



N° d'ordre : .....

**UNIVERSITE DE M'SILA**  
**FACULTE DES MATHEMATIQUES ET DE**  
**L'INFORMATIQUE**  
Département d'Informatique

**MEMOIRE de fin d'étude**  
**Présenté pour l'obtention du diplôme de MASTER**  
**Domaine : Mathématiques et Informatique**  
**Filière : Informatique**  
**Spécialité : Systèmes d'Informations Avancés**

Par: DAHOUMI Fares

**SUJET**

**Identification de la langue et catégorisation thématique des  
textes d'un corpus multilingue en utilisant les algorithmes :  
NB, SVM**

Soutenu publiquement le :    /    /2013 devant le jury composé de :

|                |                      |            |
|----------------|----------------------|------------|
| .....          | Université de M'sila | Président  |
| Mr. KADRI Said | Université de M'sila | Rapporteur |
| .....          | Université de M'sila | Examineur  |
| .....          | Université de M'sila | Examineur  |

**Promotion : 2012/2013**

# Table des matières

|  |    |
|--|----|
| <b>Introduction général</b> .....  | 1  |
| <b>Chapitre 1. L'intelligence artificielle et data mining</b> .....                | 3  |
| 1. Introduction .....  | 3  |
| 2. L'intelligence artificielle(IA) .....   | 3  |
| 2.1. Définition de l'intelligence artificielle.....                                | 3  |
| 2.2. Autre définitions.....  | 4  |
| 2.2.1. Définitions reposant sur le concept d'Intelligence.....                     | 4  |
| 2.2.2. Définitions prenant l'humain comme référence.....                           | 4  |
| 2.2.3. Définitions liées à la difficulté des problèmes visés .....                 | 4  |
| 2.3. Premiers programmes d'IA.....   | 4  |
| 2.4. Caractérisation de l'intelligence artificielle .....                          | 5  |
| 3. Data Mining.....  | 6  |
| 3.1. Définition du data mining .....   | 6  |
| 3.2. Pourquoi faire la fouille de données ?.....                                   | 7  |
| 3.3. Les enjeux de la fouille de données .....                                     | 7  |
| 3.4. Les tâches de la fouille de données .....                                     | 8  |
| 3.4.1. La classification.....  | 8  |
| 3.4.2. L'estimation.....   | 8  |
| 3.4.3. La Segmentation .....   | 8  |
| 3.4.4. La Prédiction .....   | 9  |
| 3.4.5. L'association .....   | 9  |
| 3.5. Processus de la fouille de données .....                                      | 9  |
| 4. La fouille de textes (Text Mining).....   | 11 |
| 4.1. Définition de la fouille de textes .....                                      | 11 |
| 4.2. Les objectifs de la fouille de textes .....                                   | 11 |
| 4.3. Les tâches .....  | 12 |
| 5 Conclusion.....  | 14 |
| <b>Chapitre 2 . Apprentissage automatique et classification de documents</b> ..... | 15 |
| 1. Introduction.....   | 15 |
| 2. apprentissage automatique .....   | 15 |
| 2.1. Qu'est-ce que l'apprentissage automatique .....                               | 15 |
| 2.2. Les tâches de l'apprentissage .....   | 16 |
| 2.2.1. L'apprentissage supervisé .....   | 16 |
| 2.2.2. Apprentissage non-Supervisé .....   | 16 |
| 3. Classification des documents textuels .....                                     | 16 |
| 3.1. L'objectif de la classification.....  | 16 |
| 3.2. Méthodes de la classification .....   | 17 |
| 3.2.1. Classification exclusive/non exclusive.....                                 | 18 |

