many NLP applications, most noticeably for text classification [2–5]. It has also been applied for part-of-speech tagging [6], statistical parsing [7–10], noun phrase chunking [11], Japanese word segmentation [12], and confusion set disambiguation [13].

There are few reported instances of applying active learning techniques to mention detection. [14] investigated several document, rather than sentence, selection strategies for Information Extraction. They observed that some strategies lead to improvement in recall while others improve precision, but it is difficult to get significant improvement in both recall and precision for an active learner to perform better than random selection.

[15] proposed a multi-criteria-based active learning approach for named entity recognition. The multiple criteria include informativeness, representativeness, and diversity. They proposed measures to quantify these properties and investigated different selection strategies. They showed that the labeling cost can be reduced by at least 80% without degrading the performance for their data sets.

[16] investigated a query-by-committee-based active learner for information extraction in the astronomy domain. It studied the effects of selective sampling on human annotators. Although active learning improved annotation efficiency overall, they observed lower inter-annotator agreement and higher per-token annotation times for the data selected by active learning.

3 Active Learning

Active learning techniques are usually divided into two types: uncertainty sampling for a single learner [17], or disagreement measurement between a committee of learners [18]. In each case, seed data needs to be provided to build an initial model or models. In the uncertainty-based approach, a single learner labels unlabeled examples and provides a confidence score for each predicted label. Samples that have the lowest confidence scores are chosen for manual labeling. In the query-by-committee approach, a committee of learners is built and each learner labels the unlabeled samples. Samples that have the highest disagreement among committee members are chosen for manual labeling.

In our work, we experiment with query-by-committee-based approaches as well as hybrid approaches in which we employ a weighted combination of multiple learners. We also explore the effect of using different sets of committee learners. We carry out a number of experiments to compare the selection strategies that we propose. In this section, we first describe the corpus used in the experiments

and then present the statistical classifiers used in the committees, along with the scoring metrics that we use.

3.1 The MALACH Information Extraction Corpus

The MALACH collection [19, 20] contains 116,000 hours of digitized interviews and testimonies in 32 languages from 52,000 survivors, liberators, rescuers and witnesses of the Holocaust. We are interested in automatic information extraction from this corpus of spontaneous conversational speech. For a small number of English testimonies, we manually transcribed them and manually annotated the transcripts with three types of information:

- **Mentions.** Named, nominal and pronominal mentions of 20 categories of entities.
- **Co-reference.** Sets of mentions that refer to the same real-world entity.
- **Relations.** Sets of relationships between pairs of mentions. For instance, given the sentence “I was in Auschwitz for a year”, the following relation exists: `LocatedAt(I, Auschwitz)`.

For the experiments described in this paper, we chose to focus on mention annotation only. We plan to explore co-reference and relation annotation in the future. For mention annotation, we excluded pronouns; their inclusion leads to inflated numbers since detection of pronominal mentions is easy and they occur with high frequency.

We split the annotated data into two parts: the pool from which the learners select the next batch of data for annotation, and the development set. The first pool consists of 99 documents, including 4772 sentences and 198K words. A total of 43K mentions have been annotated for the documents in the pool (16K mentions if pronouns are excluded). The development evaluation set consists of 5 complete testimonies, which cover over 10 hours of speech. It consists of 1700 sentences, 73K words and 16K total mentions (6K mentions if pronouns are excluded).

3.2 Granularity of Selection

When using active learning techniques to select samples for human annotation, we need to first decide on the granularity of sample selection. The samples can be documents, sentences, or tokens. A possible problem with selecting samples at the document level [14] is that a document may only be partially useful from a learning point of view and it is impossible to add only the interesting examples so as to maximize the effect of active learning. This problem is particularly acute in our corpus since the documents, which are testimonies of survivors, are very long and contain a lot of redundant information. Selecting tokens as samples, as in [15], has the advantage that all samples are of equal length but suffers from the problem that a human annotator has to annotate mentions by looking at only the tokens. The selected tokens may also contain

partial mentions when mention boundaries are incorrectly identified. Sentences, on the other hand, contain enough non-redundant contextual information for effective annotation. Therefore, we use sentence-level blocks for active learning, as in [16].

