

Processus de classification multi-label hiérarchique sémantique pour le Big Data





- 1. Contexte, Problématique
- 2. Prototype
 - Contraintes, objectifs
 - Architecture
- 3. Processus
- 4. Travaux futurs
 - Evaluations
 - Améliorations
- 5. Conclusion







Secteur d'activité +

Secti	teur a activité +
	Tout cocher ✔ / Tout décocher 🗙
	Administration publique ▼
	Agriculture, sylviculture, pêche
	Art et loisir
	BTP ▼
	Construction de bâtiment
	Démolition
	Génie civil
	Promotion immobilière
	Préparation de site
	Urbanisme
	Banque, finance et assurance
	Commerce et distribution
	Energie
	Enseignement & Formation
	Environnement
	Hôtellerie, restauration, tourisme
	Industrie extractive
	Industrie manufacturière ▼
	Information et communication
	Informatique et télécommunication 🔻
	Matériel informatique
	 Service et ingénierie informatique, NTIC
	Télécommunication
	Recherche & Développement
	Santé et action sociale
	Services ▼
	Transport et logistique ▼







PROBLEMATIQUE

Problème:

Réduire la quantité d'informations à traiter, tout en gardant une représentation proche des connaissances des experts.



PROBLEMATIQUE

Littérature :

- Systèmes existants restreints à un cadre précis
- Pas de système générique
- Les systèmes n'évoluent pas en fonction des données

Verrou scientifique:

Pas de processus d'analyse de données générique, à la fois adapté aux grands volumes de données et à des connaissances métier.



Contraintes

- Traitement générique
 - → Adaptation à différents cas d'utilisation
- Evolution du processus en fonction des données
 - → Apprentissage continu
 - → Montée en charge
- Recherche d'informations de Valeur pour les entreprises
 - → Utilisation de la connaissance métier
 - → Découverte d'informations dirigée par les données
- Mise en corrélation des informations
 - → Rapprochement des données similaires





Passage à l'échelle



Stocker et analyser de grands volumes de données







Hétérogénéité (Variété) des données



Contraintes

<?xml version="1.0" encoding="iso-8859-1" ?> - <rss version="2.0"> - <channel> <title>Hobbyman.se</title> <link>http://www.hobbyman.se/</link> <description>Tips och Trix för hemmet</description> <language>en-us</language> <generator>Nucleus CMS v3.23</generator> <copyright>@</copyright> <category>Weblog</category> <docs>http://backend.userland.com/rss</docs> - <image> <url>http://www.hobbyman.se//nucleus/nucleus2.gif</url>

<title>Hobbyman.se</title> dink>http://www.hobbyman.se/</link>

</image>

- <item>

<title>Så här gör du ditt eget t-shirt tryck med strykjärn...</tit <link>http://www.hobbyman.se/index.php?itemid=10</link>

- <description>

- <![CDATA[</p>

Fyller din kom inte vad du ska ge i present? Ge en personlig prese: T-shirt (vit eller färgad)

Items





A Trappes, le 28/02/2012

Communiqué de presse

La fédération des Yvelines du PCF, engagé dans le Front de Gauche et la campagne présidentielle avec Jean Luc Mélenchon, et ses candidats aux élections législatives, expriment leur entière solidarité à Maamar Kaoulal, salarié du magasin Carrefour de Chambourcy qui a tenté de mettre fin à ces jours récemment.





Intégration des connaissances métier



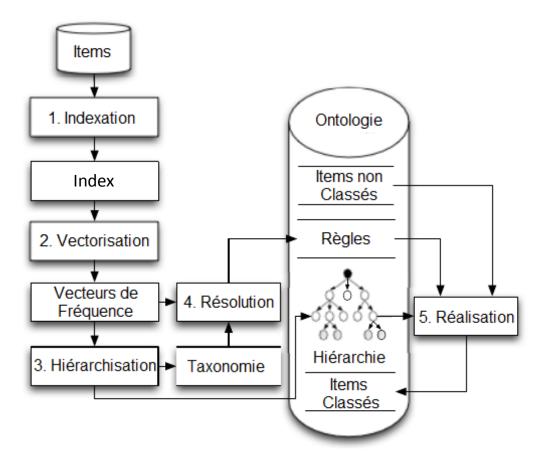
Thomas Hassan – thomas.hassan@gmail.com - Equipe de projet Checksem – Laboratoire Electronique Informatique et Image (LE21 – UMR CNRS 6306) IUT Dijon-Auxerre – Université de Bourgogne, BP 47870, 21078 Dijon Cedex, France

