

العنوان:	تصميم خوارزمية حديثة لزيادة كفاءة تصنيف النص العربي آليا باستخدام الانظمة الذكية
المؤلف الرئيسي:	سليمان، عبدالعزيز حسن خرساني
مؤلفين آخرين:	أحمد، عوض حاج علي(مشرف)
التاريخ الميلادي:	2013
موقع:	الخرطوم
الصفحات:	1 - 158
رقم MD:	855316
نوع المحتوى:	رسائل جامعية
اللغة:	Arabic
الدرجة العلمية:	رسالة دكتوراه
الجامعة:	جامعة النيلين
الكلية:	كلية علوم الحاسوب وتقانة المعلومات
الدولة:	السودان
قواعد المعلومات:	Dissertations
مواضيع:	تصميم الخوارزميات، تصنيف النص العربي، النظم الذكية، علوم الحاسبات
رابط:	http://search.mandumah.com/Record/855316



جامعة النيلين
كلية الدراسات العليا
كلية علوم الحاسوب وتقانة المعلومات

بحث مقدم لنيل درجة الدكتوراه في علوم الحاسوب بعنوان:

تصميم خوارزمية حديثة لزيادة كفاءة تصنيف النص
العربي آليا باستخدام الانظمة الذكية

إشراف الأستاذ الدكتور:

عوض حاج علي احمد

إعداد الطالب:

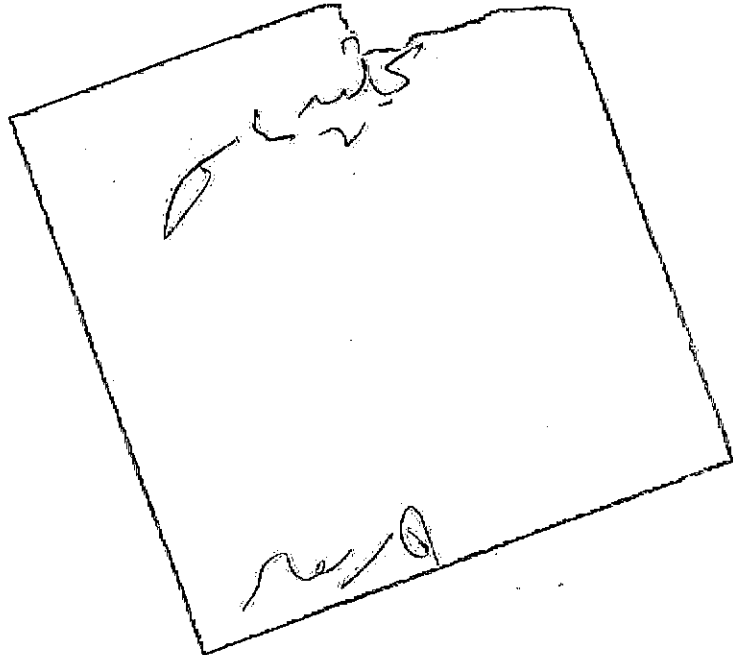
عبد العزيز حسن خرساني سليمان

(2013)

الاستهلال

(وَلَقَدْ نَعَلِمَ أَنَّهُمْ يَقُولُونَ إِنَّمَا يُعَلِّمُهُ بَشَرٌ لِّسَانُ الَّذِي
يُلْحِدُونَ إِلَيْهِ أَعْجَمِي وَهَذَا لِسَانٌ عَرَبِيٌّ مُبِينٌ)

القرآن الكريم : سورة النحل : الآية: (103)



الأهداء

إلى...

والدي...

الذنان سقياني من كأس الحب والحنان ، وأطعماني العزيمة والصبر والأمل وألبساني

الشموخ والتواضع

إلى...

زوجتي وأبنائي وإخوتي الأعزاء وجميع أفراد العائلة الكريمة....

حفظهم الله ورعاهم

إلى...