3.3 Framework Description

We propose a new active learning based framework for sample selection which provides considerable savings in the annotation task and can provide better performance for the same amount of annotation. Figure 1 shows the architecture of the framework. As mentioned before, the central idea is to use an ensemble of classifiers -each one differing from the others in some respect. Once trained, the classifiers can then be used to detect mentions in the unlabeled sentences. These detected mentions can then be compared, for each sentence, and a score assigned to this sentence indicating the agreement, or lack thereof, among the ensemble classifiers. The sentences on which the classifiers disagree the most are, intuitively, the ones that can contribute the most to the learning. The sentences are then ranked according to this metric and the required number sampled from this ranked list.

In the case of true active learning, these samples would be provided to the annotator to label and then added to the training data for the main classifier used for the actual task of mention detection. However, for this paper, we consider controlled experiments where the entire dataset has been annotated in advance. Therefore, the path between the selection of the samples to the annotator has effectively been short-circuited and the samples are added directly to the training data for the main classifier. We structured our setup in such a way since we wanted to experiment with adding successively increasing amounts of training data without incurring the overhead of annotation at each step. The same exact advantages and savings will apply to a real setting where the selected samples need to be annotated at each step of learning. At each step, the set of sampled sentences is split into equal parts and added to the seed data for the ensemble classifiers as well.

Our approach differs from previous approaches to the same problem in two significant ways. First, we measure the disagreement at the granularity of sentences rather than documents, which allows us to be much more discerning when selecting samples. The second point is that our framework allows targeting the sample selection towards specific types of mention that can be specified by the user. This is important because for some tasks, like ACE, some types of mentions carry more weight than the others.

3.4 Statistical classifiers

At each step in the active learning process, we build **two** maximum-entropy based statistical classifiers [21] using the labeled data available at that step.

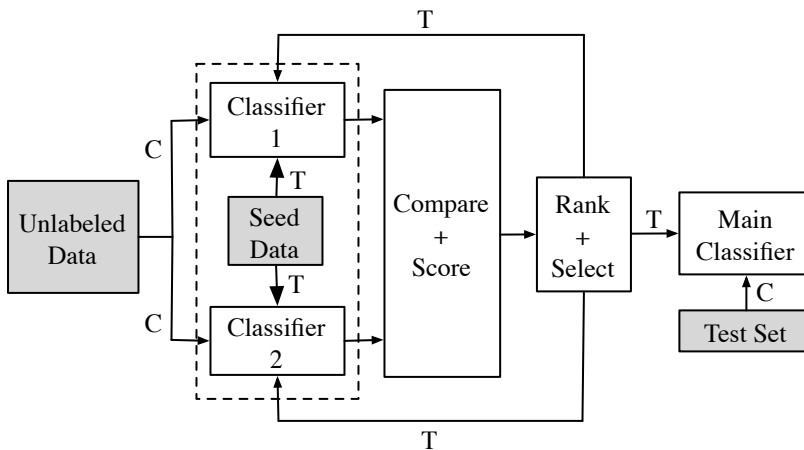


Fig. 1. The architecture of the active learning framework. The edges marked with a **T** represent training steps and the ones marked with a **C** represent classification.

There are two dimensions of classifier training - the feature set of the classifier and the data used for training. For our experiments, we use two different combinations of these dimensions:

- Each classifier is trained with the *same* features but on a *different* half of the available labeled data. We refer to this as the *data-different* (DD) setting.
- Each classifier is trained on the entire available labeled data but using a *different* feature set : the “inside” classifier uses lexical features derived solely from the current token, and the “outside” classifier uses features involving surrounding tokens only. We refer to this experimental setting as *feature-different* (FD).

