

The method for detecting plagiarism in a collection of documents

Natalya Shakhovska, Iryna Shvorob

Abstract - The development of the intelligent system for searching for plagiarism by combining two algorithms of searching fuzzy duplicate is considered in this article. This combining contributed to the high computational efficiency. Another advantage of the algorithm is its high efficiency when small-sized documents are compared. The practical use of the algorithm makes it possible to improve the quality of the detection of plagiarism. Also, this algorithm can be used in different systems text search.

Keywords - Keywords: intellectual system, searching for plagiarism, signatures, data space.

I. INTRODUCTION

Nowadays, the Internet is the biggest source of information. Now, people can easily search, get access and browse the web to get the information they need. Just imagine how difficult it would be to do scientific research without the Internet and web space. Furthermore, due to the size and digital structure of the internet, it is easy to illegally use someone else work now.

The problem of plagiarism has a direct relationship with the scientific community. The most common plagiarism is written text document which is formed by copying some or all parts of the original document, sometimes with some modifications. Identification of documents which were copied is stressful and time-consuming process to humans due to the large number of documents which have to be analyzed. The documents in digital format make the process of plagiarism quite simple, it means that such cases of plagiarism can be traced automatically.

Plagiarism detection depends on many factors[5]. The first factor is the presentation of the document, which essentially covers the characteristics of the document as a preliminary step to compare [7, 8, 9]. These representations include model of identification tags, N-grams, probabilistic models, algorithms "scales" and others. Most of these representations work well in detecting verbatim plagiarism, but are vulnerable to identify complex patterns of plagiarism.

The second factor is a similarity and a measure of proximity, which is used to calculate the similarities or differences between sentences. Given the behavior of plagiarists, which usually includes insertion, deletion or substitution of words necessary to determine which activities are best for detection of plagiarism.

With the development of information systems the number of areas to identify plagiarism text only increased. This is the area of scientific papers, various publications in the field of journalism, fiction genres [13 – 15].

Currently, there are many methods and algorithms that can detect plagiarism of text objects. But over time, there are new challenges associated with the development of information systems. These tasks require more qualitative and more accurate detection of plagiarism in text.

Purpose of this paper is to improve the efficiency and quality of plagiarism detection in text objects by the use of the combined algorithm.

II. THE INTELLIGENT SYSTEM OF DETERMINE THE DEGREE OF RESEMBLANCE OF THE TEXTS

Two algorithms of searching for fuzzy duplicates named Lex Rand and Opt Freq are used in developing the system to search for plagiarism [3].

Lex Rand algorithm implemented in the following way. At first, the dictionary for the collection is created and the words with the largest and smallest values of IDF are removed. Then based on the dictionary generated 10 additional dictionaries that contain approximately 30% fewer words than the original. The words are removed at random.

11I-Match signatures are built for each document. Documents which have at least one the same signature considered duplicate. Such approach greatly increases the fullness of duplicate detection when the relative accuracy is reduced by only 14% [1,3].

Opt Freq algorithm implements the method of "optimal search frequency" and its used to search for similar documents in a wide range of applications, from web to clustering news. The gist of it is this. Instead of classical metrics TF*IDF a modified version of it is proposed. We introduce a heuristic concept of "optimal frequency" for the word "equal" $\ln\left(\frac{10}{1000000}\right) = 11.5$ which means

"the optimal" entering of word in 10 documents from 1000000. If the real value of IDF is less than "optimal", then it slightly (by law parabola) rises to

$IDF_{opt} = \sqrt{\frac{IDF}{11.5}}$, and if it is greater it significantly (as hyperbole) reduces to

$$IDF_{opt} = \sqrt{\frac{11.5}{IDF}} \quad (1)$$

For the collection the dictionary is created. This dictionary puts every word in accordance with the number

Natalya Shakhovska, Iryna Shvorob - Lviv Polytechnic National University, S. Bandery Str., 12, Lviv, 79013, UKRAINE,
E-mail: natalya233@gmail.com, irka.shvorob@gmail.com

of documents in which this word occurs at least once (df). Then the frequency dictionary for document is built and the "weight" wt of each word is calculated by the formula:

$$wt = TF * IDF_{opt}, \quad (2)$$

where

$$TF = 0.5 + 0.5 * \frac{tf}{tf_{max}}, \quad (3)$$

$$IDF = \log \left(\frac{df}{N} \right), \quad (4)$$

$$IDF_{opt} = \begin{cases} \sqrt{\frac{IDF}{11.5}}, & IDF < 11.5 \\ \frac{11.5}{IDF}, & IDF \geq 11.5 \end{cases} \quad (5)$$

tf (term frequency) is the ratio of occurrences of a word to the total number of words of the document. Thus, the estimated importance of words within a single document:

$$tf = \frac{n_i}{\sum_k n_k} \quad (6)$$

where n_i is the number using the word in a document, and the denominator – the total number of words in this document.

