

UNIVERSITÉ LUMIÈRE LYON2
Année 2003

THÈSE
pour obtenir le grade de
DOCTEUR
en
INFORMATIQUE

présentée et soutenue publiquement par

Radwan JALAM
le 4 juin 2003

**Apprentissage automatique et
catégorisation de textes multilingues**

préparée au sein du laboratoire ERIC
Equipe de Recherche en Ingénierie des Connaissances

sous la direction de
Jean-Hugues CHAUCHAT

DEVANT LE JURY, COMPOSÉ DE :

| | |
|--|--|
| Annie MORIN, Rapporteur | Maître de conférences habilitée, IRISA, Rennes |
| Yves KODRATOFF, Rapporteur | Directeur de recherche, CNRS, LRI Orsay |
| Martin RAJMAN, Rapporteur | Professeur, Ecole Polytechnique Fédérale, Lausanne |
| Geneviève BOIDIN-LALLICH, Examineur | Professeur, Université Claude Bernard-Lyon 1 |
| Ludovic LEBART, Examineur | Directeur de recherche, CNRS, ENST Paris |
| Jean-Hugues CHAUCHAT, Directeur de thèse | Professeur, Université Lumière-Lyon 2 |

Table des matières

| | |
|---|-----------|
| Introduction | 1 |
| I. Catégorisation de textes monolingues | 5 |
| 1. Catégorisation de textes | 7 |
| 1.1. Introduction | 8 |
| 1.2. Définition de la catégorisation de texte | 8 |
| 1.3. Comment catégoriser un texte ? | 9 |
| 1.3.1. Représentation, le codage, des textes | 10 |
| 1.3.2. Choix de classifieurs | 11 |
| 1.3.3. Évaluation de la qualité des classifieurs | 12 |
| 1.4. Applications de la catégorisation de texte | 12 |
| 1.4.1. Catégorisation de textes : une fin en soi | 13 |
| 1.4.2. Catégorisation de textes : un support pour différentes appli- cations | 13 |
| 1.5. Difficultés particulières de la catégorisation de textes | 13 |
| 1.5.1. Grandes dimensions | 14 |
| 1.5.2. Imprécision des fréquences | 15 |
| 1.5.3. Déséquilibre | 15 |
| 1.5.4. Ambiguïté | 15 |
| 1.5.5. Synonymie | 15 |
| 1.5.6. Subjectivité de la décision | 16 |
| 1.6. Lien avec la recherche documentaire | 16 |
| 1.7. Jeu de données utilisé pour l'évaluation | 18 |
| 1.8. Conclusion | 19 |
| 2. Approches pour la représentation de textes | 21 |
| 2.1. Introduction | 22 |

| | |
|--|-----------|
| 2.2. Choix de termes | 22 |
| 2.2.1. Représentation en « sac de mots » | 22 |
| 2.2.2. Représentation des textes par des phrases | 23 |
| 2.2.3. Représentation des textes avec des racines lexicales et des lemmes | 24 |
| 2.2.4. Méthodes basées sur les n-grammes | 24 |
| 2.3. Codage des termes | 26 |
| 2.3.1. Codage TF \times IDF | 27 |
| 2.3.2. Codage TFC | 27 |
| 2.4. Réduction de la dimension | 28 |
| 2.4.1. Réduction locale de dimension | 29 |
| 2.4.2. Réduction globale de dimension | 29 |
| 2.4.3. Sélection de termes | 29 |
| 2.4.4. Extraction de termes | 30 |
| 2.5. Conclusion | 30 |
| 3. Sélection multivariée de termes | 33 |
| 3.1. Introduction | 34 |
| 3.2. Méthode du χ^2 univariée | 34 |
| 3.3. Méthode du χ^2 multivarié | 36 |
| 3.4. Expérimentation | 36 |
| 3.5. Conclusion | 38 |
| 4. Pourquoi les n-grammes fonctionnent | 39 |
| 4.1. Introduction | 40 |
| 4.2. Intérêt du codage en n-grammes | 40 |
| 4.3. Étapes de la recherche des mots caractéristiques | 41 |
| 4.3.1. Recherche des n-grammes caractéristiques et des mots qui les contiennent | 41 |
| 4.3.2. Filtrage des mots « parasites » | 42 |
| 4.3.3. Algorithme complet | 42 |
| 4.4. Exemple d'application | 42 |
| 4.4.1. Données indexées de Reuters | 42 |
| 4.4.2. Quelques résultats | 44 |
| 4.4.3. Discussion des résultats sur la collection Reuters | 44 |
| 4.5. Conclusion | 46 |
| 5. Techniques pour la construction de classifieurs | 51 |
| 5.1. Introduction | 52 |
| 5.1.1. Manière de construction du classifieur | 52 |
| 5.1.2. Caractéristique du modèle | 53 |
| 5.2. Méthode de Rocchio | 53 |
| 5.3. Arbres de décision | 55 |

