

REPUBLIQUE ALGERIENNE DEMOCRATIQUE ET POPULAIRE
Ministère de l'Enseignement Supérieur et de la Recherche Scientifique

Université Hadj Lakhdar – Batna
Faculté des Sciences
Département d'Informatique



MÉMOIRE
Présenté par
Tahar DILEKH

Pour obtenir le grade de
Magister
Spécialité : Système d'Information et de Connaissance (SIC)

Implémentation d'un outil d'indexation et de recherche des textes en arabe

Soutenue publiquement le 28 /09 /2011 devant le jury formé de :

Pr. Mohammed BENMOHAMMED	Professeur	Président	Université de Constantine
Dr. Abdelmadjid ZIDANI	M.C.	Rapporteur	Université de Batna
Dr. Ali BEHLOUL	M.C.	Co-Rapporteur	Université de Batna
Pr. Azzeddine BILAMI	Professeur	Examinateur	Université de Batna
Dr. Brahim BELATTAR	M.C.	Examinateur	Université de Batna

ABSTRACT

Arabic, one of the six official languages of the United Nations, is the mother tongue of more than 300 million people¹. Becoming a center of research and commercial development, the importance of the domain of information retrieval (IR) is due to the essential need for such tools to people in the Net era. The number of Internet users in 2002 was about 4.4 million which represents about 1.5 % of the population of the Arab world². Few search engines are available for the growing number of Arabic Internet users; however many effort are being deployed,

Arabic is a highly inflected language and has a complex morphological structure. The information retrieval (IR) on Arabic texts requires the basic form of the word (root or lemma); therefore stemming process is necessary. This process can be defined as a process of removing all affixes (prefixes, infixes and/or suffixes) from words in order to bring them to their roots or lemmas.

Morphological complexity of Arabic language makes particularly developing software for natural languages process difficult. In Semitic languages like Arabic, the majority of names, adjectives and verbs are derived from few thousand roots by adding new letters, for example, the words مكتبة (library), كتاب (book), كتب (pounds), نكتب (he wrote), and نكتب (writing), the root [Wig98].

Every natural language has its own characteristics and features. Thus, it is difficult to follow the same pattern of stemming and apply the same techniques for all languages. A technique of stemming could be relevant to a language, but could not be effectively applied to other languages. There are several techniques used for word stemming; these include technical dictionaries, morphological analysis, deletion of affixes, statistics, and translation.

In this work, we proposed a hybrid method that incorporates three different techniques so as the Arabic stemming process resolve problems in connection with each mentioned technique.

These three techniques are: removal of affix given by Kadri [Kad08], dictionaries, and morphological analysis

These techniques require some adjustments to be relevant for use. Each technique is fitted and adapted individually to solve practical problems related to itself.

Therefore, the main contribution of this experiment was to demonstrate the effectiveness of the *hybrid method* compared to other methods, and the choice of removing the suffix before prefix during the operation of Arabic stemming process. For example

¹ <http://www.frcu.eun.eg/www/homepage/cdc/cdc.htm>. Accédé le 12/12/2008.

² <http://sonomabusiness.com/archives/2002-09-column-levini.html>. Accédé le 05/12/2008.

Word	Root	Suffix – Préfixe	Prefix - Suffix
أَلْمٌ	أَلْمٌ	أَلْمٌ	أَلْمٌ
Their pain	pain	Pains	error

Key Words: Information retrieval, lemmatization, Arabic language.

RESUME

L'arabe, une des six langues officielles des Nations Unies, est la langue maternelle de plus de 300 millions de personnes¹. Le domaine recherche d'information (RI) arabe, devenu un centre de la recherche et du développement commercial est du à la nécessité essentielle de tels outils pour des personnes dans l'ère électronique. Le nombre d'internautes arabophones en 2002 était environ 4.4 millions, environ 1.5% de la population du monde arabe². Mais, de l'autre côté de la réalité, peu de moteurs de recherche sont mis à la disposition des utilisateurs arabophones, bien que les efforts soient en marche pour servir le nombre croissant d'utilisateurs.

