

Technical Report of *Language Analyzer*
Comprehension

Zhongxia Li

December 24, 2010

Abstract

This document is a comprehensive technical report for *Language Analyzer*, including *Chinese Morphological Analyzer*.

Language Analyzer (LA) is a essential module for text processing in SF1-R system, it's goal is to convert the input text which was filled with words and characters into more consistent *index terms*. Index terms are the representation of the content of a document that are used for indexing and searching. *LA* provides various analyzing methods for text processing, the first pass of processing is *tokenization* in which the characters are split apart into basic categories, and then other analyzing methods are applied. For Chinese language, the Chinese segment processing for detecting words and other morphological analysis, such as part-of-speech (POS), will be applied.

Changes

Date	Author	Notes
2010-12-17	Zhongxia Li	Initial version
2010-12-24	Zhongxia Li	Update the Chapter for Chinese Morphological Analyzer

Contents

1	Overview	3
1.1	What Is Language Analyzer?	3
1.2	Terms of Output	3
1.3	Architecture of Language Analyzer Modules	4
2	Tokenizer	5
2.1	Introduction	5
2.2	Options for Tokenizer	6
2.3	Setting Options in Configuration File	7
3	Analyzer	8
3.1	Introduction	8
3.2	Configuration and Options for Analyzer	8
4	Usage Of Language Analyzer	10
4.1	LA Configuration for Indexing Collection	10
4.2	Application Interface	10
5	Chinese Morphological Analyzer	12
5.1	Introduction	12
5.2	Chinese Word Segmentation	13
5.2.1	Dictionary Based Approach	13
5.2.2	Character Based Tagging	13
5.2.3	POC Tagging	13
5.3	Part-Of-Speech Tagging	15
	References	16

Chapter 1

Overview

1.1 What Is Language Analyzer?

The *Language Analyzer* is an important component for text processing in SF1-R system, it aims to convert the input text which was filled with words and characters into more consistent *index terms* by applying various analyzing methods including morpheme analyzing. The index terms are the representation of the content of a document that are used for indexing and searching. Text preprocessing is an essential pre-stage for indexing and query for a search engine.

1.2 Terms of Output

The input text of *LA* maybe consist of alphabetic, numeric or specail characters, and more often Chinese characters. Generally, *LA* generates three categories of terms as outputs.

- *Raw Term*: These are the raw terms that completely represent the original raw text.
- *Primary Term*: Within the Raw Terms, the terms that consists either alphabetic or numeric characters fall into this category.
- *Secondary Term*: The terms that are analyzed from the Primary Terms by analyzers fall into this category.

Each term is assigned a word offset corresponding to the term's relative position in the original text. The offset of expanded primary and secondary terms are set based on the raw terms.

1.3 Architecture of Language Analyzer Modules

There are two main kinds of processors for *Language Analyzer* which are Tokenizer and Analyzers, the input text is first passed through the Tokenizer and then processed using other optional analyzers such as StemAnalyzer, NGramAnalyzer, ChineseAnalyzer, etc. The basic architecture of *Language Analyzer* modules for SF1-R is show as figure 1.1. LA Manager exposes interfaces to other modules of SF1-R.

