

# Collecte et analyse d'offres d'emploi

Claire DE MARICOURT

Soutenance de stage de fin d'études

5 septembre 2018



# Plan

- 1 Présentation du projet
  - Problématique
  - Environnement
  - Indicateur de tension
  - Codes ROME
- 2 Collecte de données
  - De nombreux sites d'offres d'emploi
  - Scraping, parsing et automatisation
- 3 Classification
  - Méthodologie
  - Méthodologie - Nettoyage
  - Méthodologie - Appariement par mesure de similarité
- 4 Résultats
- 5 Conclusion

# Présentation du projet

## Problématique

Collecter automatiquement des offres d'emploi sur Internet et les classer par métier

# Présentation du projet

- Environnement : la Dares (Direction de l'animation de la recherche, des études et des statistiques)



- Département d'accueil : DMQ (Département Métiers et Qualifications)
- Objectif à terme : révision de l'indicateur de tension

# Présentation du projet

- Calcul de l'indicateur de tension

$$tension = \frac{offres}{demandes}$$

- Source de données : Pôle Emploi
- Objectif : étendre le champ des données
- Nécessité : coder les codes ROME de ces offres pour pouvoir les utiliser à des fins statistiques

# Nomenclature ROME

K			<b>SERVICES A LA PERSONNE ET A LA COLLECTIVITE</b>
K	21		<b>Formation initiale et continue</b>
K	21	08	<b>Enseignement supérieur</b>
K	21	08	Attaché / Attachée temporaire d'enseignement et de recherche -ATER-
K	21	08	Chargé / Chargée de cours
K	21	08	Chargé / Chargée d'enseignement
K	21	08	Chargé / Chargée d'enseignement du supérieur
K	21	08	Chef de département de l'enseignement supérieur
K	21	08	Enseignant-chercheur / Enseignante-chercheuse
K	21	08	Lecteur / Lectrice de langues dans l'enseignement supérieur
K	21	08	Maître / Maîtresse de conférences
K	21	08	Maître / Maîtresse de langues dans l'enseignement supérieur
K	21	08	Moniteur / Monitrice d'initiation à l'enseignement supérieur
K	21	08	Professeur / Professeure de l'enseignement supérieur
K	21	08	Professeur / Professeure des universités
K	21	08	Responsable filière enseignement supérieur
K	21	08	Responsable matière enseignement supérieur
K	21	08	Responsable programme enseignement supérieur

# Les grands domaines

Lettre	Grand domaine de Pôle Emploi
A	Agriculture et pêche, espaces naturels et espaces verts, soins aux animaux
B	Art et façonnage d'ouvrages d'art
C	Banque, assurances et immobilier
D	Commerce, vente et grande distribution
E	Communication, media et multimédia
F	Construction, bâtiment et travaux publics
G	Hôtellerie-restauration, tourisme, loisirs et animation
H	Industrie
I	Installation et maintenance
J	Santé
K	Services à la personne et à la collectivité
L	Spectacle
M	Support à l'entreprise
N	Transport et logistique

# Collecte de données

- De nombreux sites d'offres d'emploi





# Collecte de données

Conducteur de bus (H/F)

23/09/2018 09h57

### Description

**Résumé**  
CONDUCTEUR DE BUS H/F, 94190 Villeneuve Saint Georges, Intérim, Temps partiel.

**L'entreprise**  
Vous recherchez un emploi ? Faites confiance à nos différences ! RAS Intérim, réseau d'agences d'emploi de 110 agences, propose des centaines d'opportunités d'emploi dans tous les secteurs d'activité, en intérim, CDD et CDI.

**Description du poste**  
Vous recherchez un emploi ? Faites confiance à nos différences !  
R.A.S. Intérim, réseau d'agences d'emploi de 94 agences, propose des centaines d'opportunités d'emploi dans tous les secteurs d'activité, en intérim, CDD et CDI.

**Réactivité, Proximité, Qualité :** R.A.S. Intérim et Recrutement est la solution sur-mesure à votre recherche d'emploi.

Nous recherchons pour un de nos clients des conducteurs de bus H/F.

Vous aurez pour mission le transport de personnes, l'encaissement des titres de transport sur des substitutions SNCF.

Poste basé dans le 94 (Villeneuve S ... lire plus

[▲ Signaler un abus](#)

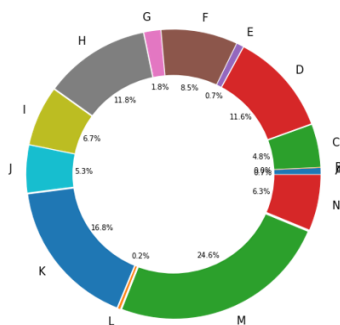
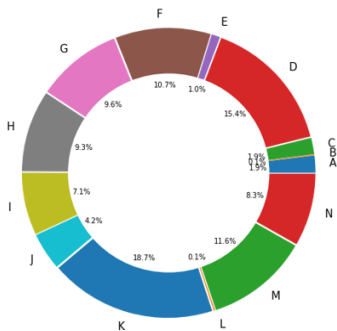
### Critères