3.5 Sentence scoring metrics

We employ a variety of metrics for measuring the degree of disagreement between two classifiers over the output labels for a given sentence. The simplest metrics look at only the output labels predicted by the classifiers. However, we propose another set of metrics that utilizes the confidence values associated with those labels as well.

- **F-measure** As the harmonic mean of precision and recall, the F-measure is often used to assess the agreement between two classifiers for the named entity recognition task. The value of F-measure is between 0 and 1, with higher values indicating greater agreement. During sentence selection, we compute the F-measure of the two classifiers on each unannotated sentence, and select sentences with the lowest F-measure values for annotation.

- **Macro-averaged F-measure** Instead of computing the F-measure over the entire set of mentions, as in the previous metric, another option is to compute the F-measure for each mention category and then take the average over all categories. This metric allows categories with a small number of mentions to be weighted equally.
- **Confidence Sum.** Our statistical classifiers can provide a probability value for each output label, indicating its confidence in choosing the label. To leverage this information, we use the normalized sum of the confidence values of the two classifiers as another metric. The higher this value, the more confident both classifiers are about a particular sentence. Therefore, that sentence will provide little or no information if added to the training set and should be ranked lower for selection.
- **Confidence Difference.** We also use the absolute difference of the two confidence values. A sentence with a higher confidence difference value indicates an explicit disagreement between the two classifiers when detecting mentions and, therefore, would prove to be more useful for annotation.

In addition to the above selection metrics, there are three additional parameters that were used in our experiments:

- **Minimum number of mentions per sentence.** This parameter is used to deal with the sparse mention problem in sentence selection. For instance, for a given sentence $S1$, if the first classifier finds only one mention in it, and the second classifier finds zero mentions, then the F-measure for the sentence is 0, which puts the given sentence at the top of the selection list. In contrast, for another sentence $S2$, if the first classifier finds 5 mentions, and the second classifier finds 3 mentions, two of which overlap with the output by the first classifier, then the F-measure is 50%. Therefore, $S2$ is erroneously ranked lower than $S1$ on the selection list. To eliminate sentences with very few mentions, a user can set the minimal number of mentions per sentence parameter to N , where N is a non-negative number. Any sentence with less than N mentions is not considered as a candidate for active learning.
- **Mention category weights.** Each mention category is associated with a weight in the above metrics. This gives us the flexibility to focus on certain categories during active learning. For instance, if a system performs weakly in the ORGANIZATION category, we may want to customize the scoring metric so that the active learner can pick up samples that are particularly useful in improving the performance in this category. The categories with higher weights are considered more important; the categories with zero weights are not considered for active learning.
- **Mention level weights.** The mention level indicates whether a mention is named, nominal, or pronominal. We can customize our learner to focus on a particular mention level by adjusting the weight associated with it.

4 Results

We carried out a number of experiments to evaluate the different sentence selection strategies presented in the previous section. All the strategies perform better than random selection, but the confidence sum metric with the feature-different classifier training setting gives the greatest improvement.

Ideally, the active learners should be retrained each time a new sample has been selected by active learning and annotated by a human. In reality, this is rarely done due to time and computational cost. Instead, active learners usually operate in “batch” mode — a set of samples, instead of a single sample, is selected at each step and the active learners are retrained by adding all the samples in the set. Given our granularity of selection, one way to define the size of this batch can be in term of number of sentences to be added at each step. However, this approach does not provide a way to differentiate between a sentence that is 100 words long and another that is only 10 words long; both are equally important. Therefore, we define the batch in terms of the number of words contained in the sentences. At each step, we add just enough complete sentences that, among themselves, contain the given number of words or as close to that number as possible. In the experiments described in this section, this size is 20000 words. First, a set of 20000 words is randomly selected to build seed models. Then at each step of active learning, an additional 20000 words are selected. We also describe the effect of this *step size* at the end of this section. All the results presented in this section are on our development evaluation set — after each step of active learning, we train a model using all the data that have been selected as of that step (including the seed data), and the model’s performance on the development set is reported.

