



REPUBLIQUE ALGERIENNE DEMOCRATIQUE ET POPULAIRE

UNIVERSITE DE M'SILA

FACULTE DES MATHÉMATIQUE  
ET DE L'INFORMATIQUE

Département de stic

جامعة المسيلة  
كلية الرياضيات والإعلام الآلي  
مكتبة الكلية  
11/15/2014/1,47

*Mémoire présenté pour l'obtention du diplôme de master*

*Domaine : mathématique et informatique*

*Filière : Informatique*

*Spécialité : technologie de l'information et la communication*

# **L'apport de la sélection des termes sur la catégorisation automatique des textes arabe**

Préparer par :

**HAFIDI Mustapha**

supervise par :

**Mr.KADRI Said**

2014/2015

## Tableau de matières

Remerciements .....	3
Introduction générale .....	9
<b>Chapitre 1 : catégorisation</b>	
<b>1.1. Introduction .....</b>	<b>11</b>
1.2. Pourquoi automatiser la classification ?.....	12
1.1. Définition de la Catégorisation de textes .....	14
1.2. Historique de la Catégorisation de textes .....	15
1.3. Les systèmes de classification et vocabulaire utilisé .....	16
1.3.1. Catégorisation (Supervisé).....	17
1.3.2. Clustering (Non supervisé) .....	17
1.4. Classification de textes et Text Mining .....	17
1.5. Démarche à suivre pour la catégorisation de textes .....	18
1.6. Le texte .....	19
1.7. Prétraitements .....	20
1.7.1. La segmentation.....	21
1.7.2. Suppression des mots fréquents.....	21
1.7.3. Suppression des mots rares.....	23
1.7.4. Le traitement syntaxique.....	23
1.7.5. Le traitement sémantique.....	23
<b>1.8. Définition de descripteurs.....</b>	<b>24</b>
1.8.1. Représentation en « sac de mots » « bag of words ».....	24
1.8.2. Représentation des textes par des phrases .....	25

1.8.3.	Représentation des textes avec des racines lexicales (stemming) .....	26
<b>1.9.</b>	<b>Pondération ou calcul de poids.....</b>	<b>27</b>
1.9.1.	Le codage TFIDF:.....	27
1.9.2.	Codage TFC : .....	28
1.9.3.	Le codage Lnu : .....	28
1.9.4.	L'entropie .....	29
1.10.	Conclusion .....	30

## **Chapitre 2: algorithmes**

2.1.	Introduction .....	32
2.1.1.	L'apprentissage automatique .....	32
2.1.2.	L'apprentissage supervisé.....	32
2.1.3.	La catégorisation est un problème de classification supervisée .....	33
2.1.4.	Comment classer ?.....	33
2.2.	Différents modèles de classifieurs.....	34
2.2.1.	Machines à Vecteurs Support – SVM .....	35
2.2.2.	Rocchio .....	36
2.2.3.	Méthode du centroïde .....	37
2.2.4.	K plus proches voisins - kPPV .....	39
2.2.5.	L'algorithme Naïve Bayes.....	40
2.2.6.	Arbres de décision .....	41
2.2.7.	Autres méthodes .....	45
2.3.	Conclusion .....	45

## **Chapitre 3: Méthodes de Sélection des termes**

3.1.	Introduction .....	48
3.2.	Le but de la Réduction de la taille du vocabulaire .....	48

3.2.1.	La fréquence ( «document frequency» ) : .....	49
3.2.2.	Le gain d'information ( «information gain» ) : .....	49
3.2.3.	L'information mutuelle ( «mutual information» ) : .....	49
3.2.4.	La statistique du $\chi^2$ : .....	50
3.2.5.	La force du terme ( «term strength» ) : .....	50
3.3.	Conclusion .....	52