<b>TYPE DE CONTRAT</b> Intérim	<b>EXPÉRIENCE</b> 2 à 5 ans	<b>TRAVAIL À</b> Temps partiel
-----------------------------------	--------------------------------	-----------------------------------

- Mise à jour du scraping et du parsing
- Ajout de nouveaux sites
- Automatisation de la collecte et gestion du stockage
- Collecte de données étiquetées sur Pôle Emploi

# Données Pôle Emploi

- Offres Pôle Emploi étiquetées manuellement
- Offres partenaires étiquetées automatiquement



Répartition des offres Pôle Emploi (G) et partenaires (D) selon les codes ROME

# Méthodologie

- Utilisation du titre

Titre  $\implies$  Titre ROME  $\implies$  Code ROME  $\implies$  Grand domaine

employé.e commercial au rayon bazar à lyon (69) cdi h/f  $\implies$   
Employé / Employée de rayon bazar  $\implies$  D1507  $\implies$  D

- Calcul de similarité après nettoyage  
sim(employé commercial rayon bazar; employée de rayon bazar) = 0.87  
sim(employé commercial rayon bazar; boulanger) = 0

# Méthodologie - Nettoyage

Étapes du nettoyage	Exemple
Tout mettre en minuscules	employé.e commercial au rayon bazar à lyon (69) cdi h/f
Enlever la ponctuation	employée commercial au rayon bazar à lyon 69 cdi hf
Enlever les chiffres	employée commercial au rayon bazar à lyon cdi hf
Enlever les <i>stopwords</i>	employée commercial rayon bazar lyon cdi hf
Enlever certains mots	employée commercial rayon bazar lyon
Enlever le nom de la ville	employée commercial rayon bazar
Lematiser	employé commercial rayon bazar

# Méthodologie - Appariement par mesure de similarité

Trois mesures de similarité + une

- Jaro-Winkler avec comparaison mot à mot
- Damerau-Levenshtein
- Tf-Idf et similarité cosinus
- Produit de la similarité de Jaro-Winkler mot à mot et de la similarité cosinus

# Méthodologie - Appariement par mesure de similarité

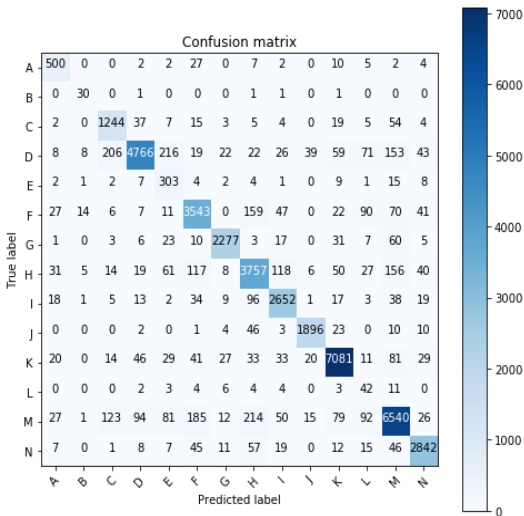
$$tok\_jw\_sim(T_1, T_2) = 0.5 \times partial\_tokens\_jw\_sim(T_1, T_2) + 0.5 \times partial\_tokens\_jw\_sim(T_2, T_1)$$

$$partial\_tokens\_jw\_sim(T_1, T_2) = \frac{\sum_{t^1 \in T_1} (\exists t^2 \in T_2, sim_{jaro-winkler}(t^1, t^2) > 0.9)}{Card(T_1)}$$

# Tableau des résultats

	Accuracy	Précision	Rappel	Score F1
<b>Jaro-Winkler mot à mot</b>	89.3%	77.4%	86.8%	80.1%
Pôle Emploi	92.5%	81.2%	88.5%	84.1%
Partenaires	86.0%	72.1%	84.4%	74.8%
<b>Damerau-Levenshtein</b>	88.6%	70.6%	81.7%	73.7%
Pôle Emploi	90.4%	78.2%	85.0%	80.7%
Partenaires	78.2%	62.3%	78.2%	65.6%
<b>Tf-idf et similarité cos</b>	88.6%	81.3%	84.9%	82.8%
Pôle Emploi	93.0%	85.5%	88.1%	86.6%
Partenaires	84.1%	75.5%	81.0%	77.6%
<b>Produit JW/Cosinus</b>	91.1%	84.8%	87.6%	85.6%
Pôle Emploi	93.8%	86.3%	90.0%	87.6%
Partenaires	88.2%	82.7%	84.4%	82.7%

# Matrice de confusion





# Perspectives et Limites

- Possibilité de définir un seuil de validation de l'étiquetage
- Nouveaux sites à scraper
- Déduplication des offres
- Codage basé sur le corps de l'offre
- Réutilisation de ce codage pour d'autres projets
- Utilisation d'autres informations pour d'autres projets (localisation, compétences, ...)