

**DOCUMENT PLAGIARISM DETECTION ALGORITHM USING
SEMANTIC NETWORKS**

AHMED JABR AHMED MUFTAH

**A project report submitted in partial fulfillment of the
requirements for the award of the degree of
Master of Science (Computer Science)**

**Faculty of Computer Science and Information Systems
Universiti Teknologi Malaysia**

NOVEMBER 2009

Dedicated to my father, who taught me that anything is possible to achieve; the only limitations in our lives are those that we impose on ourselves.

ACKNOWLEDGEMENT

Most of my gratitude goes to my supervisor Assoc. Prof. Dr. Naomie. Her patience and considerate nature made her accessible whenever I needed her assistance. I indeed thank her for showing me how to identify interesting problems and how a research can be started and finished correctly.

I acknowledge that my UTM colleagues are the greatest. My especial thank to Omar and Mohammed Hakami for their unrelenting encouragement during project-1. Also many thanks to Ali Alfaris, Amjad Esmail, Murad Rassam, Yassir, and Falah for helping me retain some sanity to get this work done.

Last but not less, I thank my beloved brothers Esam, Nasser, Hesham and Hussam for their ultimate support during the course of my study especially in my final semester. Hey guys thanks for everything.

ABSTRACT

The vast increase of available documents in the World Wide Web (WWW) and the ease access to these documents has lead to a serious problem of using other's works without giving credits. Although many methods have been developed to detect some instances of plagiarism such as changing the structure of sentences or when slightly replacing words by their synonyms, it is often hard to reveal plagiarism when the copied sentences are deliberately modified. This project proposes an algorithm for plagiarism detection over the Web using semantic networks. The corpus of this study contains 610 documents downloaded from the Web, 10 of those were selected to be the source of 20 manually plagiarized documents. The algorithm was compared to N-grams representation and the achieved results show that an appropriate semantic representation of sentences derived from WordNet's relations outperforms N-grams with different similarity measures in detecting the plagiarized sentences. It also show that a proposed method based on extracting named entities and common nouns is in-general capable for retrieving the source documents from the Web using a search engine API when sentences are being moderately plagiarized.

ABSTRAK

Peningkatan keluasan sedia dokumen-dokumen di World Wide Web (WWW) dan kemudahan akses kepada dokumen-dokumen ini telah menyebabkan masalah yang serius dengan menggunakan karya-karya lain tanpa memberikan kredit. Walaupun banyak kaedah telah dibangunkan untuk mengesan beberapa kes plagiarisme seperti menukar struktur kalimat ataupun ketika sedikit menukar kata dengan mereka sinonim, sering sukar untuk mendedahkan plagiarisme ketika menyalin kalimat-kalimat yang sengaja diubahsuai. Projek ini mencadangkan sebuah algoritma untuk mengesan plagiarisme melalui Web menggunakan rangkaian semantik. Korpus kajian ini mengandungi 610 dokumen-download dari Web, 10 daripada mereka yang terpilih untuk menjadi sumber secara manual menjiplak daripada 20 dokumen. Algoritma ini dibandingkan dengan N-gram representasi dan keputusan yang dicapai menunjukkan bahawa representasi semantik yang tepat dari kalimat yang berasal dari hubungan WordNet melebihi N-gram dengan berbagai ukuran persamaan dalam mengesan menjiplak kalimat. Hal ini juga menunjukkan bahawa kaedah yang dicadangkan berdasarkan pada ekstraksi bernama entiti dan kata benda umum adalah jeneral yang mampu untuk mengambil dokumen-dokumen sumber dari Web menggunakan enjin pencari API ketika kalimat-kalimat yang sedang sedang dijiplak.

