

Technical Report of *Chinese Morphological
Analyzer(Chen)*

Vernkin Smith

March 20, 2009

Abstract

Practical results show that performance of statistic segmentation system outperforms that of hand-crafted rule-based systems. And the evaluation shows that the accuracy drop caused by out-of-vocabulary (OOV) words is at least five times greater than that of segmentation ambiguities. The better performance of OOV recognition the higher accuracy of the segmentation system in whole, and the accuracy of statistic segmentation systems with character-based tagging approach outperforms any other word-based system. This Report is about a supervised machine-learning approach to Character-based Chinese word segmentation. A maximum entropy tagger is trained on manually annotated data to automatically assign to Chinese characters, or *hanzi*, tags that indicate the position of a *hanzi* within a word.

Contents

1	Overview	2
1.1	A Decade Review of Chinese Word Segmentation	2
1.2	Main Difficulties in Chinese Segmentation	2
2	Introduce to Character-based Chinese Segmentation	5
3	Maximum Entropy Modeling	7
3.1	The Modeling Problem	7
3.2	Parameter Estimation	8
3.3	Usage of the MaxEnt Model	9
3.3.1	Representing Features	9
3.3.2	Create a Maxent Model Instance	9
3.3.3	Adding Events to Model	10
3.3.4	Training the Model	11
3.3.5	Using the Model	11
4	POS Tagger	13
4.1	The Tagging Model	13
4.2	Feature Selection	13
A	Appendix - Project schedules and milestones	16

Chapter 1

Overview

1.1 A Decade Review of Chinese Word Segmentation

During the last decade, especially since the First International Chinese Word Segmentation Bakeoff was held in July 2003, the study is automatic Chinese word segmentation has been greatly improved. Those improvements could be summarized as following: (1) on the computation sense Chinese words in read text that have been well-defined by "segmentation guidelines + lexicon + segmented corpus"; (2) practical results show that performance of statistic segmentation system outperforms that of hand-crafted rule-based systems; (3) the evaluation shows than the accuracy drop caused by out-of-vocabulary (OOV) words is at least five times greater than that of segmentation ambiguities; (4) the better performance of OOV recognition the higher accuracy of the segmentation system in whole, and the accuracy of statistic segmentation systems with character-based tagging approach outperforms any other word-based system.

1.2 Main Difficulties in Chinese Segmentation

It is generally agreed among researchers that word segmentation is a necessary first step in Chinese language processing. However, unlike English text in which sentences are sequences of words delimited by white spaces, in Chinese text, sentences are represented as strings of Chinese characters or *hanzi* without similar natural delimiters. Therefor, the first step in a Chinese language processing task is to identify the sequence of words in a sentence and mark boundaries in appropriate places. This may sound simple enough but in reality identifying words in Chinese is a non-trivial problem that has drawn a large body of research in the Chinese language processing community.

It is easy to demonstrate that the lack of natural delimiters itself is not

the heart of the problem. In a hypothetical language where all words are represented with a finite set of symbols, if one subset of the symbols always start a word and another subset, mutually exclusive from the previous subset, always end a word, identifying words would be a trivial exercise. Nor can the problem be attributed to the lack of inflectional morphology. Although it is true in Indo-European languages inflectional affixes can generally be used to signal word boundaries, it is conceivable that a hypothetical language can use symbols other than inflectional morphemes to serve the same purpose. Therefore the issue is neither the lack of natural word delimiters nor the lack of inflectional morphemes in a language, rather it is whether the language has a way of unambiguously signaling the boundaries of a word.

The real difficulty in automatic Chinese word segmentation is the lack of such unambiguous word boundary indicators. In fact, most hanzi can occur in different positions within different words. The examples in Table 1 show how the Chinese character 产 (“produce”) can occur in four different positions. This state of affairs makes it impossible to simply list mutually exclusive subsets of hanzi that have distinct distributions, even though the number of hanzi in the Chinese writing system is in fact finite. As long as a hanzi can occur in different word-internal positions, it cannot be relied upon to determine word boundaries as they could be if their positions were more or less fixed.

