

An Example-Based Chinese Word Segmentation System for CWSB-2

Chunyu Kit Xiaoyue Liu

Department of Chinese, Translation and Linguistics
City University of Hong Kong
Tat Chee Ave., Kowloon, Hong Kong
{ctckit, xyliu0}@cityu.edu.hk

Abstract

This paper reports the example-based segmentation system for our participation in the second Chinese Word Segmentation Bakeoff (CWSB-2), presenting its basic ideas, technical details and evaluation. It is a preliminary implementation. CWSB-2 valuation shows that it performs very well in identifying known words. Its unknown word detection module also illustrates great potential. However, proper facilities for identifying time expressions, numbers and other types of unknown words are needed for improvement.

1 Introduction

Word segmentation is to identify lexical items, especially individual word forms, in a text. It involves two fundamental tasks, both aiming at minimizing segmentation errors: one is to infer out-of-vocabulary (OOV) words, also known as unknown (or unseen) word detection, and the other to identify in-vocabulary (IV) words, with an emphasis on disambiguation. OOV words and ambiguities are the two major causes to segmentation errors.

Accordingly, word segmentation approaches can be divided into the categories summarized in Table 1 in terms of the resources in use to tackle these two causes. The closed and open tracks in CWSB correspond, respectively, to the last two categories, both involving inferring OOV words

Category	Resource in use		Major Task	
	Lexicon	Tr. Corpus	OOV	Disamb.
WD ^a	-	(-) ^b	+	
WS/CL ^c	+	-	-	+
WS/IL ^d	+	-	+	+
WS/TC ^e	(+) ^f	+	+	+
WS/TC+L ^g	+	+	+	+

^aWord discovery, or unsupervised lexical acquisition

^bInput data is used for unsupervised training

^cWord segmentation with a complete lexicon

^dWord segmentation with an incomplete lexicon

^eWord segmentation with a pre-segmented training corpus

^fIt can be extracted from the given training corpus.

^gWord segmentation with a pre-segmented training corpus and an extra lexicon

Table 1: Categories of segmentation approach

beyond disambiguating IV words. Word discovery and OOV word detection pursue a similar target, i.e., inferring new words. The continuum connecting them is the size of the lexicon in use: the former assumes few words known and the latter an existing lexicon to some scale. Inferring new words is an essential task in word segmentation, for a complete lexicon is rarely a realistic assumption in practice.

This paper presents our segmentation system for participation in CWSB-2. It takes an example-based approach to recognize IV words and follows description length gain (DLG) to infer OOV words in terms of their text compression effect. Sections 2 and 3 below introduce the example-based and DLG-based segmentation respectively. Section 4 presents a strategy to combine their strength and Section 5 reports our system’s performance in CWSB-2. Following error analysis in Section 6, Section 7 concludes the paper.

2 Example-based segmentation

How to utilize as much information as possible from the training corpus to adapt a segmentation system towards a segmentation standard has been a critical issue. Kit et al. (2002) and Kit et al. (2003) attempt to integrate case-based learning with statistical models (e.g., n-gram) by extracting transformation rules from the training corpus for disambiguation via error correction; Gao et al. (2004) adopt a similar strategy for adaptive segmentation, with transformation templates (instead of case-based rules) to modify word boundaries (instead of individual words).

The basic idea of example-based segmentation is very simple: existing pre-segmented strings in training corpus provide reliable examples for segmenting similar strings in input texts. In contrast to dictionary checking for locating *possible* words in an input sentence to facilitate later segmentation operations, pre-segmented examples give *exact* segmentation to copy.

The example-based segmentation can be implemented in the following steps.

1. Find all exemplar pre-segmented fragments, with regards to a training corpus, and all possible words, with regards to a lexicon, from each character in an input sentence;
2. Identify the optimal sequence, among all possibilities, of the above items over the sentence following some optimization criterion.

If adopting the minimal number of fragments or words in a sequence as optimization criterion, we have a maximal matching approach to word segmentation. However, it differs remarkably from the previous maximal matching approaches: it matches pre-segmented fragments, instead of dictionary words, against an input sentence. It can be carried out by a best-first strategy: repeatedly select the next longest example or word until the entire sentence is properly covered. Unfortunately, the best-first approach does not guarantee to give the best answer. For CWSB-2, we implemented a program following the Viterbi algorithm to perform a complete search in terms of the number of fragments, and then words, in a sequence.