كل من يحمل جزءاً من حياتنا وأعاننا على ترتيب أفكارنا

إليكم جميعاً أهدي ثمرة جهدي المتواضع

الشكر والعرفان

الشكر لله وحده من قبل ومن بعد، الذي يسر لي إكمال هذا البحث وأنعم علي بإتمامه في هذه الصورة

الشكر إلى من أنعم علينا بنعمة العلم وجعله لنا سراجا منيرا نستنير به ونقهر الجهل

الشكر إلى من جاد علي بنصح وتوجيه وإرشاد ومتابعة أستاذي الجليل الدكتور:

أ.د/عوضه صالح محلي

الذي كان نعم المعين في إخراج هذا البحث بهذه الصورة، لما قدمه لي من نصائح ثرة جزاه الله عني كل الخير وجعل ذلك في ميزان حسناته.

وكذلك الشكر أ.د/السماني عبد المطلب احمد الذي لم يبخل عني بالتوجيه والمعلومات

كما أتقدم بالشكر لكل من جامعة النيلين وجامعة سلمان بن عبدالعزيز، ولكل من شارك وجاد علي بنصح وتوجيه وإرشاد، أو قدم إلي أي مساعدة .

ولكم الشكر وجزاكم الله خيرا ..

مستخلص الدراسة

تصنيف النص Text categorization هو عملية تصنيف المستندات والوثائق إلى مجموعات محددة مسبقا من الفئات استنادا علي مضمون ومحتوي المستندات والوثائق . وهدفت الدراسة الي تصميم خوارزمية لتحسين دقة تجريد النص العربي وتصنيفه وخاصة الافعال المضعفة والأسماء الثنائية والحرف المكرر وذلك باستخدام برامج ذكية تساهم في تصنيف النص العربي . تم استخدام حزمة برنامج visual studio2010 وحزمة #c لأنها تناسب برامج النظم الذكية المستخدمة في تصنيف النصوص . وخلصت الدراسة الي مجموعة نتائج اهمها زيادة نسبة التصنيف للنصوص العربية الي 96% . تم تحسين خوارزمية التجريد وذلك بزيادة الافعال المضعفة والاسماء الثنائية . تم تحسين خوارزمية التجريد وذلك بازالة الحرف المكرر(————) في النص العربي . تم عمل واجهة يمكن من خلالها زيادة وزن الكلمة يدويا في النظام وذلك اذا راي المستخدم اهمية الكلمة لتحديد فئة ما . وهناك حاجة ماسة لاستخدام اللغة العربية لعلم التصنيف الالي وذلك لمواكبة اللغات الأخرى وتقنيات العصر اوصت الدراسة باستخدام معجم عربي حديث يستفيد من المعالجة الآلية، نجد مثلا الكثيرين ممن لديهم خلفية في الحاسب ليس لديهم معرفة باللغة العربية وعلماء اللغة العربية ليس لديهم معرفة بالحاسب إلا قليلا، وضرورة ربط الترجمة الالية بالتصنيف الالي

للنصوص العربية المجمعة، وان يصبح تصنيف النص العربي اليا ضمن مقررات

لمناهج علوم الحاسوب، ويمكن استخدام المصنف Support Vector Machines

(SVM) لتصنيف النص العربي.

Abstract

Text categorization is the process of classifying documents into a predefined set of categories based on the content of the documents. The study aimed to design an algorithm to improve the accuracy of the Arabic text stemming and classification especially verbs multiple reward and bilateral names and duplicate character. And using intelligent software contribute to the Arabic text classification . Software package was used visual studio2010 and c # package to fit it smart software systems used to classify texts . And The study concluded that the most important result set to increase the percentage of classification of Arabic texts to 96% . Abstraction algorithm has been improved by increasing verbs multiple reward and bilateral names. abstraction algorithm has been improved and duplicate by deleting the character in the Arabic text(————) . been working interface from which you can increase Weight word manually in the system and the user so if Ray importance of the floor to determine the category. There is an urgent need for the use of the Arabic language for automated taxonomy in order to keep up with other languages, techniques, and study

recommended. There is no modern Arabic dictionary benefit from processing machine, We find, for example, many people who have a background in the computer have no knowledge of the Arabic language and Arab linguists have no knowledge of the computer only slightly, the need to link machine translation automatic with classification of Arabic texts collected, the Arabic text classification becoming within the decisions of the computer science approaches, can be used Support Vector Machines (SVM) to the Arabic text classification.

