

Overview of the NLPCC 2015 Shared Task: Chinese Word Segmentation and POS Tagging for Micro-blog Texts

Xipeng Qiu^(✉), Peng Qian, Liusong Yin, Shiyu Wu, and Xuanjing Huang

School of Computer Science, Fudan University,
825 Zhangheng Road, Shanghai, China
{xpqiu,pqian11,lsyin14,sywu13,xjhuang}@fudan.edu.cn

Abstract. In this paper, we give an overview for the shared task at the 4th CCF Conference on Natural Language Processing & Chinese Computing (NLPCC 2015): Chinese word segmentation and part-of-speech (POS) tagging for micro-blog texts. Different with the popular used newswire datasets, the dataset of this shared task consists of the relatively informal micro-texts. The shared task has two sub-tasks: (1) individual Chinese word segmentation and (2) joint Chinese word segmentation and POS Tagging. Each subtask has three tracks to distinguish the systems with different resources. We first introduce the dataset and task, then we characterize the different approaches of the participating systems, report the test results, and provide an overview analysis of these results. An online system is available for open registration and evaluation at <http://nlp.fudan.edu.cn/nlpcc2015>.

1 Introduction

Word segmentation and Part-of-Speech (POS) tagging are two fundamental tasks for Chinese language processing. Benefiting from the developments of the machine learning techniques and the large scale shared corpora, such as Chinese Treebank [9], Chinese word segmentation and POS tagging have achieved a great progress. The state-of-the-art method is to regard these two tasks as sequence labeling problem [6, 8], which can be handled with supervised learning algorithms such as Maximum Entropy (ME) [1], averaged perceptron [2], Conditional Random Fields (CRF)[4]. However, their performances are still not satisfying for the practical demands to analyze Chinese texts, especially for informal texts. The key reason is that most of annotated corpora are drawn from news texts. Therefore, the system trained on these corpora cannot work well with the informal or specific-domain texts.

In this shared task, we focus to evaluate the performances of word segmentation and POS tagging on relatively informal micro-texts.

2 Data

Different with the popular used newswire dataset, we use relatively informal texts from Sina Weibo¹. The training and test data consist of micro-blogs from various topics, such as finance, sports, entertainment, and so on. Both the training and test files are UTF-8 encoded.

The information of dataset is shown in Table 1. The out-of-vocabulary (OOV) rate is slight higher than the other benchmark datasets. For example, the OOV rate is 5.58% in the popular division [10] of the Chinese Treebank (CTB 6.0) dataset [9], while the OOV rate of our dataset is 7.25%.

Table 1. Statistical information of dataset.

Dataset	Sents	Words	Chars	Word Types	Char Types	OOV Rate
Training	10,000	215,027	347,984	28,208	39,71	-
Test	5,000	106,327	171,652	18,696	3,538	7.25%
Total	15,000	322,410	520,555	35,277	4,243	-

There are total 35 POS tags in this dataset. A detailed list of POS tags is shown in Table 2.

Table 2. Statistical information of POS tags.

词性(POS)	Labels	Occurrences
名词	NN	84,006
实体名	人名	PER 3,232
	机构名	ORG 2,578
	地名	LOC 9,701
	其他	NR 550
	邮件	EML 3
	型号名	MOD 34
	网址	URL 11
副词	疑问副词	ADQ 340
	副词	AD 26,155
形貌	形容词	JJ 9,477
	形谓词	VA 3,339
动词	动词	VV 51,294
	情态词	MV 3,700
	趋向动词	DV 781
	被动词	BEI 927
	把动词	BA 600
时间短语	NT	5,881

词性(POS)	Labels	Occurrences
代词	人称代词	PNP 4,903
	疑问代词	PNQ 492
	指示代词	PNI 834
连词	并列连词	CC 2,725
	从属连词	CS 866
数量	数词	CD 10,764
	量词	M 7,917
	序数词	OD 1,219
助词	方位词	LC 4,725
	省略词	ETC 673
	语气词	SP 1,076
	限定词	DT 3,579
	叹词	IJ 20
	标点	PU 52,922
	结构助词	DSP 13,756
	介词	P 9,488
时态词	AS 3,382	

¹ <http://weibo.com/>

2.1 Background Data

Besides the training data, we also provide the background data, from which the training and test data are drawn. The purpose of providing the background data is to find the more sophisticated features by the unsupervised way.