4.1 Effect of different scoring metrics

In the experiments described here, we consider only named and nominal mentions. The weights for all mention categories are set to 1.

Random Selection. We established a baseline performance by selecting sentences at random from the unlabeled data pool. We performed 5 runs and averaged the numbers from each run. The performance for the random selection strategy is shown in Figure 2, including each run and their average. The average performance is used as the baseline and shown in later graphs. The X-axis represents the data size in number of words, and the Y-axis presents the F-measure on the development set.

F-measure and Macro-Averaged F-measure. Figure 3 shows the curves for the F-measure and Macro-Averaged F-measure metrics compared to the baseline. As we can see, the F-measure metric works very well. Even at 40K words (after only one step of active learning), we observe an increase of about 2 absolute F-measure points with the same amount of labeled data, compared to random sentence selection. To attain the same performance, the random strategy needs 1.5 times as much labeled data. The selection strategy based on macro-averaged F-measure also performs better than random selection, but the improvement is

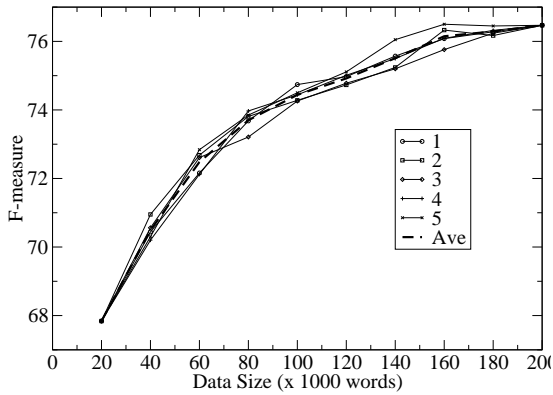


Fig. 2. Performance for the baseline strategy of selecting sentences at random.

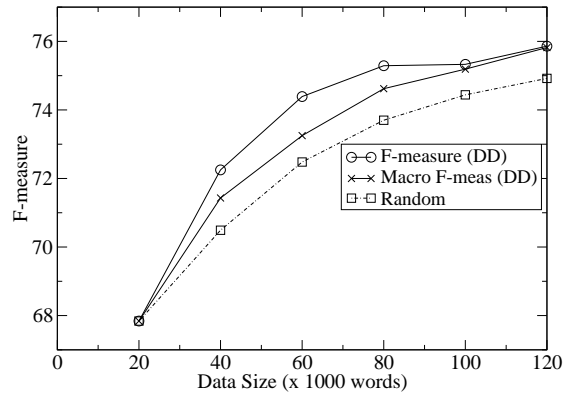


Fig. 3. Performance for the F-measure and Macro-Averaged F-measure metrics.

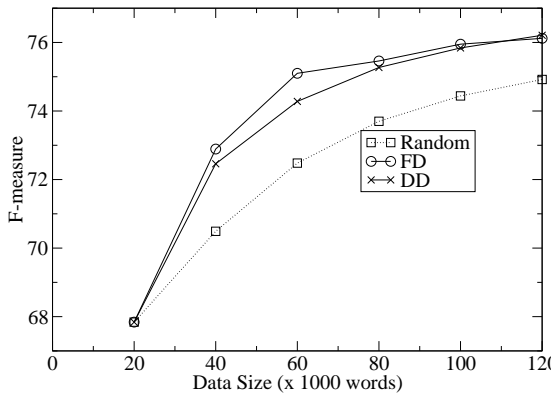


Fig. 4. Performance for the Confidence Sum metric, computed for both classifier training settings.

