

# A NEW METHOD FOR CHINESE TEXT CONTENT IDENTIFICATION

<sup>1</sup>LI WEIWEI, <sup>2</sup>ZHANG TAO, <sup>3</sup>LIN WEIMIN, <sup>4</sup>DENG SONG, <sup>5</sup>SHI JIAN, <sup>6</sup>WANG CHEN

<sup>1-6</sup> China Electric Power Research Institute (Nanjing), Jiangsu Province, 210003, China

E-mail: [<sup>1</sup>liweiwei@epri.sgcc.com.cn](mailto:liweiwei@epri.sgcc.com.cn), [<sup>2</sup>zhangtao@epri.sgcc.com.cn](mailto:zhangtao@epri.sgcc.com.cn),  
[<sup>3</sup>linweimin@epri.sgcc.com.cn](mailto:linweimin@epri.sgcc.com.cn), [<sup>4</sup>dengsong@epri.sgcc.com.cn](mailto:dengsong@epri.sgcc.com.cn),  
[<sup>5</sup>shijian2@epri.sgcc.com.cn](mailto:shijian2@epri.sgcc.com.cn), [<sup>6</sup>wangchen2@epri.sgcc.com.cn](mailto:wangchen2@epri.sgcc.com.cn)

## ABSTRACT

This article designed and implemented a method to identify the sensitive data in Chinese text. It can be used in data leakage prevention. There are two main innovation of this paper. One innovation is a method of preprocessing sensitive data based on statistical and the other is a method of determining the threshold based on self-learning. Experiments prove that the method is simple and practical.

**Keywords:** *Sensitive Data, Content Identification, Data Leakage Prevention*

## 1. INTRODUCTION

There are three main algorithms for text content identification: one is a classification algorithm based on probability and information theory such as Naive Bayes algorithm. Another is based on the TFIDF. The third is based on learning knowledge, such as support vector machine. The method of pretreatment of content identification was complex and lack of flexible method to determine thresholds in the past.

This paper proposed and implemented a method of identification of sensitive data based on Chinese text content. Through learning the sensitive text library and known-classification text library, we determined the threshold of the sensitive data. We analyzed the unknown-classification text and formed a feature vector, then calculated the similarity value between this vector and the feature vector of sensitive data by cosine formula. Finally, we compared the similarity value with the threshold to judge whether the text is sensitive or not. This article described a detailed account of the process and experiments showed that the method was simple, practical and had a high accuracy rate.

The content of this paper is divided into four main sections. The first chapter is the introduction, which analyzed the background and introduced the content of research and the structure of this article. The second chapter is architecture, which introduced the realization of architecture. The third chapter is functional

components, which introduced the details of process of pretreatment, text recognition and determination of the threshold. The fourth chapter is the experiment, which can verify the effectiveness of this method.

## 2. ARCHITECTURE

Generally, the process of text content identification can be divided into the following steps: first, we created the data sets, which including the training set and test set. Then, we established the model to show the text and selected the feature. Then, we established the classifier by learning from training set. Finally, we did the experiment and evaluated the performance.

In this article, the data set contained the training set and test set. The training set included sensitive text library and the known classification text library. The former was used to form classifier by learning. The latter was used to form threshold. The latter library contained two small libraries. One of the small libraries contained a certain number of sensitive texts and the other contained the same number of safe texts.

In this article, the text was mainly expressed through the way of vector space model (VSM) [6], which described the Chinese text content in the form of vector. We proposed a method of identification of sensitive data based on Chinese text content. The main architecture was shown as Figure.1.

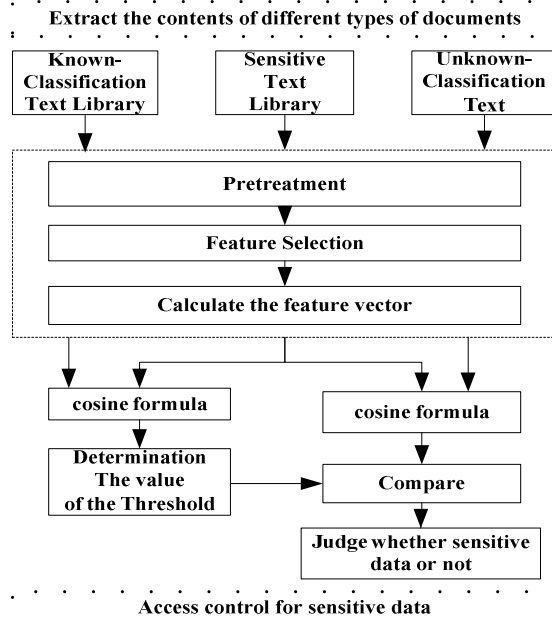


Fig.1 The Architecture Of Method To Identify Sensitive Data Based On Chinese Text Content

The main steps were shown as follow:

Step 1: The sensitive text library formed a feature vector after pretreatment, feature selection and vector space weight calculation.

