





Computer Science Department
Faculty of Computer and Information Sciences
Ain Shams University

Implementation of Enhanced Ontological Techniques for Information Retrieval Purposes

A thesis submitted as a partial fulfillment of the requirements for the degree
of Master of Science in Computer and Information Sciences.

By

Amr Aly Amin Aly ElSehemy

B.Sc. in Computer & Information Sciences

Demonstrator at Computer Sciences Department

Faculty of Computer and Information Sciences, Ain Shams University

Under the supervision of

Prof. Dr. Taymoor Nazmy

Vice Dean for Research and Higher Studies

Faculty of Computer and Information Sciences, Ain Shams University

Dr. Mohamed Abdeen

Computer Science Department

Faculty of Computer and Information Sciences, Ain Shams University

Cairo, Egypt May 2015

Acknowledgement

All praise to Allah, the gracious and merciful, for helping me complete this work.

No words can describe my gratitude towards my parents, pushing me forward every step of the way. I could not have asked for a more supportive mother, doing everything she can to help. My father put no limits to the sacrifices he was willing to make to see me succeed. My sister and brother are the greatest gift my parents gave me, always putting me first. To them I owe the accomplishment of this work, and to them I put it forward as the fruit of their sacrifices.

I also present my sincerest thanks to my thesis advisors: Prof. Dr. Taymoor Nazmy, and Dr. Mohamed Abdeen. I thank them for the hours and hours spent brainstorming our way over obstacles. I'm lucky to have had such supportive advisors.

I sincerely appreciate the help of my friends, supporting me with great and insightful advice. I thank all my colleagues for their limitless support and encouragement.

Abstract

Keyword-based search engines were used to solve such problem, however, in many cases; keyword-based search engines either miss some documents or retrieve non-relevant documents. New technologies were presented to enhance the search results. One of these technologies is based on semantic web. Many semantic search techniques and models were made to enhance the traditional keyword-based search. Including conceptual knowledge such as ontologies in the information retrieval process contributes to the solution of major problems found in keyword-based search. Advances were made in languages such as English, German, French and Spanish. Although Arabic language is spoken by as many as 422 million native speakers, the literature does not fully cover it yet.

In this work, an architecture for an ontology-based information retrieval for Arabic language is presented. The architecture presented consists of four main modules, the query parser, the indexer, the search and the ranking modules. This work included building a semantic index; providing weighted links between the ontology concepts and the documents. Furthermore, a document categorizer to the architecture was built. The document categorizer was used as an additional process to enhance the overall document ranking.

To test our work, three Arabic domain ontologies were built. These ontologies are Sports, Economics and Politics. A knowledge base was built that consisted of 79 classes and 1456 instances. The document categorizer was evaluated using the WATAN-2004 corpus. We introduced some modifications merging the categories;

to be relevant to the work in here. The corpus contains around 20,000 articles in five categories (Culture, Religion, Economy, News and Sports). Twenty retrieval operations were manually chosen to assess different cases. The operations were applied on a sample of 40,316 documents with a size 320 Mega Bytes of pure text. Each operation was applied three times, once using the keyword-based search, and the other two using the presented ontology-based search with and without the classification module. The documents were downloaded from www.aljazeera.net news website.

This work was evaluated using the precision and recall metrics. The presented ontology-based search was compared against the traditional keyword-based search. The results showed that the ontology-based system performs better than the keyword-based system whenever ontology instances existed. The classification module was evaluated using the 10-fold evaluation technique. It achieved an average accuracy 78%, with a minimum error rate of 4%.

Further evaluation to the system was made by comparing between the ontology-based search with and without integrating the automatic document categorizer. The results showed an enhancement by 5% after integrating the categorizer module.

Table of Contents

Chapter 1:	Introduction	3
1.1	Overview	3
1.2	Previous Work	4
1.3	Challenges of Arabic Ontology and Semantic Web.....	8
1.4	Objectives	10
1.5	Thesis Outline	10
Chapter 2:	Information Retrieval and Semantic Web.....	13
2.1	Introduction.....	13
2.2	The Information Retrieval Process	14
2.3	The Retrieval Models	16
2.4	Information Retrieval Evaluation Techniques	26
2.5	Knowledge Engineering Ontologies	31
2.6	Semantic Web.....	34
2.7	Semantic Knowledge Representation	39
2.8	Semantic Knowledge Acquisition	48
2.9	Semantic Knowledge Annotation	49
2.10	Classification Models	49