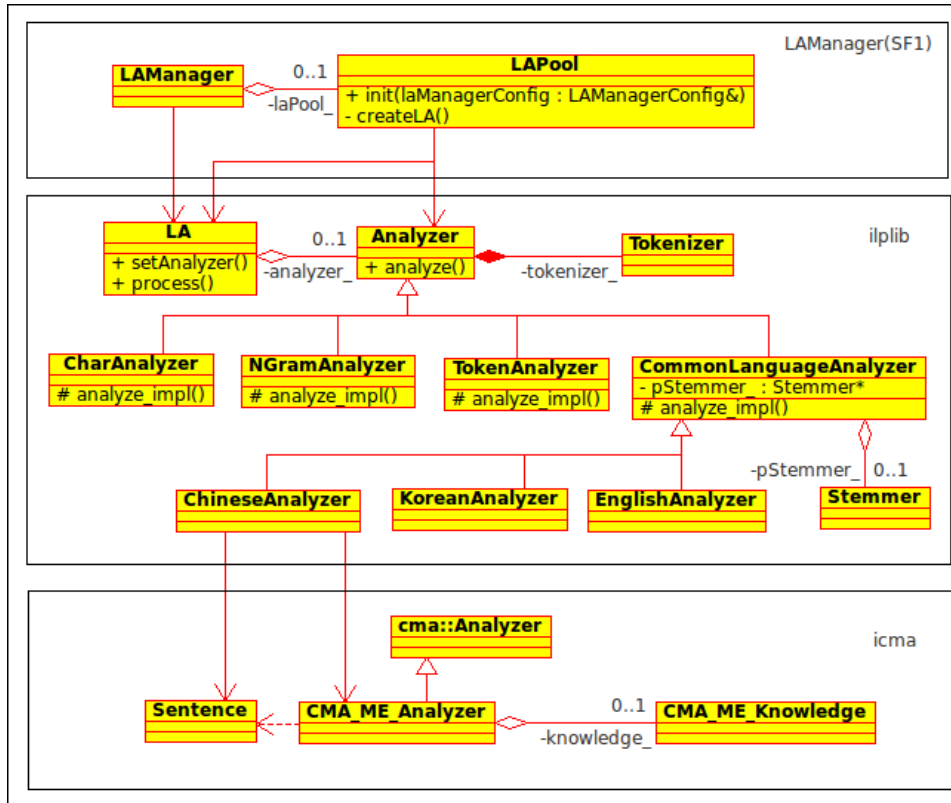


Figure 1.1: Basic Architecture of Language Analyzer Modules

For CJK family of languages, the key problem is word segmentation, where the breaks corresponding to words or terms must be identified in the continuous sequence of characters. The *Chinese Morphological Analyzer* (CMA) module implemented several Chinese segmentation algorithms which are optional, including dictionary based algorithms and char-based tagging approach. The *Maximum Entropy Module* is employed for building Chinese segment tagger, POC (Position of Character), and building part-of-speech (POS) tagger.

Chapter 2

Tokenizer

2.1 Introduction

The tokenizing is done by categorizing characters into two basic categories: delimiters (white space and special characters by default) and content characters (alphabetic, numeric and other language characters). Continuous characters in the same category are grouped into a single term.

Each term is assigned with a word offsets which shows it's relative position in the original text. Take input text "It's like A@B" for example, the output terms and their offsets are show in table 2.1.

Table 2.1: Raw Term Example

TERM	WORD OFFSET
It	0
'	1
s	2
like	3
A	4
@	5
B	6

Notice that the white space doesn't occupy a position. So the white spaces that are between "s" and "like" and between "like" and "A" are ignored.

As we can see in table 2.1, these terms are Raw Terms. The Raw Terms are used to rebuild the original text, since they exactly represented the original text when put together. The Primary Terms, in this example, are 'It' 's'

'like' 'A' and 'B' whose offsets are based on Raw Terms.

2.2 Options for Tokenizer

Table 2.1 showed an output sequence of terms for our example, however, there are several options for us to decide how the special characters (such as '[', '@']) will be parsed. The options are: *Allow*, *Divide* and *Unite*, as described in the table below.

Method	Description
Allow	The characters set as <i>allow</i> will be removed from the list of delimiters. E.g. "A@B" => "A@B"
Divide	The characters will be set to be a delimiter. E.g. "A@B" => "A", "B"
Unite	The character will be used to concatenate the two adjacent tokens on either side of the character. E.g. "A@B" => "AB"

The following shows the primary terms when different options are set for Tokenizer.

- Allow: The character is recognized as a content character.

TERM	WORD OFFSET
It's	0
like	1
A@B	2

- Divide: The character is recognized as a delimiter. Acts like the default Tokenizer.

TERM	WORD OFFSET
It	0
s	2
like	3
A	4
B	6

- Unite: The character is used to concatenate two terms on both side.