| | | |
|--|--|-----------|
| 5.3.1. | Phase d'apprentissage | 56 |
| 5.3.2. | Phase de classification | 61 |
| 5.3.3. | Critiques de la méthode | 61 |
| 5.4. | Classifieurs à base d'exemples | 62 |
| 5.4.1. | K-plus proches voisins | 63 |
| 5.5. | Fonctions à bases radiales | 67 |
| 5.6. | Machine à Vecteurs de Support | 68 |
| 5.6.1. | Cas des classes linéairement séparables | 69 |
| 5.6.2. | Cas des classes non séparables | 70 |
| 5.7. | Évaluation de classifieurs de textes | 71 |
| 5.7.1. | Évaluation des classifieurs, l'approche « binaire » | 72 |
| 5.7.2. | Évaluation des classifieurs, l'approche « multi-classes » | 76 |
| 5.8. | Contributions personnelles | 77 |
| 5.8.1. | Nouvelle utilisation des SVM | 77 |
| 5.8.2. | Nouvelle utilisation des réseaux RBF | 77 |
| 5.8.3. | Nos expérimentations | 77 |
| 5.9. | Conclusion : quel est le meilleur classifieur ? | 79 |
| II. Catégorisation de textes multilingues | | 83 |
| 6. Catégorisation multilingue : les solutions proposées | | 85 |
| 6.1. | Introduction | 86 |
| 6.2. | Intérêt accru aux traitements multilingues | 86 |
| 6.2.1. | Davantage de collections numériques | 86 |
| 6.2.2. | Plus de personnes connectées en ligne | 87 |
| 6.2.3. | Plus de globalisation et de pays unifiés | 87 |
| 6.2.4. | Réseau plus rapide et plus souple | 88 |
| 6.3. | Recherche documentaire multilingues | 88 |
| 6.3.1. | Approches basées sur la traduction automatique | 89 |
| 6.3.2. | Thésaurus multilingues | 90 |
| 6.3.3. | Utilisation de dictionnaires | 90 |
| 6.4. | Nos solutions pour catégoriser des textes multilingues | 91 |
| 6.4.1. | Premier schéma : le schéma trivial | 92 |
| 6.4.2. | Deuxième schéma : choisir une seule langue d'apprentissage | 93 |
| 6.4.3. | Troisième schéma : mélanger les ensembles d'apprentissage | 93 |
| 6.5. | Conclusion | 94 |
| 7. Identification de la langue | | 97 |
| 7.1. | Introduction | 98 |
| 7.2. | Approches linguistiques | 99 |
| 7.2.1. | Présence de certains chaînes de caractères spécifiques | 99 |
| 7.2.2. | Présence de certains mots | 100 |