However, a serious problem with this example-based approach is the *sparse data* problem. Long exemplar fragments are more reliable but small

in number, whereas short ones are large in number but less reliable. In the case of no exemplar fragment available for an input sentence, this approach draws back to the maximal match segmentation with a dictionary. How to incorporate statistical inference into example-based segmentation to infer more reliable optimal segmentation beyond string matching remains a critical issue for us to tackle.

3 DLG-based segmentation

DLG is formulated in Kit and Wilks (1999) and Kit (2000) as an empirical measure for the compression effect of extracting a substring from a given corpus as a lexical item. DLG optimization is applied to detect OOV words for our participation in CWSB-2. It works as follows in two steps.

1. Calculate the DLG for all known words and all new word candidate (i.e., substrings with frequency ≥ 2 , preferably, in the test corpus), based on frequency information in the training and the test corpora;
2. Find the optimal sequence of such items over an input sentence with the greatest sum of DLG.

Step 2 above in our system re-implements only the first round of DLG-based lexical learning in Kit (2000). It is implemented by the same algorithm as the one for example-based segmentation, with DLG as optimization criterion. Evaluation results show that this learning-via-compression approach discovers many OOV words successfully, in particular, person names.

4 Integration

The example-based segmentation is good at identifying IV words but incapable of recognizing any new words. In contrast, the DLG-based segmentation performs slightly worse but has potential to detect new words. It is expected that the strength of the two could be combined together for performance enhancement.

However, because of inadequate time we had to take a shortcut in order to catch the CWSB-2 deadline: DLG segmentation is only applied to recognize new words among the sequences of mono-character items in the example-based segmentation output.

Track	P	R	F	OOV	R _{OOV}	R _{IV}
AS _c	.944	.902	.923	.043	.234	.976
PKU _c	.929	.904	.916	.058	.252	.971
MSR _c	.965	.935	.950	.026	.189	.986

Table 2: System performance in CWSB-2

5 Performance

Our group took part in three closed tracks in CWSB-2, namely, AS_c, PKU_c and MSR_c, with a preliminary implementation of the example-based word segmentation presented above. Our system’s performance in terms of CWSB-2’s official scores is presented in Table 2. Its R_{OOV} scores look undesirable, showing that applying the first round of DLG-based segmentation to sequences of mono-character items is inadequate for the OOV word discovery task. Nevertheless, its R_{IV} scores are, in general, quite close to the top systems in CWSB-2, although it does not have a disambiguation module to polish its maximal matching output.

However, this is not to say that the DLG-based segmentation deserves no credit in unknown word detection. It does recognize many OOV words, as shown in Table 3. The low R_{OOV} rate has to do with our system’s incapability in handling time expressions, numbers, and foreign words.

6 Error analysis

Most errors made by our system are due to the following causes: (1) no knowledge, overt or implicit, in use for recognizing time expressions, numbers and foreign words, as restricted by CWSB-2 rules, (2) a premature module for OOV word detection, (3) no further disambiguation besides example application, and (4) significant inconsistency in the training and test data.

The inconsistency exists not only between the training and test corpora for each track but, more surprisingly, also within individual training corpora. Some suspected cases are illustrated in Tables 4, 5 and 6. They are observed to be in a large number in the CWSB-2 corpora. Scoring with a golden standard involving so many of them appears to be problematic, for it penalizes the systems for handling such cases right and rewards the others for producing “correct” answers. What

AS_c: 小森(106) 左詩雅(45) 馮禮(31) 霍金(29) 瑪雅(21) 丹紐(20) 張庭(18) 英特爾(17) 黑盤鮑(16) 沈時華(15) 證期會(13) 厲莉(12) 黃銘坤(11) 卡門(11) 佛指舍利(11) 陳慶浩(10) 璩女(9) 范曉萱(8) 郭女(8) 邱太三(7) 網站(7) 狄亞尼(7) 張惠妹(6) 米田真澄(6) 鄭優(5) 丁克華(5) 經發會(5) 蘇盈貴(5) 陳秀玲(5) 八美圖(5) 房貧(5) 蔡信弘(5) 郭玉鈴(4) 地型(4) 御膳(4) 白藜蘆醇(4) 羅麗泰(3) 張珏(3) 米倉涼子(3) 陳兆伸(2) 羅勒(2) 幼發拉底河(2) 溥傑(2)