TERM	WORD OFFSET
Its	0
like	1
AB	2

2.3 Setting Options in Configuration File

We can set the options for Tokenizing in the configuration file of SF1-R as below.

```
<Tokenizing>
  <Tokenizer id="tok_divide" method="divide" value="@# $" code="" />
  <Tokenizer id="tok_unite" method="unite" value="/" code="" />
</Tokenizing>
```

In the *Tokenizing* Element of the XML file of configuration, we can edit the attributes of *Tokenizer* elements to set options. The *id* attribute is the name of Tokenizer, the *method* attribute is the option we chose, the *value* attribute indicates what character(s) to apply this method, the character(s) also can be represented as UCS2 *code* in *code* attribute.

Chapter 3

Analyzer

3.1 Introduction

When text are tokenized into a list of Raw Terms and Primary Terms, the Primary Terms can be further processed by several Filters and Analyzers. The terms created during this process are Secondary Terms. Secondary Terms have the same word offsets as the Primary Terms that they expanded from.

Various analyzing methods are available to be used in SF1-R. Some of these methods are character-based, like token, ngram, and matrix, while some are language based methods. Analyzers will process the tokens returned from the Tokenizer and extract terms from it. These terms are used for indexing and searching.

3.2 Configuration and Options for Analyzer

We can define (set) different Analyzer Methods by setting different options for each Method, the defined Methods can be applied in SF1-R. The following are the core options that are required for each **<Method>** element (except those that marked with "optional"). Other options vary depending on the *analysis*.

Element	Attribute	Description
Method	id	The name of the analyzer instance. The ID is used like variables in <Property> configurations in Collections settings, along with Tokenizer settings.

analysis		This attribute decides the analyzer type. There are two analysis categories, one in which is language independent, which are <i>token</i> , <i>ngram</i> , and <i>emphmatrix</i> . The other category is language dependent, it includes <i>English</i> , <i>Chinese</i> , <i>Korean</i> and other languages. The analysis attribute is followed by options that correspond to the option.
casesensitive	(optional)	Default is "yes". This attribute is used to set the case-sensitivity of the document field property. If the setting is turned on, the indexing and searching on the document will be done case sensitivity.
idxflag	(optional)	Default is <i>all</i> . Indexing flag to indicates which type of terms should return in the indexing. It has four values: 1) <i>all</i> returns both primary and second terms; 2) <i>prime</i> returns only primary terms; 3) <i>second</i> returns only secondary terms; 4) <i>none</i> returns neither primary nor second terms. <i>Primary Term</i> is tokenized by tokenizer and <i>Secondary Term</i> is analyzed by Specific Language Analyzer basing on the <i>Primary Term</i> .
schflag	(optional)	Default is <i>second</i> . Searching flag to indicates which type of terms should return in the searching. The detail see <i>idxflag</i> .

The following shows an example of Analyzer Methods defined in configuration file of SF1-R.

```
<LanguageAnalyzer dictionarypath="..." >
  <Method id="la_token" analysis="token"/>
  <Method id="la_ngram" analysis="ngram" min="2" max="3" maxno=
    ="2194967296" apart="n" idxflag="second" schflag="second"/>
  <Method id="inner_la_korall_mia" analysis="korean" casesensitive="<
    yes" >
    <settings mode="label" option="R1H-S-" specialchar="#" <
      dictionarypath=""/>
  </Method>
  <Method id="la_mia" analysis="multilang" advooption="default,<
    inner_la_korall_mia;cn,char" casesensitive="yes" lower="no"/>
  <Method id="inner_la_cnall_sia" analysis="chinese" casesensitive="<
    yes" >
    <settings mode="label" option="R+H+S+T3" specialchar="#" <
      dictionarypath="/home/zhongxia/codebase/icma/db/icwb/utf8"/>
  </Method>
  <Method id="la_sia" analysis="multilang" advooption="default,<
    inner_la_korall_sia;cn,ma,inner_la_cnall_sia"/>
</LanguageAnalyzer>
```

