

masling 22

République algérienne démocratique populaire
Ministère De L'Enseignement Supérieur et Recherche scientifique

Université de m'sila

Département de l'informatique et mathématique



Mémoire de fin étude

Pour l'obtention de diplôme de Master informatique

Option : système d'information avancée

Thème

Etude des algorithmes de classification supervisés
KNN et Naïve Bayes
pour la classification des documents textuels

Réalisé par :

YOUCEFI Asma

Encadré par :

Gasmi Abdel Al kader

Bouamama salim

Promotion : 2010/ 2011

Sommaire

Introduction	1
Chapitre 1 – Data Mining et Text Mining.....	3
1. Data Mining.....	3
1.1 Définition de Data Mining.....	3
1.2 pourquoi la fouille de données.....	4
1.3 Les différents types de données.....	4
1.3.1 Les données discrètes.....	4
1.3.2 Les données continues.....	5
1.3.3 Les données textuelles.....	5
1.4 Le processus de fouille de données.....	5
1.4.1 CRISP-DM, le Cross Industry Standard Processfor Data Mining.....	5
1.4.2 Les phases de CRISP-DM.....	7
1.5 Les tâches de fouille de données.....	8
1.5.1 La description	8
1.5.2 L'estimation	8
1.5.3 La prédiction	9
1.5.4 La Classification	9
1.5.5 La segmentation	11
1.5.6 L'association	11
1.6 Exemples d'application possibles.....	13
2. Text Mining.....	13
2.1 Les tâche de Text Mining.....	13
2.2 Les fonctions.....	14
2.3 Méthodes utilisées pour la fouille de texte.....	15
2.4 Les étapes de la fouille de textes.....	16
2.5 Les applications de fouille de textes.....	16
2.5.1 Les études	16
2.5.2 Intelligence économique.....	17
2.5.3 La recherche médicale.....	17

2.5.4 La recherche académique.....	17
2.5.5 Le triage automatisé.....	18
2.5.6 Catégorisation des textes	18
Conclusion	19

Chapitre 2 - La Classification.....21

1 préface.....	21
2 Définition de la classification.....	21
2.1 La classification non supervisée (Clustering).....	22
2.2 La classification supervisée (catégorisation).....	22
2.2.1 Définition	22
2.2.2 Définition formelle.....	22
2.3 Représentation de textes objet de catégorisation.....	23
2.3.1 Représentation en « sac de mots »	23
2.3.2 Représentation par phrases.....	24
2.3.3 Représentation avec des racines lexicales.....	24
2.3.4 Représentation avec des lemmes.....	24
3 Processus de catégorisation.....	24
4 Les applications de la catégorisation.....	25
5 Les algorithmes de catégorisation de documents.....	26
5.1 K plus proche voisin (KNN).....	26
5.2 Naïve Bayes.....	26
5.3 Arbres de décision.....	26
5.4 Les réseaux de neurone.....	27
5.5 Machine à vecteur de support (SVM).....	27
6 Critères d'évaluation des classificateurs.....	28
7 Difficultés particuliers de la catégorisation de textes.....	29
7.1 La subjectivité de la décision.....	30
Conclusion	31

Chapitre 3 - Algorithmes de classification: KNN et NAIVE BAYES.....32

1 Préface.....	32
2 L'Algorithme KNN (K plus proche voisin).....	32

2.1	Présentation de l'algorithme	32
2.2	L'algorithme KNN en détail	32
2.3	Définition de la distance.....	33
2.4	Choix de la fonction de similarité.....	35
2.5	Choisir k	35
2.6	Sélection de la classe.....	35
2.7	Utilisation de l'algorithme K plus proches voisins pour la classification des textes..	36
2.8	Discussion.....	37
3.	L'Algorithme Naïve Bayes	38
3.1	Classificateur bayésien naïf	38
3.2.	L'algorithme NAIVE BAYES en détail.....	39
3.3.	Caractéristiques de l'apprentissage bayésien.....	40
3.4	Discussion.....	40
3.4.1	Avantages.....	41
3.4.2	Inconvénients.....	41
	Conclusion.....	42
	Chapitre 4 - Implémentation et expérimentation.....	43
1	Outil de mise en œuvre.....	43
1.1	Langage de programmation.....	43
1.2	Environnement de développement.....	43
2.	Présentation du corpus d'expérimentation.....	44
2.1	Construction de la liste des mots outils (Stop Word).....	44
3.	Le prétraitement de texte.....	45
4.	La classification	47
4.1	KNN.....	47
4.2	NAIVE BAYES.....	47
4.2.1	La phase d'apprentissage.....	47
4.2.1	La phase de classification.....	48
5.	Résultats expérimentaux.....	48
	Conclusion.....	52
	Références Bibliographique.....	53

Introduction

La recherche 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 au plan international, le développement de l'infrastructure de communication et de l'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.

