

Le Machine Learning pour coder dans une nomenclature : l'enquête Associations

Florian Lécivain (Insee)

Séminaire Big Data de l'Insee (14 janvier 2020)





• Contexte et objectif



• Les données utilisées



• Apprentissage supervisé



• Limites et pistes d'amélioration

01 Contexte et objectif

- Objectifs : dénombrer le nombre d'associations actives en France, connaître leur(s) domaine(s) d'activité, évaluer la participation associative, connaître leurs ressources...
- Champ : ensemble des associations relevant de la loi de 1901
- Première édition réalisée en 2014 (sur l'année de référence 2013).
- La seconde édition a été lancée fin 2019 (sur l'année de référence 2018).

- Base de sondage constituée à partir :
 - du Répertoire National des Associations (RNA)
 - du répertoire Sirène
- Échantillon stratifié selon :
 - **Le domaine d'activité**
 - L'âge de l'association
 - Le nombre de salariés (pour les associations employeuses)

- 3 nomenclatures différentes pour l'activité entre le répertoire Sirène, le RNA et l'enquête.
- On dispose avec le RNA de l'objet social de l'association, texte libre qui décrit ses activités :

Notre association a pour objet l'enseignement spécialisé de la musique, danse et théâtre.

Réaliser produire diffuser et promouvoir des projets culturels et des créations artistiques.

Organiser les loisirs de la collectivité dans son ensemble. Favoriser la pratique de l'éducation physique et sportive. Renforcer la solidarité morale des habitants.

Développement connaissance pratique du sport en général et en particulier des sports de combat ainsi que le maintien des traditions propres à chaque discipline.

Nous sommes une association de commerçants elle sert à animer le centre ville par des festivals commerciaux qui fait gagner à nos clients des lots.

- **Objectif** : exploiter l'objet social afin de coder le domaine d'activité de l'association selon la nomenclature propre à l'enquête.
- Expérimentation réalisée en collaboration avec l'unité SSP Lab de l'Insee de janvier à mai 2019.
- Tirage de l'échantillon en mai, en mettant à profit le domaine d'activité codé avec ces travaux.