الفهرس العام

الصفحة	الموضوع	الرقم
.i	الآية	1
.ii	الإهداء	2
.iii	الشكر و العرفان	3
.iv	مستخلص الدراسة	4
.vi	Abstract	5
.viii	الفهرس العام	6
.xv	فهرس الاشكال و المخططات	7
.xvi	فهرس الجداول	8
الفصل الاول - أساسيات البحث		
2	تمهيد	1 - 1
3	مشكلة البحث	2 - 1
4	اهداف البحث	3 - 1
5	اهمية البحث	4 - 1
5	حدود وفترة الدراسة	5 - 1
5	منهجية البحث	6 - 1
6	تقسيمات البحث	7 - 1
6	مصادر جمع البيانات	8 - 1
7	الدراسات السابقة	9 - 1
15	الفجوة البحثية	10 - 1

الفصل الثاني - الذكاء الاصطناعي Artificial Intelligence		
17	تمهيد	1 - 2
18	تاريخ بداية الذكاء الصناعي	2 - 2
19	تعريف الذكاء الصناعي Artificial Intelligence	3 - 2
21	الذكاء الإنساني	4 - 2
23	الفرق بين الذكاء الصناعي والذكاء الانساني	5 - 2
24	اساليب الذكاء الصناعي	6 - 2
24	اسلوب استخدام القوانين	1-6 - 2
25	اسلوب شبكات المعاني	2-6 - 2
25	اسلوب تمثيل الاطارات	3- 6 - 2
25	اسلوب الرؤية الالكترونية	4- 6 - 2
27	اسلوب معالجة اللغات الطبيعية	5 - 6 - 2
27	الكلام Speech	1- 5 - 6 - 2
27	النظر Vision	2- 5 - 6 - 2
27	الروبوت Roboties	3 - 5 6 - 2
28	التعليم Learning	4-5 - 6 - 2

28	مجالات الذكاء الصناعي	7 - 2
28	الذراع الآلية الذكية	1 - 7 - 2
30	الأنظمة الخبيرة (Expert Systems)	2 - 7 - 2
35	تفوق النظام الخبير / الذكاء الاصطناعي على برامج الحاسبة التقليدية	8 - 2
36	أهمية استخدام الذكاء الصناعي	9 - 2
38	تطبيقات علم الذكاء الاصطناعي	10 - 2
40	مكونات الذكاء الصناعي	11 - 2
43	فروع علم الذكاء الاصطناعي	12 - 2
44	منطق الذكاء	13 - 2
44	التمييز النمطي والنموذجي	14 - 2
46	لغة البرمجة المستخدمة لإنتاج برامج الذكاء الصناعي	15 - 2
الفصل الثالث - معالجة اللغات الطبيعية Natural Language Processing		
48	تمهيد	1 - 3
49	البداية والتاريخ	2 - 3
50	The term and concept المصطلح والمفهوم	3 - 3

53	Objectives and reasons الأهداف والأسباب	4 - 3
54	تواصل أفضل مع الحاسب	1 - 4 - 3
55	تواصل أفضل بين البشر	2 - 4 - 3
55	الوصول الفعال للمعلومات	3 - 4 - 3
56	تحديات تواجه المعالجة الآلية للغة	5 - 3
56	تقطيع الكلام والأصوات والوحدات المعجمية	1 - 5 - 3
57	فك الغموض أو اللبس	2 - 5 - 3
57	العبارات الطلبية	3 - 5 - 3
58	تطبيقات المعالجة الآلية للغة applications of machine language processing	6 - 3
58	Machine translation الترجمة الآلية	1 - 6 - 3
60	Auto summarization التلخيص الآلي	2 - 6 - 3
61	التوليد الآلي للغة	3 - 6 - 3
62	extracting information استخراج المعلومات	4 - 6 - 3
62	Information Retrieval استرجاع المعلومات	5 - 6 - 3
62	answer the questions الإجابة على الأسئلة	6 - 6 - 3
64	texts mining التنقيب في النصوص	7 - 6 - 3