Step 2: The known-classification text formed a feature vector after pretreatment, feature selection and vector space weight calculation. Then we calculated the similarity between this vector and the feature vector of sensitive data by cosine formula. By gathering statistics, determine the threshold of sensitive data.

Step 3: The unknown-classification text formed a feature vector after pretreatment, feature selection and vector space weight calculation. Then we calculated the similarity between this vector and the feature vector of sensitive data by cosine formula. We compared the value of similarity with the threshold of sensitive data to judge whether the text is sensitive or not.

### 3. FUNCTIONAL COMPONENTS

#### 3.1 Pretreatment

The first step to identify sensitive data based on Chinese text content is pretreatment. Document collection which expressed as  $T_{pre} = \{T_1, T_2, \dots, T_i\}$  was divided into individual phrases and identified the parts of speech by the interface of ICTCLAS and compiled statistics of word length and word frequency. For

example, nouns was shown as (n), verb was shown as (v), adjective was shown as (a) and other.

Text  $T_i$  was shown as follow after pretreatment:

$$T_i = ((a_{i1}, l_{i1}, p_{i1}), (a_{i2}, l_{i2}, p_{i2}), \dots, (a_{in}, l_{in}, p_{in}))$$

Where,  $T_i$  is a Chinese text,  $a_{in}$  is the phrase which is got by the interface of ICTCLAS,  $l_{in}$  is the length of  $a_{in}$  and  $p_{in}$  is the parts of speech of  $a_{in}$ .

#### 3.2 Feature Selection

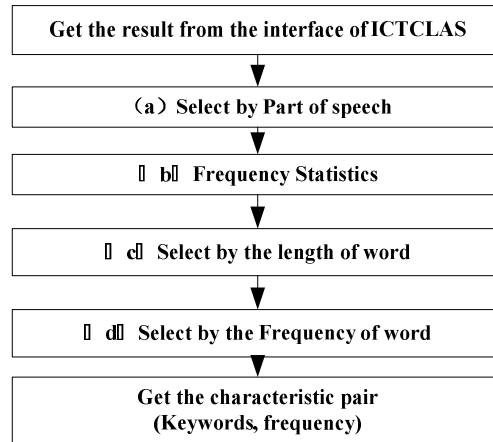


Fig.2 The Process Of Feature Selection

In the process of learning and identifying the Chinese text library, if all the parts of speech of the words were regarded as keyword, this would lead a lot of calculation and too much redundant information. Therefore, we extracted part of the result which get form the interface of ICTCLAS to reduce the dimension of vector space. This made keywords more representative and calculation more simple and effective.

The process of extraction is as follows:

(a) Select by part of speech

In the process of extraction, we only keep the keywords which can most strongly express the content of the article to eliminate redundancy. According to the statistical, the noun can express better than other parts of speech. Therefore, we only keep the noun.  $T_i$  was shown as  $Ta_i$  as follow after been selected by part of speech:

$$Ta_i = ((a_{i1}, l_{i1}), (a_{i2}, l_{i2}), \dots, (a_{in}, l_{in}))$$

Where,  $Ta_i$  is the text which was extracted to keep noun and  $(a_{in}, l_{in}) \in T_i$  in which  $a_{in}$  is noun.

(b) Frequency statistics



We gathered statistics the frequency for each keyword and added the frequency of keyword on the basis of  $Ta_i$  and show as  $Tb_i$  as follow:

$$Tb_i = ((a_{i1}, f_{i1}), (a_{i2}, f_{i2}), \dots, (a_{in}, f_{in}))$$

Where,  $Tb_i$  is the text after gathering statistics of the frequency of keyword and  $f_{in}$  is the frequency of  $a_{in}$ .