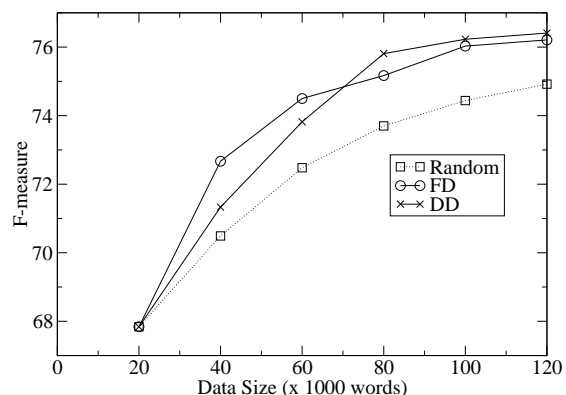


Fig. 5. Performance for the Confidence Difference metric.

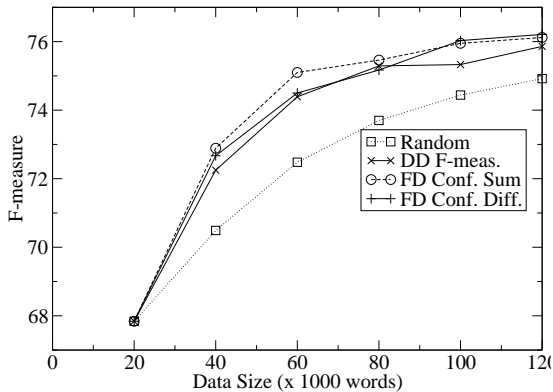


Fig. 6. Comparing the 3 best selection strategies.

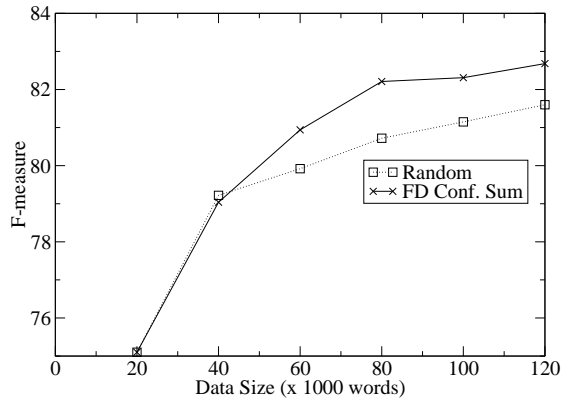


Fig. 7. Performance of the *PERSON* category.

smaller than that for the F-measure metric. As indicated, these metrics were computed only for the data-different classifier setting.

Confidence Sum. Figure 4 shows the curves for the confidence sum metric. This metric is computed for both the data-different and the feature-different setting. As shown in the figure, the classifiers trained on different features (FD) lead to slightly higher improvement in the first two steps than the classifiers with the same feature set (DD). Overall, the confidence sum performs much better than random selection: at 60K words (after only two steps of active learning), the confidence sum strategy attains better performance than the baseline does with twice this amount of labeled data.

Confidence Difference. Figure 5 shows the results for the confidence difference metric. Similar to the confidence sum metric, using the classifiers trained on different features (FD) leads to better improvements in the first two steps. At 60K words, the system achieves the same performance as the random strategy at 100K words.

The Best 3 Strategies. Figure 6 compares the curves for the 3 best strategies - FD confidence sum, FD confidence difference, and DD F-measure. The FD confidence sum strategy outperforms the other two strategies in the first three active learning steps. At the fourth step, the FD confidence difference strategy catches up. At 60K words, the FD confidence sum strategy achieves the performance of the random strategy at about 130K words. This indicates a reduction of over 50% in the amount of annotated data and demonstrates the effectiveness of active learning for making annotation more efficient.

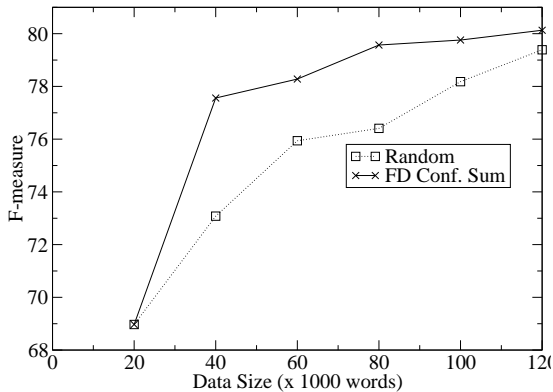


Fig. 8. Targeting named mentions.