Table 1. A hanzi can occur in multiple word-internal positions

Position	Example
Left	产生 ‘to come up with’
Word by itself	产小麦 ‘to grow wheat’
Middle	生产线 ‘assembly line’
Right	生产 ‘to produce’

The fact that a hanzi can occur in multiple word-internal positions leads to ambiguities of various kinds. For example, 文 can occur in both word-initial and word-final positions. It occurs in the word-final position in 日文 (“Japanese”) but in the word-initial position in 文章 (“article”). In a sentence that has a string “日文章”, as in (1a), an automatic segmenter would face the dilemma whether to insert a word boundary marker between 日 and 文, thus grouping 文章 as a word, or to mark 日文 as a word, to the exclusion of 章. The same scenario also applies to 章, since like 文, it can also occur in both word-initial and word-final positions.

1. (a): Segmentation 1

日文章 鱼 怎么说?

Japanese octopus how say

“How to say octopus in Japanese?”

1. (b): Segmentation 2

日 文章 鱼 怎么说?

Japan article fish how say

Ambiguity also arises because some hanzi should be considered to be just word components in certain contexts and words by themselves in others. For example, 魚 can be considered to be just a word component in 章魚. It can also be a word by itself in other contexts. Presented with the string 章魚 in a Chinese sentence, a human or automatic segmenter would have to decide whether 魚 should be a word by itself or form another word with the previous hanzi. Given that 日, 文章, 章魚, 魚 are all possible words in Chinese, how does one decide that 日文章魚 is the right segmentation for the sentence in (1) while 日文章 魚 is not? Obviously it is not enough to know just what words are in the lexicon. In this specific case, a human segmenter can resort to world knowledge to resolve this ambiguity, knowing that 日文章 魚 would not make any kind of real-world sense.

Chapter 2

Introduce to Character-based Chinese Segmentation

Follow the *Overview* and *Abstract*, this chapter would introduce Character-based Chinese Segmentation. The Maximum Entropy Modeling (see chapter-3) used in the segmentation and more detail see the following chapters. In this charper, we first formalize the idea of tagging *hanzi*(Chinese Character) based on their word-internal positions and describe the tag set we used.

First we convert the manually segmented words in the corpus into a tagged sequence of Chinese characters. To do this, we tag each character with one of the four tags, LL, RR, MM and LR depending on its position within a word. It is tagged LL if it occurs on the left boundary of a word, and forms a word with the character(s) on its right. It is tagged RR if it occurs on the right boundary of a word, and forms a word with the character(s) on its left. It is tagged MM if it occurs in the middle of a word. It is tagged LR if it forms a word by itself. We call such tags position-of-character (POC) tags to differentiate them from the more familiar part-of-speech (POS) tags. For example, the manually segmented string in (2a) will be tagged as (2b):

Example Senetence 2:

(a) 上海 计划 到 本 世纪 末 实现 人均 国内 生产 总值 五千 美元

(b) 上/LL 海/RR 计/LL 划/RR 到/LR 本/LR 世/LL 纪/RR 末/LR 实/LL 现/RR 人/LL 均/RR 国/LL 内/RR 生/LL 产/RR 总/LL 值/RR 五/LL 千/RR 美/LL 元/RR

(c) Shanghai plans to reach the goal of 5,000 dollars in per capita GDP by the end of the century.