Chapter 4

Usage Of Language Analyzer

4.1 LA Configuration for Indexing Collection

For each *Property* of *Collection*, we can set *Indexing* Element to indicate which Analyzer or Tokenizer(s) will be applied to parse the *Property* content. As showed below, is an example of cofiguration for data collection *ChnWiki* indexed yb SF1-R.

```
<Collection name="ChnWiki" encoding="UTF-8" rankingmodel="ranker1" ←
  wildcardtype="unigram" generalclassifierlanguage="chinese">
  <DocumentSchema>
    <Property name="DOCID" type="string" index="no" />
    <Property name="DATE" type="string" index="no" />
    <Property name="Title" type="string" index="yes" mining="yes">
      <Indexing filter="no" analyzer="la_sia" tokenizer="" rankweight←
        ="heavy"/>
    </Property>

    <Property name="Content" type="string" index="yes" mining="yes" ←
      ">
  <Display length="160" snippet="yes" summary="no" >
    <settings summarynum="2"/>
  </Display>
  <Indexing filter="no" analyzer="la_sia" tokenizer="tok_divide, ←
    tok_unite" rankweight="normal"/>
  </Property>
</DocumentSchema>
</Collection>
```

4.2 Application Interface

The LAManager exposes interfaces of *Language Analyzer* to other modules of SF1-R, associated with the configuration setting described above. It's

easy to use, initialize LA and parse input text with specified configuration (*AnalysisInfo*), the output term list will be returned.

Application Interfaces:

```
bool LAManager::getTermList(  
    const izenelib::util::UString & text,  
    const AnalysisInfo& analysisInfo,  
    la::TermList& termList );  
  
bool LAPool::init(const sfiv5::LAManagerConfig & laManagerConfig);
```

Usage:

```
LAManagerConfig config;  
AnalysisInfo analysisInfo;  
String input_text = "content to be parsed";  
  
// setting config  
  
LAPool::getInstance()->init( config );  
LAManager laMgr_ = new LAManager();  
  
la::TermList& termList; // output  
laMgr_->getTermList(input_text, analysisInfo, termList )
```

Chapter 5

Chinese Morphological Analyzer

5.1 Introduction

For CJK family of languages, a key problem is word segmentation, where the breaks corresponding to words or terms must be identified in the continuous sequence of characters. The *Chinese Morphological Analyzer (CMA)* module implemented several Chinese segmentation algorithms which are optional, including dictionary based algorithms and char-based tagging approach. The *Maximum Entropy Module* is employed for building Chinese segment tagging, POC (Position of Character), and building part-of-speech (POS) tagger.

There are two main difficulties in Chinese segmentation, segmentation ambiguities and out-of-vocabulary (OOV) words. Practical results show that performance of statistic segmentation system outperforms that of hand-crafted rule-based systems. And 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. 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, tags that indicate the position of Chinese character within a word.

5.2 Chinese Word Segmentation

5.2.1 Dictionary Based Approach

The dictionary base algorithms for Chinese Word Segmentation is totally dictionary based matching, the algorithms we implemented are Forward Maximum Matching and Forwards Minimum Matching algorithm. For example, for Forwards Minimum Matching based approach, all the ambiguities would be contained in the result. Take the input is ABCDEF (A to H represents a character respectively), and the dictionary contains AB, ABC, DEF. The output segments will be AB, C, DEF.

5.2.2 Character Based Tagging

This section will introduce Character-based Chinese Segmentation by applying the Maximum Entropy Modeling. 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 two tags, *B* or *E* depending on its position within a word. It is tagged *B* if it occurs on the left boundary of a word, and forms a word with the character(s) on its right. It is tagged *B* if it occurs on the beginning position of a word, and tagged with *E* in other cases (In the non-beginning position of a word). We call such tags as position-of-character (POC) tags to differentiate them from the more familiar part-of-speech (POS) tags.

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; Hockenmaier and Brew, 1998; Xue, 2001] 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.