L'Arabe est une langue fortement flexionnelle qui a une structure morphologique complexe. La recherche d'information sur le texte arabe exige la forme de base du mot (racine ou lemme) pour être la plus pertinente, donc le processus de lemmatisation est nécessaire. La lemmatisation peut être définie comme un processus qui consiste à retirer tous les affixes (préfixes, infixes, ou/et suffixes) des mots pour ramener ces derniers à leurs lemmes ou racines.

La complexité morphologique de la langue arabe rend particulièrement difficile le développement des applications pour le traitement en langue naturelle. Dans les langues sémitiques comme l'arabe, la plupart des lemmes de nom, d'adjectif, et de verbe sont dérivés de quelques mille racines par l'insertion de nouvelles lettres, par exemple, les mots مكتبة (bibliothèque), كتاب (livre), أیل (Ivres), 写 (Il a écrit), و نكتب (Nous écrivons), de la racine كتب [Wig98].

Chaque langue naturelle a ses propres caractéristiques et dispositifs. Ainsi, il est difficile de suivre la même configuration de lemmatisation et d'appliquer les mêmes techniques pour toutes les langues. Une technique de lemmatisation pourrait être pertinente à une langue, mais ne peut pas être effectivement appliquée à d'autres langues. Il existe plusieurs techniques utilisées pour la lemmatisation des mots. Celles-ci incluent, des techniques de dictionnaires, d'analyse morphologique, de suppression des affixes, de statistiques, et de traduction.

Dans ce travail, nous avons proposé une méthode hybride qui incorpore trois techniques différentes pour que la lemmatisation arabe résolve les problèmes liés à chaque technique précédente.

Ces trois techniques sont: suppression d'affixe proposée par Kadri [Kad08], dictionnaires, et analyse morphologique.

¹ <http://www.frcu.eun.eg/www/homepage/cdc/cdc.htm>. Accédé le 12/12/2008.

² <http://sonomabusiness.com/archives/2002-09-column-levini.html>. Accédé le 05/12/2008.

Ces techniques ont besoin d'une certaine adaptation pour être pertinentes pour l'utilisation. Chaque technique est adaptée individuellement pour résoudre les problèmes pratiques liés à elle-même.

La contribution principale de ce travail concerne la démonstration de l'efficacité de la *méthode hybride* comparée aux autres méthodes, et le choix de l'enlèvement des suffixes avant les préfixes pendant l'opération de lemmatisation Arabe. Par exemple

Mot	Racine	Suffixe – Préfixe	Préfixe - Suffixe
الله	لله	لله	لله
Leurs douleurs	douleur	douleurs	Erreur

Mots clés : Recherche d'information, lemmatisation, langue arabe.

TABLE DES MATIERES

ABSTRACT	1
RESUME	3
REMERCIEMENTS	5
TABLE DES MATIERES.....	6
LISTE DES TABLEAUX	10
LISTE DES FIGURES	11
INTRODUCTION GENERALE	12
1. Objectifs	13
1.1. Objectif général	13
1.2. Objectifs spécifiques	13
2. Méthodologie	13
3. Organisation du mémoire	15
CHAPITRE 1 : LA RECHERCHE D'INFORMATION.....	16
1. Introduction:.....	16
2. Processus de recherche d'information.....	16
2.1. Modèles de RI.....	18
2.1.1. Le modèle booléen	18
2.1.1.1. Avantages	19
2.1.1.2. Inconvénients.....	19
2.1.2. Le modèle probabiliste.....	20
2.1.2.1. Avantages	21
2.1.2.2. Inconvénients.....	21
2.1.3. Le modèle LSI (Latent Semantic Indexing)	22
2.1.4. Le Modèle vectoriel	24
2.1.4.1. Avantages	25
2.1.4.2. Inconvénients.....	25
2.2. Critères d'évaluation des SRI	25
2.2.1. Évaluation.....	25
2.2.1.2. Précision	26
2.2.1.3. Rappel	26
2.2.2. La courbe de Rappel/Précision	26
2.2.3. Mesures globales	27
2.2.3.1 La précision moyenne interpolée IAP (Interpolated Average Precision).....	27