TABLE OF CONTENTS

CHAPTER	CONTENT	PAGE
	TITLE	i
	DECLARATION	ii
	DEDICATION	iii
	ACKNOWLEDGEMENT	iv
	ABSTRACT	v
	ABSTRAK	vi
	TABLE OF CONTENTS	vii
	LIST OF TABLES	x
	LIST OF FIGURES	xii
	LIST OF APPENDICES	xv
1	INTRODUCTION	1
	1.1 Introduction	1
	1.2 Problem Background	4
	1.3 Problem Statement	5
	1.4 Project Objectives	5
	1.5 Project Scope	6
	1.6 Project Justification	6
	1.7 Report Organization	7
2	LITERATURE REVIEW	8
	2.1 Introduction	8

2.2	Document Plagiarism	9
2.3	Plagiarism Detection Methods	10
	2.3.1 Detection based on Stylometry Analysis	11
	2.3.2 Detection based on Documents Comparison	12
	2.3.2.1 Semantic-Based Detection	12
	2.3.2.2 Syntactic-Based Detection	14
2.4	Existing Web-Based Plagiarism Detection Tools	16
2.5	Semantic Networks	20
2.6	Document Preprocessing	23
	2.6.1 Tokenization	24
	2.6.2 Stop-word Removal	24
	2.6.3 Stemming	24
	2.6.4 Document Chunking	25
2.7	Document Representations & Similarity Measures	28
	2.7.1 Semantic Based-Representation	29
	2.7.2 Syntactic Based-Representation	34
	2.7.2.1 Fingerprinting	34
	2.7.2.2 Term Weighting Schemes	37
	2.7.2.3 N-Grams	38
2.8	Algorithms for Approximate Similarity	40
	2.8.1 Signature Scheme Algorithms	40
	2.8.2 Inverted Index-Based Algorithms	47
2.9	Discussion and Summary	51
3	METHODOLOGY	53
3.1	Introduction	53
3.2	Operational Framework	54
	3.2.1 Initial Study and Literature Review	55
	3.2.2 Corpus Preparation	55
	3.2.3 Document Preprocessing	57
	3.2.4 Applying Plagiarism Detection Techniques	58
	3.2.4.1 Semantic Relatedness Approach	58

3.2.4.2	N-grams Approach	66
3.2.5	Web Document Retrieval	70
3.2.6	Implementation	73
3.2.7	Findings Evaluation	75
4	EXPERIMENTAL RESULTS	78
4.1	Introduction	78
4.2	Information about the Corpus	79
4.3	Sentence-to-Sentence similarity	82
4.3.1	N-grams Approach	82
4.3.2	Semantic Relatedness Approach	85
4.4	Results and Comparisons	97
4.4.1	Results of Corpus Sentence Retrieval	97
4.4.2	Results of Web Document Retrieval	102
4.4.3	Comparison with Existing Tools	111
4.5	Discussion and Summary	115
5	CONCLUSION	117
5.1	Introduction	117
5.2	Achievements and Constraints	118
5.3	Future work	119
5.4	Summary	121
	REFERENCES	122
	APPENDICES A-E	128-150

LIST OF TABLES

TABLE NO.	TITLE	PAGE
2.1	Properties of some existing plagiarism detection tools based on [59]	19
2.2	Some of the relations between concepts in WordNet (N=noun, V=Verb, Adj=adjective, Adv=adverb)	21
2.3	Statistics about WordNet 2.1	22
2.4	Common similarity measures between binary vectors	39
2.5	Common similarity measures between sets	39
2.6	Common factors that influence the performance of Inverted index algorithms	48
2.7	Signature-based versus Inverted index-based algorithms	52
3.1	Integrated libraries in the project and their roles	73
4.1	Number of plagiarized sentences in documents pairs (query-vertical)/(source -horizontal)	80
4.2	Statistics about the corpus and query documents	81
4.3	Statistics about part-of-speech tagging	81
4.4	Part-of-speech tagging of s_1 and s_2	86

4.5	Shortest path between word pairs in the joint set and $T2$ (“-1” no path exists, “=” equals, “?” not of the same part of speech)	91
4.6	Subsumer depth between word pairs in the joint set and $T2$ (-1 no depth exists, = equals, ? not of the same part of speech)	91
4.7	Word-to-word similarity between the joint set and $T2$	92
4.8	Raw semantic and order vectors for $T1$	93
4.9	Raw semantic and order vectors for $T2$	94
4.10	Information contents of word in the joint set	95
4.11	N-grams recall rate in 610 corpus documents with 0.5 cutoff threshold	97
4.12	Recall rate when increasing number of documents with 0.5 cutoff threshold	99
4.13	Precision, Recall, and Harmonic Mean (F-measure) in 110 corpus documents with 0.5 cutoff threshold	100
4.14	Recall rate across similarities in 110 corpus documents	100
4.15	Semantic-R recall rate in 110 corpus documents with Alpha=0.2, Beta=0.45 and 0.8 cutoff threshold	102
4.16	Results of using 3-grams searching with 64 results/query limit	104
4.17	Results of using weighted 3-grams with 64 results/query limit	105
4.18	Results of using selective searching with 64 results/query limit.	108
4.19	Results of using 3-grams searching with 8 results/query limit.	109
4.20	Results of using weighted 3-grams with 8 results/query limit.	110
4.21	Results of using selective searching with 8 results/query limit.	111

