

TECHNIQUES D'APPRENTISSAGE

IFT 603

## **FeedUS**

Classificateur de flux RSS

Travail présenté à

M. Shengrui Wang

Par

Marc-Alexandre Côté-Harnois      07 166 997

Julien Filion      07 177 770

Simon Renaud-Deputter      07 149 640

Université de Sherbrooke

Hiver 2010

# Table des matières

Tables de figures :	ii
Introduction	1
Notre projet	1
Introduction au flux RSS	2
Difficultés	2
L'application FeedUS	3
Comparaison	4
Fonctionnement	5
Ajouter un Flux RSS	5
Pré-traitement des articles	5
Affichage des articles	5
Entraînement des algorithmes	6
Algorithmes	7
Algorithme non-supervisé	7
Algorithme supervisé	7
K-NN	7
Évaluation	8
K-NN	8
Naive Bayes	8
Random Forest	8

## **Tables de figures :**

Figure 1 – Interface principale de l'application .....	3
Figure 2 – Tableau de comparaison des algorithmes de classification .....	4

## **Introduction**

Le projet présenté dans ce document a été réalisé dans le cadre du cours IFT603 de l'Université de Sherbrooke par l'équipe Hocus, composée de Marc-Alexandre Côté, Julien Fillion et Simon Renaud, à l'hiver 2010.

### ***Notre projet***

Notre projet est une application permettant de classifier des articles récupérés à partir de flux RSS. Le but est de permettre à l'utilisateur de définir des catégories dans lesquelles seront classés les documents. Une interface simple permet la visualisation des articles directement dans le logiciel ainsi que leur classification. Cette application simplifie la vie de l'utilisateur en répartissant automatiquement les nouveaux documents dans les différentes catégories prédéfinies par celui-ci.

## Introduction au flux RSS

Les flux RSS sont des fichiers dont le contenu est généré automatiquement en fonction des mises à jour d'un site web. On en retrouve sur les sites d'actualité ou les blogs pour présenter les titres des dernières informations consultables en ligne.

Le terme RSS signifie que le contenu du fichier respecte le format RSS qui s'appuie lui-même sur le langage XML. Généralement, chaque flux contient les informations suivantes :

- Titre : Le titre de l'article.
- Description : Brève description de l'article. (Qualité selon l'émetteur du flux)
- Source : Lien Internet vers l'article intégrale.
- Date : Date d'issue de l'article.

### **Difficultés**

La classification de flux RSS présente quelques difficultés. Premièrement, le flux en tant que tel ne contient pas beaucoup d'information. En effet, la majorité des émetteurs de flux RSS n'indique aucune description ou bien simplement la première phrase de l'article. Pour cette raison, il est très difficile de seulement s'y fier.

Afin de contourner le manque d'information dans le flux, on peut suivre le lien menant à la page web où se trouve le document en entier. Cela amène une deuxième difficulté : récupérer uniquement le contenu de l'article.

Troisièmement, l'utilisateur s'attend à voir des résultats rapidement, c'est-à-dire sans avoir à classifier des centaines d'articles avant d'obtenir un classement potable. Pour cette raison, il faut que les algorithmes de classification et de segmentation performant même en présence de données peu nombreuses.

En raison de ces difficultés, un excellent traitement des données avant de les utiliser dans les algorithmes est nécessaire. Également, une part de responsabilité revient à l'utilisateur lorsqu'il classe les documents. Les articles d'une même classe devront être nombreux et présenter de fortes similitudes.

## L'application FeedUS

L'application offre plusieurs fonctionnalités permettant de gérer les flux, les catégories, les documents et les algorithmes. L'arborescence dans le cadre de gauche permet de voir les articles disponibles dans l'application. Les items avec une étoile représentent les catégories, créés soit par l'utilisateur, soit par l'algorithme de segmentation. L'icone de cadenas indique que le document en question a été assigné à une catégorie manuellement par l'utilisateur et qu'il est verouillé. L'enveloppe ouverte indique un document qui a été lu, tandis que l'enveloppe fermée, un nouvel article.