3 Description of the Task

In recent years, word segmentation and POS tagging have undergone great development. In this shared task, we wish to investigate the performances of Chinese word segmentation and POS tagging for the micro-blog texts.

3.1 Subtasks

This task focus the two fundamental problems of Chinese language processing: word segmentation and POS tagging, which can be divided into two subtasks:

1. **CWS** subtask: The first subtask is Chinese word segmentation (CWS). Word is the fundamental unit in natural language understanding. However, Chinese sentences consists of the continuous Chinese characters without natural delimiters. Therefore, Chinese word segmentation has become the first mission of Chinese natural language processing, which identifies the sequence of words in a sentence and marks the boundaries between words.
2. **S&T** subtask: The second subtask is joint Chinese word segmentation and POS tagging.

3.2 Tracks

Each participant will be allowed to submit the three runs for each subtask: **closed track** run, **semi-open track** run and **open track** run.

1. In the **closed** track, participants could only use information found in the provided training data. Information such as externally obtained word counts, part of speech information, or name lists was excluded.
2. In the **semi-open** track, participants could use the information extracted from the provided background data in addition to the provided training data. Information such as externally obtained word counts, part of speech information, or name lists was excluded.
3. In the **open** track, participants could use the information which should be public and be easily obtained. But it is not allowed to obtain the result by the manual labeling or crowdsourcing way.

4 Participants

Sixteen teams have registered for this task. Finally, there are 27 qualified submitted results from 10 teams. A summary of qualified participating teams are shown in Table 3.

Table 3. Summary of the participants.

	CWS			S&T		
	closed	open	semi-open	closed	open	semi-open
NJU	✓	✓	✓			
BosonNLP	✓	✓		✓	✓	
CIST	✓		✓	✓		✓
XUPT	✓			✓		
CCNU	✓	✓				
ICT-NLP	✓					
BJTU	✓	✓	✓	✓	✓	✓
SZU		✓			✓	
ZZU			✓			
WHU				✓		✓

5 SubTask 1: Chinese Word Segmentation

The evaluation measures are precision, recall, and an evenly-weighted F1.

5.1 Baseline Systems

Currently, the mainstream method of word segmentation is discriminative character-based sequence labeling. Each character is labeled as one of {B, M, E, S} to indicate the segmentation. {B, M, E} represent *Begin*, *Middle*, *End* of a multi-character segmentation respectively, and S represents a *Single* character segmentation.

For the joint word segmentation and POS tagging, the state-of-the-art method is also based on sequence learning with cross-labels, which can avoid the problem of error propagation and achieve higher performance on both subtasks[5]. Each label is the cross-product of a segmentation label and a tagging label, e.g. {B-NN, I-NN, E-NN, S-NN, ...}. The features are generated by position-based templates on character-level.

Sequence labeling is the task of assigning labels $\mathbf{y} = y_1, \dots, y_n$ to an input sequence $\mathbf{x} = x_1, \dots, x_n$. Given a sample \mathbf{x} , we define the feature $\Phi(\mathbf{x}, \mathbf{y})$. Thus, we can label \mathbf{x} with a score function,

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} F(\mathbf{w}, \Phi(\mathbf{x}, \mathbf{y})), \quad (1)$$

where \mathbf{w} is the parameter of function $F(\cdot)$.

For sequence labeling, the feature can be denoted as $\phi_k(y_i, y_{i-1}, \mathbf{x}, i)$, where i stands for the position in the sequence and k stands for the number of feature templates.

Here, we use two popular open source toolkits for sequence labeling task as the baseline systems: FNL² [7] and CRF++³. Here, we use the default setting of CRF++ toolkit with the feature templates as shown in Table 4. The same feature templates are also used for FNL.

Table 4. Templates of CRF++ and FNL.

unigram feature	$c_{-2}, c_{-1}, c_0, c_{+1}, c_{+2}$
bigram feature	$c_{-1} \circ c_0, c_0 \circ c_{+1}$
trigram feature	$c_{-2} \circ c_{-1} \circ c_0, c_{-1} \circ c_0 \circ c_{+1}, c_0 \circ c_{+1} \circ c_{+2}$

5.2 Participant Systems

In CWS subtask, the best F1 performances are 95.12, 95.52 and 96.65 for closed, semi-open and open tracks respectively. The best system outperforms the baseline systems on closed track. The best system on semi-open track is better than that on closed track. Unsurprisingly, the performances boost greatly on open track.

Table 5. Performances of word segmentation.