LIST OF FIGURES

FIGURE NO.	TITLE	PAGE
1.1	Hierarchical semantic knowledge base[57]	3
2.1	Taxonomy of plagiarism detection methods	10
2.2	An example HTML report generated by DocCop	16
2.3	A sample report with timeline returned from Plagium	17
2.4	The interface of EVE2 for Web searching	18
2.5	An example of a synset in WordNet	20
2.6	Bipolar adjective structure (\rightarrow = similarity, $--\rightarrow$ =antonymy)	22
2.7	N-Unit non-overlapped chunking strategy with N=5	26
2.8	N-Unit chunking with K-overlap where N=5 and K=2	26
2.9	An example given by [55] to illustrate the different between most specific subsumers in WordNet.	32
2.10	Framework for signature-based algorithms [7]	41
2.11	Two documents represented as records	41

2.12	Prefix Filter scheme of Figure 2.5 with 80% Overlap similarity threshold.	42
2.13	Two vectors with hamming distance = 4	43
2.14	Two vectors with hamming distance $\leq k$ 4 must agree on one of the $k + 1$ partitions	44
2.15	The two vectors in Figure 2.14 with hamming distance = 8 and agree on one partition	44
2.16	Enumeration scheme for two vectors with hamming distance=3	45
2.17	Formal specification of PartEnum[7]	46
2.18	Formal specification of All-Pairs[6]	49
3.1	Operational Framework	54
3.2	The procedure used in obtaining the semantic attributes between two concepts	63
3.3	The algorithm for semantic relatedness between a pair of sentences	65
3.4	Binary vector representation of a sentence	66
3.5	An inverted index implementation for Cosine similarity[6]	67
3.6	An inverted index implementation for Jaccard Similarity	70
3.7	An inverted index implementation for Dice coefficients	70
3.8	The procedure of evaluating Web document retrieval techniques	72
3.9	Response format from querying the Google API	74
4.1	The inverted index for document Q	84
4.2	The hypernym trees for the words: <i>Idea</i> (5 senses), <i>Learning</i> (2 senses), <i>Adaptation</i> (3 senses), <i>Evolution</i> (2 senses) and <i>Biology</i> (3 senses).	89

4.3	The hypernym trees for the words: <i>Idea</i> (5 senses), <i>Learning</i> (2 senses), <i>Adaptation</i> (3 senses), <i>Evolution</i> (2 senses) and <i>Aspect</i> (3 senses)	90
4.4	Recall rate (y-axis) across similarities (x-axis) in 110 corpus documents	101
4.5	Recall rate in one-to-one exact copies	113
4.6	Recall rate in one-to-one plagiarized by synonym replacing	114
4.7	Recall rate in one-to-many plagiarized by synonym replacing	114

LIST OF APPENDICES

APPENDIX	TITLE	PAGE
A	Stop-words and their corresponding frequencies in the Brown corpus	128
B	Information about ScienceDirect source documents	129
C	Information about Wikipedia corpus documents	130
D	The PENN TreeBank English POS tag set and their mappings	144
E	Examples of original/plagiarized sentence pairs and the corresponding similarities based on equation 3.6	145

CHAPTER 1

INTRODUCTION

1.1 Introduction

The World Wide Web (Web) is the biggest source of information these days. People now can easily search for, access, and browse Web pages to get the information they need, one can imagine how difficult the academic research would be without the Internet and the Web. It is also now easy, and again because the scale and the digital structure of the Web, to use someone else's work illegally.