|  |           |
|--|-----------|
| 3.2.2. Classification supervisée/non supervisée .....  | 18        |
| 3.2.3. Classification hiérarchique/partitionnement.....  | 18        |
| 4. La catégorisation (classification automatique supervisée) automatique de textes.....            | 18        |
| 4.1. L'historiqu .....   | 18        |
| 4.2. Définition .....  | 19        |
| 4.3. Définition formelle.....  | 19        |
| 4.4. Comment catégoriser un texte ?.....   | 20        |
| 5. Approches pour la représentation de textes .....  | 21        |
| 5.1. Choix des termes .....  | 21        |
| 5.1.1. Représentation en «sac de mots».....  | 22        |
| 5.1.2. Représentation des textes par des phrase .....  | 23        |
| 5.1.3. Représentation des textes avec des racines lexicales et des lemmes.....                     | 23        |
| 5.1.4. Méthodes basées sur les n-grammes.....  | 23        |
| 5.2. Codage des terme .....  | 24        |
| 5.2.1. Le codage TF_IDF .....  | 24        |
| 5.2.2. Codage TFC .....  | 25        |
| 5.3. Réduction de la dimension .....   | 25        |
| 5.4. Extraction de termes.....   | 26        |
| 6. Conclusion .....  | 26        |
| <b>Chapitre 3. Identification de la langue et Les algorithmes de catégorisation de textes.....</b> | <b>27</b> |
| 1. Introduction .....  | 27        |
| 2. Méthodes d'identification de la langue.....   | 27        |
| 2.1. Approches linguistiques .....   | 27        |
| 2.1.1. Présence de certains chaînes de caractères spécifiques .....                                | 27        |
| 2.1.2. Présence de certains mots.....  | 28        |
| 2.2. Approches statistiques et probabilistes .....   | 28        |
| 2.2.1. Mots les plus fréquents.....  | 29        |
| 2.2.2. Méthodes basées sur les n-grammes.....  | 30        |
| 3. Les algorithmes de catégorisation de textes .....   | 30        |
| 3.1 Naïve Bayes .....  | 30        |
| 3.1.1. Les limites .....   | 31        |
| 3.2 K plus proches voisins .....   | 32        |
| 3.2.1. Définition de la distance.....  | 32        |
| 3.2.2. Choisir K .....   | 33        |
| 3.2.3. Les domaines d'application .....  | 33        |
| 3.2.4. Limites.....  | 33        |
| 3.3 Les arbres de décision .....   | 34        |
| 3.3.1. Principe.....   | 35        |
| 3.3.2. Les domaines d'application .....  | 35        |
| 3.3.3. Les limites .....   | 36        |
| 3.4 Les réseaux de neurones.....   | 36        |
| 3.4.1. L'auto-apprentissage .....  | 37        |
| 3.4.2. Les domaines d'application .....  | 37        |

|  |           |
|--|-----------|
| 3.4.3. Les limites .....   | 38        |
| 4. Mesure d'erreur en apprentissage .....  | 38        |
| 5. conclusion .....  | 40        |
| <b>Chapitre 4 . Machines à support de vecteurs.....</b>                                | <b>41</b> |
| 1. Introduction.....   | 41        |
| 2. Machines à support de vecteurs (ou SVM).....  | 41        |
| 3. Principe.....   | 43        |
| 3.1 Séparateur linéaire.....   | 43        |
| 3.1.1. Forme primale .....   | 45        |
| 3.1.2. Forme duale .....   | 45        |
| 3.2 Séparateur non linéaire.....   | 46        |
| 4. La régression.....  | 50        |
| 5. Domaines d'application des SVM .....  | 51        |
| 5. conclusion .....  | 51        |
| <b>Chapitre 5. Implémentation et expérimentation .....</b>                             | <b>52</b> |
| 1. Introduction.....   | 52        |
| 2. Présentation des caractéristiques techniques du système .....                       | 52        |
| 2.1. Visual Studio 2012.....   | 52        |
| 2.2. Langage de programmation .....  | 53        |
| 3. Présentation du corpus d'expérimentation.....                                       | 53        |
| 4. Processus de catégorisation entrepris par notre application .....                   | 55        |
| 5. Présentation de l'application réalisée(le classifieur) .....                        | 56        |
| 5.1. Interface principale.....   | 56        |
| 5.2. Prétraitement sur les textes .....  | 57        |
| 5.2.1. Punctuation.....  | 57        |
| 5.2.2. Racinisation .....  | 57        |
| 5.2.3. Mots outils ou mots vides (stop-words) .....                                    | 58        |
| 5.3. Identification des langues.....   | 60        |
| 5.4. Représentation de texte par tf*Idf et apprentissage .....                         | 61        |
| 5.5. Choix de l'algorithme d'apprentissage .....                                       | 62        |
| 6. Évaluation des résultats du classifieur .....                                       | 63        |
| 7. Interprétation des résultats obtenus.....   | 67        |
| 7.1. Identification de langue.....   | 67        |
| 7.2. Classification des documents selon leur contenu (classification thématique) ..... | 67        |
| 7.2.1. Pour les documents anglais .....  | 67        |
| 7.2.2. Pour les documents français .....   | 68        |
| 7.2.3. Pour les documents arabe .....  | 68        |
| 8. conclusion .....  | 69        |
| <b>Conclusion générale .....</b>   | <b>70</b> |
| <b>Bibliographie .....</b>   | <b>72</b> |