64	تحويل النص إلى كلام منطوق to speech Spoken	8 - 6 - 3
65	فهم الصوت understand the audio	9 - 6 - 3
66	التعرف الضوئي على الحروف Optical Character Recognition	10 - 6 - 3
66	مستويات تحليل اللغات الطبيعية	7 - 3
66	التحليل الصرفي	1 - 7 - 3
67	التحليل النحوي	2 - 7 - 3
67	التحليل الدلالي	3 - 7 - 3
الفصل الرابع = الإطار التحليلي		
70	تمهيد	1 - 4
70	اهمية الصرف في اللغة العربية	2 - 4
71	نظرة عامة إلى العلاقة بين اللغة العربية وتقنيات الحاسوب	3 - 4
73	تعريفات مصطلحات	4 - 4
74	معالجة الصرف العربي آلياً	5 - 4
75	خصائص الصرف العربي	6 - 4
76	خصائص المعالج الصرفي الآلي	7 - 4

77	مراحل تطور معالجة الصرف العربي آلياً	8 - 4
77	بناء أدوات التحليل والتوليد الصرفي الآلي	9 - 4
78	تعريف المحلل الصرفي الآلي	10 - 4
79	الفعل الثلاثي المضعف	11 - 4
81	خوارزمية المحلل الصرفي	12 - 4
الباب الخامس - الإطار التصميمي		
86	مقدمة عن لغة ال C#	1 - 5
86	مميزات C#	2 - 5
89	اهم المصطلحات في لغة C#	3 - 5
89	تصميم الخوارزمية المقترحة لعملية تجريد النص العربي	4 - 5
91	تفاصيل خطوات الخوارزمية المقترحة للتصنيف	5 - 5
الفصل السادس - الإطار التطبيقي		
107	المدرّب	1 - 6
107	الشاشة الرئيسية للنظام التدريب	1 - 1 - 6
108	خطوات انشاء فئة جديدة بالبرنامج	2 - 1 - 6
112	خطوات تدريب الفئة الجديدة	3 - 1 - 6

110	خطوات فتح الفئات المدربة	4 - 1 - 6
111	خطوات الدخول علي اوزان الكلمات في فئة محددة	5 - 1 - 6
112	خطوات اضافة الافعال المضعفة والاسماء الثنائية	6 - 1 - 6
112	المصنف	2 - 6
112	خطوات تصنيف مجموعة نصوص عربية مجمعة	1 - 2 - 6
113	خطوات تصنيف المستندات العربية	2 - 2 - 6
114	نتائج التصنيف	3 - 2 - 6
الفصل السابع = النتائج و التوصيات و المراجع و المصادر		

118	الخاتمة	1 - 7
119	النتائج	2 - 7
122	التوصيات	3 - 7
123	المصادر و المراجع	4 - 7
123	المراجع العربية	1 - 4 - 7
125	المصادر الأجنبية	2 - 4 - 7
126	مواقع الانترنت	3 - 4 - 7
127	المنشورات	4 - 4 - 7
128	الأوراق العلمية	5 - 4 - 7
130	الملاحق	5 - 7

فهرس الأشكال

الصفحة	الشكل	الرقم
82	مخطط التحليل الصرفي	1 - 4
90	المخطط الانسيابي لعملية التجريد للخوارزمية المقترحة	1 - 5
100	مخطط اختيار الكلمة	2 - 5
101	معادلة وزن الكلمة	3 - 5
103	تدريب الفئات	4 - 5
104	شاشة التدريب	5 - 5
107	شاشة التصنيف	6 - 5
110	الشاشة الرئيسية للنظام التدريب	1 - 6
111	انشاء فئة جديدة بالبرنامج	2 - 6
112	تدريب الفئة الجديدة	3 - 6
113	فتح الفئات المدربة	4 - 6
113	اضافة الافعال المضعفة والاسماء الثنائية	5 - 6
115	تصنيف مجموعة نصوص عربية مجمعة	6 - 6
116	تصنيف المستندات العربية	7 - 6