| | | |
|-----------|--|------------|
| 7.2.3. | Approche lexicale | 101 |
| 7.2.4. | Approche plus linguistique | 101 |
| 7.3. | Approches statistiques et probabilistes | 102 |
| 7.3.1. | Mots les plus fréquents | 102 |
| 7.3.2. | Méthodes basées sur les n -grammes | 103 |
| 7.4. | Expériences pour la reconnaissance de langue | 107 |
| 7.5. | Conclusion | 109 |
| 8. | Traduction automatique | 111 |
| 8.1. | Introduction | 112 |
| 8.2. | Premières approches historiques de la traduction automatique | 112 |
| 8.2.1. | Décryptage | 112 |
| 8.2.2. | Analyse par micro-contexte | 112 |
| 8.2.3. | Imiter la traduction humaine | 113 |
| 8.3. | Nouvelles approches, plus modestes | 113 |
| 8.3.1. | Mémoire de traduction | 113 |
| 8.3.2. | Sous-langages et langages contrôlés | 114 |
| 8.4. | Évaluer la traduction automatique | 114 |
| 8.5. | Conclusion | 115 |
| 9. | Cadre pour la catégorisation de textes multilingues | 117 |
| 9.1. | Introduction | 118 |
| 9.2. | Méthodes pour la catégorisation de textes multilingues | 118 |
| 9.2.1. | Nouveau cadre pour la catégorisation multilingue | 118 |
| 9.2.2. | Détection de la langue du texte à classer | 119 |
| 9.2.3. | Traduction du texte à classer | 119 |
| 9.3. | Application sur les corpus CLEF | 120 |
| 9.3.1. | Constitution du corpus | 120 |
| 9.3.2. | Représentation des textes | 122 |
| 9.3.3. | Algorithmes d'apprentissage | 123 |
| 9.3.4. | Reconnaissance de la langue | 123 |
| 9.3.5. | Catégorisation des articles | 123 |
| 9.4. | Discussion | 128 |
| 9.5. | Conclusion | 129 |
| | Conclusion et perspectives | 133 |
| | Index des auteurs cités | 137 |
| | Bibliographie | 141 |

Introduction

La recherche accorde, ces dernières années, beaucoup d'importance au traitement des données textuelles [Lebart, 2001, Kodratoff, 2001] et en particulier aux données multilingues. Ceci pour plusieurs raisons : un nombre croissant de collections mises en réseau et distribuées au plan international, le développement de l'infrastructure de communication et de l'Internet, la progression constante du nombre de personnes connectées au réseau mondial et dont la langue maternelle n'est pas l'anglais [Peters and Sheridan, 2001]. Ceci a créé de nouveaux besoins pour organiser et traiter ces immenses volumes de données. Les traitements manuels de ces données (systèmes experts, ou à base de connaissances) s'avèrent très coûteux en temps et en personnel, ils sont peu flexibles et leur généralisation à d'autres domaines sont quasi impossibles ; c'est pourquoi on cherche à mettre au point des méthodes automatiques [Moulinier, 1996, Sebastiani, 2002].

Dans ce travail nous nous intéressons à l'utilisation de l'apprentissage automatique pour catégoriser (ou classer) les textes. L'apprentissage automatique est un processus d'induction général qui permet la construction automatique de classificateurs [Mitchell, 1997]. Ici, il s'agit d'affecter une ou plusieurs catégories à des documents : l'objectif est de trouver une liaison fonctionnelle, que l'on appelle également **modèle de prédiction**, entre les textes à classer et l'ensemble des catégories. Pour estimer le modèle de prédiction, il faut disposer d'un ensemble de textes préalablement étiquetés, dit **ensemble d'apprentissage**, à partir duquel on estime les paramètres du modèle de prédiction le plus performant possible, c'est-à-dire qui produit le moins d'erreurs en prédiction.

L'objectif de la catégorisation de texte est donc d'associer automatiquement une étiquette à tout nouveau texte à classer. Une application typique est le classement automatique de dépêches de presse en différents thèmes (par exemple, "les actualités du monde", "l'économie", "les sports", *etc.*), à l'aide de leurs contenus textuels [Lewis, 1992a].

Dans ce travail, nous étendons la catégorisation de textes à la catégorisation de textes *multilingues*. La nouveauté apportée est la possibilité d'inférer pour un texte ré-

digé dans une langue quelconque. En reprenant l'exemple précédent, nous apprenons à sélectionner, parmi l'ensemble des dépêches, celles qui concernent l'*économie* ; ceci quelque soit leur langue.

Comme dans le cas monolingue habituel, la phase d'apprentissage s'effectue à partir d'un corpus d'apprentissage étiqueté. Bien entendu, comme pour la catégorisation de textes monolingue, le texte à classer doit appartenir aux mêmes domaines que ceux utilisés lors de l'apprentissage. On ne saurait, par exemple, essayer de classer un article scientifique à partir d'un modèle construit sur un ensemble d'apprentissage constitué d'articles de journaux de mode [Sebastiani, 2002].

Cette extension au cas multilingue introduit des contraintes supplémentaires. Il faut adapter le processus habituellement mis en œuvre pour classer les nouveaux textes ; et certaines techniques à base linguistique, utilisées en monolingue, deviennent alors inopérantes [Biskri and Delisle, 2001].

