

LTP: A New Active Learning Strategy for CRF-Based Named Entity Recognition

Mingyi Liu, Zhiying Tu, Tong Zhang, Tonghua Su, Zhongjie Wang

School of Computer Science and Technology, Harbin Institute of Technology, Harbin, China
{liumy, tzy_hit, thsu, rainy}@hit.edu.cn, muyejunzi@163.com

Abstract

In recent years, deep learning has achieved great success in many natural language processing tasks including named entity recognition. The shortcoming is that a large amount of manually-annotated data is usually required. Previous studies have demonstrated that active learning could elaborately reduce the cost of data annotation, but there is still plenty of room for improvement. In real applications we found existing uncertainty-based active learning strategies have two shortcomings. Firstly, these strategies prefer to choose long sequence explicitly or implicitly, which increase the annotation burden of annotators. Secondly, some strategies need to invade the model and modify to generate some additional information for sample selection, which will increase the workload of the developer and increase the training/prediction time of the model. In this paper, we first examine traditional active learning strategies in a specific case of BiLstm-CRF that has widely used in named entity recognition on several typical datasets. Then we propose an uncertainty-based active learning strategy called Lowest Token Probability (LTP) which combines the input and output of CRF to select informative instance. LTP is simple and powerful strategy that does not favor long sequences and does not need to invade the model. We test LTP on multiple datasets, and the experiments show that LTP performs slightly better than traditional strategies with obviously less annotation tokens on both sentence-level accuracy and entity-level F1-score.

Keywords:

active learning, named entity recognition, CRF

1. Introduction

Over the past few years, papers applying deep neural networks (DNNs) to the task of named entity recognition (NER) have achieved noteworthy success [1, 2, 3]. However, under typical training procedures, the advantages of deep learning are established mostly relied on the huge amount of labeled data. When applying these methods on domain-related tasks, their main problem lies in their need for considerable human-annotated training corpus, which requires tedious and expensive work from domain experts. Thus, to make these methods more widely applicable and easier to adapt to various domains, the key is how to reduce the number of manually annotated training samples.

Active learning are designed to reduce the amount of data annotation. Unlike the supervised learning setting, in which samples are selected and annotated at random, the process of active learning employs one or more human annotators by asking them to label new samples

that are supposed to be the most informative in the creation of a new classifier. The greatest challenge in active learning is to determine which sample is more informative. The most common approach is uncertainty sampling, in which the model preferentially selects samples whose current prediction is least confident.

Quite a lot of works have been done to reduce the amount of data annotation for NER tasks through active learning. However, these state-of-the-art approaches mainly face two problems. One of the problems is that they tend to choose the long sequences explicitly or implicitly, which will be an undesirable behavior when someone seeks to maximize performance for minimal cost annotation. Another problem is they may need to invade and modify the original model, which will increase the workload of the developer and increase the computing cost. **In this work**, we try to propose a simple but effective active learning strategy that does not prefer long sequence and does not need to invade origi-

nal model.

When evaluating the effect of NER, most of the works only use the value of entity-level F_1 score. However, in some cases, this could be misleading, especially for languages that do not have a natural separator, such as Chinese. And the NER task is often used to support downstream tasks (such as QA, task-oriented dialogue), which prefer that all entities in the sentence are correctly identified. So **in this work**, we not only evaluate the entity-level F_1 score but also the sentence-level accuracy.

We first experiment with the traditional uncertainty-based active learning algorithms, and then we proposed our own active learning strategy based on the lowest token probability with the best labeling sequence. Experiments show that our selection strategy is superior to traditional uncertainty-based active selection strategies in multiple Chinese datasets and English datasets both in entity-level F_1 score and overall sentence-level accuracy. Finally, we make empirical analysis with different active selection strategies.

The remainder of this paper is organized as follows. In Section 2 we summarize the related works in named entity recognition and active learning. In section 3 we brief introduced the data representation and CRF. Section 4 describes in details the active learning strategies we propose. Section 5 describes the experimental setting, the datasets, and discusses the empirical results. And the last section is the conclusion.

2. Related Work

2.1. Named entity recognition

The framework of NER using deep neural network can be regarded as a composition of encoder and decoder. For encoders, there are many options. Collobert et al.[4] first used convolutional neural network (CNN) as the encoder. Traditional CNN cannot solve the problem of long-distance dependency. In order to solve this problem, RNN[5], BiLSTM[6], Dilated CNN[7] and bidirectional Transformers[8] are proposed to replace CNN as encoder. For decoders, some works used RNN for decoding tags [5, 9]. However, most competitive approaches relied on CRF as decoder[2, 10].

2.2. Active learning

Active learning strategies have been well studied [11, 12], [13]. These strategies can be grouped into following categories: *Uncertainty sample* [14, 15, 16, 17], *query-by-committee*[18, 19], *information density*[20], *fisher information*[?]. There were some works that

compared the performance of different types of selection strategies in NER/sequence labeling tasks with CRF model [21, 22, 23]. These results show that, in most case, uncertainty-based methods perform better and cost less time.

However, we found that these studies are mainly based on English datasets, and don't pay much attention to Chinese datasets. Additionally, traditional uncertainty-based strategies always choose long sequence explicitly or implicitly, which significantly increases the burden on the annotators. And some strategies [24] invade the model and let the model perform additional tasks for sample selection. So, in this work we proposed a new active learning strategy that does not favor long sequences and does not need to invade the model.