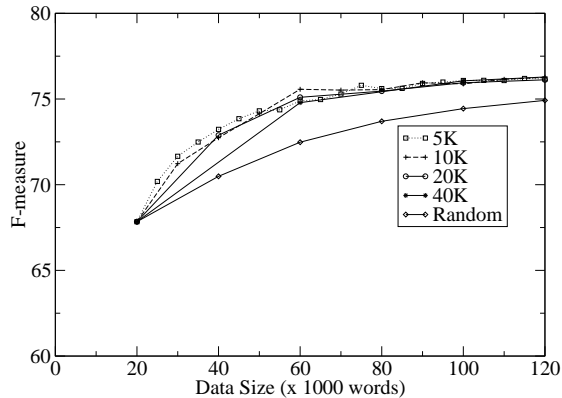


Fig. 9. Effect of step size on confidence sum metric.

4.2 Focusing on specific types of mentions

As we have described in the previous section, our system allows us to weight the different types and levels of mentions differently in order to focus our learning on specific mention types and levels. We provide the results of two such experiments: one focusing on the PERSON category and one focusing on named mentions.

First, we target the learning to mentions of type PERSON. To achieve this, we weight the PERSON category twice as much as any of the other mention categories. The results are shown in Figure 7. Note that the F-measure along the y-axis in the graph is the performance for the PERSON category.

The second experiment was for named mentions. For this experiment, we set the weights to all other levels to 0. The result is an active learner specifically tuned to maximizing the performance of the system on named mentions. Figure 8 shows the results for this experiment. This plot clearly illustrates the advantages of active learning. We see a 5-point F-measure increase with only one step of learning, which could otherwise have only been achieved by annotating almost 2.5 times as much data.

4.3 Effect of step size

The size of the step used in the learning needs to be optimized. On the one hand, getting higher performance gains from smaller amounts of labeled data to be added at each step represents a huge savings in the annotation task; on the other hand, retraining active learners at each step takes time and resources. We need to balance the cost of retraining active learners and the gains from smaller batches.

We experimented with different step sizes for the feature-different (FD) confidence-sum metric. The results are shown in Figure 9. From that graph, it seems that a step size of 10000 words might have provided the best balance between annotation and performance. Step sizes higher than 20000 do not seem to have the right granularity for effective learning.

5 Conclusions

We conducted several active learning experiments for the task of detecting mentions of entities in human transcripts of spontaneous conversational speech. Specifically, we proposed and compared a variety of sentence selection strategies for active learning. The best strategy uses the sum of the confidence values of a pair of classifiers trained with different feature sets. Compared to random sentence selection, this strategy required 50% of the data to achieve the same performance. For named mentions, the strategy required only 42% of the data needed by random selection to achieve the same performance.

In the future, we would like to test our active learning strategies on data from a different domain and data in languages other than English, such as the ACE mention annotation corpus. We would also like to explore active learning for co-reference and relation annotation.

Acknowledgment

This project is part of an on-going effort funded by NSF under the Information Technology Research (ITR) program, NSF IIS Award No. 0122466. Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the NSF.