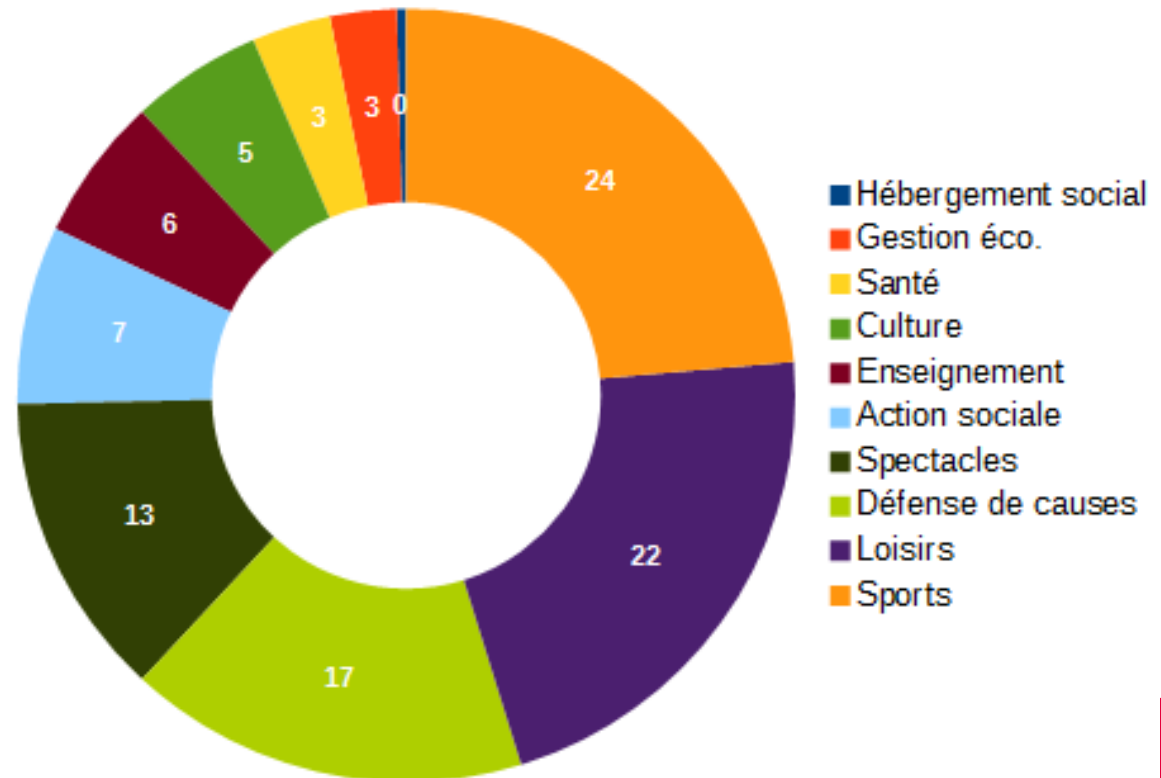
02 Les données utilisées

- Les réponses de la première enquête Associations fournissent, après appariement avec le RNA, un jeu de données labellisées de **17 200 observations**.
 - usage possible de techniques d'apprentissage supervisé

Figure 1 – Répartition des associations selon leur principal domaine d'activité

(en 2013, source : Enquête Associations 2014)

- 64 modalités regroupées en **10 grands domaines**



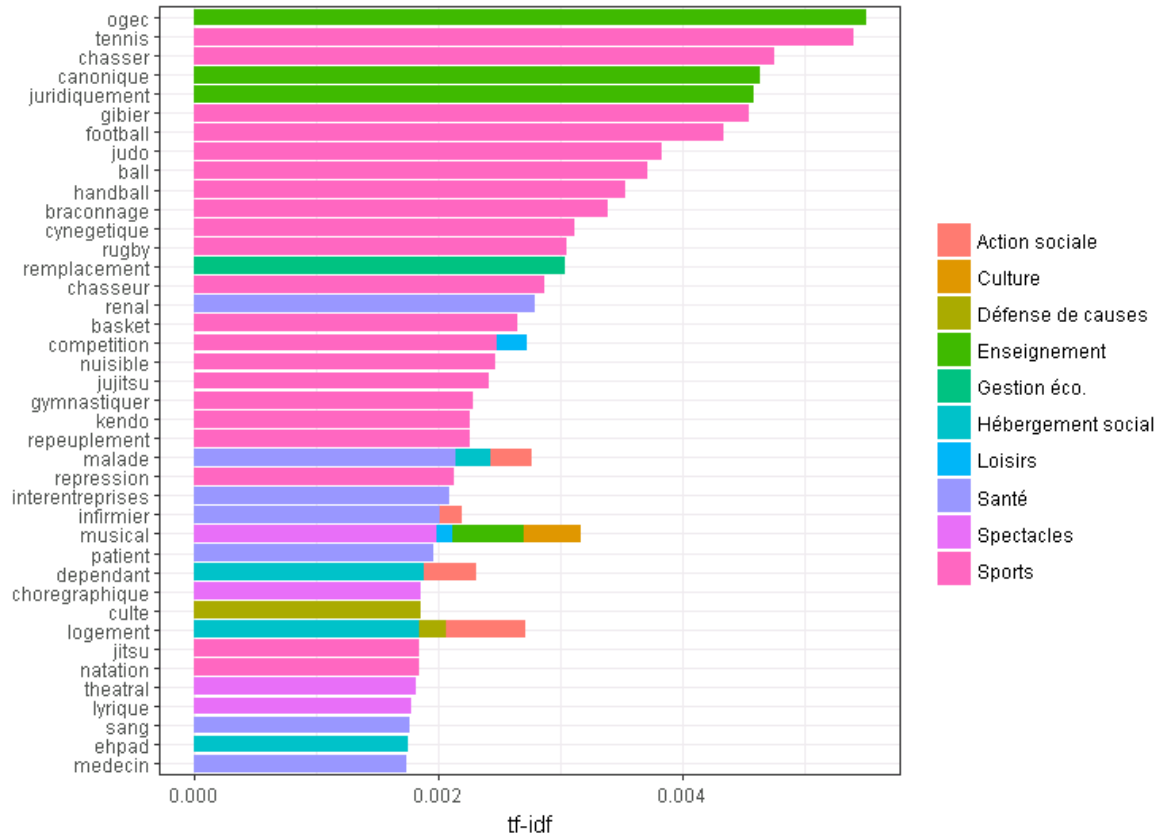
- En plus de l'**objet social**, 3 autres variables corrélées au domaine d'activité sont utilisées :
 - Le statut employeur ou non de l'association
 - La taille de la commune d'implantation de l'association
 - L'âge de l'association