3. NER Model

3.1. Data Representation

We represent each input sentence following Bert format; Each token in the sentence is marked with BIO scheme tags. Special $[CLS]$ and $[SEP]$ tokens are added at the beginning and the end of the tag sequence, respectively. $[PAD]$ tokens are added at the end of sequences to make their lengths uniform. The formatted sentence in length N is denoted as $\mathbf{x} = \langle x_1, x_2, \dots, x_N \rangle$, and the corresponding tag sequence is denoted as $\mathbf{y} = \langle y_1, y_2, \dots, y_N \rangle$.

3.2. CRF Layer

CRF are statistical graphical models which have demonstrated state-of-art accuracy on virtually all of the sequence labeling tasks including NER task. Particularly, we use linear-chain CRF that is a popular choice for tag decoder, adopted by most DNNs for NER.

A linear-chain CRF model defines the posterior probability of \mathbf{y} given \mathbf{x} to be:

$$P(\mathbf{y}|\mathbf{x}; A) = \frac{1}{Z(\mathbf{x})} \exp \left(P(y_1; \mathbf{x}_1) + \sum_{k=1}^{n-1} P(y_{k+1}; \mathbf{x}_{k+1}) + A_{y_k, y_{k+1}} \right) \quad (1)$$

where $Z(\mathbf{x})$ is a normalization factor over all possible tags of \mathbf{x} , and $P(y_k; \mathbf{x}_k)$ indicates the probability of taking the y_k tag at position k which is the output of the previous DNN layer, such as bilstm, softmax. A is a parameter called a transfer matrix, which can be set manually or by model learning. In our experiment, we let the model learn the parameter by itself. $A_{y_k, y_{k+1}}$ means the probability of a transition from tag states y_k to y_{k+1} . We use \mathbf{y}^* to represent the most likely tag sequence of \mathbf{x} :

$$\mathbf{y}^* = \arg \max_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}) \quad (2)$$

Table 1: Example of data representation. [PAD] tag are not shown.

Sentence	Trump	was	born	in	the	United	States
Tag	[CLS]	B-PER	O	O	O	B-LOC	I-LOC I-LOC [SEP]

The parameters A are learnt through the maximum log-likelihood estimation, that is to maximize the log-likelihood function ℓ of training set sequences in the labeled data set \mathcal{L} :

$$\ell(\mathcal{L}; A) = \sum_{l=1}^L \log P(\mathbf{y}^{(l)} | \mathbf{x}^{(l)}; A) \quad (3)$$

where L is the size of the tagged set \mathcal{L} .

4. Active Learning Strategies

The biggest challenge in active learning is how to select instances that need to be manually annotated. A good selection strategy $\phi(\mathbf{x})$, which is a function used to evaluate each instance \mathbf{x} in the unlabeled pool \mathcal{U} , will select the most informative instance \mathbf{x} .

Algorithm 1 Pool-based active learning framework

Require: Labeled data set \mathcal{L} ,
unlabeled data pool \mathcal{U} ,
selection strategy $\phi(\cdot)$,
query batch size B

while not reach stop condition **do**
// Train the model using labeled set \mathcal{L}
 $train(\mathcal{L})$;
for $b = 1$ to B **do**
//select the most informative instance
 $\mathbf{x}^* = \arg \max_{\mathbf{x} \in \mathcal{U}} \phi(\mathbf{x})$
 $\mathcal{L} = \mathcal{L} \cup \langle \mathbf{x}^*, label(\mathbf{x}^*) \rangle$
 $\mathcal{U} = \mathcal{U} - \mathbf{x}^*$
end for
end while

Algorithm 1 illustrate the entire pool-based active learning process. In the remainder of this section, we describe various query strategy formulations of $\phi(\cdot)$ in detail.

4.1. Token-based (Local) Strategies

The token-based strategy treats the labeling sequence as a set of isolated tokens, and evaluates uncertainty by aggregating the information of these tokens.

Minimum Token Probability (MTP) selects the most informative tokens, regardless of the assignment

performed by CRF. This strategy greedily samples tokens whose highest probability among the labels is lowest:

$$\phi^{MTP}(\mathbf{x}) = 1 - \min_i \max_j P(y_i = j | \mathbf{x}_i; A) \quad (4)$$

where $P(y_i = j)$ is the probability that j is the label at position i in the sequence.

Entropy is a popular measure of informativeness. The entropy of a discrete random variable Y can be represented by $H(Y) = -\sum_i P(y_i) \log P(y_i)$, and means the information needed to "encode" the distribution of outcomes for Y . **Token Entropy (TE)** is a way to use the entropy of model's posteriors over its labeling:

$$\phi^{TE} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M P(y_i = j | \mathbf{x}_i; A) \log P(y_i = j | \mathbf{x}_i; A) \quad (5)$$

where N is the length of \mathbf{x} without [PAD], j ranges over all possible token labels.

Settles [?] argue that querying long sequences should not be explicitly discouraged, if in fact they contain more information. They extend TE into **Maximum Token Entropy (MTE)**:

$$\phi^{MTE}(\mathbf{x}) = N \times \phi^{TE}(\mathbf{x}) \quad (6)$$

4.2. Sentence-based (Global) Strategies

Different from token-based strategies, sentence-based strategies treat labeling sequence \mathbf{y} as whole. Most of these strategies have high complexity or require intrusive models.