The problem of plagiarism has its direct association to academia. Maurer et al. [3] defined it as “the unacknowledged used of someone else's work”. The most common type is written-text plagiarism in which the plagiarized document is formed by copying some or all parts of the original document(s) possibly with some alterations. Plagiarism is classified into *intra-corp*al and *extra-corp*al with respect to the location of the source document(s) [1]. The former happens when both the copy and source documents are within the same corpus, such as within a collection of students' submissions or within a digital library. While in the latter, the copy and

source documents are not of the same corpus. Here the source documents could be from textbooks or most commonly Web documents. Unless the problem of locating the source documents is solved, it is hard to prove this kind of plagiarism. Identifying Web documents from which copying has occurred is stressful and time consuming for a human inspector given the large number of documents that need to be compared. As the digital structure of Web documents made it easy to plagiarize, fortunately it means that such instances of plagiarism could be traced in an automated manner.

There are two methods to provide an access to a large number of Web documents. The first method is by indexing documents through Web crawling; this has the inherent problems of Web documents that face any Web retrieval system such as bulk size, heterogeneity, and duplication [2], however the system could be tuned for the retrieval purposes, for example if the purpose is to detect plagiarism, the system can be employed to return the most syntactically or semantically similar documents to the query document. The other method, which this project will use, is utilizing general-purposes search engines (such as Google, Yahoo, and Bing) as they provide access services to their systems. The suspected document can be considered as a sequence of queries submitted to the search engine, the result are then compared with the input document.

Intuitively it is required to partition the query document into more primitive units plausible for querying the search engine and for documents comparisons. Sentences are suitable for both cases since they carry ideas and also plagiarism patterns (e.g., insertion, deletion, and/or substitution).

Similarity between sentences (or more generally objects) can be captured numerically using similarity measures such as Jaccard similarity, Overlap similarity, Cosine similarity. These measures are called symmetric functions and widely used

in many Information Retrieval applications. Each measure returns a value indicating the degree of similarity between pairs of objects usually between 0 and 1.

Beside the similarity measures, another aspect is the document (or sentence) representation. There are many representations that have been developed including document fingerprinting [17], bag-of-words model [10], N-grams (consecutive words of length N). Another important representation comes from semantic networks. A semantic network or net “is a graphic notation for representing knowledge in patterns of interconnected nodes and arcs” [50]. Concepts in semantic networks are usually organized in hierarchical structure as illustrated in Figure 1.1