PKU_c: 罢免(38) 世清(23) 任免(21) 拉姆斯菲尔德(20) 哈苏(19) 海合会(17) 军级(16) 加纳(15) 水心村(12) 网友(11) 小阜村(10) 宋双(10) 吐逊江(9) 申奥(9) 库福尔(8) 金稅(8) 亚布力(8) 金门(7) 三阳镇(7) 教主(6) 罢免书(6) 法輪功(6) 香客(6) 普京(6) 福清市(5) 柯尔克孜(5) 华能(5) 罗林斯(5) 松诺斯(5) 洋务(5) 马祖(5) 甲午(5) 哈斯达(4) 奥申委(4) 高唐县(4) 天津(4) 门市部(4) 妈祖(4) 刁民(4) 团会(4) 商机(4) 临武县(4) 湄洲島(4) 圈套(4) 法塔赫(3) 流星雨(3) 堆龙德庆县(3) 阿族(3) 哥德堡(3) 信众(3) 钱柜(3) 衡水(3) 内坦亚(3) 双丰村(3)

MSR_c: 博古(26) 玲英(19) 游景玉(19) 任尧森(17) 秦机(15) 亚仿(14) 进占(14) 刘积仁(13) 东宝(13) internet(12) 猴王(12) 双保(11) 张肇群(10) 抚州(10) 南丁格尔(10) 彭珮云(10) 海塘(10) 王常力(9) 穆守家(9) 大关村(8) 樊皇(8) 张钧(8) 嶂石岩(7) 三老四严(7) 局域网(7) 透支(6) 秦家山(6) 东软(6) 后金(6) 八旗(6) 弄堂(5) 黎秀芳(5) 提灌(5) 大关(5) 王丙乾(5) 棉铃虫(5) 百亿次(4) 西文(4) 米夫(4) 陆冰(4) 郑守仁(4) 秦邦宪(4) 关小瑛(3) 二进制(3) 顺延(3) 李朋朋(3) 瞿秋白(3) 黄华平(3) 汪赛进(3) 汪延(3) 何元亮(3) 陆佑楣(3) 阜平县(3) 中信所(3) 导流(3) 徐殿龙(3) 重男轻女(2)

Table 3: Illustration of new words successfully detected, with frequency in parentheses

is even more worth noting is that (1) an inconsistent case involves more than one word, and (2) the difference between a correct and an erroneous judgment of a word is 1, in a sense, but the difference between one system that loses it for doing right and another that earns it by doing wrong is surely greater.

7 Conclusions

In the above sections we have reported the example-based word segmentation system for our participation in CWSB-2, including its basic ideas, technical details and evaluation results. It has illustrated an excellent performance in IV word identification and nice potential in OOV word discovery. However, its weakness in handling time expressions, numbers and other types of unknown words has hindered it from performing better. We are expecting to implement a full-fledged version of the system for improvement.

Acknowledgements

The work described in this paper was supported by the RGC of HKSAR, China, through the CERG grant 9040861. We wish to thank Alex Fang and Robert Neather for their help.

Training & Answer	f_T/f_A	Golden Standard	f_T/f_A
繁殖場	4/8	繁殖場	0/0
老歌	28/7	老歌	0/0
準備率	5/7	準備率	0/0
小火車	11/6	小火車	0/0
本文	186/5	本文	0/0
第三者	41/4	第三者	0/0
新一代	29/4	新一代	0/0
稱之為	129/4	稱之為	0/0
狗不夠	23/3	狗不夠	0/0
短時間	47/3	短時間	0/0
中正機場	33/2	中正機場	0/0
台北市長	32/2	台北市長	0/0
民意代表	85/2	民意代表	0/0
清醒過來	10/2	清醒過來	0/0
統治者	62/2	統治者	0/0
新生命	23/2	新生命	0/0
也就是說	192/1	也就是說	0/0
不景氣	149/1	不景氣	0/0
中央政府	66/1	中央政府	0/0
挑戰性	31/1	挑戰性	0/0
看電影	80/1	看電影	0/0
研究者	68/1	研究者	0/0
面無表情	13/1	面無表情	0/0
曾文水庫	13/1	曾文水庫	0/0
無意識	20/1	無意識	0/0
無話不談	6/1	無話不談	0/0
另一面	29/1	另一面	0/0
紅透半邊天	4/1	紅透半邊天	0/0
民間團體	24/7	民間團體	25/0
女性主義	17/3	女性主義	53/0
是不是	1201/2	是不是	2/0