Culotta and McCallum [15] employ a simple uncertainty-based strategy for sequence models called least confidence (LC), which sort examples in ascending order according to the probability assigned by the model to the most likely sequence of tags:

$$\phi^{LC}(\mathbf{x}) = 1 - P(\mathbf{y}^* | \mathbf{x}; A) \quad (7)$$

This confidence can be calculated using the posterior probability given by Equation 1. Preliminary analysis revealed that the LC strategy prefer selects longer sentences:

$$P(\mathbf{y}^* | \mathbf{x}; A) \propto \exp \left(P(y_1^*; \mathbf{x}_1) + \sum_{k=1}^{n-1} P(y_{k+1}^*; \mathbf{x}_{k+1}) + A_{y_k^* y_{k+1}^*} \right) \quad (8)$$

Since Equation 8 contains summation over tokens, LC method naturally favors longer sentences. Although the LC method is very simple and has some shortcomings, many works have proved the effectiveness of the method in sequence labeling tasks.

Scheffer et al. [16] propose a method called **Margin**, which queries samples with the smallest margin between the posteriors for its two most likely annotations:

$$\phi^M(\mathbf{x}) = -(P(\mathbf{y}_1^*|\mathbf{x}; A) - P(\mathbf{y}_2^*|\mathbf{x}; A)) \quad (9)$$

where, \mathbf{y}_1^* and \mathbf{y}_2^* are the first and second most likely tag sequence of \mathbf{x} . **Margin** requires the model to calculate the unnecessary second most likely tag sequence.

Different from **TE** and **TTE**, **Sequence Entropy (SE)** considers the entropy of the sequence instead of the entropy of the token:

$$\phi^{SE}(\mathbf{x}) = - \sum_{\hat{\mathbf{y}}} P(\hat{\mathbf{y}}|\mathbf{x}; A) \log P(\hat{\mathbf{y}}|\mathbf{x}; A) \quad (10)$$

where $\hat{\mathbf{y}}$ ranges all over possible tag sequences for \mathbf{x} . This calculation cost will increase exponentially with the length of \mathbf{x} and the number of tag categories.

The most recent uncertainty-based selection strategy is called **Bayesian Active Learning by Disagreement (BALD)**[24, 25]. BALD measures the uncertainty of the sample by observing the changes in the forward propagation result of the sample through multiple random dropouts[26]. Let $\tilde{\mathbf{y}}^1, \tilde{\mathbf{y}}^2, \dots, \tilde{\mathbf{y}}^T$ represent the result from apply T independently sampled dropout masks:

$$\phi^{BALD}(\mathbf{x}) = 1 - \frac{\max_{\tilde{\mathbf{y}}} \text{count}(\tilde{\mathbf{y}})}{T} \quad (11)$$

where $\text{count}(\tilde{\mathbf{y}})$ means the number of occurrences of $\tilde{\mathbf{y}}$ in $\tilde{\mathbf{y}}^1, \tilde{\mathbf{y}}^2, \dots, \tilde{\mathbf{y}}^T$. Normally, the value of T is 100. BALD will cost a lot of time on repeating forward propagation when the data pool is large.

4.3. Lowest Token Probability (LTP)

Unlike existing strategies, we believe that local information and global information have their own advantages, and the two can complement each other. We look for the most probable sequence assignment (global), and hope that each token (local) in the sequence has a high probability.

$$\phi^{LTP}(\mathbf{x}) = 1 - \min_{y_i^* \in \mathbf{y}^*} P(y_i^*|\mathbf{x}_i; A) \quad (12)$$

We proposed our select strategy called **Lowest Token Probability (LTP)**, which selects the tokens whose probability under the most likely tag sequence \mathbf{y}^* is lowest. It is not difficult to infer from the formulation that

LTP utilizes global and local information, and implicitly implements **Margin** but does not require additional calculations¹.

Table 2 compares all the uncertainty-based active learning strategies mentioned in this section. Strategies that do not need to invade the model and do not require additional calculations are selected as the comparison method of our strategies.

5. Experiments

5.1. Datasets

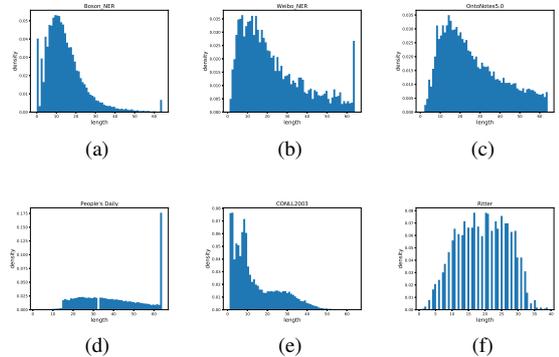


Figure 1: Distribution of sample lengths on different datasets

We have experimented and evaluate the active learning strategies mentioned on Section 4 on four Chinese datasets and two english datasets. *People's Daily* is a collection of newswire article annotated with 3 balanced entity types; *Boson_NER*² is a set of online news annotations published by bosonNLP, which contains 6 entity types; *Weibo_NER*[27, 28] is a collection of short blogs posted on Chinese social media Weibo with 8 extremely unbalanced entity types; *OntoNotes-5.0*[29] Chinese dataset used in this paper is a collection of broadcast news articles, which contains 18 entity types. *CONLL2003*[30] is a well known english dataset consists of Reuters news stories between August 1996 and August 1997, which contains 4 different entity types; *Ritter*[31] is a english dataset consist of tweets annotated with 10 different entity types. All datasets are formatted in the "BIO" sequence representation. In order to be able to perform batch training, the length of all

