

Récupération de flux de corpus RSS et traitements linguistiques sur le texte

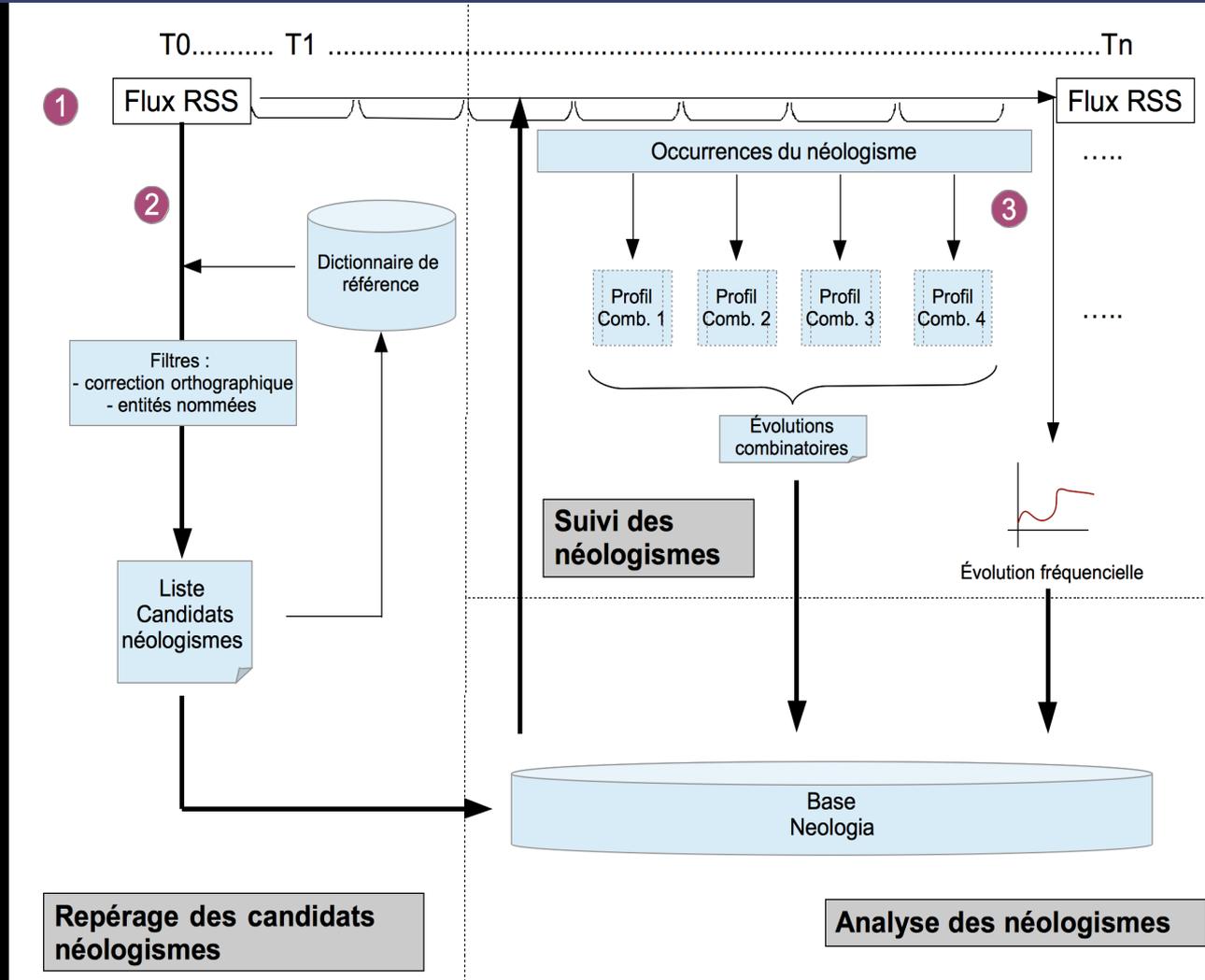
Emmanuel Cartier – LIPN - RCLN

Sommaire

- Contexte : projet de veille néologique sur grand corpus
- Objet de la mission
- Contraintes et calendrier

Contexte

Projet de veille néologique sur grand corpus : veille des mots nouveaux dans la langue et de l'évolution du sens des mots