Table 4: Some inconsistent cases in AS corpus

Training & Answer	f_T/f_A	Golden Standard	f_T/f_A
鄉政府	14/26	鄉政府	0/0
區政府	6/1	區政府	0/0
鎮政府	5/21	鎮政府	0/0
世紀之交	24/19	世紀之交	0/0
改為	23/18	改為	0/0
生產總值	66/15	生產總值	0/0
個人所得稅	10/9	個人所得稅	0/0
無黨派人士	10/5	無黨派人士	0/0
新年伊始	45/5	新年伊始	0/0
組成部分	42/5	組成部分	0/0
固定資產	27/4	固定資產	0/0
世界各地	21/4	世界各地	0/0
鄉鎮企業	126/4	鄉鎮企業	0/0
至關重要	20/4	至關重要	0/0
綜合國力	15/4	綜合國力	0/0
總參謀長	25/4	總參謀長	0/0
經濟社會	25/3	經濟社會	0/0
勞動生產率	13/3	勞動生產率	0/0
去年底	32/3	去年底	0/0
世界紀錄	30/3	世界紀錄	0/0
無黨派	11/3	無黨派	0/0
議事日程	15/3	議事日程	0/0
中小企業	22/3	中小企業	0/0
不同於	11/2	不同於	0/0
此事	25/2	此事	0/0
火樹銀花不夜天	3/1	火樹銀花不夜天	0/0
雨夾雪	13/1	雨夾雪	0/0
縣政府	24/5	縣政府	1/0
意味着	49/4	意味着	1/0
很多	112/3	很多	14/0
發達國家	48/1	發達國家	1/0

Table 5: Some inconsistent cases in PKU corpus

Training & Answer	f_T/f_A	Golden Standard	f_T/f_A
營業部	12/7	營業部	0/0
高性能	16/6	高性能	0/0
“八五	29/5	“八五	0/0
大型企業集團	6/3	大型企業集團	0/0
徑流量	3/3	徑流量	0/0
海信電器	1/2	海信電器	0/0
投資銀行?	4/2	投資銀行	0/0
一版	10/2	一版	0/0
浙東	3/2	浙東	0/0
2000多人	7/1	2000多人	0/0
採訪團	2/1	採訪團	0/0
產油國家	1/1	產油國家	0/0
第一、二、三?	4/1	第一、二、三	0/0
葛洲壩大江	1/1	葛洲壩大江	0/0
國際互聯網	1/1	國際互聯網	0/0
計算機應用軟件	1/1	計算機應用軟件	0/0
家門口?	10/1	家門口	0/0
兩點?	16/1	兩點	0/0
呂梁?	4/1	呂梁	0/0
一口氣?	16/1	一口氣	0/0
中國人民	122/1	中國人民	0/0
中外合資經營	3/1	中外合資經營	0/0

Table 6: Some inconsistent cases in MSR corpus

References

- E. Brill. 1993. *A Corpus-Based Approach to Language Learning*. PhD thesis, University of Pennsylvania, Philadelphia.
- J. Gao, A. Wu, M. Li, C. Huang, H. Li, X. Xia and H. Qin. 2004. Adaptive Chinese word segmentation. In *ACL-04*. Barcelona, July 21-26.
- C. Kit and Y. Wilks. 1999. Unsupervised learning of word boundary with description length gain. In M. Osborne and E. T. K. Sang (eds.), *CoNLL-99*, pp.1-6. Bergen, Norway, June 12.
- C. Kit 2000. *Unsupervised Lexical Learning as Inductive Inference*. PhD thesis, University of Sheffield.
- C. Kit, H. Pan and H. Chen. 2002. Learning case-based knowledge for disambiguating Chinese word segmentation: A preliminary study. *SIGHAN-1*, pp.33-39. Taipei, Sept. 1, 2002.
- C. Kit, Z. Xu and J. J. Webster. 2003. Integrating n-gram model and case-based learning for Chinese word segmentation. In Q. Ma and F. Xia (eds.), *SIGHAN-2*, pp.160-163. Sapporo, 11 July, 2003.
- D. Palmer. A trainable rule-based algorithm for word segmentation. In *ACL-97*, pp.321-328. Madrid.