¹If there is a small probability token in the best sequence, then there is a high probability that the margin between 1st best sequence and 2nd best sequence is small

²https://bosonnlp.com/resources/BosonNLP_NER_6C.zip

Table 2: Qualitative comparison of uncertainty-based active learning strategies

	MTP	LC	TE	TTE	LTP	Margin	SE	BALD
Local(Token) Information	√		√	√	√			
Global(Sentence) Information		√			√	√	√	√
Favor long sequence explicitly		√		√				
Invade model						√	√	
Additional compute						√	√	√

samples is limited to 64. Those samples that were originally longer than 64 will be split according to commas or directly truncated to meet the length requirement.

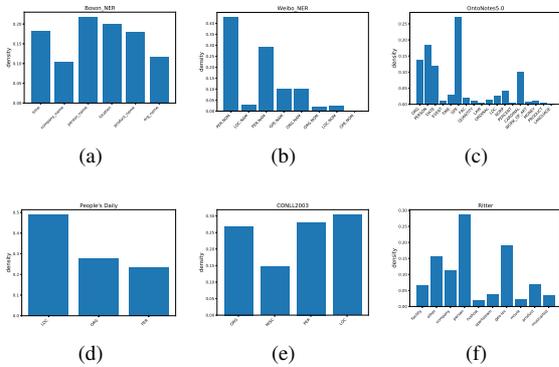


Figure 2: Distribution of entity types on different datasets

Table 3 shows some statistics of the datasets in terms of dimensions, number of entity types, distribution of the labels, etc. Figure 1 gives the distribution of sample lengths on different datasets. Figure 2 presents the distribution of entity types on different datasets. According to the description and statistical information of these datasets, we can conclude that these datasets are 6 datasets with obvious differences in language, text style, entity distribution, length distribution, data magnitude, etc.

5.2. Experimental Setting

For each dataset, we random choose 1% warmstart samples as initial training set \mathcal{L}_1 . We train initial model on this data, then we apply active learning strategy to choose additional 2% samples based on model’s uncertainty estimates and train a new model based on this data. In each iteration, we train from scratch to avoid negative effects accumulated from previous training. We train each model to convergence in each iteration. We fix the number of active learning iterations at 25 because of each algorithm does not improve obviously after 25 iteration.

In the NER model, we use a 300d word embedding pre-trained on the Chinese Wikipedia corpus[32] for the Chinese datasets, and a 100d glove word embedding pre-trained on the English Wikipedia corpus[33] for the English datasets. We uniformly set the learning rate as 0.001 and the training batch size as 64. The transition matrix A in CRF is left to let the model learn by itself. It must be noted that the goal of this article is not to obtain SOTA of NER, but to compare the performance of different active learning strategies under same conditions. So, the NER model itself and its parameters may not be the best but fair.

We empirically compare the selection strategy proposed in Section 4, as well as the uniformly random baseline (**RAND**) and long baseline (**LONG**). We evaluate each selection strategy by constructing learning curves that plot the overall F_1 -score (for entities) and *accuracy* (for sentences). In order to prevent the contingency of experiments, we have done 5 independent experiments for each selection strategy on each dataset using different random initial training set \mathcal{L}_1 . All results reported in this paper are averaged across these experiments.

5.3. Results

Entity-level F_1 -scores are shown in Figure 4, it is clear that all active learning strategies (except **TE**) perform better than the random baseline on 4 Chinese datasets. Our strategy is not weaker than other strategies on all datasets, slightly better than other strategies on *Boson_NER*, *Weibo_NER*, and *CONLL2003*, and significantly surpasses other strategies on *Ritter*.

Figure 5 shows the results of **sentence-level accuracy** on six datasets. The results exceeded our expectations and are very interesting. Firstly, the results confirm that entity-level F_1 -score is sometimes misleading (two social media datasets, *Weibo_NER* and *Ritter*) as what we mention in Section 1. Secondly, our strategy **LTP** is better than the rest of methods, while it not obvious on the large data set of canonical text, which is similar to text for pre-trained word embedding.

