

Foille de données TP 2

M2 Informatique, Université Paris Diderot

Anne-Claire Harter

2016/2017

L'objectif de ce second TP est de vous essayer à classer la classification par exemple. À partir d'un jeu de données que vous allez transformer, vous allez pouvoir essayer différents modèles et comprendre comment choisir le meilleur.

Voire aussi dans cette affaire est [la page d'aide de l'API sklearn](#). Toutes les fonctions utilisées ici sont expliquées. Vérifiez la dernière version de sklearn est installée sur votre ordinateur :

a
a

Récapitulatif de la transformation des données

Récapitulatif de la donnée

Pour ce TP nous utiliserons un jeu de données composées de messages SMS que vous pouvez

charger sur le site de cours :

<https://sites.google.com/site/da-miningp7/s-propos-de-cours/SMS Spam Collection>.

Il s'agit de classer en anti-spam à partir de ces données.

Chaque ligne du fichier représente un message et son label (spam/ham). Nous utiliserons ces données pour entraîner un algorithme de prédiction qui fournira pour un message une réponse entre 0 (ham) et 1 (spam).

À partir de là, nous allons utiliser ce jeu de données de manière non supervisée (sans regarder les labels) et appliquer des algorithmes de classification.

Etape 1 - Transformation de données

1. écrire une fonction `read_dataset(filename)` qui prend en argument le chemin vers le fichier et retourne une liste de paires (pe/e). Le `pe` est un entier 0 pour les hams (non-spams) et 1 pour les spams.

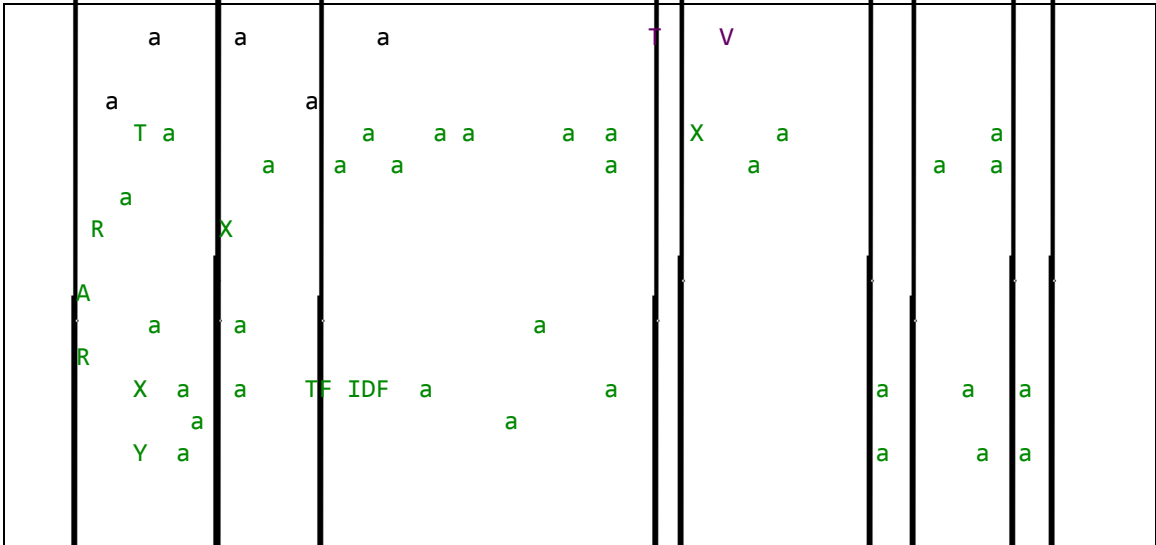
```
a a a a a a a
R a a a a a
a a a b a ab a a
A a a a a
R a a a a a
a a a a a
```

2. écrire une fonction `spams_count(comments)` qui prend en argument la liste renvoyée par la fonction précédente et retourne le nombre de spams dans le jeu de données.

```
a a b a a
R a a a
A a a a
R a a b a
a a a a
```

3. Importer la classe `sklearn.feature_extraction.text.TfidfVectorizer`. Celle-ci nous permet d'appliquer l'algorithme TF-IDF en codeurs. Le constructeur accepte un certain nombre de paramètres optionnels que nous pourrions donner pour adapter l'algorithme à nos besoins et retourner un objet capable d'appliquer ce algorithme. Vous trouverez [ici](#) la documentation de cette classe. Observer aussi les arguments et choisir les valeurs pertinentes. Il n'y a pas une seule bonne réponse, mais il faut être capable de justifier nos choix.

4. Créez une fonction `transform_text(pairs)` qui retourne une matrice `X` (les lettres sont à la forme TF-IDF) et un vecteur `y` (0 si le message est un ham, 1 si il s'agit d'un spam).



Apprentissage non supervisé (clustering)

Etape 2 - Un simple KMean

Trouvez dans la bibliothèque `sklearn` la fonction qui vous sera utile pour appliquer l'algorithme de KMeans. Lancez sur l'ensemble `X` un KMeans avec `K = 10`. Vous aurez un exemple de script sur la page d'aide de l'algorithme. Inspirez-vous en !

Évaluez ce algorithme en regardant la valeur de la fonction `sklearn.metrics.silhouette_score`. Expliquez ce que représente cette valeur. Bon si vous êtes en avance : créez une fonction qui retourne les 10 mots les plus utilisés dans chaque cluster.

Etape 3 - Choisir le meilleur K

Cette fois, vous ne connaissez pas la valeur de `K`. Comment choisir la meilleure valeur de `K` ? Implémentez cette solution.

E xercice 4 - Comparaison de algorithmes de clustering

Vous n'avez presque rien fait pour lancer ce 4^{ème} algorithme de classification hiérarchique sur ces données. L'algorithme s'appelle `AgglomerativeClustering` dans la bibliothèque `sklearn`.

Apprentissage supervisé (classification)

Ce 4^{ème} fois, on cherche à prédire le type d'un message (spam ou ham). On va donc utiliser y qui n'a été laissé de côté dans la partie précédente.

E xercice 5 - Classification par K-Nearest Neighbors

1. Utilisez la fonction `sklearn.model_selection.train_test_split` pour obtenir des ensembles d'entraînement et de test.
2. Lancez un algorithme des plus proches voisins avec `k = 1` sur l'ensemble d'entraînement. Évaluez les résultats sur l'ensemble de test.
3. Faites de même avec `k = 3`. Lequel donne les meilleurs résultats ?
4. Vous n'avez rien fait de chose d'intermédiaire. Pourquoi les étapes 2 et 3 sont-elles interdites ?
5. Utilisez `sklearn.model_selection.StratifiedKFold` pour choisir le meilleur `k` en réalisant 100 (ou 50 en les nombres impairs) splits sur l'ensemble d'entraînement. Une fois que vous avez trouvé le meilleur `k`, relancez l'algorithme avec ce `k` sur l'ensemble d'entraînement et l'ensemble de test.

E xercice 6 - Classification par K-Nearest Neighbors

Reprenez la dernière question de l'exercice 5 pour choisir les meilleurs paramètres d'un arbre de décision, puis d'une forêt aléatoire. Quel algorithme donne le meilleur résultat ?