2.2.3.2. La R-précision.....	27
2.2.3.3. La F-mesure [Van79].....	27
3. RI en langue arabe	28
3.1. Les ressources arabes :	28
3.1.1. Corpus	28
3.1.2. Dictionnaire	29
3.1.3. Outils	29
3.1.3.1. Analyseurs morphologiques.....	29
3.1.3.2. Les concordanciers	29
3.1.3.3. Racineurs	29
3.2. Lemmatisation	30
4. Conclusion.....	30
CHAPITRE 2 : PROPRIETES MORPHOLOGIQUES DE L'ARABE	31
1. Introduction:	31
2. Particularité de la langue arabe	31
3. Morphologie arabe	34
4. Structure d'un mot	35
4.1. Les antéfixes:	35
4.2. Les préfixes:.....	36
4.3. Les suffixes:.....	36
4.4. Les post fixes:.....	37
5. Les catégories des mots	39
5.1. Verbe	39
5.2. Nom	39
5.3. Particule.....	40
6. Problèmes du traitement automatique de l'arabe	42
7. Conclusion.....	42
CHAPITRE 3 : PRETRAITEMENTS NECESSAIRES	44
1. Introduction:	44
2. Encodage	44
2.1. L'Unicode:	44
2.2. UTF-8	44
2.3. Produits Unicode supportant l'écriture arabe	45
2.4. L'encodage de corpus et requêtes:	45

3. Segmentation:.....	45
3.1. Définition	46
3.2. Le système d'écriture arabe:.....	46
3.4. Les types de segmentation	47
3.5. Les clitics	48
3.6. Segments arabes	48
3.6.1. Segments principaux.....	49
3.6.2. Segments secondaires	49
3.7. Les solutions de segmentation.....	50
3.7.1. Le modèle de segmentation: Guesser [Bes03]	50
3.7.1.1. Le Guesser (pronostiqueur) de clitics.....	51
3.7.1.2. Capteur de Clitics.....	51
4. Les mots vides.....	51
5. Normalisation.....	53
6. Conclusion.....	54
CHAPITRE 4 : LEMMATISATION.....	56
1. Introduction:.....	56
2. Définition.....	57
3. Difficultés de la lemmatisation des mots arabes.....	57
4. Les Techniques de lemmatisation.....	58
4.1. La technique de dictionnaire.....	58
4.2. Suppression d'affixe	58
4.3. Techniques d'analyse morphologique	60
4.4. Techniques statistiques.....	61
4.5. Techniques de traduction	62
5. La méthode proposée.....	62
5.1. Suppression d'affixe	63
5.2. La technique de dictionnaire.....	66
5.3. Techniques d'analyse morphologique	67
6. Conclusion.....	68
CHAPITRE 5 : IMPLEMENTATION ET EXPÉRIMENTATION	70
1. Introduction:.....	70
2. Le corpus de test:	70
3. Implémentation:	73

3.1. Indexation.....	73
3.2. Recherche d'information	73
3.3. Architecture du système.....	74
3.3.1 Encodage.....	75
3.3.2. Normalisation	76
3.3.3. Segmentation	76
3.3.4. Élimination des mots vides.....	76
3.3.5. Lemmatisation:.....	77
3.3.5.1. La méthode PS-M :.....	77
3.3.5.2. La méthode SP-M :.....	77
3.3.5.3. La méthode PS+M (Préfixe Suffixe Avec Modèle):	78
3.3.5.4. La méthode SP+M :	79
3.3.5.5. La méthode HY (Hybride):	79
3.3.6. Pondération des termes d'indexation.....	79
3.3.7. Techniques de création des index.....	80
3.3.8. Méthode de recherche	80
3.3.8.1. L'appariement document-requête	81
4. Expérimentation et évaluation	81
5. Conclusion.....	88
CONCLUSIONS ET PERSPECTIVES	90
1. Conclusion.....	90
2. Perspective	91
2.1. Lemmatisation des mots.....	91
2.2. Approche Sémantique	91
2.3. Approche Hybride	92
BIBLIOGRAPHIE	93