Table 3: Training(Testing) Data Statistics. #S is the number of total sentences in the dataset, #T is the number of tokens in the dataset, #E is the number of entity types, ASL is the average length of a sentence, ASE is the average number of entities in a sentence, AEL is the average length of an entity, %PT is the percentage of tokens with positive label,%AC is the percentage of a sentences with more than one entity, %DAC is the percentage of sentences that have two or more entities. English datasets are marked in bold.

corpus	#S	#T	#E	ASL	ASE	AEL	%PT	%AC	%DAC
Boson_NER	27350 (6825)	409830 (99616)	6 (6)	14.98 (14.59)	0.67 (0.67)	3.93 (3.87)	17.7% (17.8%)	41.8% (41.8%)	14.7% (14.8%)
Weibo_NER	3664 (591)	85571 (13810)	8 (8)	23.35 (23.36)	0.62 (0.66)	2.60 (2.60)	6.9% (7.3%)	33.6% (36.3%)	14.8% (17.7%)
OntoNotes5.0 (bn-zh)	13798 (1710)	362508 (44790)	18 (18)	26.27 (26.19)	1.91 (1.99)	3.14 (3.07)	22.8% (23.4%)	72.5% (75.4%)	48.0% (51.5%)
People’s Daily	50658 (4620)	2169879 (172590)	3 (3)	42.83 (37.35)	1.47 (1.33)	3.23 (3.25)	11.1% (11.6%)	58.3% (54.4%)	35.8% (29.1%)
CONLL2003	13862 (3235)	203442 (51347)	4 (4)	14.67 (15.87)	1.69 (1.83)	1.44 (1.44)	16.7% (16.7%)	79.9% (80.4%)	44.2% (48.8%)
Ritter	1955 (438)	37735 (8733)	10 (10)	19.30 (19.93)	0.62 (0.60)	1.65 (1.62)	5.3% (4.9%)	38.1% (39.2%)	15.3% (15.5%)

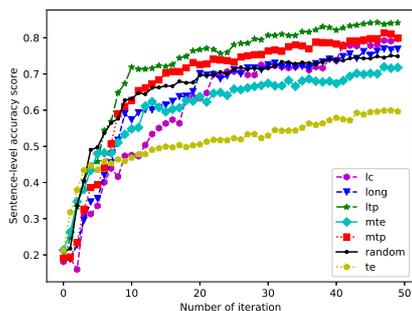


Figure 3: Sentence-level accuracy score results on CONLL2003 with 1% samples selected each iteration.

We know that the most obvious effect of active learning is to select one sample at a time, although this is not realistic due to the cost of retraining. The more samples selected each time, the worse the active learning effect. Therefore, in the case of a large data pool, selecting 2% of the samples in each round cannot clearly reflect the differences between different strategies. In order to clearly reflect the differences between strategies, we constructed an additional experiment on *CONLL2003* with 1% samples selected each iteration. Results are given in Figure 3.

Figure 6 shows **average length** of the samples selected by different active learning strategies. Unlike other active learning strategies, **LTP** does not have a obvious bias towards sample length. From another aspect, **LTP** use less annotation cost to achieve better performance than other strategies.

5.4. Discussion

In this section, we will briefly discuss *possible reasons* for the gap between different selection strategies.

The core of active learning is to select "informative" samples, but there is no unified standard to measure "informative". One thing is certain, the samples that are not correctly labeled by the model are informative samples for the model. Therefore, we use the proportion of samples in each iteration of selection the model is not correctly labeled as the effectiveness of each iteration of selection. Figure 7 shows the results. We can find that **LTP** can more effectively select samples that are incorrectly predicted by the model.

6. Conclusion

We proposed a new active learning strategy for CRF-based named entity recognition. The experiment shows that compared with the traditional active selection strategies, our strategy has better performance, but lower annotation cost.

7. Acknowledgments

Research in this paper is partially supported by the National Key Research and Development Program of China (No 2018YFB1402500), the National Science Foundation of China (61832004, 61772155, 61802089, 61832014).