Plan de la thèse

Notre travail s'intéresse à l'application de méthodes issues de l'apprentissage automatique à la catégorisation de textes multilingues. Il comporte deux parties.

Une première partie donne une présentation générale de la catégorisation de textes :

- les définitions, objectifs généraux et domaines d'application (chapitre 1 page 7).
- l'adaptation des algorithmes d'apprentissage aux spécificités des textes ; représentation vectorielle d'un texte ; choix des unités d'information (mots, lemme, phrase, n-grammes, *etc.*) ; choix des valeurs numériques associées (présence/absence, fréquences, information mutuelle, *etc.*) ; réduction de l'espace de représentation, par sélection de termes, ou bien par extraction (chapitre 2 page 21).
- la méthode de sélection de termes multivariée (chapitre 3 page 33).
- le codage en n-grammes (chapitre 4 page 39) :
 - les avantages et inconvénients des n-grammes, en particulier dans le contexte multilingue,
 - le lien entre les n-grammes et les mots.
- les méthodes d'apprentissage et la mesure de leurs performances (chapitre 5 page 51).
- les tests réalisés pour comparer les algorithmes d'apprentissage sur les textes monolingues : les méthodes des "*k*-plus proches voisins", des machines à vecteurs de support (SVM), les fonctions à bases radiales (RBF) et les graphes d'induction (sections 5.8.3 page 77 et 7.4 page 107).

La deuxième partie s'intéresse à l'apprentissage de textes multilingues en comparant deux chaînes possibles (chapitres 6 page 85 et 9 page 117) :

- **chaîne 1** : reconnaissance de la langue, puis utilisation de règles de classement construites pour chaque langue ; il faut alors avoir construit un modèle adapté à chacune des langues.

-
- **chaîne 2** : utilisation de la traduction automatique dans le processus de catégorisation ; cette solution permet d'utiliser un seul ensemble de règles de classement. Ici, il y a deux options :
 1. construire un modèle unique sur l'ensemble d'apprentissage d'une langue donnée ; ensuite, pour classer un nouveau texte, (i) reconnaissance de sa langue, (ii) traduction de ce texte vers la langue d'apprentissage, (iii) application du modèle de prédiction sur le texte traduit ; ici la phase de traduction n'intervient que dans la phase de classement.
 2. faire intervenir la traduction automatique dès la phase d'apprentissage : à partir d'un ensemble étiqueté de textes en différentes langues, traduction automatique de tous ces textes vers une langue cible et apprentissage sur cet ensemble de textes traduits ; ensuite, pour classer un nouveau texte, la procédure est la même qu'en 1.

Ces deux chaînes dépendent de la qualité du traducteur automatique (chapitre 8 page 111).

Nous testons nos algorithmes sur des corpus multilingues pour pouvoir évaluer et comparer la performance de ces deux chaînes de traitement (section 9.3 page 120).

Nos principales contributions portent sur :

- l'exploration des raisons de l'efficacité de la représentation des textes par le codage en n-grammes ; ce travail a été présenté aux JADT-2002 [Jalam and Chauchat, 2002] ; le chapitre 4 page 39 développe ce travail.
- la proposition d'une nouvelle méthode de sélection de termes multivariée pour l'apprentissage basées sur la statistique χ^2 , publiée dans SFdS-2003 [Clech et al., 2003] ; cette proposition est détaillée dans le chapitre 3 page 33.
- une nouvelle utilisation des SVM et des RBF avec le χ^2 comme noyau, publiée dans EGC-2001 et IJCNN-2001 [Jalam and Teytaud, 2001, Teytaud and Jalam, 2001] ; les sections 5.8.3 page 77 et 7.4 page 107 présentent les résultats auxquels nous sommes parvenus avec ces méthodes.
- la proposition d'un nouveau cadre pour la catégorisation de textes multilingues, dans les chapitres 6 page 85 et 9 page 117.
- l'adaptation du corpus CLEF pour disposer d'un échantillon de textes multilingues étiquetés en vue de l'apprentissage ; ce corpus adapté est, à notre connaissance, le premier jeu de données dans ce domaine. Il est présenté au chapitre 9 page 117.