Given a manually segmented corpus, a POC-tagged corpus can be derived trivially with perfect accuracy. The reason why use such POC-tagged sequences of characters instead of applying n -gram rules to segmented corpus directly[Palmer, 1997[5]; Hockenmaier and Brew, 1998[3]; Xue, 2001[6]] is that they are much easier to manipulate in the training process. In addition, the POC tags reflect our observation that the ambiguity problem is

due to the fact that a hanzi can occur in different word-internal positions and it can be resolved in context. Naturally, while some characters have only one POC tag, most characters will receive multiple POC tags, in the same way that words can have multiple POS tags. Table 2 shows how all four of the POC tags can be assigned to the character 产 (“produce”):

Table 2. A character can receive as many as four tags

Position	Tag	Example
Left	LL	产生 ’to come up with’
Word by itself	LR	产小麦 ’to grow wheat’
Middle	MM	生产线 ’assembly line’
Right	RR	生产 ’to produce’

If there is ambiguity in segmenting a sentence or any string of hanzi, then there must be some hanzi in the sentence that can receive multiple tags. For example, each of the first four characters of the sentence in (1) would have two tags. The task of the word segmentation is to choose the correct tag for each of the hanzi in the sentence. The eight possible tag sequences for (1) are shown in (3a), and the correct tag sequence is (3b).

Example Sentence 4:

(a) 日/LL,LR 文/RR,LL 章/LL,RR 鱼/RR,LR 怎/LL 么/RR 说/LR ?

(b) 日/LL 文/RR 章/LL 鱼/RR 怎/LL 么/RR 说/LR ?

Also like POS tags, how a character is POC-tagged in naturally occurring text is affected by the context in which it occurs. For example, if the preceding character is tagged LR or RR, then the next character can only be tagged LL or LR. How a character is tagged is also affected by the surrounding characters. For example, 关(“close”) should be tagged RR if the previous character is 开(“open”) and neither of them forms a word with other characters, while it should be tagged LL if the next character is 心(“heart”) and neither of them forms a word with other characters. This state of affairs closely mimics the familiar POS tagging problem and lends itself naturally to a solution similar to that of POS tagging. The task is one of ambiguity resolution in which the correct POC tag is determined among several possible POC tags in a specific context. Our next step is to train a maximum entropy model on the perfectly POC-tagged data derived from a manually segmented corpus to automatically POC-tag unseen text.

Chapter 3

Maximum Entropy Modeling

This chapter provides a brief introduction to Maximum Entropy Modeling. The Advantage of maximum entropy model includes: Based on features, allows and supports feature induction and feature selection; offers a generic framework for incorporating unlabeled data; only makes weak assumptions; gives flexibility in incorporating side information; natural multi-class classification.

Maximum Entropy (ME or maxent for short) model is a general purpose machine learning framework that has been successfully applied in various fields including spatial physics, computer vision, and Natural Language Processing (NLP). This introduction will focus on the application of maxent model to NLP tasks. However, it is straightforward to extend the technique described here to other domains.

3.1 The Modeling Problem

The goal of statistical modeling is to construct a model that best accounts for some training data. More specific, for a given empirical probability distribution \tilde{p} , we want to build a model p as close to \tilde{p} as possible.

Of course, given a set of training data, there are numerous ways to choose a model p that accounts for the data. It can be shown that the probability distribution of the form 3.1 is the one that is closest to \tilde{p} in the sense of Kullback-Leibler divergence, when subjected to a set of feature constraints:

$$p(y | x) = \frac{1}{Z(x)} \exp \left[\sum_{i=1}^k \lambda_i f_i(x, y) \right] \quad (3.1)$$

Here $p(y | x)$ denotes the conditional probability of predicting an *outcome* y on seeing the *context* x . $f_i(x, y)$'s are feature functions (described in detail later), λ_i 's are the weighting parameters for $f_i(x, y)$'s. k is the number of features and $Z(x)$ is a normalization factor (often called partition function) to ensure that $\sum_y p(y|x) = 1$.

ME model represents evidence with binary functions¹ known as *contextual predicates* in the form:

$$f_{cp,y'}(x, y) = \begin{cases} 1 & \text{if } y = y' \text{ and } cp(x) = true \\ 0 & \text{otherwise} \end{cases} \quad (3.2)$$

Where cp is the contextual predicate that maps a pair of *outcome* y and *context* x to $\{true, false\}$.