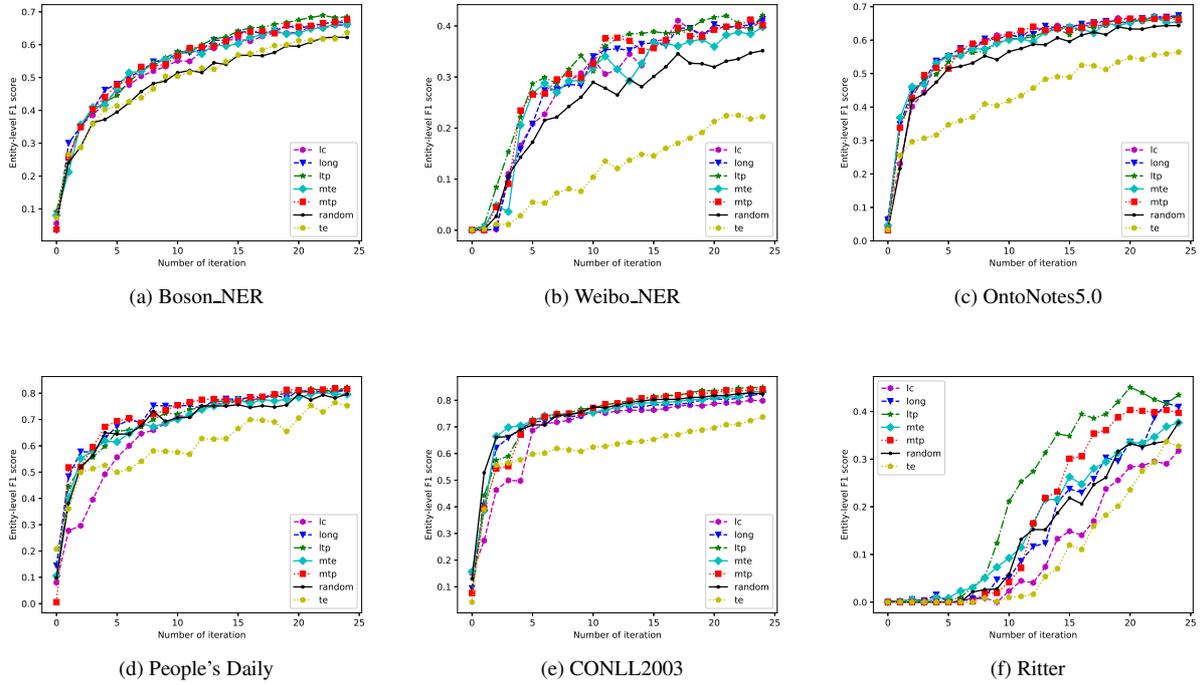


Figure 4: Entity-level F_1 -score results on different datasets

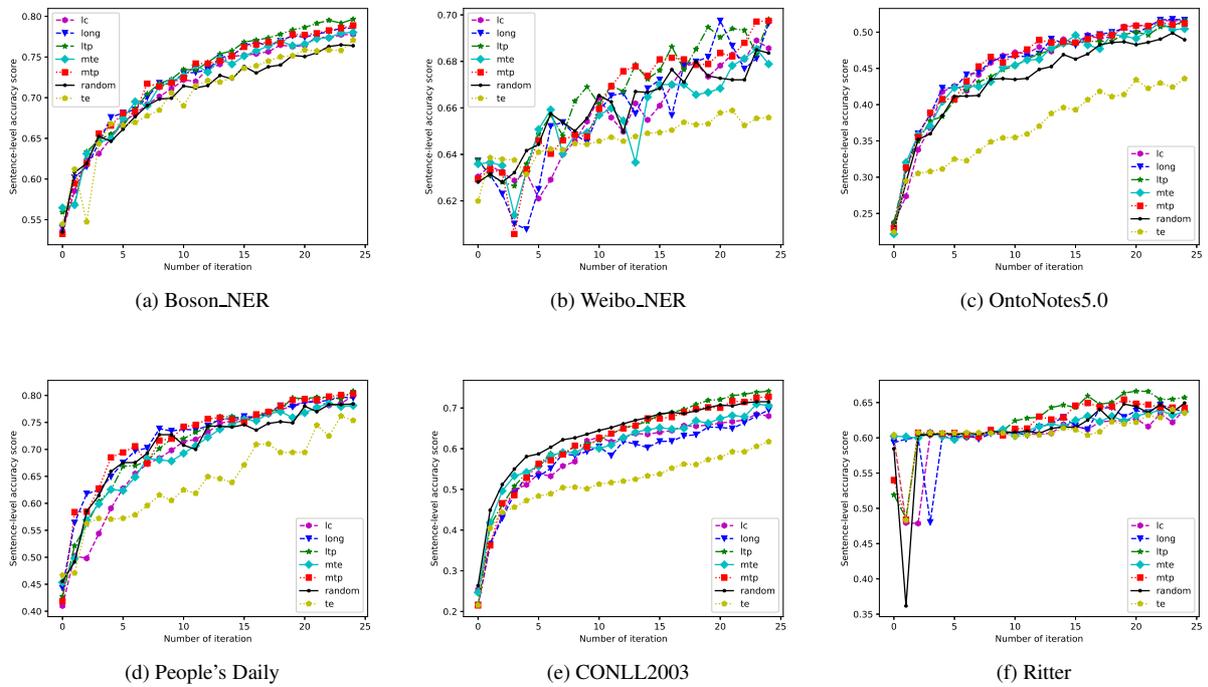


Figure 5: Sentence-level accuracy score results on different datasets