5.2.3 POC Tagging

The POC tagger here uses the same probability model as the POS tagger. The probability model is defined over $H \times T$, where H is the set of possible

contexts or "*histories*" and T is the set of possible tags. The model's joint probability of a history h and a tag t is defined as

$$p(h, t) = \pi \mu \prod_{j=1}^k \alpha_j^{f_j(h, t)} \quad (5.1)$$

Where π is a normalization constant, $\{\mu, \alpha_1, \dots, \alpha_k\}$ are the model parameters and $\{f_1, \dots, f_k\}$ are known as features, where $f_j(h, t) \in \{0, 1\}$. Each feature f_j has a corresponding parameter α_j , that effectively serves as a "*weight*" of this feature. In the training process, given a sequence of characters $\{c_1, \dots, c_k\}$ and their POC tags $\{t_1, \dots, t_k\}$ as training data, the purpose is to determine the parameters $\{\mu, \alpha_1, \dots, \alpha_k\}$ that maximize the likelihood of the training data using p :

$$L(P) = \prod_{i=1}^n P(h_i, t_i) = \prod_{i=1}^n \pi \mu \prod_{j=1}^k \alpha_j^{f_j(h_i, t_i)} \quad (5.2)$$

The success of the model in tagging depends to a large extent on the selection of suitable features. Given (h, t) , a feature must encode information that helps to predict t . The features used are instantiations of the feature templates. Feature templates (1) to (3) for common characters and (1) for special characters represent character features while (2) for special characters represents tag features. $C_{-2} \dots C_1$ are characters and $T_{-1} \dots T_1$ are POC tags.

There are two groups of feature templates basing on the context. The C represents character array, and $C[0]$ is the current character, $C[1]$ is the previous character and so on. The T represents the characters' types array (type see section-??), and $T[0]$ is the type of the current character.

The default feature set (eight features) is:

1. The single character group: (C_i) , $i \in [-2, 1]$;
2. The adjacent two characters group: (C_i, C_{i+1}) , $i \in [-2, 0]$;
3. The previous and next characters (C_{-1}, C_1) .

The feature set for special characters (like digits, letters) (four features) is:

1. The next character C_1 ;
2. The Single Type group: (T_i) $i \in [-1, 1]$.

In general, given (h, t) , these features are in the form of co-occurrence relations between t and some type of context h , or between t and some properties of the current character. For example,

$$f_i(h_i, t_i) = \begin{cases} 1 & \text{if } t_{i-1} = B \text{ \& } t_i = E \\ 0 & \text{otherwise} \end{cases}$$

This feature will map to 1 and contribute towards $p(h_i, t_i)$ if c_{i-1} is tagged B and c_i is tagged E .

The feature templates encode three types of contexts. First, features based on the current and surrounding characters are extracted. Given a character in a sentence, this model will look at the current character, the previous two and next characters. For example, if the current character is *Men* (plural marker), it is very likely that it will occur in the non-begin position of a word, thus receiving the tag E . On the other hand, for other characters, they might be equally likely to appear on the beginning of a word. In those cases where it occurs within a word depends on its surrounding characters. For example, if the current character is (“love”), it should perhaps be tagged B if the next character is *hu* (“protect”). However, if the previous character is (“warm”), then it should perhaps be tagged E . Second, for special characters, features based on the previous tags (2) are extracted. Information like this is useful in predicting the POC tag for the current character just as the POS tags are useful in predicting the POS tag of the current word in a similar context. When the training is completed, the features and their corresponding parameters will be used to calculate the probability of the tag sequence of a sentence when the tagger tags unseen data. Given a sequence of characters c_1, \dots, c_n , the tagger searches for the tag sequence t_1, \dots, t_n with the highest probability

$$P(t_1, \dots, t_n \mid C_1, \dots, C_n) = \prod_{i=1}^n P(t_i \mid h_i) \quad (5.3)$$

And the conditional probability of for each POC tag t given its history h is calculated as

$$P(t \mid h) = \frac{p(h, t)}{\sum_{t' \in T} p(h, t')} \quad (5.4)$$

5.3 Part-Of-Speech Tagging

[TODO]

Bibliography