## INTRODUCTION GENERALE

La recherche en informatique accorde ces dernières années, beaucoup d'importance au traitement des données textuelles. Ceci pour plusieurs raisons : un nombre croissant de collections mises en réseau et distribuées sur le plan international, le développement des infrastructures de communication et d'Internet. Les traitements manuels de ces données s'avèrent très coûteux en temps et en personnel, ils sont peu flexibles et leur généralisation à d'autres domaines est presque impossible ; c'est pour cela que l'on cherche à mettre au point des méthodes automatiques.

Le domaine de la fouille de textes (Texte Mining) s'est développé pour répondre à volonté à la gestion par contenu des sources volumineuses de textes. A l'heure actuelle, de nombreux logiciels de classification de textes sont disponibles, ils ont fait l'objet de publications et leurs champs d'application s'élargit de jour en jour. En général, ces systèmes sont basés sur des algorithmes d'apprentissage automatique (approche statistique, approche syntaxique et approche connexionniste).

Nous nous intéressons ici plus particulièrement aux algorithmes d'apprentissage et nous avons utilisé l'algorithme de Support Vector Machines (SVM). Pour pouvoir utiliser de tels algorithmes, il est nécessaire de transformer les données, initialement en format texte, en une représentation numérique. Nous avons choisi pour ce faire, la méthode de sélection des termes les plus pertinents. Une fois ce prétraitement terminé, nous pouvons effectuer la classification à l'aide de ces algorithmes.

Ce travail s'articule ainsi, autour cinq chapitres organisés comme suit :

Le premier chapitre donne une présentation générale des disciplines suivantes : L'intelligence artificielle (IA), le Data Mining, la fouille de textes ainsi que quelques techniques qui leur sont liées.

Le deuxième chapitre donne une présentation générale sur la apprentissage automatique et la catégorisation des textes, ainsi que les différentes méthodes utilisées dans la littérature.

Le troisième chapitre est consacré à l'explication détaillée des différentes d'approches d'identification de la langue et Les algorithme de catégorisation de textes.

Le quatrième chapitre donne en détails la méthode de classification de textes appelée SVM ou Support Vector Machines.

Le cinquième chapitre décrit le fonctionnement de l'application suivi par une phase d'évaluation des résultats obtenus à l'aide du classifieur implémenté afin de mesurer sa performance.

Enfin, nous clôturerons par une conclusion qui mettra le point sur l'essentiel de ce travail et ses perspectives.

Nous disons de nous que nous sommes des *Hommes supérieurs*, autrement dit des *sages*, en raison de l'importance que nous attribuons à notre intelligence. Pendant des millénaires, nous avons essayé de comprendre le processus de la pensée, à savoir comment un simple amas de chair peut percevoir, comprendre, prévoir et manipuler un monde bien plus étendu et complexe que lui-même. Le domaine de l'intelligence artificielle, ou l'IA, va encore plus loin :

Il tente non seulement de comprendre des entités intelligentes, mais aussi d'en construire.

L'IA est un des champs les plus récents parmi les sciences et l'ingénierie. Les travaux ont véritablement débuté juste après la seconde guerre mondiale et le terme a été forgé en 1956.

À l'heure actuelle, l'IA est composée d'une grande diversité de sous-disciplines allant des plus générales (apprentissage, perception) aux plus spécifiques (jouer aux échecs, démontrer des théorèmes mathématiques, écrire des poèmes, conduire un véhicule au milieu de la circulation et diagnostiquer des maladies, etc...).

