

## Tables des matières

<b>Table des matières</b> .....	1
<b>Introduction générale</b> .....	4
<b>Les comptes rendus hospitaliers</b> .....	5
<b>Chapitre I. la catégorisation automatique des textes</b> .....	6
<b>1. Introduction</b> .....	7
<b>2. Formalisation du problème</b> .....	7
<b>3. Définition</b> .....	8
<b>4. Application de la catégorisation des textes</b> .....	9
4.1 Indexation automatique des textes.....	9
4.2 Organisation des documents.....	9
4.3 Filtrage des textes.....	9
4.4 Informatique linguistique.....	10
4.5 La catégorisation hiérarchique des pages web.....	10
<b>5. Les méthodes de classification</b> .....	10
5.1 Regroupement hiérarchique.....	10
5.2 Regroupement base sur une partition.....	11
5.3 Méthodes basé sur la densité ou une grille.....	11
5.4 Méthode basé sur un modèle.....	11
<b>6. Les types de classificaion</b> .....	11
6.1 Classification supervisé (catégorisation).....	11
6.2 Classification non supervisé (clustering).....	11
<b>7. Les différents contextes de classificaion</b> .....	12
7.1 La classification bi-classe (Filtrage).....	12
7.2 La classification multi-classe (Routage).....	13
7.1 La classification multi-classe disjointe.....	13
7.1 La classification ordonnée.....	13
<b>8. Lien avec la recherche documentaire</b> .....	13
<b>9. Comment catégoriser un texte</b> .....	13
<b>9.1 Le prétraitement</b> .....	16
9.1.1 La segmentation.....	17
9.1.2 Suppression des mots fréquents et élimination des "Mots vides" .....	17
9.1.3 Suppression des mots rares.....	19
9.1.4 Le traitement morphologique.....	19
9.1.5 Le traitement syntaxique.....	20
9.1.6 Le traitement sémantique.....	20
9.1.7 Réduction de la dimension.....	21
<b>9.2 Codage des termes</b> .....	21
9.2.1 Codage Term frequency × inverse document frequency (TF × IDF).....	21
9.2.2 Le Codage LNU.....	22

<b>9.3 Sélection des termes</b> .....	23
9.3.1 L’algorithme de sélection des termes.....	23
<b>10. Méthodes de représentation des documents</b> .....	24
10.1 La représentation en sac mots .....	24
10.2 La représentation des textes par des phrases .....	25
10.3 Représentation des textes avec des racines lexicales et des lemmes.....	26
10.4 Représentation des textes avec la méthode des n-grammes.....	26
<b>11. Difficultés particulières de la catégorisation des textes</b> .....	28
11.1 Redondance(Synonymie) .....	28
11.2 Polysémie (Ambiguïté) .....	29
11.3 L’homographie .....	29
11.4 Déséquilibre .....	29
11.5 La graphie .....	29
11.6 Les variations morphologiques .....	30
11.7 Les mots composés .....	30
11.8 Présence-Absence de termes .....	30
11.9 Complexité de l’algorithme d’apprentissage.....	30
11.10 Sur-apprentissage.....	30
<b>I.12 Les différents corpus utilisés dans la catégorisation des textes</b> .....	31
<b>Conclusion</b> .....	31
<b>Chapitre II. L’état de l’art</b> .....	32
<b>II.1 Introduction</b> .....	33
<b>II.2 L’état de l’art</b> .....	33
<b>II.2.1 Les projets effectués dans le domaine de la classification des rapports textuels</b> .....	33
<b>II.2.2 Choix du classifieurs</b> .....	34
II.2.2.1 Classificateur bayésien .....	34
II.2.2.2 Les arbres de décision.....	37
a. Définition .....	37
b. Algorithme général d’apprentissage par arbres de décision.....	37
II.2.2.3 Algorithme des k-voisins les plus proches .....	38
a. Définition.....	38
b. Algorithme de classification par k-PPV.....	39
II.2.2.4 Machines à support vectoriel (svm) .....	39
a. Définition.....	40
<b>II.2.3 Comparaison des algorithmes</b> .....	40
<b>II. 2.4 Analyse du corpus Reuters</b> .....	41
<b>II. 2.5 Quelle est la meilleure méthode ?</b> .....	41
<b>Conclusion</b> .....	43
<b>Chapitre III. Conception et implémentation</b> .....	44
<b>III.1 Introduction</b> .....	45

<b>III.2 Le processus du travail</b> .....	45
<b>III.3 Le corpus utilisé dans la phase d'apprentissage</b> .....	47
<b>III.4 La phase d'apprentissage</b> .....	48
III.4.1 présentation de l'environnement de prétraitement Weka.....	48
III.4.2 Charger le corpus d'apprentissage .....	48
III.4.3 Eliminer les mots vides.....	49
III.4.4 techniques utilisés dans la sélection des termes.....	49
III.4.5 sauvegarder les données.....	50
III.4.6 Les résultats obtenus dans la phase d'apprentissage .....	51
<b>III. 5 La phase de test</b> .....	51
III.5.1 L'algorithme de la classification d'un nouveau document.....	51
<b>III. 6 Configuration</b> .....	53
<b>Conclusion</b> .....	54
<b>Conclusion générale et perspectives</b> .....	55
<b>Bibliographie</b> .....	56
<b>List des figures</b> .....	61
<b>Glossaire</b> .....	62