Voici l'interface principale de l'application :

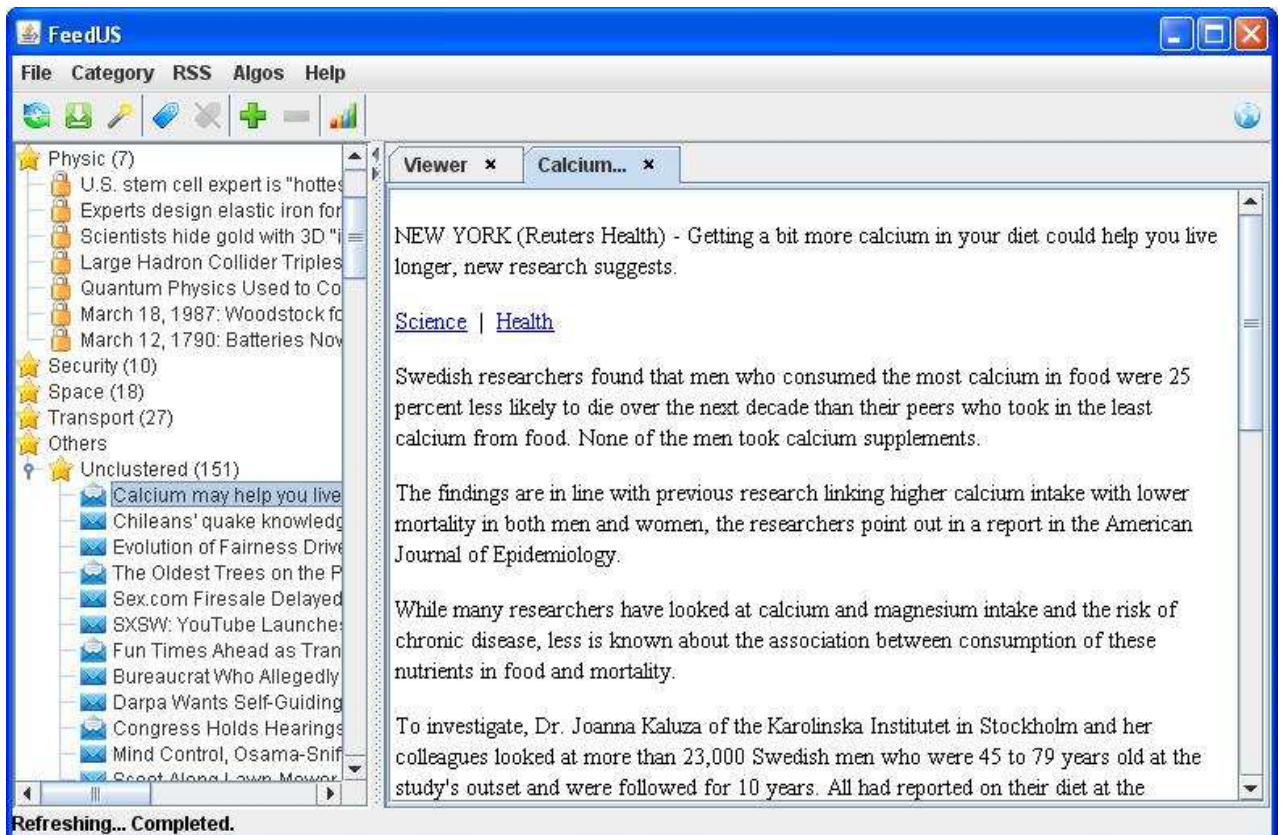


Figure 1 – Interface principale de l'application

## Comparison

Un outil très utile dans l'application est la possibilité de comparer visuellement la justesse des algorithmes de classifications.