The modeler can choose arbitrary feature functions in order to reflect the characteristic of the problem domain as faithfully as possible. The ability of freely incorporating various problem-specific knowledge in terms of feature functions gives ME models the obvious advantage over other learn paradigms, which often suffer from strong feature independence assumption (such as naive bayes classifier).

For instance, in part-of-speech tagging, a process that assigns part-of-speech tags to words in a sentence, a useful feature may be (DET is DETERMINER for short):

$$f_{prev_tag_DET,NOUN}(x, y) = \begin{cases} 1 & \text{if } y = NOUN \text{ and } prev_tag_DET(x) = true \\ 0 & \text{otherwise} \end{cases}$$

which is *activated* when previous tag is DETERMINER and current word's tag is NOUN. In Text Categorization task, a feature may look like (RO is ROMANTIC for short):

$$f_{doc_has_RO,love_story}(x, y) = \begin{cases} 1 & \text{if } y = love_story \text{ and } doc_has_RO(x) = true \\ 0 & \text{otherwise} \end{cases}$$

which is *activated* when the term ROMANTIC is found in a document labeled as type:love_story.

Once a set of features is chosen by the modeler, we can construct the corresponding maxent model by adding features as constraints to the model and adjust weights of these features. Formally, We require that:

$$E_{\tilde{p}} < f_i > = E_p < f_i >$$

Where $E_{\tilde{p}} < f_i > = \sum_x \tilde{p}(x, y) f_i(x, y)$ is the empirical expectation of feature $f_i(x, y)$ in the training data and $E_p < f_i > = \sum_x p(x, y) f_i(x, y)$ is the feature expectation with respect to the model distribution p . Among all the models subjected to these constraints there is one with the Maximum Entropy, usually called the Maximum Entropy Solution.

3.2 Parameter Estimation

Given an exponential model with n features and a set of training data (empirical distribution), we need to find the associated real-value weight for each of the n feature which maximize the model's log-likelihood:

$$L(p) = \sum_{x,y} \tilde{p}(x, y) \log p(y | x) \quad (3.3)$$

Selecting an optimal model subjected to given constraints from the exponential (log-linear) family is not a trivial task. There are two popular iterative scaling algorithms specially designed to estimate parameters of ME models of the form 3.1: *Generalized Iterative Scaling* [Darroch and Ratcliff, 1972][1] and *Improved Iterative Scaling* [Della Pietra et al., 1997][2].

Recently, another general purpose optimization method *Limited-Memory Variable Metric* (L-BFGS for short) method has been found to be especially effective for maximum entropy parameters estimating problem [Malouf, 2003][4]. L-BFGS is the default parameter estimating method in the implementation.

3.3 Usage of the MaxEnt Model

This section covers the basic steps required to build and use a Conditional Maximum Entropy Model.

3.3.1 Representing Features

The mathematical representation of a feature used in a Conditional Maximum Entropy Model can be written as:

$$f_{cp,y'}(x, y) = \begin{cases} 1 & \text{if } y = y' \text{ and } cp(x) = true \\ 0 & \text{otherwise} \end{cases} \quad (3.4)$$

where cp is the *contextual predicate* which maps a pair of *outcome* and *context* into $\{true, false\}$.

This kind of math notation must be expressed as features of literal string in order to be used in this toolkit. So a feature in part-of-speech tagger which has the form (DET is DETERMINER for short):

$$f_{prev_tag_DET, NOUN}(x, y) = \begin{cases} 1 & \text{if } y = NOUN \text{ and } prev_tag_DET(x) = true \\ 0 & \text{otherwise} \end{cases}$$

can be written as a literal string: “tag-1=DETERMINER_NOUN”.

3.3.2 Create a Maxent Model Instance

A *maxent* instance can be created by calling its constructor:

```
#include <maxent/maxentmodel.hpp>
using namespace maxent;
MaxentModel m;
```

This will create an instance of MaxentModel class called m. Please note that all classes and functions are in the namespace maxent. For illustration purpose, the include and using statements will be ignored intentionally in the rest of this section.