On retrouve des domaines d'étude similaires en intelligence artificielle.

## 2. L'intelligence artificielle (IA)

### 2.1. Définition de l'intelligence artificielle

L'Intelligence Artificielle (IA) est la science dont le but est de faire par une machine des tâches que l'homme accomplit en utilisant son intelligence. La terminologie « malheureuse » d'Intelligence Artificielle est apparue en 1956. On peut lui préférer celle d'informatique heuristique. On ne parle pas donc de cours de machine intelligente, ni de programme intelligent [34].

## Bibliographie

- [1] Smyth, 2000, Bref historique sur le data mining , 2000.
- [3] Christopher D.Manning & Hinrich Schutze, Foundation of statistical natural language processing. MIT press, 1999.
- [4] Pierce, J.R, an introduction to information theory symbols, signals and noise, Dover publications, 1980.
- [2] Baeza\_yates R. & B Ribeiro-Neto, Modern information retrieval, ACM press books, 1999.
- [5] M.Boudjemia et A.bekri, Application d'un algorithme évolutif au problème de la classification non supervisée. Mémoire Ingénieur, 2007.
- [6] M.J.A. Berry et G. Linoff, Inter Editions, Data mining: techniques appliqués au marketing, à la vente et au service client, 1997.
- [7] P. Naïm et M. Bazsalicza, Editions Eyvolles, Data mining pour le web, 2001.
- [8] Ben Messaoud, Data Mining, institut universitaire de technologie lumière, licence C.E.STAT, laboratoire ERIC, 5 avenue Pierre Mendés France 69676 Bron Cedex, 2007.
- [9] Gilbert, Data Mining : une nouvelle façon de faire de la statistique, chaire de statistique appliquée, conservatoire national des arts et métiers, 292 rue Saint Martin, 75003 Paris, 2006.
- [10] Pascal Vincent, Modèles à noyaux à structure locale, Thèse présentée à la Faculté des études supérieures en vue de l'obtention du grade de Philosophiæ Doctor (Ph.D.) en informatique , Université de Montréal, 2003.
- [11] Jalam, Radwan , Apprentissage automatique et catégorisation de textes multilingues, Thèse de doctorat, Université Lumière Lyon 2, 2003.
- [12] ABIDI Karima, la catégorisation de texte multilingue, mémoire de magistère d'informatique, école supérieur d'informatique, 2011.
- [13] Sebastián Pena Saldarriaga, Approches textuelles pour la catégorisation et la recherche de documents manuscrits en-ligne, thèse de doctorat, Université de Nantes, 2010 .
- [14] Nicolas Béchet, Extraction et regroupement de descripteurs morphosyntaxiques pour des processus de Fouille de Textes, thèse de doctorat, Université des Sciences et Techniques du Languedoc, 2008.
- [15] Gardarin Georges, Huaizhong Kou, Karine Zeitouni DocCat : un composant logiciel de catégorisation de documents et de marquage sémantique , XML Laboratoire PRiSM, Université de Versailles-Saint-Quentin 45 , 2003.

## Conclusion générale

L'utilisation des langues comme moyen de communication à travers le support informatique a été longtemps appréhendée avec beaucoup d'hésitation par la communauté scientifique du monde où cet outil trouvera beaucoup d'utilisations importantes. En effet, les langues rencontrent plusieurs difficultés qui s'y rattachent, notamment : le problème de l'ambiguïté des expressions, le problème de reconnaissance des formes fléchies, par exemple la langue arabe étant fortement flexionnelle, le problème d'absence de travaux publiés sur l'extraction de l'information exprimée en langue arabe à travers l'utilisation de modèles statistiques du langage. S'ajoute à tout cela, la diversité des techniques et méthodes relatives au processus de classification qui pose un problème de choix, tout cela pose un énorme défi difficile à surmonter.

Malgré tous ces problèmes, nous avons abordé ce domaine en espérant qu'on arrive à trouver des solutions faisables à chacun des problèmes précités.