Voici le tableau de comparaison :

Documents	KNN	Naive Bayes	Random Forest Tr...	User
Star Wars Lightsaber Bookends	Cinema	Cinema	Gadget	
iBooks Store Loaded with Project Gutenberg Tit...	Gadget	Gadget	Gadget	Gadget
Cheap Pocket Video Camera Shoots for Hours	Gadget	Gadget	Gadget	Gadget
Hipsters Grieve: The \$150 Walmart Fixie	Transport	Transport	Transport	Transport
Netflix Streaming Comes to the Wii	Gaming	Gaming	Gaming	Gaming
3-D Tabletop Display Gets Rid of the Glasses	Gadget	Gadget	Gadget	Gadget
BlizzCon Returns To Anaheim in October		Gaming	Gaming	Gaming
PAX: Shank, The Gory 'Cinematic Brawler'	Gaming	Gaming	Gaming	Gaming
This Week in The Clone Wars: Anakin on the Att...	Cinema	Cinema	Cinema	Cinema
Dork Tower Friday		Cinema	Gaming	
Australian LEGO Champion In Full Flight		Cinema	Military	
8 Things Parents Should Know About How to Tr...	Cinema	Cinema	Cinema	Cinema
Ten Geeky Places to Visit in Portland	Cinema	Cinema	Space	
Need Some Help Making That iPad Decision?	Gadget	Gadget	Gadget	Gadget
March 26, 1999: 'Melissa' Wreaks Havoc on Net	Cinema	Security	Security	Security
TJX Hacker Gets 20 Years in Prison	Security	Security	Security	Security
WIPO: Dope-Vaporizer Seller Not Bogarting Do...		Space	Computer Science	
Retro Tron Intro Mimics '60s Cool of Saul Bass	Cinema	Gaming	Gaming	Gaming
Playlist: Top-Shelf Tracks From SXSW Standouts	Gadget	Gadget	Space	
Alt Text: How Will Nintendo 3DS Work? 5 Eye-P...	Gaming	Gaming	Gaming	Gaming
SXSW: Byrne Doc Ride, Rise, Roar Burns With ...		Gadget	Space	
Calling All Pill-Poppers! Who's Your Alice?	Cinema	Cinema	Cinema	Cinema
Photos: Inside Harry Potter and the Forbidden J...		Cinema	Cinema	Cinema
NASA sets next shuttle launch for April 5	Space	Space	Space	
New iPad Top-Guides: The End of the iPad? Photo...		Gadget	Gaming	

Figure 2 – Tableau de comparaison des algorithmes de classification

## Fonctionnement

### **Ajouter un Flux RSS**

Au premier démarrage de l'application, l'utilisateur doit inscrire des flux RSS à partir desquels les articles désirés seront récupérés. Pour ajouter un flux, il suffit d'avoir son adresse source.

Malheureusement, certains sites ne suivent pas de standards en ce qui concerne l'emplacement du contenu des articles. C'est pourquoi, l'utilisateur peut spécifier certains filtres (sur le nom des balises ainsi que sur les attributs ID, class et name) qui seront appliqués à tous les textes récupérés à partir de ce flux RSS. Il arrive parfois (souvent) que les articles d'un même flux ne soient pas cohérents entre eux. Dans de tels cas, l'application ignore les textes fautifs.

### **Pré-traitement des articles**

Par la suite, l'utilisateur demande au programme d'aller chercher les nouveaux textes sur Internet. Pour chaque document récupéré plusieurs étapes doivent être effectuées afin de le rendre compatible avec les différents algorithmes de classifications. Tout d'abord, la page Internet référencée par le document RSS est récupérée. Pour n'obtenir que le contenu de l'article, les règles associées au flux RSS sont utilisées. Une fois les informations récupérées, elles doivent être épurées de toutes les balises HTML. Ensuite, les mots se trouvant dans la liste de « StopWords », c'est-à-dire la liste des mots les plus fréquents de la langue, sont retirés. Ensuite, une lexémisation (stemming) est appliquée aux mots restants. Finalement, le vecteur TF-IDF de ce document, sera généré, à l'aide de la fréquence de chaque mot et le nombre de documents dans lesquels ce même mot apparaît.