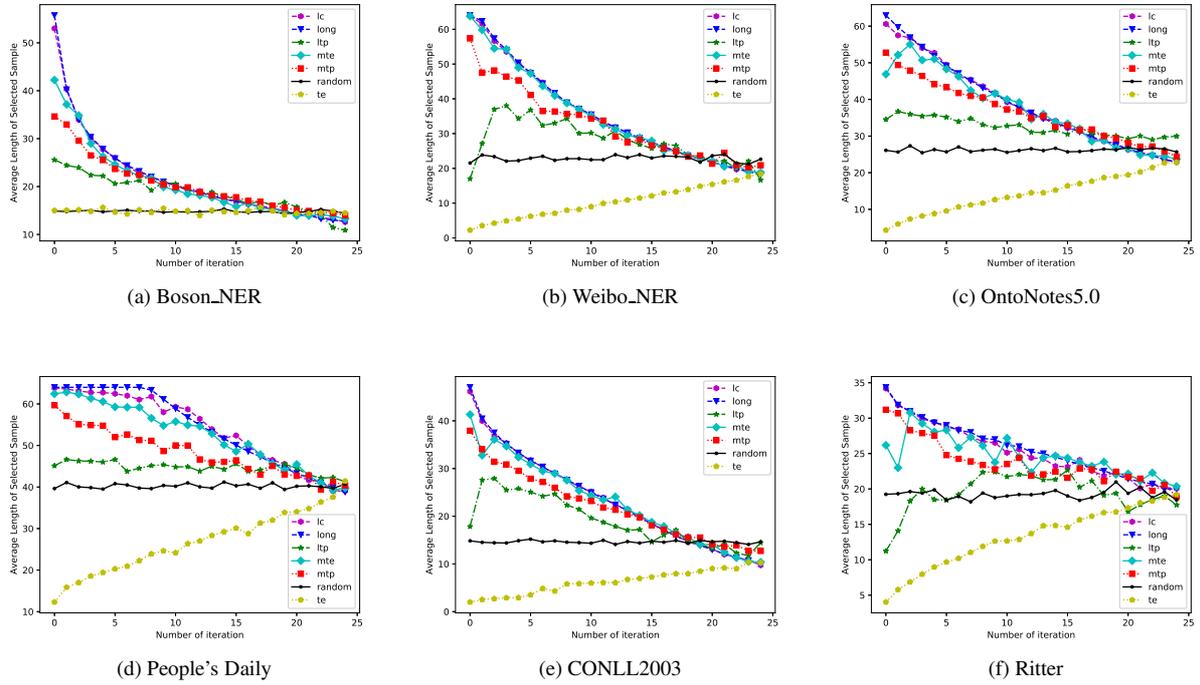


Figure 6: Average length of the samples selected by active learning strategies.

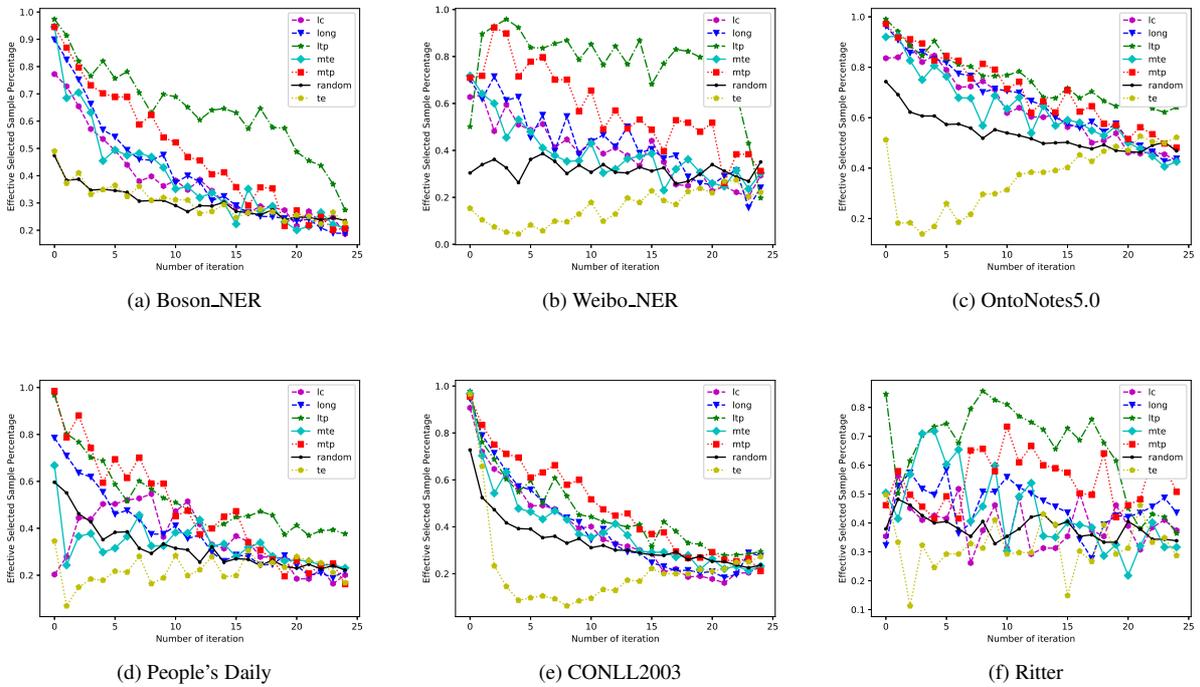


Figure 7: The results of effective selected sample percentage on different datasets