- Nettoyage des données textuelles
 - Suppression de la casse, de la ponctuation, des caractères spéciaux, etc.
 - Suppression des termes qui n'apportent aucune information (*Stopwords*)
 - Substitution des chiffres par le terme « number »
- Lemmatisation (dictionnaire Morphalou) → réduction de la dimension
- Améliorations apportées au fil de l'eau :
 - Prise en compte des fautes d'orthographe (ou de frappe) les plus courantes
 - Ajouts et corrections de certains lemmes

- Après mise en forme des N objets sociaux, on compte l'occurrence de chaque terme (vocabulaire de taille V) dans chaque descriptif → on obtient une **matrice document-terme de taille $N \times V$**
- Cette matrice est complétée de la même manière avec les **bigrammes**.
- Au final on obtient 16 200 lemmes et 55 800 bigrammes différents, soit 72 000 variables (de type indicatrices) → sélection de variables
- Les travaux sont menés en utilisant une **représentation simplifiée de l'objet social dite *bag of words***, où seule l'occurrence d'un terme est prise en compte.

Aperçu des termes les plus discriminants

Figure 2 - Termes les plus discriminants (selon le score tf-idf)



03

Apprentissage supervisé

- **3 familles de modèles testées** : Forêts aléatoires, *Extreme Gradient Boosting* (XGBoost) et machines de vecteurs à support (SVM)
- Optimisation des hyper-paramètres de chaque algorithme
- 5 types de sélection de variables testées (entre 28 et 6000 variables)
- Répartition des données entre jeu d'apprentissage et jeu test (80/20)

- Indicateurs de performance retenus :
 - Au niveau de chaque domaine
 - Précision
 - Rappel
 - **Score F1** : moyenne harmonique entre la précision et le rappel
 - Au niveau global
 - Précision globale et rappel global : moyenne pondérées des indicateurs au niveau domaine
 - **Score F1 global** : moyenne harmonique entre la précision globale et le rappel global
- Meilleures performances obtenues avec XGBoost (score F1 global de 68.6 %), mais les Forêts aléatoires (67 %) et les SVM (67.1 %) obtiennent des performances proches

Tableau 1 – Performances du modèle XGBoost

| | <i>F1-score</i> | <i>Précision</i> | <i>Rappel</i> |
|--------------------|-----------------|------------------|---------------|
| Ensemble | 68.6 | 68.7 | 68.5 |
| Action sociale | 67.8 | 63.9 | 72.1 |
| Culture | 46.6 | 61.1 | 37.7 |
| Défense de causes | 58.9 | 64.7 | 54 |
| Enseignement | 70.5 | 68.2 | 72.9 |
| Gestion éco. | 67.2 | 72.7 | 62.5 |
| Hébergement social | 60.7 | 61.7 | 59.8 |
| Loisirs | 58.5 | 54.5 | 63.1 |
| Santé | 64.1 | 71 | 58.5 |
| Spectacles | 79.6 | 78.4 | 80.7 |
| Sports | 82.3 | 82.7 | 81.9 |

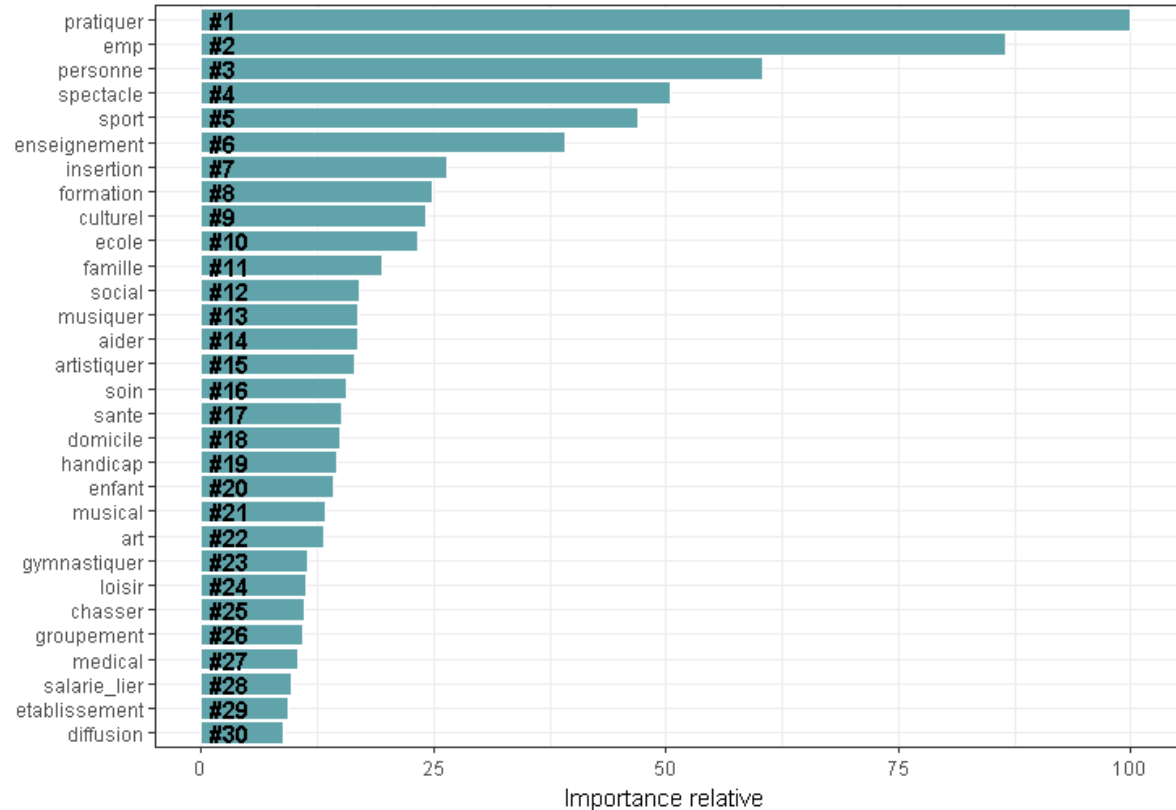
- Sélection des lemmes et des bigrammes présents dans au moins 5 descriptifs → **score F1 global de 68.6 %.**