df (inverse document frequency) is inversion frequency with which a certain word is found in the documents collection. Consideration df reduces weight widely used words:

$$df = \log \frac{|T|}{|T_i \supset a_i|} \quad (7)$$

where $|T|$ is count of text documents in collection; $|T_i \supset a_i|$ is count of text documents, where word a_i occurs (where $n_i \neq 0$).

Then the 6 words with the largest values of wt are selected and concatenated in alphabetical order into the string. The check sum of the resulting line is calculated as the signature of document [3, 6].

Also, it is very important, where part of text is arisen [10, 11].

First of all, we introduce the concept of weight sentence.

$$Location = \frac{1}{n \cdot m} \quad (8)$$

where $n = \overline{1..3}, m = \overline{1..3}$ – the place calls to the main part and paragraph respectively. Begin and end of text or paragraph estimated value of 1, the middle is as 3. Coefficient key phrase is determined by entering the sentence U of elements of a set of significant sentences from A membership function:

$$Cuephrase = \mu_A(U) \quad (9)$$

$A = \{ \langle \text{«Conclusion»}, \langle \text{«In the end»}, \langle \text{«By the way»} \dots \} \}$.

Index of statistical significance is formed on the basis of visiting sentence key-words specified by the author of the article:

$$Statterm = \mu_K(U) \quad (10)$$

The value added is defined as the presence of terms related words sentences that appear in the article's headline to the total number of words in a sentence (words) except for words whose length is less than 3 characters:

$$Addterm = \frac{\text{word}}{\text{words}} \quad (11)$$

The weight of text block U is:

$$Weight(U) = Location(U) + Cuephrase(U) + Statterm(U) + Addterm(U) \quad (12)$$

So after being allowed to study all the documents necessary to accomplish the following: to exclude a statement that its content has hit the consolidated data repository and perform the final sorting sentences. For the task of bringing to the final ranking factor "information novelty" use the following method:

– Let we have two sets of sentences $B = \emptyset$ and $A = \{A_i | i = 1, 2, \dots, N\}$, N is count of sentences in text. For every sentence A_i the usefulness $P(i)_i$ q_i : $P(i)_i = q_i, i = 1, 2, \dots, N$ (13)

– The sentences from set A sort Descending $P(i)_i$

– If A_i has the biggest $P(i)_i$, we take it in B. The usefulness for sentences in A set s

$$P(i) = \frac{P(i)}{k q_i} \quad (14)$$

where $k > 0$ – factor clipping similar sentences.

– Is A empty? If NOT, go to 1.

The next problem is information estimating from different sources [15- 16]. For semi-structured data type text file with a known format - dictionary data types defined formatting released the text of the formatting, copying its contents:

$$object \rightarrow Find \left(\pi_{formattype} \left(\sigma_{object} (Dic) \right) \right) \quad (15)$$

```
foreach object
Selection
. ParagraphFormat.Alignment = Left (1,
formattype)
. Font.type = Mid (formattype, 3, 1)
. Font.Caps = Right (formattype, 1)
InStr (1, . Text, Right (formattype, 2);
Copy.Selection
```

III. THE SYSTEM ARCHITECTURE

To build an information system model is used CASE-tool AllFusion Erwin Data Modeler, which enables model based infological model of information system build its datalogical model and create a database in any database management system. The development of the summarization system provides in the notation IDEF1X.

During the implementation of systems analysis for this area following charts were developed[4]:

1) IDEF0-diagram for subtasks of the main business process (figure 1);

2) IDEF3-diagram for the block "Choosing the algorithm of working with words" (figure 2).