## Chapitre 4: réalisation

4.1.	Introduction .....	54
4.2.	Méthodologie.....	54
4.2.1.	Corpus de textes.....	54
4.2.2.	Caractéristiques des classificateurs utilisés .....	55
4.3.	Environnement matériel et logiciel : .....	56
4.3.1.	Le langage de programmation (c#) : .....	56
4.3.2.	L'environnement de programmation : .....	57
4.4.	Structure et fonctionnement de l'application : .....	57
4.4.1.	Interface principale .....	57
4.4.2.	La Prétraitement des textes et calcule les fréquences des termes : .....	58
4.4.3.	Représentation des textes par la fréquence TF*IDF avec l'étape du l'apprentissage : .....	59
4.4.4.	L'algorithme Naïve bayes : .....	60
4.4.5.	Le résultat de catégorisation .....	61
4.4.6.	Réduction de la taille du vocabulaire : .....	64
4.4.7.	Résultat de catégorisation après la rédaction : .....	64
4.4.8.	Evaluation de résultat final : .....	65
4.5.	Evaluation des résultats .....	66
4.5.1.	Le résultat avant la rédaction : .....	66

3.5.2. le méthode de « la méthode $x^2$ » :.....	66
4. Conclusion :.....	67
Conclusion générale .....	68
Bibliographie .....	69

## Introduction générale

La quantité d'information disponible sous format électronique sur Internet ou dans l'intranet des entreprises croît de façon extrêmement rapide. Par exemple, le 22 juillet 2002, Google avait recensés 2.073.418.204 pages ; le 16 mars 2003, ce nombre est passé à 3.083.324.652 pages, soit moitié plus en moins d'un an (voir <http://www.google.fr>). L'utilisateur, submergé par cette masse d'informations, ne se pose plus la question d'accéder à l'information. Son problème devient : comment trouver l'information dont il a besoin, parmi toutes celles qui est accessible. Dans ce contexte, La classification de textes a pour objectif de regrouper les textes similaires, c'est à dire thématiquement proches, au sein d'un même ensemble. L'intérêt d'une telle démarche est d'organiser les connaissances de façon à pouvoir effectuer, par la suite, une recherche ou une extraction d'information efficace. Le volume de documents numériques s'accroissant, des besoins en classification automatique se sont fait ressentir aussi bien sur internet (moteurs de recherche), qu'au sein des entreprises (classement de documents internes, dépêches d'agences, etc.). On distingue dans le domaine de la classification automatique deux types d'approches : la classification supervisée et la classification non supervisée. Ces deux méthodes diffèrent sur la façon dont les classes sont générées. En effet dans le cas de la classification non supervisée, les groupes de documents (classes) sont calculés automatiquement par la machine [SAL 83, IWA 95], tandis qu'ils sont, dans l'approche supervisée [JOA 98b, SEB 02, YAN 99a], définis par un expert. Dans ce dernier cas, il est intéressant de représenter les documents et les classes à l'aide d'un même formalisme et celui généralement utilisé est un espace vectoriel [SEB 02, BES 02]. Dans cet article, nous nous intéresserons à la catégorisation, c'est à dire aux algorithmes d'apprentissage supervisés et plus particulièrement aux méthodes basées sur une représentation vectorielle de documents

## **Conclusion générale**

La classification de textes s'est avérée au cours des dernières années comme un domaine majeur de recherche pour les entreprises comme pour les particuliers.

Ce dynamisme est en partie dû à la demande importante des utilisateurs pour cette technologie. Elle devient de plus en plus indispensable dans de nombreuses situations où la quantité de documents textuels électroniques rend impossible tout traitement manuel. La catégorisation de textes a essentiellement progressé ces dix dernières années grâce à l'introduction des techniques héritées de l'apprentissage automatique qui ont amélioré très significativement les taux de bonne classification. Il reste néanmoins difficile de fournir des valeurs chiffrées sur les performances qu'un système de classification peut actuellement atteindre. Les travaux de recherche dans le domaine se focalisent surtout sur deux aspects : l'efficacité et l'amélioration de performances. Dans ces deux optiques nous avons entamé notre Project de catégorisation et sélection des termes.