References

1. ACE: Automatic Content Extraction. <http://www.nist.gov/speech/tests/ace/> (2005)
2. Lewis, D., Catlett, J.: Heterogeneous uncertainty sampling for supervised learning. In: Proceedings of ICML 1994, 11th International Conference on Machine Learning, New Brunswick, NJ (1994) 148–156
3. McCallum, A.K., Nigam, K.: Employing EM in pool-based active learning for text classification. In: Proceedings of ICML 1998, 15th International Conference on Machine Learning, Madison, US (1998) 350–358
4. Schohn, G., Cohn, D.: Less is more: Active learning with support vector machines. In: Proceedings of ICML 2000, 17th International Conf. on Machine Learning. (2000) 839–846
5. Tong, S., Koller, D.: Support vector machine active learning with applications to text classification. *Journal of Machine Learning Research* **2** (2001) 45–66
6. Argamon-Engelson, S., Dagan, I.: Committee-based sample selection for probabilistic classifiers. *Journal of Artificial Intelligence Research* **11** (1999) 335–360

7. Thompson, C.A., Califf, M.E., Mooney, R.J.: Active learning for natural language parsing and information extraction. In: Proceedings of ICML 1999, 16th International Conference on Machine Learning. (1999) 406–414
8. Tang, M., Luo, X., Roukos, S.: Active learning for statistical natural language parsing. In: Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, Philadelphia, PA, USA (2002) 120–127
9. Steedman, M., Hwa, R., Clark, S., Osborne, M., Sarkar, A., Hockenmacier, J., Ruhlen, P., Baker, S., Crim, J.: Example selection for bootstrapping statistical parsers. In: Proceedings of HLT-NAACL 2003, Edmonton, Canada (2003) 236–243
10. Osborne, M., Baldrige, J.: Ensemble-based active learning for parse selection. In: Proceedings of HLT-NAACL 2004, Boston, Massachusetts, USA (2004) 89–96
11. Ngai, G., Yarowsky, D.: Rule writing or annotation: Cost-efficient resource usage for base noun phrase chunking. In: Proceedings of 38th Annual Meeting of the Association for Computational Linguistics, Hong Kong (2000)
12. Sassano, M.: An empirical study of active learning with support vector machines for japanese word segmentation. In: Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, Philadelphia, PA, USA (2002) 505–512
13. Banko, M., Brill, E.: Scaling to very very large corpora for natural language disambiguation. In: Proceedings of 39th Annual Meeting of the Association for Computational Linguistics, Toulouse, France (2001) 26–33
14. Finn, A., Kushmerick, N.: Active learning selection strategies for information extraction. In: ECML-03 Workshop on Adaptive Text Extraction and Mining, Croatia (2003)
15. Shen, D., Zhang, J., Su, J., Zhou, G., Tan, C.L.: Multi-criteria-based active learning for named entity recognition. In: Proceedings of the 42nd Meeting of the Association for Computational Linguistics (ACL'04), Main Volume, Barcelona, Spain (2004) 589–596
16. Hachey, B., Alex, B., Becker, M.: Investigating the effects of selective sampling on the annotation task. In: Proceedings of the 9th Conference on Computational Natural Language Learning (CoNLL), Ann Arbor (2005) 144–151
17. Cohn, D.A., Ghahramani, Z., Jordan, M.I.: Active learning with statistical models. *Journal of Artificial Intelligence Research* **4** (1995) 129–145
18. Seung, H.S., Opper, M., Sompolinsky, H.: Query by committee. In: *Computational Learning Theory*. (1992) 287–294
19. Gustman, S., Oard, D.S.D., Byrne, W., Picheny, M., Ramabhadran, B., Greenberg, D.: Supporting access to large digital oral history archives. In: Proceedings of the Joint Conference on Digital Libraries. (2002) 18–27
20. Oard, D., Soergel, D., Doermann, D., Huang, X., Murray, G., Wang, J., Ramabhadran, B., Franz, M., Gustman, S., Mayfield, J., Kharevych, L., Strassel, S.: Building an information retrieval test collection for spontaneous conversational speech. In: Proceedings of SIGIR'04, Sheffield, U.K. (2004)
21. Florian, R., Hassan, H., Ittycheriah, A., Jing, H., Kambhatla, N., Luo, X., Nicolov, N., Roukos, S.: A statistical model for multilingual entity detection and tracking. In: Proceedings of HLT-NAACL 2004, Boston, Massachusetts, USA (2004) 1–8