Contexte

- Architecture séquentielle :
 - Récupération de corpus sur le web en flux continu (fils RSS)
 - Stockage sur un serveur en format txt (ou xml) UTF-8
 - Traitements linguistiques 1 : recherche des mots nouveaux (formes nouvelles) – par comparaison avec un lexique de référence
 - Traitements linguistiques 2 : analyse statistique des séquences de caractères, de mots, de séquences répétées **en vue d'établir des profils combinatoires des usages des lexies**
 - Comparaison des profils combinatoires en diachronie => différences signalent une évolution de sens
 -

Objet de la mission

Intervention principale sur le crawler de corpus RSS

Mise en place d'un crawler de fils RSS fonctionnant en continu et effectuant les traitements suivants :

- récupération des flux de quotidiens nationaux et régionaux français
- récupération des fichiers xml RSS et des méta-informations => stockage
- récupération des articles liés (url) : nettoyage, conversion et stockage dans une base de corpus.

Nettoyage fichiers (boiterplate cleaning)

TECHNO CAMPUS
Europe

INTERNATIONAL EUROPE Européennes 2014 Allemagne Belgique Espagne Grèce Italie Royaume-Uni

Le rouble subit une spectaculaire chute de 8 %

Le Monde.fr avec AFP | 15.12.2014 à 17h43 • Mis à jour le 15.12.2014 à 20h53

Réagir Classer Partager

Recommander Partager 224 personnes le recommandent.



Le rouble subissait lundi 15 décembre un quasi krach, s'effondrant d'environ 8 % face à l'euro et au dollar, à des niveaux jamais vus. Cette dégringolade intervient alors que la banque centrale a dressé un tableau noir de l'économie russe pour l'année prochaine.

Déjà en chute libre, la monnaie russe a encore décroché et l'euro a dépassé 78 roubles, contre 72,28 roubles la veille au soir. Le dollar a quant à lui atteint 63 roubles, soit 5 de plus que la veille. A ces niveaux, elle a perdu 42 % de sa valeur face à la monnaie européenne et 48 % face au dollar.

Dans un rapport trimestriel de politique monétaire, la Banque centrale a indiqué que dans un scénario où le baril de pétrole resterait à 60 dollars, son niveau actuel, le produit intérieur brut de la Russie pourrait chuter de 4,5 % à -0,9 %.

"La Duce vita", l'Italie d'hier et d'aujourd'hui

Cliquez sur un pays pour accéder à ses dernières actualités

Vidéo

Chute du rouble : la société russe sidérée

Le Monde.fr ÉDITION ABONNÉS

Trouvez le cadeau idéal

1 abonnement Le Monde.fr + 1 surprise au choix

DEEZER | univers|ciné
amazon | iTunes

Je choisis !

Titre,
Méta-informations,
Texte à l'exclusion
de tout le reste de la
page

Conversion, stockage

- Conversion : format texte, UTF-8
- Stockage – décision en cours entre:
 - Stockage des fichiers texte sur disque
 - Base de données (relationnelle ou XML)
 - Outil dédié à la gestion et l'interrogation de corpus

Traitements linguistiques 1

- Objectif : repérage des mots-formes inconnus
- Segmentation en “mots” du texte
- Comparaison avec dictionnaire de référence
 - => mots “inconnus”, avec tri des noms propres
- **Attention : tâche bien plus complexe qu'il n'y paraît**
 - Segmentation fautive
 - Erreurs typographiques
 - Noms propres...

Traitements linguistiques 2

- Objectif : repérage des nouveaux emplois par une étude quantitative des unités lexicales visant à établir sur corpus la combinatoire des mots
 - Segmentation en “mots” du texte, normalisation
 - Calcul du vocabulaire et de la fréquence de chaque élément
 - Calcul des n-grams répétés (fenêtre variable)
 - Calcul des n-grams sur informations morpho-syntaxiques
 - Calcul n-grams sur informations mixtes (formes et informations morpho-syntaxiques)

Contraintes et calendrier

- Langage de programmation : (de préférence) Python, Perl ou Java
- Livraison d'une API + batch (interface graphique facultative)
- accompagnement par un expert pour les problématiques TAL
- Possibilité d'utiliser outils préexistants opensource (analyseur morphosyntaxique notamment)
- Livraison obligatoire en mai-juin 2015!