(c) Select by the length of word

In the Chinese text, the word has a stronger ability to express than the single word. We calculated the length of the keywords and delete single word, and then further expressed  $Tb_i$  as  $Tc_i$  as follow:

$$Tc_i = ((a_{i2}, f_{i2}), (a_{i3}, f_{i3}), \dots, (a_{in}, f_{in}))$$

Where,  $Tc_i$  is the text after delete single word and  $a_{i1}$  is the keyword which is not a single word.

(d) Select by the frequency of the word

In the Chinese text, the keyword which only appears once is unrepresentative and accidental. Therefore, delete these keywords. The final characteristics was expressed as follow:

$$Td_i = ((a_{i2}, f_{i2}), (a_{i3}, f_{i3}), \dots, (a_{in}, f_{in}))$$

Where,  $Td_i$  is the text after gathering statistics of the frequency of keyword and for which  $f_{in} > 1$ .

### 3.3 Calculate The Feature Vector

(1) Calculate the feature vector of sensitive data.

After pretreatment and feature selection, the collection of sensitive text was shown as:  $T = \{Td_1, Td_2, \dots, Td_n\}$ . Where,  $Td_i$  was shown as follow:

$$Td_i = ((a_{i2}, f_{i2}), (a_{i3}, f_{i3}), \dots, (a_{in}, f_{in}))$$

Where,  $Td_i$  is the text after gathering statistics of the frequency of keyword and for which  $f_{in} > 1$ .

Calculate the weight of the keyword is the effective way to measure the characteristic value. Currently, the TF-IDF formula is the most widely used which based on statistical methods. This formula proved to be feasible and effective in a large number of actual uses. TF-IDF formula which is commonly used is expressed as follow:

$$d_{ij} = t_{ij} * \log(N/n_j)$$

Where,  $t_{ij}$  is the number of occurrences of  $a_{ij}$  in  $T^i$ ,  $N$  is the total number of documents and  $n_j$  is the number of  $a_{ij}$  include in the document library.

After calculate the weight of the keyword, feature vector of sensitive data library is expressed as:

$$V = ((a_{11}, d_{11}), (a_{12}, d_{12}), \dots, (a_{1m}, d_{1m}), \dots, (a_{n1}, d_{n1}), (a_{n2}, d_{n2}), \dots, (a_{nm}, d_{nm}))$$

Where,  $d_{nm}$  is the weight of  $a_{nm}$  calculated by TF-IDF formula for sensitive data library. Simply expressed as:

$$V = (d_{11}, d_{12}, \dots, d_{1m}, \dots, d_{n1}, d_{n2}, \dots, d_{nm})$$

(2) Calculate the feature vector of known-classification text and unknown-classification text in library.

According to the sensitive data feature vector  $V$ , we calculated the weight of the keywords  $a_{ij}$  respectively and get the feature vector as follows:

$$V' = ((a_{11}, d'_{11}), (a_{12}, d'_{12}), \dots, (a_{1m}, d'_{1m}), \dots, (a_{n1}, d'_{n1}), (a_{n2}, d'_{n2}), \dots, (a_{nm}, d'_{nm}))$$

Where,  $a_{mn}$  in  $V'$  is equal to  $a_{nm}$  in  $V$  and  $d'_{nm}$  is the weight of  $a_{nm}$  calculated by TF-IDF formula for known-classification text. Simply expressed as:

$$V' = (d'_{11}, d'_{12}, \dots, d'_{1m}, \dots, d'_{n1}, d'_{n2}, \dots, d'_{nm})$$

We used the same method to get the feature vector of unknown classification document and expressed as:

$$V'' = (d''_{11}, d''_{12}, \dots, d''_{1m}, \dots, d''_{n1}, d''_{n2}, \dots, d''_{nm})$$

### 3.4 Calculation Of Cosine

We used the cosine formula to calculate the similarity between two feature vectors, which is shown as follows:

$$\cos \theta = \frac{V_1 \cdot V_2}{\|V_1\| \|V_2\|}$$

Where,  $V_1$  and  $V_2$  is feature vector of the document;  $V_1 \cdot V_2$  dot product of standard vector and define as  $\sum_{i=1}^n V_{1i} V_{2i}$  and the norm in the denominator  $\|V_1\|$  is defined as  $\sqrt{V_1 \cdot V_1}$ .

### 3.5 Determination Of The Threshold

The threshold of sensitive data was important to determine whether the document was sensitive or not by comparing it with the results of the cosine similarity calculation.