## Bibliographie

- BÉCHET, N. (2008). « Extraction et regroupement de descripteurs morphosyntaxiques pour des processus de Fouille de Textes ». doctorat.
- BÉCHET, N. (2008, dec 8). « Extraction et regroupement de descripteurs morphosyntaxiques pour des processus de Fouille de Textes ». doctorat.
- C.D.Loupy. (2000). « *Évaluation de l'Apport de Connaissances Linguistiques en Désambiguïsation Sémantique et Recherche Documentaire* » .
- F.Sebastiani. (2002). « *Machine learning in automated text categorization* » .
- G.Brown, H. (1998). « *The Guru System in TREC-6* » .
- I.Moulinier. (1996). « *Une approche de la catégorisation de textes par l'apprentissage symbolique* ». *interstices*. (s.d.). Récupéré sur [interstices.info](http://interstices.info):  
[https://interstices.info/encart.jsp?id=c\\_41867&encart=3&size=600,500](https://interstices.info/encart.jsp?id=c_41867&encart=3&size=600,500)
- J. Brank, M. G.-F. (2002). Interaction of Feature Selection Methods and Linear Classification Models. *Workshop on Text Learning*.
- J.Clech. (2004). « *Contribution méthodologique à la fouille de données complexes* » .
- J.Clech, D. (2004). « *Une technique de réétiquetage dans un contexte* .
- Karima, A. (2011). Consulté le 2011, sur  
<http://share.esi.dz/55/1/La%20cat%C3%A9gorisation%20de%20texte%20multilingue.pdf>
- Lang, K. (1995). « *NewsWeeder : Learning to Filter Netnews* » .
- Lewis. (1992). « *An evaluation of phrasal and clustered representations on a text* .
- M.McGill, G. &. (1983). « *Introduction to Modern Information Retrieval* » .
- P.Hayes, S. (1990). « *Construe/Tis : A system for content-based* .
- Pedersen., Y. Y. (1997). Comparative Study on Feature Selection in Text Categorization.
- R.Armstrong, D. T. (1995). « *WebWatcher : a Learning apprentice for the World Wide Web* » .
- Sebastiani, F. (1999). A Tutorial on Automated Text Categorisation. *ASAI-99*.
- Sebastiani, F. (1999). A Tutorial on Automated Text categorisation .
- Sebastiani, F. (2002). « *machine learning in automted texte categorization* ». Récupéré sur  
<http://nmis.isti.cnr.it/sebastiani/Publications/ACMCS02.pdf>
- Stricker, M. (2000). « Réseaux de neurones pour le traitement automatique du langage : conception et réalisation de filters d'informations ». paris, doctorat.
- theses*. (s.d.). Récupéré sur [theses.ulaval.ca](http://theses.ulaval.ca):  
<http://theses.ulaval.ca/archimede/fichiers/22376/ch02.html>
- C.D.Loupy. (2000). « *Évaluation de l'Apport de Connaissances Linguistiques en Désambiguïsation Sémantique et Recherche Documentaire* » .
- F.Sebastiani. (2002). « *Machine learning in automated text categorization* » .
- G.Brown, H. (1998). « *The Guru System in TREC-6* » .
- I.Moulinier. (s.d.). *Une approche de la catégorisation de textes par l'apprentissage symbolique* .
- I.Moulinier. (1996). « *Une approche de la catégorisation de textes par l'apprentissage symbolique* » .
- J.Clech. (2004). « *Contribution méthodologique à la fouille de données complexes* » .
- J.Clech, D. (2004). « *Une technique de réétiquetage dans un contexte* .
- Lang, K. (1995). « *NewsWeeder : Learning to Filter Netnews* » .
- Lewis. (1992). « *An evaluation of phrasal and clustered representations on a text* .
- M.McGill, G. &. (1983). « *Introduction to Modern Information Retrieval* » .
- P.Hayes, S. (1990). « *Construe/Tis : A system for content-based* .
- R.Armstrong, D. T. (1995). « *WebWatcher : a Learning apprentice for the World Wide Web* » .
- Sebastiani, F. (1999). A Tutorial on Automated Text categorisation .
- Stricker, M. (2000). « Réseaux de neurones pour le traitement automatique du langage : conception et réalisation de filters d'informations ». paris, doctorat.