PROTOTYPE

Objectifs

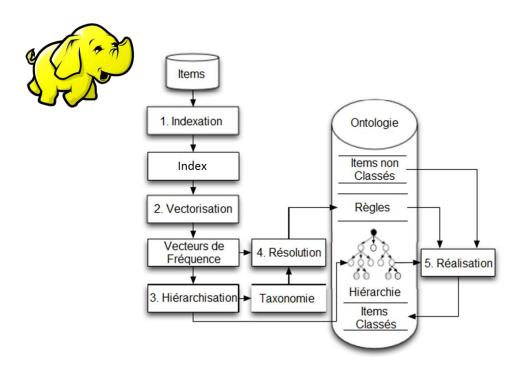
- Traitement des données brutes, réduction de la dimension
- Recherche d'informations pertinentes
- Création d'un modèle de classification évolutif (ontologie)
- Classification des nouveaux items

Architecture: processus à 5 phases distinctes



Montée en charge :

Utilisation du framework Hadoop et du modèle MapReduce

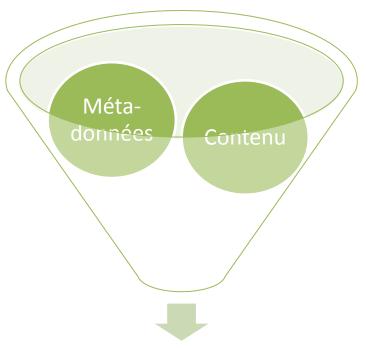


→ Comment distribuer chacune des phases du processus ?





Génération d'un index des termes



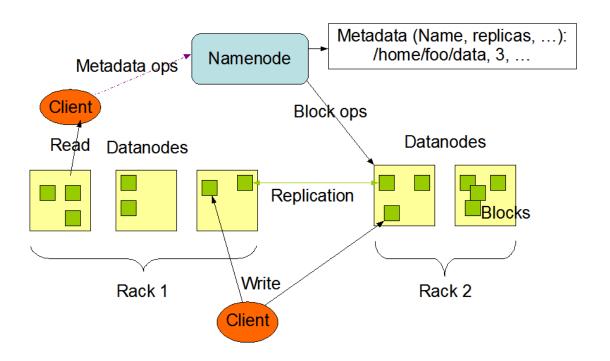
Parsing, Traitement du Langage Naturel*

- Réduction de la dimension
- Suppression des données inutiles

Indexation

Distribution de l'index sur un système de fichiers HDFS

HDFS Architecture



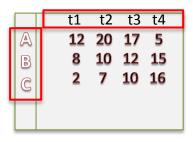




Détection des termes pertinents basée sur l'index.

Mesures statistiques telles que TF-IDF:

- Uni-grams (termes)
- N-grams (collocations)
 - → Définition des concepts basée sur les termes les plus fréquents

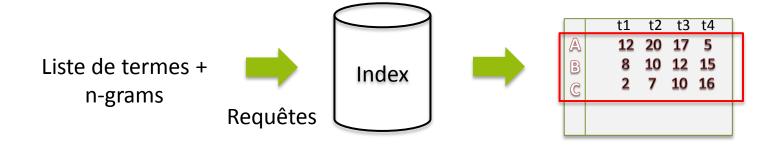






Utilisation d'implémentations existantes de l'algorithme de collocations (MapReduce)

Requêtes sur l'index pour générer les vecteurs de fréquence (1 vecteur par concept) :

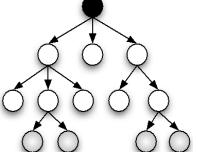


Hiérarchisation

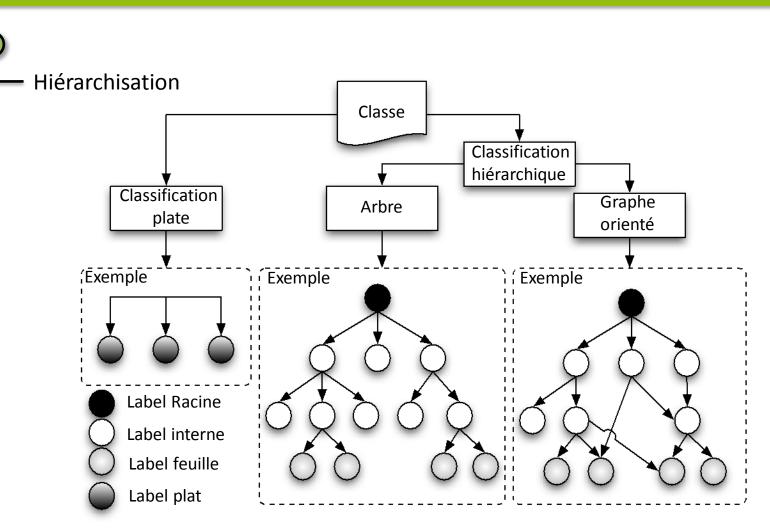
Utilisation de la matrice de fréquence pour générer une hiérarchie de concepts

t1	t2	t3	t4	
12	20	17	5	
8	10	12	15	
2	7	10	16	
				1
	12 8	12 20 8 10	12 20 17 8 10 12	12 20 17 5 8 10 12 15