Parmi les nouvelles technologies qu'on tend à appliquer aux textes, la *fouille de données*, qui ne se limite pas au traitement des données structurées sous forme de tables numériques mais offre des moyens d'investigation des corpus en langage naturel.

Le domaine de la fouille de textes (text 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.

La préparation des données est une étape importante, si ce n'est primordial, du processus d'extraction de connaissances à partir de données. En schématisant, il s'agit de définir au mieux les éléments traités et la représentation utilisée pour l'apprentissage. La qualité (du modèle) de prédiction dépend donc grandement de la qualité de la préparation effectuée en amont.

La préparation consiste à homogénéiser les données et à les disposer en tableau lignes/colonnes, car il s'agit presque toujours de la structure la mieux adaptée à l'exploitation des données. Formellement, chaque ligne/colonne peut être considérée comme un objet vecteur ayant un nombre fixe de composantes. Ce vecteur ligne/colonne sera vu comme un objet mathématique que l'on pourra manipuler selon ses propriétés.

Nous nous intéressons ici plus particulièrement aux algorithmes d'apprentissage et nous avons utilisé et comparé deux algorithmes : les K plus proches voisins (KNN) et

Introduction

l'algorithme NAIVE Bayes. Pour pouvoir utiliser de tels algorithmes, il est nécessaire de transformer les données, initialement en format texte, en une représentation numérique. Une fois ce prétraitement terminé, nous pouvons effectuer la classification à l'aide de nos algorithmes ; ensuite nous ferons une étude comparative entre les deux algorithmes.

Notre travail s'intéresse à l'analyse de données textuelles pour les préparer à la classification par les techniques de fouille de données et faire une comparaison entre les deux algorithmes. Il comporte quatre chapitres et une brève conclusion avec une perspective.

Le premier chapitre est une présentation générale de la fouille de données (Data Mining), et la fouille de textes(text mining).

Le deuxième chapitre présente la notion de la classification et la catégorisation des textes.

Le troisième chapitre détaille les deux algorithmes KNN et NAIVE BAYES.

Le quatrième chapitre est consacré à l'implémentation.

Conclusion et perspectives

Les techniques du data mining ont pour objectif de satisfaire la nécessité de classer des nouveaux documents. Ce projet a pour objectif principal d'étudier et d'implémenter deux algorithmes de Data Mining. Ces algorithmes dits des algorithmes de classification supervisée.

Après avoir réalisé les mesures expérimentales, et interprété les résultats donnés par KNN et ceux donnés par NAIVE Bayes. On montre tantôt que Naïve Bayes est plus performant et plus efficace que KNN tantôt c'est l'inverse.

L'algorithme NAIVE Bayes est simple, en vue des résultats obtenus, on peut conclure que l'algorithme Naïve donne de bons résultats avec un bon taux de précision, l'apprentissage est rapide, permettant de traiter des données volumineuses, il convient donc particulièrement bien au problème de la fouille de textes.

Il serait maintenant intéressant de poursuivre cette recherche, qui a permis de découvrir les avantages de ces algorithmes dans le cadre de la catégorisation de texte. Nous espérons que cette contribution pourra ouvrir de nouvelles perspectives à d'autres collègues et d'ouvrir un champ de recherche.

A mon avis je prépose de faire encore une étude approfondie et utilise un grand nombre de test pour trancher et décider.