3.3.3 Adding Events to Model

Typically, training data consists of a set of events (samples). Each event has a *context*, an *outcome*, and a *count* indicating how many times this event occurs in training data.

Remember that a *context* is just a group of *contextpredicates*. Thus an event will have the form:

$$[(predicate_1, predicate_2, \dots, predicate_n), outcome, count]$$

Suppose we want to add the following event to our model:

$$[(predicate_1, predicate_2, predicate_3), outcome1, 1]$$

We need to first create a context:

```
std::vector<std::string> context;
context.append("predicate1");
context.append("predicate2");
context.append("predicate3");
. . .
```

It's possible to specify feature value (must be non-negative) in creating a context:

```
std::vector<pair<std::string, float> > context;
context.append(make_pair("predicate1", 1.0));
context.append(make_pair("predicate2", 2.0));
context.append(make_pair("predicate3", 3.0));
. . .
```

Before any event can be added, one must call `begin_add_event()` to inform the model the beginning of training.

```
m.begin_add_event();
```

Now we are ready to add events:

```
m.add_event(context, "outcome1", 1);
```

The third argument of `add_event()` is the count of the event and can be ignored if the count is 1.

One can repeatedly call `add_event()` until all events are added to the model.

After adding the last event, `end_add_event()` must be called to inform the model the ending of adding events.

```
m.end_add_event();
```

3.3.4 Training the Model

Train a Maximum Entropy Model is relatively easy. Here are some examples:

```
// train the model with default training method
m.train();
// train the model with 30 iterations of L-BFGS method
m.train(30, "lbfgs");
// train the model with 100 iterations of GIS method and apply
// Gaussian Prior smoothing with a global variance of 2
m.train(100, "gis", 2);
// set terminate tolerance to 1E-03
m.train(30, "lbfgs", 2, 1E-03);
```

The training methods can be either “gis” or “lbfgs” (default). Also, if `m.verbose` is set to 1 (default is 0), training progress will be printed to `stdout`.

You can save a trained model to a file and load it back later:

```
m.save("new_model");
m.load("new_model");
```

A file named `new_model` will be created. The model contains the definition of context predicates, outcomes, mapping between features and feature ids and the optimal parameter weight for each feature.

If the optional parameter `binary` is true and the library is compiled with `zlib` support, a compressed binary model file will be saved which is much faster and smaller than plain text model. The format of model file will be detected automatically when loading:

```
m.save("new_model", true); //save a (compressed) binary model
m.load("new_model");      //load it from disk
```

3.3.5 Using the Model

The use of the model is straightforward. The `eval()` function will return the probability $p(y|x)$ of an *outcome* given some *context*:

```
m.eval(context, outcome);
```

`eval_all()` is useful if we want to get the whole conditional distribution for a given context:

```
std::vector<pair<std::string, double> > probs;
m.eval_all(context, probs);
```

eval_all() will put the probability distribution into the vector *probs*. The items in *probs* are the outcome labels paired with their corresponding probabilities. If the third parameter *sort_result* is true (default) *eval_all()* will automatically sort the output distribution in descendant order: the first item will have the highest probability in the distribution.

Chapter 4

POS Tagger

This Chapter discusses the steps involved in building a Part-of-Speech (POS) tagger using Maximum Entropy Model in detail.

4.1 The Tagging Model

The task of POS tag assignment is to assign correct POS tags to a word stream (typically a sentence). The following table lists a word sequence and its corresponding tags:

To attack this problem with the Maximum Entropy Model, we can build a conditional model that calculates the probability of a tag y , given some contextual information x :

$$p(y|x) = \frac{1}{Z(x)} \exp \left[\sum_{i=1}^k \lambda_i f_i(x, y) \right]$$

Thus the possibility of a tag sequence $\{t_1, t_2, \dots, t_n\}$ over a sentence $\{w_1, w_2, \dots, w_n\}$ can be represented as the product of each $p(y|x)$ with the assumption that the probability of each tag y depends only on a limited context information x :

$$p(t_1, t_2, \dots, t_n | w_1, w_2, \dots, w_n) \approx \prod_{i=1}^n p(y_i | x_i)$$

Given a sentence $\{w_1, w_2, \dots, w_n\}$, we can generate K highest probability tag sequence candidates up to that point in the sentence and finally select the highest candidate as our tagging result.