117	نتائج التصنيف	8 - 6
-----	---------------	-------

فهرس الجداول

الصفحة	الموضوع	الرقم
83	امثلة الاسماء الثنائية والافعال المضعفة	1-4
91	الاسماء الثنائية والافعال المضعفة	1-5
93	الكلمات المستثناة	2 - 5
96	الكلمات المستثناة للكلمات العربية المشتركة الشائعة جدا في الأخبار	3 - 5
97	الكلمات المستثناة للكلمات المستخدمة كثيرا في الانترنت	4 - 5
98	السوابق واللواحق	5 - 5
120	بيانات الفئات	1 - 7
121	النتائج	1 - 7

الفصل الأول

Chapter one

أساسيات البحث

The essentials of search

أولا: خطة البحث

1.1 تمهيد :

لم يكن من السهل على القطاع الخاص أن ينجز مخططا شاملا ذا رؤية مستقبلية لإعداد المجتمع العربي لعصر المعلومات، واستخدام الحواسيب، نظرا لوضعه الداخلي اجتماعيا واقتصاديا وسياسيا، ولأن هذا المخطط بكل المقاييس يرتبط ارتباطا عضويا بالدول العربية في إطار التسابق الدولي للتمكن من الإعلام والانترنت، الشيء الذي يسمح بالتحكم من جهة في دواليب الأجهزة القومية، ومسايرة التطور العلمي من جهة أخرى للسيطرة آليا على كل التوجهات العلمية على الصعيد الدولي. يمكن القول في هذا الصدد إن وعي الأنظمة العربية كان غائبا بأهمية الطفرة الإعلامية التي تحققت دوليا في مجال التخطيط لعصر المعلومات، أي تصاميم الحواسيب، وأنظمتها، ووضع برامجها، ونظم تشغيلها، وتطوير تقنياتها، حيث تلعب اللغات الوطنية حبر الزاوية، إما أن تكون قوة وذات فعالية في التخطيط القومي من أجل تنمية مستدامة والدفع بعجلة الصناعة الوطنية التي أصبحت خاضعة للتقنيات الحديثة في مضمار ما أضحت تفرضه العولمة من آليات جديدة للتعامل مع التطور الصناعي العالمي، وإما ألا تكون، فتغرق حينذاك في تبعية مطلقة، وتفقد هويتها وقدرتها على اللحاق بالركب الحضاري الذي أصبحت سمته الأساسية قائمة على وسائل الاتصال.

تعد اللغة في هذا المضمار الوجه اللامع والبارز في الخريطة الإعلامية، وأساس

كل المشاريع العلمية، لا باعتبارها أداة التخاطب والتواصل والتعليم والثقافة، بل لكونها

أضحت معياراً قائماً على التحدي لارتباطها بالتخطيط والتصنيع والتنمية، وأي لغة لم تدخل مجال التقنيات الإلكترونية ولم تستوعب التطور الحاصل في مجال الإعلاميات فإن مصير أهلها سيعرف تدهوراً وانحطاطاً.

إن أهم ما يميز عصرنا هو تضخم المعلومات، وضرورة الإحاطة بها لمعالجتها وتحليلها وتنظيمها وفهرستها وتصنيفها واسترجاعها، الأمر الذي لم يعد بإمكان الطاقات البشرية القيام به، هذا بالإضافة إلى أن شيوع استخدام الحواسيب في عملية التعليم والتعلم ومعالجة اللغات الطبيعية.

أضحى حجر الزاوية للتزود بالبحث العلمي ومواكبة الفيض الهائل من المعلومات، مما يدعو حتماً ضرورة التعامل مع الحاسوب، وتطوير أدواته على قاعدة اللغة الوطنية، أي إدخالها في غمار التقنيات الحديثة، واستخدامها في ضوء الوسائل العلمية المتاحة لكي تصبح أداة عملية في مجال النشر الإلكتروني.

2.1 مشكلة البحث:

اللغة العربية هي اللغة الأم لأكثر من 300 مليون نسمة ونجد أنه من الصعب انتشار تقنية حديثة دون أن تستوعب إمكانية التعامل مع اللغة العربية. ومن أهم التقنيات الحديثة واسعة الانتشار التي تستوعب في التعامل معها إمكانيات اللغة العربية، تقنية إرسال الرسائل عبر شبكة الانترنت والتي تسمى البريد الإلكتروني، والتي تتيح التواصل بين كثير من فئات المجتمع عبر الرسائل الإلكترونية وكذلك تبادل المعلومات .