## Résumé/ملخص / Abstract

Avec l'avènement de l'informatique et l'accroissement du nombre de documents électroniques stockés sur les divers supports électroniques et sur le Web, particulièrement les données textuelles, le développement d'outils d'analyse et de traitement automatique des textes, notamment la classification automatique de textes, est devenu indispensable, pour assister les utilisateurs, de ces collections de documents, à explorer et à répertorier toutes ces immenses banques de données textuelles. Ainsi la catégorisation automatique de textes, qui consiste à assigner un document à une ou plusieurs catégories, s'impose de plus en plus comme une technologie clé dans la gestion de l'intelligence, les résultats obtenus sont utiles aussi bien pour la recherche d'information que pour l'extraction de connaissance soit sur internet (moteurs de recherche), qu'au sein des entreprises (classement de documents internes, dépêches d'agences, etc.). À l'égard des différentes approches de classification automatique de textes, décrites dans l'état de l'art, se reposant sur une architecture classique basée sur un seul point de vue, nous avons introduit une nouvelle utilisation du classifieur « Naïve Bayes » avec la méthode de sélection des termes " chi2", L'objectif principal de nos travaux, est d'améliorer les performances et l'efficacité du modèle de classification.

**Mots clés :** classification , catégorisation automatique , textes , Naïve Bayes , méthode de sélection , chi2.

مع تطور المعلوماتية وتزايد الوثائق الالكترونية ، تاتي ادوار التحليل و المعالجة الاوتوماتيكية لهذه البيانات النصية كضرورة حتمية لمساعدة المستخدمين للاستكشاف وفهرسة هذه القواعد النصية الضخمة . و في هذا الإطار ، اصبح التصنيف الاوتوماتيكي للنصوص وسيلة من السائل التكنولوجية الرئيسية لإدارة هذا النوع من الإشكاليات ، و النتائج المتحل عليها مفيدة سواء في البحث عن المعلومات او استخراجها . على غرار التقنيات المختلفة المستعملة في هذا المجال قمنا بتصنيف النصوص عن طريق خوارزمية « Naïve Bayes » مع استعمال طريقة " chi2 " لتقليل عد كلمات النص بدون التأثير في التصنيف و الهدف من هذه المذكرة هو امكانية تصنيف النصوص باستخدام طرق رياضية و مقارنة النتائج .

**المصطلحات الرئيسية :** التحليل ، و المعالجة الاوتوماتيكية ، النصية ، خوارزمية ، تصنيف ، طريقة.

*With the advent of computers and the increasing number of electronic documents stored on various electronic media and web, especially text data, development of analysis tools and automatic processing of texts, including automatic text classification has become essential to assist users of these document collections, to explore and identify all these huge banks of textual data. And automatic categorization of text, which is to assign a document to one or more categories, is becoming increasingly recognized as a key technology in the management of intelligence, the results are useful both for the search information to extract knowledge or on the Internet (search engines), and at the company (ranking of internal documents, news agencies, etc.). In respect of different approaches to automatic text classification, described in the prior art, relying on a conventional architecture based on a single point of view, we introduced a novel use of the classifier "Naïve Bayes" with the method of selection of terms « chi2 » The main objective of our work is to improve the performance and efficiency of the classification model. The reference corpus Reuters will be used to conduct a comparative study of results.*

**Kay words :** analysis , automatic processing , texts , classification , categorization , Naïve Bayes , chi2.