Figure 3 – Matrice de confusion associée à la précision

- Mauvaises prédictions : 28 % des unités classées par le modèle dans l'hébergement social appartiennent en réalité à l'action sociale.

| Domaine prédit \ Vrai domaine | Sports | Spectacles | Santé | Loisirs | Hébergement social | Gestion éco. | Enseignement | Défense de causes | Culture | Action sociale |
|-------------------------------|--------|------------|-------|---------|--------------------|--------------|--------------|-------------------|---------|----------------|
| Sports | 0.83 | 0.02 | 0.01 | 0.07 | 0 | 0.01 | 0.02 | 0.01 | 0.02 | 0.02 |
| Spectacles | 0.04 | 0.78 | 0 | 0.05 | 0 | 0 | 0.07 | 0.01 | 0.04 | 0.01 |
| Santé | 0.01 | 0 | 0.71 | 0 | 0.05 | 0.02 | 0.02 | 0.08 | 0.01 | 0.11 |
| Loisirs | 0.11 | 0.05 | 0.01 | 0.55 | 0.01 | 0.02 | 0.03 | 0.1 | 0.05 | 0.08 |
| Hébergement social | 0.01 | 0 | 0.04 | 0.03 | 0.62 | 0 | 0.01 | 0.01 | 0.01 | 0.28 |
| Gestion éco. | 0.04 | 0 | 0.01 | 0.02 | 0 | 0.73 | 0.03 | 0.13 | 0.02 | 0.04 |
| Enseignement | 0.04 | 0.03 | 0.02 | 0.04 | 0.02 | 0.04 | 0.68 | 0.06 | 0.03 | 0.05 |
| Défense de causes | 0.02 | 0.02 | 0.02 | 0.1 | 0.01 | 0.04 | 0.08 | 0.65 | 0.02 | 0.06 |
| Culture | 0.04 | 0.03 | 0 | 0.11 | 0 | 0.01 | 0.06 | 0.11 | 0.61 | 0.03 |
| Action sociale | 0.01 | 0.01 | 0.04 | 0.05 | 0.08 | 0.04 | 0.04 | 0.05 | 0.03 | 0.64 |

Figure 4 – Variables les plus influentes



- Variables les plus importantes :
« pratiquer », statut employeur ou non,
« personne »,
« spectacle » et « sport »

04 Limites et pistes d'amélioration

- Des performances réduites pour certains domaines, en particulier la culture
- Une lemmatisation qui peut encore être améliorée, avec par exemple la prise en compte de la fonction grammaticale qui aurait permis de traiter les homonymies
- Utilisation possible de représentations plus complexes, telles que le *word embedding*
- Utilisation possible d'autres modèles d'apprentissage supervisé (réseaux de neurones par exemple), voir de combinaisons de différents algorithmes

Merci de votre attention !

florian.lecrivain@insee.fr

Retrouvez-nous sur :

[insee.fr](https://www.insee.fr)