Cependant, et après avoir mené l'étude à sa fin, nous pensons qu'on a quand même pu relever le défi. et par la même occasion, apprendre beaucoup de nouvelles connaissances tout au long de la réalisation de ce travail telles que :

- ✓ La programmation objet sous le langage C #.
- ✓ Des nouveaux concepts dans des domaines d'actualité tels que : l'intelligence artificielle, l'apprentissage automatique, le Data Mining, le Text Mining, et les SRI.
- ✓ La classification automatique supervisée de documents qui se situe dans l'intersection des domaines précités.
- ✓ Le traitement automatique du langage naturel TALN (en particulier la langue arabe).

Toutefois, le sujet abordé étant très vaste, il reste beaucoup à faire pour améliorer notre système. A cet effet, nous proposons comme perspectives :

- ✓ Ajouter d'autres langues latines et non latines pour rendre le système multilingue dans le sens propre du mot.
- ✓ Utiliser des corpus d'apprentissage et de test de grandes tailles pour donner plus de crédibilité aux résultats obtenus.
- ✓ Appliquer d'autres algorithmes d'apprentissage (les arbres de décision, les RNA, AdaBoost, K-NN, Rocchio, les algorithmes génétiques) et effectuer des comparaisons entre eux en terme performance et temps d'exécution.
- ✓ Utiliser d'autres formes de documents (HTML, XML, ...).

- [16] ✓ Etudier le cas de la classification multi-label et la classification non supervisée (clustering).
- [17] ✓ Toute autre idée jugée utile, réalisable et bénéfique dans ce domaine.
- [18] Romana VINCI, classification automatique de textes dans la catégorie non thématiques, these pour obtenir le grade de docteur à Ecole nationale supérieure des télécommunications, Présentée et soutenue le 09 Février 2007.
- [19] Sadik Benson, Analyse de Dictionnaire Termiques pour la Classification Automatique par les Techniques de Text Mining, application à la Langue Arabe, Pour l'Obtention du Diplôme de magister, université sabbat-abasa - soif, 2007.
- [20] <http://www.grappa.univ-lille1.fr/polyx/india/sortir005.html> 24/04/2013
- [21] Loïc Grivel, Outils de classification et de catégorisation pour la fouille de textes, Equipe ISIS, Université de Marne-La-Vallée, 2006.
- [22] Ronen Feldman & James Sanger, The text mining handbook, advanced approaches in analyzing unstructured data, Cambridge University Press, New York, USA, 2007.
- [23] Zighed, D.A. and Rakotonmalala R., Graphes d'Induction, Apprentissage et Data Mining, Hermes Science Publication, Paris, 2008.
- [24] René Lefebvre, Gilles Venturi, Data mining, Gestion de la relation client, Personnalisation de sites web, Eyrolles, 2001.
- [25] Stéphane Tuffery, Fouille de données et scoring, bases de données et gestion de la relation client, Dunod, Paris, 2002.
- [26] Philippe Besse, Cours d'Apprentissage Statistique & Data mining, Institut National des Sciences Appliquées de Toulouse-- 31077 - Toulouse cedex 4, 2006.
- [27] Ludovic Mercier, Les machines à vecteurs support pour la classification en imagerie hyperspectrale: implémentation et mise en œuvre, Travail d'Etude et de Synthèse Technique en informatique, 2010.
- [28] Roda Jourani, Reconnaissance de voyelles pour l'obtention du diplôme des études supérieures approchées, université Mohammed V-royal, Rabat, 2006.
- [29] Delorme, L'aspect de la fouille de données dans l'analyse de texte, Union (Conservatoire national des arts et métiers) 2002, Université de Montpellier, 2002.
- [30] Alexander Beaulieu, 1997, A Support Vector Machine, Modern information retrieval, ACM press book, 2002.