Fig. 1. IDEF0-diagram for subtasks of the main business process

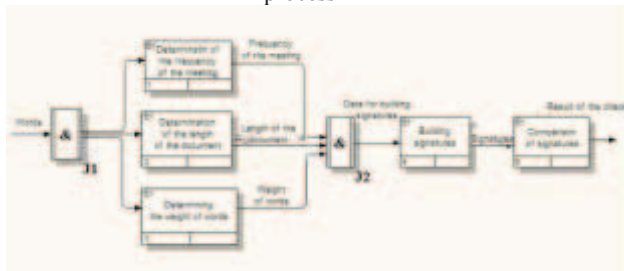


Fig. 2. IDEF3-diagram for the block "Choosing the algorithm of working with words"

The work "Checking for plagiarism in the text" is divided into 6 works: "Preparing of the text document", "Supplement Knowledge Base", "Formation of requirements to the search process", "Introduction additional set of data", "Selecting the search method", "Introduction to expert estimates". These works are carried out in the system sequentially, one after another. A text document which gets into the system due to user actions is applied to the input to the work "Preparation of the text document". Text information namely data entered in the system after this work is the result and therefore the input information for the work "Supplement knowledge base". This information is converted into data format suitable for the system in which they are ready for further processing. Checking words of text is carried out as a result of the "Supplement Knowledge Base". The result of this work is a set of sentences which will be applied to the input of the "Formation of requirements to the search process" for further processing and the input of "Introduction additional set of data". These works are carried out the text processing. Words and formed signatures are the result of the work "The introduction additional set of data". They apply to the input of "Selecting the search method" and then searching for plagiarism is carried out. The next work "Introduction to expert estimates" provides the end result — a numeric value of the searching for plagiarism.

The IDEF3-diagram for the block "Choice of algorithm with the words" is consists of such units of work: "Determination of the frequency of the meeting" (determines the number of meeting of words in the text, returns the number of meeting of words in the document), "Determination of the length of the document" (determines the length of the document), "Determining the weight of words" (the data obtained in previous studies and knowledge base are used and keywords of the text are assigned of weight), "Building signatures" (connecting words into signatures), "Comparison of signatures" (checking of signatures, a collision of hash codes takes place).

Figures 3a and 3b shows the software implementation of the developed system and the results of its implementation.

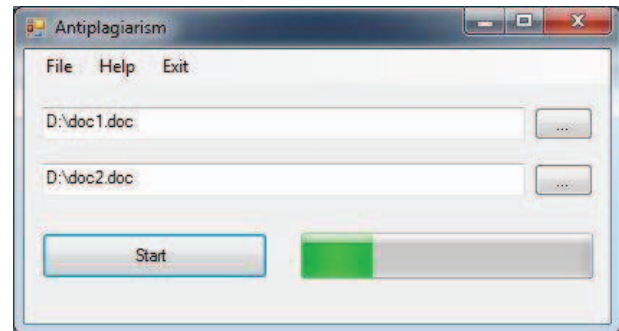


Fig. 3a. The test example of the program

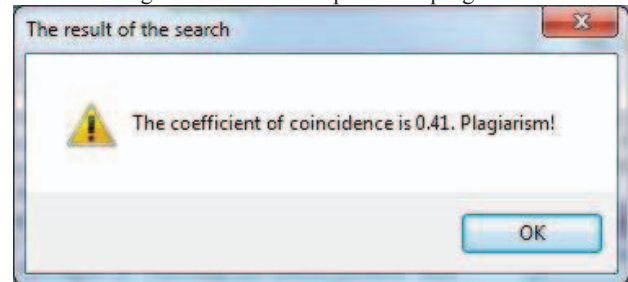


Fig. 3b. The results of the program

Before checking for plagiarism the system must makes its preliminary treatment.

This stage applies to all requested documents, as well as primary documents. There are four steps at this stage:

- remove stop words;
- fragmentation of the text;
- tokenization of the text;
- selection of roots of words.

IV. CONCLUSION

The problem of plagiarism has been established and discussed. Two algorithms for finding fuzzy duplicates were considered and incorporated. The verification of the considered algorithms and combined algorithm was made. The received data are shown in Table 1. It is worth noting, as a result of combining it has improved its performance and result of searching of fuzzy duplicates.

System analysis for the intelligent system of determines the degree of resemblance of the texts was carried out and two charts were developed.

Basic steps for preprocessing the text were identified.