### **Affichage des articles**

Une fois la mise à jour complétée, l'utilisateur peut faire un rafraîchissement de l'arborescence afin de voir les nouveaux articles. L'affichage dépend des algorithmes sélectionnés dans le menu « *Algos* ».

Par défaut, l'algorithme de classification est « *User* » et celui de regroupement est « *None* », ce qui résulte en l'affichage de tous les articles classés par catégorie que l'utilisateur a assignés. Le reste des articles seront regroupés sous une même branche « *Unclustered* ».

L'utilisateur peut changer l'algorithme de classification ce qui affichera le classement des articles selon la méthode utilisée à l'exception des articles verrouillés. Il peut également modifier l'algorithme de segmentation causant ainsi le regroupement des articles non-classifiées sous la branche « *Others* ».



### ***Entraînement des algorithmes***

Après que l'utilisateur ait classé quelques articles dans les catégories qu'il aura préalablement créées, l'application peut (ré)entraîner les algorithmes en se basant sur les décisions de l'utilisateur ainsi que sur les paramètres propres à chaque méthode qu'il aura définies.

## **Algorithmes**

### ***Algorithme non-supervisé***

Nous avons implanté un algorithme de clustering pour aider l'utilisateur à sélectionner les articles qui l'intéressent. Nous appliquons cet algorithme sur tous les documents non classés. Nous avons choisi d'implémenter Bisecting K-Mean comme algorithme non-supervisé. Nous utilisons la distance du cosinus comme mesure de distance.

### ***Algorithme supervisé***

Les algorithmes supervisés nous permettent de proposer à l'utilisateur des documents qu'il risque d'aimer. Nous avons choisi d'utiliser trois algorithmes de classification : K-NN, Random Forest et Naive Bayes. Nous avons utilisé l'implantation du logiciel Rapid Miner pour Random Forest et Naive Bayes alors que nous avons notre propre implémentation de K-NN.

### ***K-NN***

Nous avons rencontré un problème avec le K-NN de base, car cet algorithme est fait pour trouver une classe pour tous les documents. Or, dans notre application, certains documents ne devraient se retrouver dans aucune classe. Pour cette raison, nous avons ajouté un seuil de similarité comme paramètre. Ce seuil intervient lorsque nous calculons les K plus proches voisins. Un document ne peut faire partie des K plus proches voisins que si la similarité entre les documents est plus grande que le produit du seuil et de la similarité moyenne de la classe du document. Pour calculer la similarité moyenne, nous calculons la similarité entre chaque document et son K<sup>ième</sup> plus proche voisin de la même classe. On fait ensuite la moyenne de ces similarités pour chaque classe. Nous ne classons pas un document lorsque nous ne pouvons pas lui associer K voisins. Il est donc important que le seuil se situe entre 0 et 1. Une valeur de 0 est le cas général de K-NN. De plus, nous utilisons la similarité du cosinus comme mesure de similitude.

## **Évaluation**

Nous avons estimé l'erreur de généralisation de notre classification à l'aide de la technique Leave-One-Out (variante de Cross Validation). L'ensemble d'entraînement est composé de 228 documents divisés en 12 catégories où chaque catégorie possède entre 7 et 35 documents.

### **K-NN**

Erreur de généralisation = 27%

Erreur de classification = 15 %

Paramètres :      Nombre de voisins requis = 3

### **Naive Bayes**

Erreur de généralisation = 38%

Erreur de classification = 0 %

Paramètres :      N/A

### **Random Forest**

Erreur de généralisation = 63%

Erreur de classification = 0 %

Paramètres :      Nombre d'arbres utilisés = 10

                 Profondeur limite de l'arbre = infini

                 Nombre d'attributs utilisés =  $\log(M+1)$  , M étant le nombre total d'attributs.