Systems	Precision	Recall	F1	Track
CRF++	93.3	93.2	93.3	baseline, closed
FNL	94.1	93.9	94.0	
NJU	95.14	95.09	95.12	closed
BosonNLP	95.03	95.03	95.03	
CIST	94.78	94.42	94.6	
XUPT	94.61	93.85	94.22	
CCNU	93.95	93.45	93.7	
ICT-NLP	93.96	92.91	93.43	
BJTU	89.49	93.55	91.48	
CIST	95.47	95.57	95.52	semi-open
NJU	95.3	95.31	95.3	
BJTU	90.91	94.46	92.65	
ZZU	85.36	85.25	85.31	
BosonNLP	96.56	96.75	96.65	open
NJU	96.03	96.15	96.09	
SZU	95.52	95.64	95.58	
CCNU	93.68	93.09	93.38	
BJTU	91.79	94.92	93.33	

² <https://github.com/xpqi/fnl/>

³ <http://taku910.github.io/crfpp/>

The participant systems are briefly described as follows.

- The **ZZU** system uses sequence labeling for CWS with CRF model. Besides the traditional discrete feature templates, the dense representation for every character (called character vector) is also used as feature of the character. They report that the features that learned from background data automatically is weaker than the features artificially designed.
- The **CIST** system opts for the 6-tag set and corresponding six n-gram character features. Accessor variety (AV) [3] features are used to measure the possibility of whether a substring is a Chinese word. They report that the ability of OOV detection can be improved by integrating unsupervised global features extracted from the provided background data.
- The **SZU** system exploits multiple heterogeneous data to boost performance of statistical models. The system considers three sets of heterogeneous data, i.e., Weibo (WB, 10K sentences), Penn Chinese Treebank 7.0 (CTB7, 50K), and People’s Daily (PD, 280K). With the additional datasets, the F1 score is boosted from 93.76% (baseline model trained on only WB) to 95.58% (+1.82%).
- The **BJTU** system uses the CRF model to integrate several features, including normal features, dictionary features (named entity, hot micro-blog words and symbols) and branch entropy (BE) features. They also use the error-driven rule learning method to expand the training data set in order to improve the accuracy of the system, and improve the adaptability of the system.
- The **ICT-NLP** system adopts the character sequence labelling model and is trained with the averaged perceptron algorithm. To further improve the performance, they combine rules to deal with numbers and English strings and use internal dictionary (extracted from the training data) to do post-processing.
- The **NJU** system applies a word-based perceptron algorithm to build the base segmenter. They also use a bootstrap aggregating model of bagging which improves the segmentation results consistently on the three tracks of closed, semi-open and open test. Besides the basic features, they also use mutual information and accessor variety features.
- The **BosonNLP** system adopts an ensemble approach by combining both discriminative and generative methods. They find that 5-tag labeling consistently provides the best results. They also use several common patterns that might be useful for segmentation and tagging, such as “NN+ $\{|\}$ ”. For the open track, they use several models (HMM, CRF) trained on the other corpora (People’s daily corpus and Chinese TreeBank 7.0).

6 SubTask 2: Joint Chinese Word Segmentation and POS Tagging

The evaluation measures are precision, recall, and an evenly-weighted F1.

In the joint word segmentation and POS tagging, the best performances are 88.93, 88.69 and 91.55 for closed, semi-open and open tracks respectively.

Table 6. Performances of joint word segmentation and POS tagging.

Systems	Precision	Recall	F1	Track
BosonNLP	88.91	88.95	88.93	closed
XUPT	88.54	87.83	88.19	
BJTU	88.28	87.67	87.97	
CIST	88.09	87.76	87.92	
BJTU	80.64	85.1	82.81	
CIST	88.64	88.73	88.69	semi-open
WHU	88.59	87.96	88.27	
BJTU	81.76	85.82	83.74	
BosonNLP	91.42	91.68	91.55	open
SZU	88.93	89.05	88.99	
BJTU	79.85	83.51	81.64	

- The **CIST** system takes segmented inputs which are produced by the word segmenter used in CWS task, and we assign POS tags on a word-by-word basis, making use of features in the surrounding context (word-based). POS tagger for closed and semi-open track differ only in the segmentation step, closed and semi-open tagger receives word sequences from closed and semi-open segmenter respectively.
- The **SZU** system also adopts a cascaded approach for POS tagging. They use an ensemble approach combining coupled sequence labeling and the guide-feature based method. Same to CWS subtask, they use three datasets to boost the performance. Finally, the tagging F1 score is improved from 87.93% to 88.99% (+1.06%).
- The **BosonNLP** system for the S&T subtask is same to CWS, which obtains the best results in closed and open track of the S&T subtask.