TABLE 1

THE RESULTS OF VERIFICATION OF ALGORITHMS

	LEX RAND ALGOR ITHM	ALG ORITHM OPT FREQ	COM BINED ALGOR ITHM
THE ACCURACY OF SEARCHING OF PLAGIARISM	39	59	61
A VALIDATION (SEC.)	35	41	39

Obviously, the algorithm is not perfect in solving the problem of determining fuzzy duplicates. In order to improve options for combining multiple algorithms.

For example, using the method of "descriptive words" can determine what class includes documents are scanned as each generated vector uniquely identifies this class. Then identify duplicates in a particular class of documents, signatures using methods based on the analysis of special characters. In this case, the possible increase effectiveness duplicate determination in a particular class of documents.

Duplication of texts in information flows is not always a negative phenomenon in terms of the user who uses the Internet for business purposes. An example of such an exception, for example, ranking brand when republication counts the number of press releases. Also you can use a number of overlapping signs "measure of importance" of a message and more.

REFERENCES

- [1] Park S.-T. Analysis of Lexical Signatures for Finding Lost or Related Documents / S.-T. Park, D. Pennock, C. Lee Giles, R. Krovetz. — Finland, 2002. — 8p.
- [2] Nikol'skij Ju.V. Sistemi shchynogo intelektu / Ju.V. Nikol'skij, V. V.Pasichnik, Ju. M. Shherbina. — L'viv: Vidavnytvo «Magnolija – 2006», 2010. — 279 s.
- [3] Zelenkov, Ju. G, Segalovich, I. V. Sravnitel'nyj analiz metodov opredelenija nechetkih dublikatov dlja Web-dokumentov / Ju. G. Zelenkov, I. V. Segalovich // Devjataja konferencija KSB". 2007.
- [4] Katrenko A. V. Systemnyj analiz: pidruchnyk z gryfom MON / Katrenko A. V. — L'viv : «Magnolija-2006», 2009. — 352 s.
- [5] Maurer H., F. Kappe, B. Zaka. Plagiarism – A Survey. Journal of Universal Computer Sciences, vol. 12, no. 8, pp. 1050 – 1084, 2006.
- [6] The Open Archives Initiative Protocol for Metadata Harvesting Protocol Version 2.0 of 2002-06-14. <http://www.openarchives.org/OAI/2.0/openarchivesprotocol.htm>
- [7] Гриценко В.И. Информационные технологии: тенденции, пути развития / В.И.Гриценко, А.А.Урсатьев //Управляющие системы и машины, № 5, С.3-20 (2001) (in Ukrainian)
- [8] Крайовський В.Я. Основні підходи до розроблення програмного комплексу автоматичного реферування текстових документів // Крайовський В.Я., Литвин В.В., Шаховська Н.Б. // Збірник наукових праць НАН України/ Інститут проблем моделювання в енергетиці. — №51. — Київ, 2009. — С. 178-186 (in Ukrainian)
- [9] Park S.-T. Analysis of Lexical Signatures for Finding Lost or Related Documents / S.-T. Park, D. Pennock, C. Lee Giles, R. Krovetz. — Finland, 2002. — 8p.
- [10] Kolcz A. Improved Robustness of Signature-Based Near-Replica Detection via Lexicon Randomization / A. Kolcz, A. Chowdhury, J. Alspector. — KDD,2004.
- [11] Andrei Z. Broder. Identifying and Filtering Near-Duplicate Documents, COM'00 // Proceedings of the 11th Annual Symposium on Combinatorial Pattern Matching, 2000. — P. 1-10.
- [12] Berson Thomas A. Differential Cryptanalysis Mod 232 with Applications to MD5. EUROCRYPT. — <http://dl.acm.org/citation.cfm?id=1754956>.
- [13] Hahn U. The Challenges of Automatic Summarization/ U. Hahn, I. Mani // Computer.- 2000.- vol.33.- №11.- P. 29–36.
- [14] Document Understanding Conferences (DUC) : Web site, 2008. -Режим доступу: <http://duc.nist.gov>. 15.10.2011.
- [15] Алыгулиев Р.М. Автоматическое реферирование документов с извлечением информативных предложений // Вычислительные технологии. — 2007. — Т. 12, № 5. — С. 5–15. (in Russian)
- [16] Yang Ch.C. Fractal Summarization for Mobile Devices to Access Large Documents on the Web/ Ch.C. Yang, F.L. Wang // Proc. of the WWW2003, May 20-24, 2003, Budapest, Hungary. P. 26–31 .