Bibliographe

- [1]. Rachel Konrad, 2001, *Data Mining: Digging user info for gold*, ZDNET News, February 7, http://news.zdnet.com/2100-9595_22-528032.html?legacy=zdn
- [2]. The Technology Review Ten, *MIT Technology Review*, January/February 2001.
- [3]. <http://www.grappa.univ-lille3.fr/polys/fouille/sortie005.html>
- [4]. <http://cedric.cnam.fr/~saporta/DM.pdf>
- [5]. http://eric.univ-lyon2.fr/~ricco/cours/slides/Introduction_au_Data_Mining.pdf
- [6]. <http://www.univ-paris1.fr/diplomes/master-droit-du-numerique/bibliotheque-numerique-du-droit-de-ladministration-electronique/tic/informatique/data-mining/data-mining-definition/>
- [7]. <http://www.ultra-fluide.com/ressources/datamining/presentation.htm>
- [8]. Sadik Bessou, *Analyse de Données Textuelles pour la Classification Automatique par les Techniques de Text Mining, application à la Langue Arabe*, Thèse de magister, université Ferhat Abbas SETIF (2006 /2007).
- [9]. Rémi Gilleron & Marc Tommasi, 2000, *Découverte de connaissances à partir de données*, <http://www.grappa.univ-lille3.fr/polys/fouille/sortie005.html>
- [10]. http://www.madchat.fr/reseau/ids%7Cnids/data_mining-fr.htm
- [11]. <http://www.commentcamarche.net/forum/affich-7099472-data-mining-et-text-minig>
- [12]. <http://www.dataalgo.com/60-text-mining.htm>
- [13]. http://lia.univ-avignon.fr/fileadmin/documents/Users/Intranet/fich_art/Patrice/Bellot.pdf
- [14]. <http://translate.google.com/translate?hl=fr&langpair=en%7Cfr&u=http://www.crisp-dm.org/Process/index.htm>
- [15]. <http://translate.google.com/translate?hl=fr&langpair=en%7Cfr&u=http://www.crisp-dm.org/Process/index.htm>
- [16]. Radwan jalam, *Apprentissage automatique et catégorisation de textes multilingues*, thèse de doctorat, université lumière lyon2, 2003
- [17]. Sayad M et Abbari A, *Modèle discriminant pour la classification de documents XML à l'aide des réseaux bayésiens et le noyau de Fisher*, thèse ingénieure, INI, 2009/2010.
- [18]. Belacel nabil, *Méthodes de Classification Multicritère Méthodologie et Applications à l'Aide au Diagnostic Médical*, Thèse doctorat, Université LIBRE DE BRUXELLES, 1999/2000.
- [19]. Saeed Raheel, *L'apprentissage Artificiel pour la Fouille de Données Multilingues : Application à la Classification Automatique des Documents Arabes*, thèse doctorat, Université Lumière Lyon2, 22 octobre 2010.

- [20]. Abidi karima, La catégorisation de texte multilingue, thèse de magistère, INI, 2010/2011.
- [21]. http://en.wikipedia.org/wiki/K-nearest_neighbor_algorithm
- [22]. Catégorisation automatique de textes et cooccurrence de mots provenant de documents non étiquetés fichiers.

التصنيف الاولي للوثائق اصبح ضروريا بسبب حجم الوثائق المتبادل والمخزين إلكترونيا. ان تعدد الوثائق وتزايد عددها المستمر ، نتج الصعوبة في التخطيط مقدما لقواعد قرار لتحديد فئة مستند جديد. ونحن هنا قدمنا طرق التعليم المتمثلة في (Naïve Bayes et KNN) التي، من خلال مستند مصنف مسبقا، تسمح لنا بتصنيف مستند جديدة. معظم خوارزميات التعلم المراقب تسعى لإيجاد نموذج يفسر الارتباط بين المعطيات الاولي والاصناف. بعد التجارب و ترجمه النتائج المحصل عليها اصبح جلي ان NAIVE Bayes احيانا اكثر فعالية من KNN احيانا العكس .

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 ((KNN) K plus proche voisin et Naïve Bayes) qui, à partir de documents déjà classés, permettent de classer de nouveaux documents.

Après avoir réalisé les mesures expérimentales, et interprété les résultats donnés par KNN et ceux donnés par NAIVE Bayes. Tantôt que Naïve Bayes est plus performant et plus efficace que KNN tantôt c'est l'inverse.

Mots clés : documents textuels, catégorisation, sac de mot, KNN, Naïve Bayes.

Abstract

Automatic classification supervised document becomes necessary because of the volume of documents exchanged and stored electronically. As the documents are numerous or their number is growing, it would be difficult to plan in advance of decision rules for determining the class of a new document. We therefore present methods of learning ((KNN) K nearest neighbor and Naive Bayes) which, from documents already classified for classifying new documents. Conducted experiments showed that NAÏVE BAYES classifier had better performance than KNN classifier sometimes KNN had better than Naïve Bayes.

Keywords: textual documents, categorization, Bag-of-words, KNN, Naive Bayes.