The process is shown as follows:

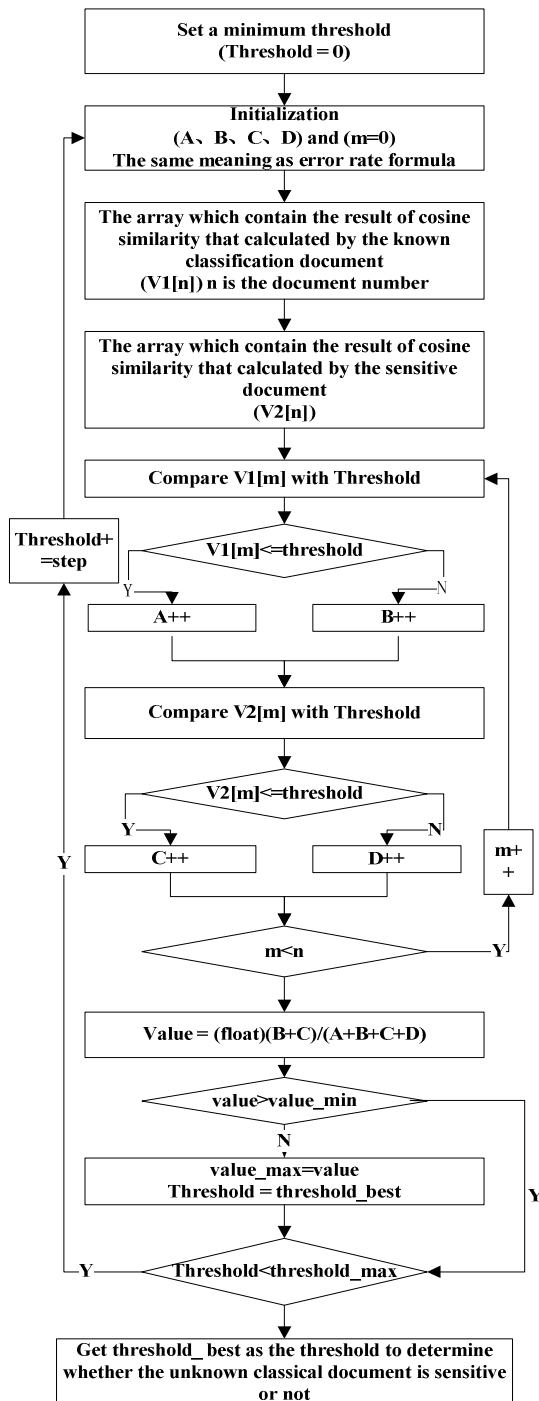


Figure 3 The Code To Determine The Threshold

In this article, we learned from the known classification documents to determine the threshold.

We collected hundred and fifty security documents and the same number of sensitive documents. And then, we calculated cosine similarity with the sensitive feature vector and got the value respectively. By setting the same step interval to threshold, we got the threshold which had the lowest error rate. The threshold would be used to determine whether the unknown classification document was sensitive or not.

The error rate is shown as follows:

$$rate = \frac{(B+C)}{(A+B+C+D)}$$

Where, A is the number of documents that are correctly identified as secure document; B is the number of documents that are mistakenly identified as sensitive document; C is the number of documents that are mistakenly identified as security document and D is the number of documents that are correctly identified as sensitive documents.

### 4. FUNCTIONAL COMPONENTS

(1) The data sets

This system can achieve identification of sensitive data in different environments by changing the training set.

We created a sensitive data identification system about education and selected 1600 documents about education from SogouC.reduced.20061102 Corpus as sensitive documents library in the experiment. We selected 100 education-related documents and 100 documents about other area randomly from the corpus for the known-classification documents library.

(2) Pretreatment and feature selection

In the process of feature selection, statistics analysis shows that: the proportion of the remaining keyword is about 32% after selecting by the part of speech, about 28% after selecting by the word length, about 11% after selecting by the word frequency. The process of feature selection greatly reduces the redundancy of the keywords and simplifies the calculation.

(3) The process of calculating the feature vector

According to the keywords, we get the feature vector of sensitive data  $V$  by the TFIDF algorithm. By the same way, we got the known classification text feature vector  $V'$  and the unknown classification text feature vector  $V''$ .