4.2 Feature Selection

We select features used in the tagging model by applying a set of feature templates to the training data. The features are:

1. Prefix / suffix characters of the word;

2. Whether it is numeric;
3. Whether it is hyphen;
4. The word and the POS of the last/previous word;
5. The word and the POS of the last/previous 2th word;
6. the words and the POSs combination of the last/previous 1st the 2th words.

The following table is the features which are selected from the actual sentence:

Table 4.1 The Features Example

curword=years
tag-1=CD
word-2=,
tag-1,2=,CD
word+1=old
curword=old
word-1=years
tag-1=NNS
tag-1,2=CD,NNS
word+2=will
word-1=old
prefix=E
prefix=El
suffix=r
suffix=er
suffix=ier

Table 4.2 The Contextual Predicates Example

Condition	Contextual Predicates
w_i is not rare	$w_i = X$
w_i is rare	X is prefix of w_i , $ X \leq 4$
	X is suffix of w_i , $ X \leq 4$
	X contains number
	X contains uppercase character
	X contains hyphen
$\forall w_i$	$t_{i-1} = X$
	$t_{i-w}t_{i-1} = XY$
	$w_{i-1} = X$
	$w_{i-2} = X$
	$w_{i+1} = X$
	$w_{i+2} = X$

Please note that if a word is rare (occurs less than 5 times in the training set) several additional contextual predicates are used to help predict the tag based on the word's form. A useful feature might be:

$$f(x, y) = \begin{cases} 1 & \text{if } y=\text{VBG} \text{ and } \textit{current_suffix_is_ing}(x) = \textit{true} \\ 0 & \text{otherwise} \end{cases}$$

and is represented as a literal string: "suffix=ing_VBG".

Appendix A

Appendix - Project schedules and milestones

APPENDIX A. APPENDIX - PROJECT SCHEDULES AND MILESTONES17

Milestone			Start	Finish	In Charge	Description	Status
1	Survey on the project		2008-02-01	2008-02-13	Vernkin	1. Have a general overview of current Chinese Segmentation techniques and their differences. 2. Select a suitable one (Character-based) and do more research on it.	Finished
2	Design the Architecture of CMA		2008-02-16	2008-02-20	Vernkin	1. Basen on WISE KMA Orange, design the Architecture for the MA Systems (CJK). 2. Moreover, design the common Architecture for the CMA. 3. Finally focus on the Architecture for the approach I selected.	Finished
3	Implement CMA and Unit Test		2008-02-23	2008-03-25	Vernkin	1. Implement the Maximum Entropy Model. 2. Implement character-based Segmentation. 3. Integrate all the components into CMA.	80% Finished

Bibliography

- [1] J.N. Darroch and D. Ratcliff. Generalized iterative scaling for log-linear models. *The Annals of Mathematical Statistics*, Vol. 43:pp 1470–1480, 1972.
- [2] Stephen Della Pietra, Vincent J. Della Pietra, and John D. Lafferty. Inducing features of random fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(4):380–393, 1997.
- [3] Julia Hockenmaier and Chris Brew. Error-driven segmentation of chinese. *Communications of COLIPS*, 1(1):69–84, 1998.
- [4] Robert Malouf. A comparison of algorithms for maximum entropy parameter estimation, 2003.
- [5] David Palmer. A trainable rule-based algorithm to word segmentation. *Proceedings of the 35th Annual Meeting of the Association of Computational Linguistics*, 1997.
- [6] Nianwen Xue. Defining and automatically identifying words in chinese. *Ph.D. thesis, University of Delaware*, 2001.