7 Analysis

The analyses of the participant systems are as follows.

1. All participant systems adopt both pre-processing and post-processing. The major purpose is to remove noises and regular patterns in micro text, such as username, URL, Email, expression symbols and the other special symbols. These processings improve the final performances greatly.
2. All the top systems adopt the ensemble based method. The improvement is obtained by combining several models which are trained on other heterogeneous annotated datasets.

3. The background data can improve the performance of CWS by about 1%. Some statistical features are extracted from the background data, such as branch entropy and mutual information.
4. For the POS tagging subtask, most systems adopt the pipeline method: first word segmentation, then POS tagging. The reason behind this is that they do their utmost to optimize the performances of CWS. These optimization cannot be applied for joint segmentation and POS tagging.

Since this is the first time for us to organize the shared task, there are a few points needed to be improved.

1. Most of the annotated data are obtained from the micro-blogs of the official news accounts. These texts are relatively more formal than the real micro-texts. Therefore, more informal micro-texts should be added in the future evaluation.
2. The noises in micro-texts should be removed with the same standard, which can reduce the differences of the participant systems with different pre-processings and post-processings.

8 Conclusion

After years of intensive researches, Chinese word segmentation and POS tagging have achieved a quite high precision. However, the performances of the state-of-the-art systems are still relatively low for the informal texts, such as micro-blogs, forums. The NLPCC 2015 Shared Task on Chinese Word Segmentation and POS Tagging for Micro-blog Texts focuses on the fundamental research in Chinese language processing. It is the first time to use the micro-texts to evaluate the performance of the state-of-the-art methods.

In future work, we hope to run an online evaluation system to accept open registration and submission. Currently, a simple system is available at <http://nlp.fudan.edu.cn/nlpcc2015>. The system also gives the leaderboards for the up-to-date results under the different tasks and tracks. Besides, we also wish to extend the scale of corpus and add more informal texts.

Acknowledgement. We are very grateful to the students from our lab for their efforts to annotate and check the data. We would also like to thank the participants for their valuable feedbacks and comments. This work was partially funded by the National Natural Science Foundation of China (61472088), National High Technology Research and Development Program of China (2015AA015408), Shanghai Science and Technology Development Funds (14ZR1403200).

References

1. Berger, A.L., Della, V.J.: Pietra, and S.A. Della Pietra. A maximum entropy approach to natural language processing. *Computational Linguistics* **22**(1), 39–71 (1996)
2. Collins, M.: Discriminative training methods for hidden markov models: theory and experiments with perceptron algorithms. In: *Proceedings of the 2002 Conference on Empirical Methods in Natural Language Processing* (2002)
3. Feng, H., Chen, K., Deng, X., Zheng, W.: Accessor variety criteria for chinese word extraction. *Computational Linguistics* **30**(1), 75–93 (2004)
4. Lafferty, J.D., McCallum, A., Pereira, F.C.N.: Conditional random fields: probabilistic models for segmenting and labeling sequence data. In: *Proceedings of the Eighteenth International Conference on Machine Learning* (2001)
5. Ng, H.T., Low, J.K.: Chinese part-of-speech tagging: one-at-a-time or all-at-once? word-based or character-based. In: *Proceedings of EMNLP*, vol. 4 (2004)
6. Peng, F., Feng, F., McCallum, A.: Chinese segmentation and new word detection using conditional random fields. In: *Proceedings of the 20th International Conference on Computational Linguistics* (2004)
7. Qiu, X., Zhang, Q., Huang, X.: FudanNLP: a toolkit for Chinese natural language processing. In: *Proceedings of Annual Meeting of the Association for Computational Linguistics* (2013)
8. Xue, N.: Chinese word segmentation as character tagging. *Computational Linguistics and Chinese Language Processing* **8**(1), 29–48 (2003)
9. Xue, N., Xia, F., Chiou, F.-D., Palmer, M.: The Penn Chinese TreeBank: Phrase structure annotation of a large corpus. *Natural language engineering* **11**(2), 207–238 (2005)
10. Yang, Y., Xue, N.: Chinese comma disambiguation for discourse analysis. In: *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers*, vol. 1, pp. 786–794. Association for Computational Linguistics (2012)