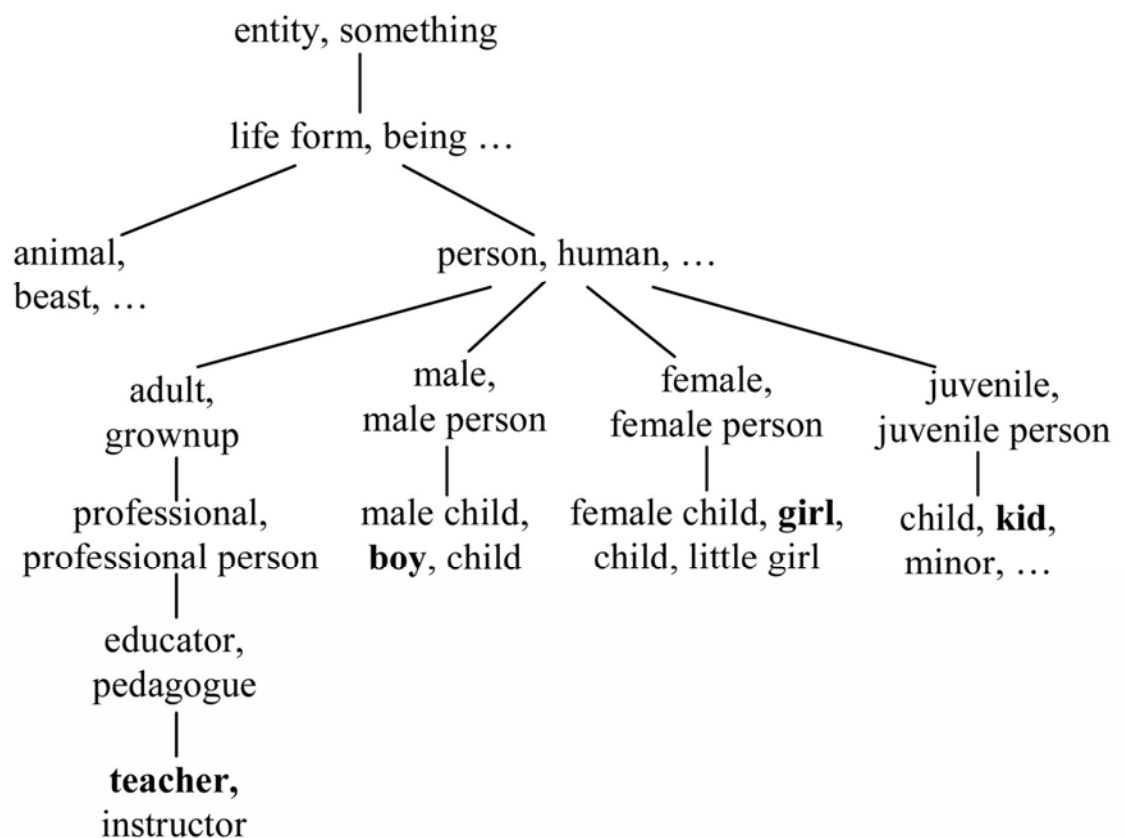


Figure 1.1 Hierarchical semantic knowledge base[57].

Usually words at upper layers of hierarchical semantic nets have more general concepts and less semantic similarity between words than words at lower layers [57].

1.2 Problem Background

In any application that involve measuring the similarity between textual contents there are two important factors that influence the accuracy of plagiarism detection. The first factor is the document representation which essentially captures the characteristics of the document as a preceding step to the comparison stage. These representations include the “Bag-of-Word” model, document Fingerprints, N-grams, probabilistic models. Most of these representations work well in detecting verbatim (word-to-word) plagiarism but have vulnerabilities in detecting complicated plagiarism patterns.

The second factor is the similarity measure that is used to calculate the similarity or dissimilarity between sentences. Considering the plagiarists behavior that usually involves insertions of words deletions and/or substitutions it is necessary to determine which measure is the best for detecting instances of plagiarism.

Retrieving the source documents from the Web using a search engine is another challenge given the fact that some plagiarism patterns are hard to locate in the setting of the Web even for a human inspector.

In this project we investigate the effectiveness of semantic net-based techniques for detecting plagiarized sentences and find out whether the achieved performance is justified comparing to other approaches. Then we determine which technique is the best for retrieving the source documents from the Web.

1.3 Problem Statement

To cater the problems introduced in section 1.2, this project is carried out to answer the following questions:

- i- Which N-gram representation is the best for sentence-based plagiarism detection?
- ii- Which similarity measure is the best for sentence-based plagiarism detection?
- iii- How can semantic networks be used to improve the detection?

1.4 Project Objectives

The main objectives of this project are stated as follows:

- i- To compare the effectiveness of different N-gram with different similarity measures in detecting plagiarized documents over the Web.
- ii- To find out whether the use of semantic networks can improve the detection of plagiarized documents.

1.5 Project Scope

- i- This project will cover plagiarism detection in English scripts.
- ii- WordNet [4] is the general semantic network that will be used in this study.
- iii- N-grams will be used with three symmetric measures; Cosine, Jaccard, and Dice coefficients.
- iv- Porter algorithm [60] will be applied in the stemming process.

1.6 Project Justification

The problem of document plagiarism detection is not new and several methods have been applied to overcome this problem over a small collection of documents or digital libraries, however, the scale of the problem has increased dramatically due to the Web.

It is also widely acceptable that traditional methods for measuring the similarity between documents are vulnerable to fail in some complex plagiarism patterns and hence it is necessary to incorporate semantic-based techniques for more accurate plagiarism detection.

1.7 Report Organization

This report is organized as follows:

Chapter 1 formulates the problem and outlines the framework and main objectives of the project.

Chapter 2 consists of four main parts; the first part introduces some terminologies of document plagiarism detection and briefly outlines some plagiarism detection methods. The second part focuses on semantic networks, in particular WordNet and its semantic relations. The third part is then devoted to document pre-processing and representation techniques and their effect in the applications of plagiarism detection, it also reviews the main approaches for semantic relatedness between concepts. The last part reviews efficient exact set similarity algorithms and discusses how they can be adopted in the case of N-grams.

Chapter 3 illustrates the methodology that will be used to fulfill the objectives of this project.

Chapter 4 presents the experimental results of this project, and finally chapter 5 concludes this research.