تميز اللغة العربية عن بقية اللغات العامة في توجيه كتابتها من اليمين إلى اليسار وان عدد حروفها 28 حرفاً، وإنها تحتوي على صيغة مختلفة لكل من المفرد والمثنى والجمع، وأيضاً المؤنث والمذكر وأكثر ما يميزها أن معظم كلماتها أصلها فعل وبسبب هذا التعقيد في التركيب الصرفية للغة العربية فإنها تحتاج إلى مجموعة من الإجراءات والخوارزميات لتجربتها وتصنيفها.

تكمن مشكلة الدراسة في صعوبة تجريد وتصنيف النص العربي آلياً بناءً على التركيب الصرفية للغة العربية.

3.1 أهداف البحث :

تهدف الدراسة (باستخدام تقنيات الأنظمة الذكية) إلى إيجاد حلول ناجعة لمشكلة تجريد

وتصنيف النص العربي وذلك بالاتي :

- بناء خوارزمية ذكية لتصنيف النص العربي.
- بناء خوارزمية لتجريد النص العربي لدراسة الأفعال المضعفة والأسماء الثنائية.

• بناء خوارزمية لتجريد النص العربي لدراسة الحرف المكرر.

• تحديد أوزان للكلمات المفتاحية آلياً ويدوياً إذا لزم الأمر .

• بناء قاموس إلكتروني للغة العربية.

4.1 أهمية البحث :

هناك حاجة ماسة إلى منظور عربي للمحتوى على الشبكة العنكبوتية العالمية، بوصفه أحد اهتمامات علم المعلومات، وبما يسهم في الوقوف على الوضع الراهن لطبيعة وخصائص هذا المحتوى، والتعامل معه حفظاً وبحثاً وتصنيفاً واسترجاعاً وأمناً.

في هذا الصدد يتم تسليط الضوء على مجالات المحتوى العربي على الشبكة العنكبوتية العالمية، تجريده وتصنيفه وفئاته، وأنواعه، والسمات العامة التي تميزه عن غيره من المحتويات الأخرى. ولما كانت اللغة العربية ذات خصائص بنائية ونحوية وإملائية قد تتشابه مع (أو تختلف عن) غيرها من اللغات، فقد أصبح لزاماً التطرق لهذه الخصائص، وتسليط الضوء عليها، بحيث تكون نقطة الانطلاق عند تناول المحتوى العربي على الشبكة العنكبوتية العالمية.

5.1 حدود وفترة الدراسة:

تتطرق الدراسة لمجموعة نصوص عربية مجمعة، لعدد 1250 مستند تم تجميعها من صحف (قوون، الراي العام، الأهرام) في الفترة (2010م-2012م)

6.1 منهجية البحث:

سوف تتناول الدراسة بالبحث والتطبيق تصميم خوارزمية ومن ثم تنفيذ برنامج ذكي لتجريد وتدريب وتصنيف النص العربي اليا وذلك باستخدام أدوات النظم الذكية المستخدمة في مجال علوم الحاسوب، النموذج العملي فقد اعتمد الباحث علي حزمة

visual studio2010 وحزمة #c لتكوين النظام الذكي.

7.1 تقسيمات البحث:

تم تقسيم الدراسة الى سبعة فصول موضحة كالآتي:

الفصل الأول: أساسيات البحث ويحتوي علي تمهيد، مشكلة البحث، أهداف البحث،

أهمية البحث، حدود البحث، منهجية البحث، تقسيمات البحث، مصادر جمع

البيانات.

الفصل الثاني : الذكاء الاصطناعي Artificial-Intelligence

الفصل الثالث: معالجة اللغات الطبيعية والمفاهيم المرتبطة بمعالجة اللغات الطبيعية

و اللغة العربية .

الفصل الرابع : الاطار التحليلي:

الفصل الخامس : الاطار التصميمي.

الفصل السادس : الاطار التطبيقي.

الفصل السابع : النتائج والتوصيات والمراجع والمصادر.

8.1 مصادر جمع البيانات:

مصادر رئيسية : الكتب والمراجع والانترنت.

مصادر ثانوية : المجلات، الدوريات، البحوث والدراسات السابقة.