Table of Contents

Chapter 3: The Proposed Approach for Ontology-Based Arabic Search Engine with Document Classification.....	58
3.1 System Architecture	58
3.2 The Preprocessing Module	62
3.3 The Indexer Module.....	63
3.4 The Query Engine Module	65
3.5 The Ranking Module	68
3.6 The Automatic Document Classifier	69
3.7 Implementation Details	72
Chapter 4: Results and Discussion	76
4.1 Introduction.....	76
4.2 Dataset Description.....	76
4.3 Test Cases.....	79
4.4 Results	79
4.5 Discussion	91
Chapter 5: Conclusion and Future Work	95
5.1 Conclusion.....	95
5.2 Future Work	97

References 100

List of Figures

Figure 2-1 Overall process of IR systems	13
Figure 2-2 The three conjunctive components of the query $q = ta \wedge (tb \vee \neg tc)$	18
Figure 2-3 Differences between humans and machines view of the same Web page.....	36
Figure 2-4 Current Web content structure vs. Semantic Web content structure	37
Figure 2-5 SPARQL example	48
Figure 2-6 Naïve Bayesian classifier algorithm.....	52
Figure 2-7 Naïve Bayesian classifier learning algorithm	54
Figure 3-1 Proposed System Architecture	59
Figure 3-2 Information Retrieval Overall Process	61
Figure 3-3 Text Preprocessing	62
Figure 3-4 Approach Database Diagram	64
Figure 3-5 Search and Retrieval Process	68
Figure 3-6 System Architecture including Classifier	71
Figure 3-7 Revised database diagram with enhanced weight	

List of Figures

added.....	72
Figure 3-8 Screenshot of Proposed System User interface.....	74
Figure 4-1 Part of the Sports Domain Ontology	78
Figure 4-2 Precision vs. Recall for query (A).....	81
Figure 4-3 Precision vs. Recall for query (B)	82
Figure 4-4 Precision vs. Recall for query (C).....	83
Figure 4-5 Precision vs. Recall for average of twenty queries	90
Figure 4-6 Comparative precision histogram for ontology- based search vs. keyword-based search.....	91

List of Tables

Table 4-1 Precision, Recall and F1 for the naive Bayesian classifier	8
Table 4-2 Precision, Recall and F1 for the KNN classifier....	8
Table 4-3 Error rates of NB over all categories	8
Table 4-4 Confusion matrix for cross validation, with full document space	8

List of Abbreviations

IR	Information Retrieval
SPARQL	SPARQL Protocol and RDF Query Language
SW	Semantic Web
PWN	Princeton Word Net
AWN	Arabic Word Net
SSE	Semantic Search Engine
RDF	Resource Description Framework
OWL	Web Ontology Language
RDFS	Resource Description Framework Schema
NLP	Natural Language Processing
NN	Neural Networks
DCMI	Dublin Core Metadata Initiative
NER	Name Entity Recognizer
VSM	Vector Space Model
W3C	World Wide Web Consortium
KB	Knowledge Base
DARPA	Defense Advanced Research Projects Agency
DAML	DARPA Agent Markup Language
NB	Naïve Bayesian
KNN	K-Nearest Neighbor

Chapter 1

Introduction

Chapter 1: Introduction

1.1 Overview

Information Retrieval (IR) is one of the well-established research areas in Information Science. The goal of the IR discipline is to search and retrieve the most relevant documents to the information needs of the user. Therefore a good IR system should retrieve only those documents that satisfy the user needs, not other unnecessary data.

The increase in the amount and complexity of reachable information in the World Wide Web caused an excessive demand for tools and techniques that can handle data semantically. The current practice in information retrieval mostly relies on keyword-based search over full-text data. The keyword-based search is modeled with bag-of-words. Such a model introduces an issue, it does not present the actual semantic information embodied in text. To deal with this issue, ontologies are proposed for knowledge representation [1]. Ontologies are considered the backbone of semantic web applications. Both the information extraction and retrieval processes can benefit from such metadata, which gives semantics to plain text. It is also used to enhance the final ranking of the retrieved documents.

Having obtained the semantic knowledge and represented them via ontologies, the next step is querying the semantic data, also known as semantic search. There are several query languages designed for semantic querying. Currently, SPARQL Protocol and