(4) The process of calculating cosine similarity value

We calculated the feature vector  $V'$  of the known classification documents respectively and calculated the similarity between  $V'$  and  $V$ . They were sorted after statistics in the table below in which the horizontal axis represented the text number and the vertical axis represented the cosine similarity value.

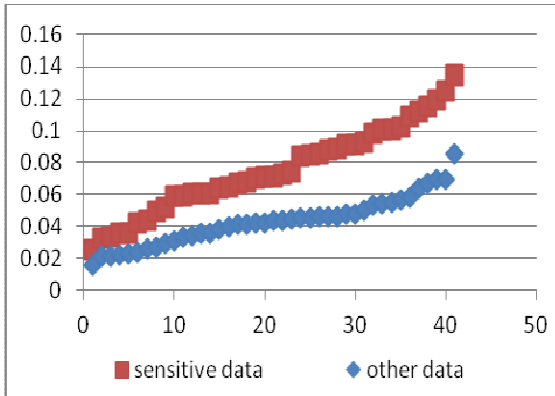


Table 1 Count Up Cosine Similarity Value

(5) The process of determining the threshold

From the bottom of the value to the top of the value, we added an interval length step by step and calculated the error rate of the threshold. Through learning the error rate of threshold, we got the threshold which has the lowest error rate. Through the experimental data, when the value of threshold is 0.060, we got the lowest error rate of 22.1%. The statistical results were illustrated below in which the horizontal axis represented the value of threshold and the vertical axis represented the error rate, recall rate and accuracy.

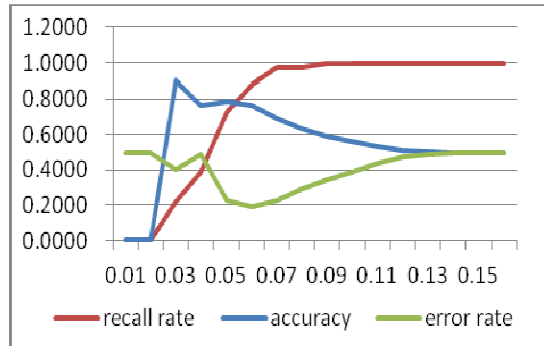


Table 2 Through learning the error rate of threshold get the threshold which has the lowest error rate

(6) The experimentation to identify the unknown classification document

We took forty unknown classified documents from the library and judged whether the document was sensitive or not by comparing cosine similarity value with the threshold. We pre-processed and analyzed the unknown classified documents library. And then, we get feature vectors based on the keywords of sensitive feature vectors. We identified the unknown classified document by comparing the cosine similarity value with threshold 0.061 which was the threshold with lowest error rate 16.52%. The experimental results show that this method is effective and feasible.

5. CONCLUSION

This article designed and implemented a method of preprocessing sensitive data based on statistical and a smart self-learning method to determine the threshold which was used to judge whether the document is sensitive or not. We described the methods and processes to achieve in detail and gave the results of the experiment.

6. ACKNOWLEDGEMENT

This work was supported in part by the information security technology research team of state grid corporation of china. In addition, I would like to offer special thanks to Zhang Tao, Lin Weimin and Deng Song for writing and experiment guidance.

REFERENCES:

[1] Lin Zhenbiao. Research on Key Technologies of File Network Leakage Prevention Based on Data stream Analysis [D]. The PLA Information Engineering University, 2009.



- [2] Yiming Yang. An Evaluation of Statistical Approaches to Text Categorization[J]. Information Retrieval, 1999,1(1-2) .
- [3] Yuan Wensheng,Wang Xiaofeng. Research on Chinese Maritime Multi-class Text Classifier Based on Na ve Bayes[J]. computer and modernization, 2011,(05)
- [4] Shi Congying, Xu Chaojun,Yang Xiaojiang. Study of TFIDF algorithm[J] Journal of Computer Applications,2009,29:168-180.
- [5] Zhang Xiaoyan, Li Qiang. The Summary of Text Classification Based on Support Vector Machines,2008,28:344-345.
- [6] Christopher J.C. Burges. A Tutorial on Support Vector Machines for Pattern Recognition[J]. Data Mining and Knowledge Discovery, 1998,2(2) .
- [7] Pang Jianfeng , Bu Dongbo , Bai Shuo. Research and Implementation of Text Categorization System Based on VSM[J]. Application Research of Computers , 2001, (9) :23-26