References

- [1] J. P. Chiu, E. Nichols, Named entity recognition with bidirectional lstm-cnns, *Transactions of the Association for Computational Linguistics* 4 (2016) 357–370.
- [2] G. Lample, M. Ballesteros, S. Subramanian, K. Kawakami, C. Dyer, Neural architectures for named entity recognition, in: *Proceedings of NAACL-HLT*, 2016, pp. 260–270.
- [3] N. Limsopatham, N. H. Collier, Bidirectional lstm for named entity recognition in twitter messages (2016).
- [4] R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, P. Kuksa, Natural language processing (almost) from scratch, *Journal of machine learning research* 12 (Aug) (2011) 2493–2537.
- [5] T. H. Nguyen, A. Sil, G. Dinu, R. Florian, Toward mention detection robustness with recurrent neural networks, *arXiv preprint arXiv:1602.07749* (2016).
- [6] Z. Huang, W. Xu, K. Yu, Bidirectional lstm-crf models for sequence tagging, *arXiv preprint arXiv:1508.01991* (2015).
- [7] E. Strubell, P. Verga, D. Belanger, A. McCallum, Fast and accurate entity recognition with iterated dilated convolutions, *arXiv preprint arXiv:1702.02098* (2017).
- [8] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, Bert: Pre-training of deep bidirectional transformers for language understanding, *arXiv preprint arXiv:1810.04805* (2018).
- [9] G. Mesnil, X. He, L. Deng, Y. Bengio, Investigation of recurrent-neural-network architectures and learning methods for spoken language understanding., in: *Interspeech*, 2013, pp. 3771–3775.
- [10] Z. Yang, R. Salakhutdinov, W. Cohen, Multi-task cross-lingual sequence tagging from scratch, *arXiv preprint arXiv:1603.06270* (2016).
- [11] S. Dasgupta, A. T. Kalai, C. Monteleoni, Analysis of perceptron-based active learning, in: *International Conference on Computational Learning Theory*, Springer, 2005, pp. 249–263.
- [12] P. Awasthi, M. F. Balcan, P. M. Long, The power of localization for efficiently learning linear separators with noise, in: *Proceedings of the forty-sixth annual ACM symposium on Theory of computing*, ACM, 2014, pp. 449–458.
- [13] Y. Shen, H. Yun, Z. C. Lipton, Y. Kronrod, A. Anandkumar, Deep active learning for named entity recognition, *arXiv preprint arXiv:1707.05928* (2017).
- [14] D. D. Lewis, J. Catlett, Heterogeneous uncertainty sampling for supervised learning, in: *Machine learning proceedings 1994*, Elsevier, 1994, pp. 148–156.
- [15] A. Culotta, A. McCallum, Reducing labeling effort for structured prediction tasks, in: *AAAI*, Vol. 5, 2005, pp. 746–751.
- [16] T. Scheffer, C. Decomain, S. Wrobel, Active hidden markov models for information extraction, in: *International Symposium on Intelligent Data Analysis*, Springer, 2001, pp. 309–318.
- [17] S. Kim, Y. Song, K. Kim, J.-W. Cha, G. G. Lee, Mmr-based active machine learning for bio named entity recognition, in: *Proceedings of the Human Language Technology Conference of the NAACL, Companion Volume: Short Papers*, Association for Computational Linguistics, 2006, pp. 69–72.
- [18] H. S. Seung, M. Opper, H. Sompolinsky, Query by committee, in: *Proceedings of the fifth annual workshop on Computational learning theory*, ACM, 1992, pp. 287–294.
- [19] J. Vandoni, E. Aldea, S. Le Hégarat-Mascle, Evidential query-by-committee active learning for pedestrian detection in high-density crowds, *International Journal of Approximate Reasoning* 104 (2019) 166–184.
- [20] K. Wei, R. Iyer, J. Bilmes, Submodularity in data subset selection and active learning, in: *International Conference on Machine Learning*, 2015, pp. 1954–1963.
- [21] Y. Chen, T. A. Lasko, Q. Mei, J. C. Denny, H. Xu, A study of active learning methods for named entity recognition in clinical text, *Journal of Biomedical Informatics* 58 (2015) 11 – 18. doi:<https://doi.org/10.1016/j.jbi.2015.09.010>. URL <http://www.sciencedirect.com/science/article/pii/S1532046415000111>.
- [22] D. Marcheggiani, T. Artières, An experimental comparison of active learning strategies for partially labeled sequences, in: *EMNLP*, 2014.
- [23] V. Claveau, E. Kijak, Strategies to select examples for active learning with conditional random fields, in: A. Gelbukh (Ed.), *Computational Linguistics and Intelligent Text Processing*, Springer International Publishing, Cham, 2018, pp. 30–43.
- [24] A. Siddhant, Z. C. Lipton, Deep bayesian active learning for natural language processing: Results of a large-scale empirical study, in: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2018, pp. 2904–2909.
- [25] Y. Gal, R. Islam, Z. Ghahramani, Deep bayesian active learning with image data, in: *International Conference on Machine Learning*, 2017, pp. 1183–1192.
- [26] Y. Gal, Z. Ghahramani, A theoretically grounded application of dropout in recurrent neural networks, in: *Advances in neural information processing systems*, 2016, pp. 1019–1027.
- [27] N. Peng, M. Dredze, Named entity recognition for chinese social media with jointly trained embeddings, in: *Processings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2015, pp. 548–554.
- [28] N. Peng, M. Dredze, Improving named entity recognition for chinese social media with word segmentation representation learning, in: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (ACL)*, Vol. 2, 2016, pp. 149–155.
- [29] R. Weischedel, S. Pradhan, L. Ramshaw, J. Kaufman, M. Franchini, M. El-Bachouti, N. Xue, M. Palmer, J. D. Hwang, C. Bohnial, et al., *Ontonotes release 5.0* (2012).
- [30] E. F. Sang, F. De Meulder, Introduction to the conll-2003 shared task: Language-independent named entity recognition, *arXiv preprint cs/0306050* (2003).
- [31] A. Ritter, S. Clark, Mausam, O. Etzioni, Named entity recognition in tweets: An experimental study, in: *EMNLP*, 2011.
- [32] S. Li, Z. Zhao, R. Hu, W. Li, T. Liu, X. Du, Analogical reasoning on chinese morphological and semantic relations, in: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, Association for Computational Linguistics, 2018, pp. 138–143. URL <http://aclweb.org/anthology/P18-2023>
- [33] J. Pennington, R. Socher, C. D. Manning, Glove: Global vectors for word representation, in: *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, 2014, pp. 1532–1543.