Tsoumakas, G., Katakis, I. & Vlahavas, I., 2010. Mining multi-label data. In *Data mining and knowledge discovery handbook*. Springer, pp. 667–685.

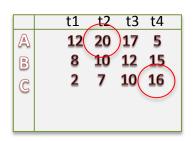
Cerri, R., Barros, R.C. & de Carvalho, A.C.P.L.F., 2014. Hierarchical multi-label classification using local neural networks. *Journal of Computer and System Sciences*, 80(1), pp.39–56. Available at: http://www.sciencedirect.com/science/article/pii/S0022000013000718.



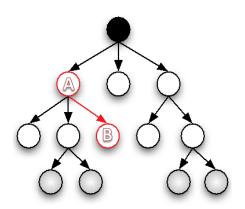


Algorithme de subsomption

Liens de parenté :







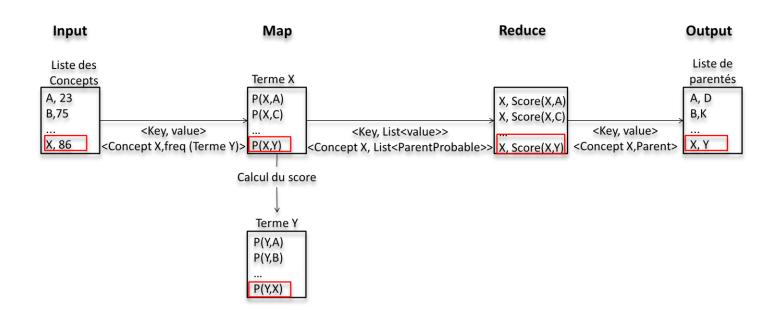
$$P(p|x) \ge t, P(x|p) < t$$

De Knijff, Jeroen and Frasincar, Flavius and Hogenboom, Frederik, 2013. Domain taxonomy learning from text: The subsumption method versus hierarchical clustering, Data & Knowledge Engineering, pp.54-69





Algorithme de subsomption MapReduce



De Knijff, Jeroen and Frasincar, Flavius and Hogenboom, Frederik, 2013. Domain taxonomy learning from text: The subsumption method versus hierarchical clustering, Data & Knowledge Engineering, pp.54-69





Génération de règles SWRL simples, basées sur la matrice de fréquence :

	t1	t2	t3	t4	
A	12	20	17	5	
B	8	10	12	15	
C	2	7	10	16	

Item(?d)^Word(?w1)^hasWord(?d,?w1) → Concept(?d1)

Enrichissement de l'ontologie

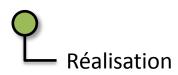


Génération des règles pour un concept est indépendante

→ Approche « diviser pour régner »

		t1	t2	t3	t4
	A	12	20	17	5
Г	B	8	10	12	15
_	C	2	7	10	16

Distribution des opérations par vecteur

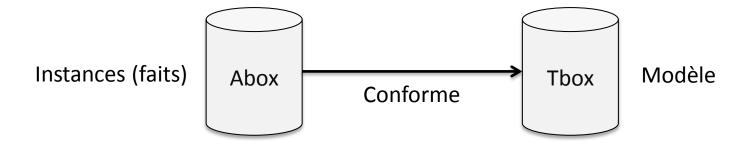


Utilisation des règles SWRL pour classer les nouveaux items

 $Item(?d)^Word(?w1)^hasWord(?d,?w1) \rightarrow Concept(?d1)$

→ Moteur d'inférence déduit les concepts les plus appropriés pour chaque item.

Assimilation des items en tant qu'individus de l'ontologie (Abox) :



Réalisation

Coût calculatoire important : inférence

→ Utilisation d'un moteur d'inférence basé sur les règles

Approches pour distribuer le raisonnement à l'échelle du Big Data se développent

Chevalier, J. (2013). A Linked Data reasoner in the Cloud. In *The Semantic Web: Semantics and Big Data* (pp. 722-726). Springer Berlin Heidelberg.