- [16] Kotsiantis S. Supervised Machine Learning , A Review of Classification ,2007
- [17] Bouckaert, Naive Bayes Classifiers That Perform Well with Continuous Variables, Lecture Notes in Computer Science, Volume 3339,Pages 1089 –1094, 2004.
- [18] Romain VINOT, classification automatique de textes dans la catégorie non thématique, thèse pour obtenir le grade de docteur d'école nationale supérieur des télécommunication , Présentée et soutenue le 09 février 2007
- [19] Sadik Bessou, Analyse de Données Textuelles pour la Classification Automatique par les Techniques de Text Mining, application à la Langue Arabe, Pour l'Obtention du Diplôme d magister ,université ferhat abbas – setif, 2007.
- [20] <http://www.grappa.univ-lille3.fr/polys/fouille/sortie005.html> 24/04/2013
- [21] Luc Grivel, Outils de classification et de catégorisation pour la fouille de textes, Equipe ISIS, Université de Marne-La-Vallée, 2006.
- [22] Ronen Feldman & James Sanger, The text mining handbook, advanced approaches in analyzing unstructured data, Cambridge University Press, New York, USA, 2007.
- [23] Zighed, D.A. and Rakotomalala R, Graphes d'induction. Apprentissage et Data Mining. Hermes Science Publication, Paris, 2000.
- [24] René Lefébure, Gilles Venturi, Data mining. Gestion de la relation client, Personnalisation de sites web, Eyrolles, 2001.
- [25] Stéphane Tufferry, Fouille de données et scoring, bases de données et gestion de la relation client, Dunod, Paris, 2002.
- [26] Philippe Besse, Cours d'Apprentissage Statistique & Data mining , Institut National des Sciences Appliquées de Toulouse— 31077 – Toulouse cedex 4, 2006.
- [27] Ludovic Mercier, Les machines à vecteurs support pour la classification en imagerie hyperspectrale implémentation et mise en œuvre, Travail d'Etude et de Synthèse Technique en informatique, 2010.
- [28] Reda Jourani, Reconnaissance de visages, pour l'obtention du diplôme des études supérieures approfondies , université Mohammed v-agdal, Rabat, 2006.
- [29] Delorme, L'apport de la fouille de données dans l'analyse de texte, Cnam (Conservation national des arts et métiers), Centre Régional de Montpellier , 2002.
- [30] **Ricardo** Baeza\_yates. & Berthier Ribeiro-Neto, Modern information retrieval. ACM press books, 1999.

[31] Stuart Russell & Peter Norvig ,Intelligence artificielle ,Pearson , 3<sup>rd</sup> Edition , 2010

[32] Dominique Pastre, cours d' intelligence artificielle, Université Paris 5, 2000

[33] Meliouh.A, cours d'intelligence artificielle, 2éme master informatique system d'information, université de M'sila , 2012 .

[34] <http://social.msdn.microsoft.com> 24/04/2013 .

## ملخص :

التصنيف الآلي للوثائق أصبح ضروريا بسبب حجم الوثائق المتبادلة والمخزنة إلكترونيا. إن تعدد الوثائق وتزايد عددها المستمر نتج عنه صعوبة في وضع منهاج ونحن هنا قدمنا طرق التعليم المتمثلة في SVM و Naïve Bayes و التي تسمح لنا بتصنيف مستند جديد انطلاقاً من مستند مصنف مسبقاً.  
كلمات مفتاحية : التصنيف الآلي، طرق التعليم، SVM، Naïve Bayes.

---

## **Abstract**

Automatic classification supervised document becomes necessary due to the volume of documents exchanged and stored electronically. As there are many documents or their number is constantly increasing. It would be difficult to program in advance of the decision rules to determine the class of a new document. We present learning methods ((SVM) Support Vector Machines And Naïve Bayes) that from documents already classified, used to classify new documents.

**Keyword :** Automatic classification, supervised, Learning methods, SVM , Naïve Bayes.

---

## **Résumé :**

La classification automatique supervisée de document devient nécessaire à cause du volume de documents échangés et stockés sur support électronique. Comme les documents sont nombreux ou que leur nombre augmente sans cesse. Il serait difficile de programmer à l'avance des règles de décision pour déterminer la classe d'un nouveau document. Nous présentons donc des méthodes d'apprentissage ((SVM) Support Vector Machine et Naïve Bayes) qui à partir de documents déjà classés ; permettent de classer de nouveaux documents.

**Mots clés :** La classification automatique, supervisée, méthodes d'apprentissage , SVM, Naïve Bayes.