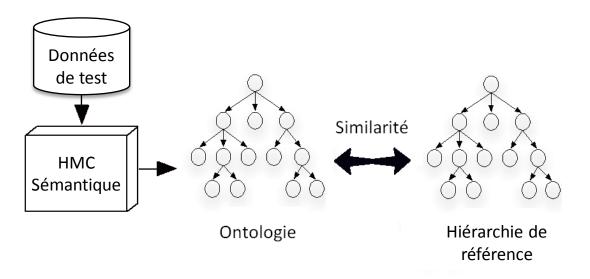
Urbani, J., Kotoulas, S., Maassen, J., Van Harmelen, F., & Bal, H. (2012). WebPIE: A Web-scale parallel inference engine using MapReduce. *Web Semantics: Science, Services and Agents on the World Wide Web*, 10, 59-75.



Evaluations

Qualité de la hiérarchie : utilisation de jeux de données standards

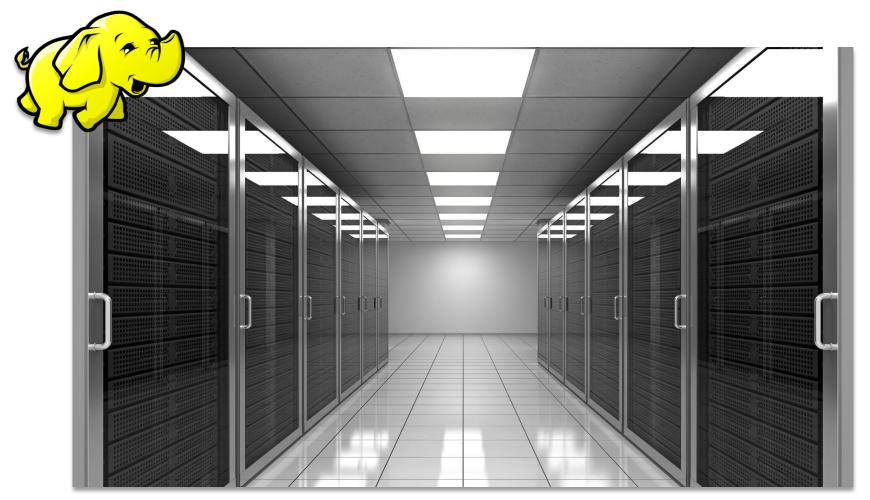
→ Comparaison entre la hiérarchie générée et la hiérarchie de référence (étalon)







Performance du processus



Améliorations

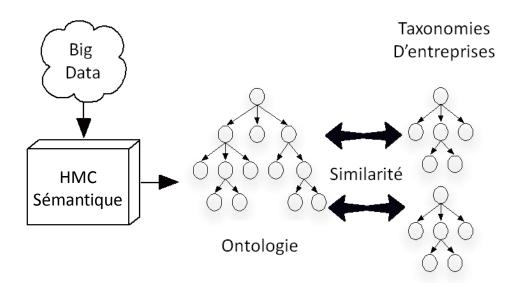
Extraction des concepts :

- Lien avec des bases de connaissance lexicales (type Wordnet)
- Lien avec des bases de connaissances générales
- Lien avec la taxonomie d'entreprise.
- Extraction d'événements complexes basés sur les bases de connaissances.

Améliorations

Rapprochement avec la connaissance métier :

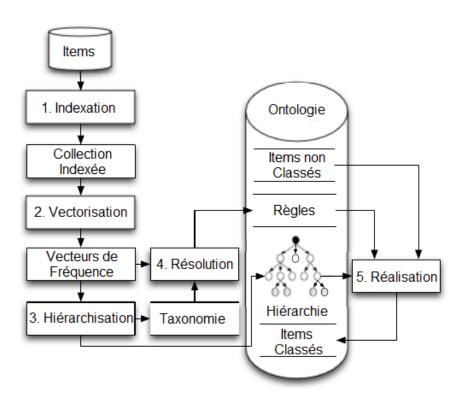
Comparaison entre les concepts de l'ontologie générée et la taxonomie d'entreprise.



Validation des concepts générés par un expert du domaine.



CONCLUSION



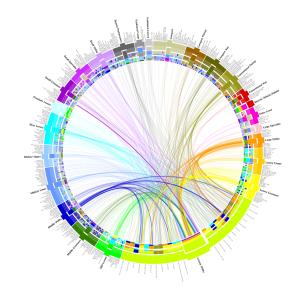
Processus d'analyse des données générique, adapté à des données et des connaissances variées

CONCLUSION

A terme:

Création d'outils d'analyse et d'aide à la décision à destination des documentalistes

Recoupement d'informations



Informations déjà parues







RESEARCH & BUSINESS APPROACH



