

## DESCRIPTION DES FICHIERS

- Le fichier **Preprocessing.ipynb** contient l'ensemble de fonctions pour traiter les données brutes des fichiers CSV fournies par les enseignants.
- Le fichier **load\_data.ipynb** est une version plus compacte de preprocessing avec les fonctions nécessaires pour faire la lecture des données nettoyées et en obtenir des features pour l'apprentissage. JE CONSEILLE VIVEMENT D'UTILISER CE FICHIER PLUTÔT QUE PREPROCESSING
- Data.csv: contient les données du fichier transtion\_ecologique traitées et filtrées. Les colonnes inutiles ont été enlevées. Les questions ouvertes ont été converties en dictionnaires avec comme clé le mot et comme valeur le nombre de répétitions. Les réponses oui/non ont été transformées en 1 pour 'oui', -1 pour 'non' et 0 si la question n'a pas été répondue. Les QCM ont été traitées comme les yes/no sauf que le tag s'incrémente à partir de 1.
  - Ordre des colonnes:
    - Colonne 0: Titre
    - Colonne 1: Catégorie de l'utilisateur
    - Colonnes 2-17: Questions
      - yes /no questions: colonnes 4,6,10
      - Questions choix multiples: 12
    - Colonnes 18: Code Postal
- Word\_count\_by\_question.csv: ce fichier contient un dictionnaire par question avec tous les mots qui apparaissent dans ce dernier.
- Word\_count\_total.csv: il contient la combinaison des tous les dictionnaires par question. C'est un dictionnaire général des données
- Les fichiers city\_information.tsv et correspondance-code-insee-code-postal.csv contiennent des informations utiles sur les communes/villes (population, aire, code zip)

## DESCRIPTION DES FONCTIONS

- **read\_data**('data.csv', numerical\_answers\_array): charge le fichier data.csv en indiquant les questions qui ont un nombre comme réponse.
- **read\_word\_count**('word\_count\_by\_question.csv'): charge soit le fichier word\_count\_by\_question.csv ou soit word\_count\_total.csv et en génère un dictionnaire ou un array de dictionnaires
- **read\_dictionary**(dictionary\_string): lire un string qui représente un dictionnaire et en génère un.
- **Sort\_dictionary**: trie un dictionnaire et renvoie une liste avec comme premier élément les mots avec le plus d'occurrences.
- **Normalize\_counts**: renvoie un dictionnaire des mots ou la valeur et sa fréquence dans le texte.

- `find_zip_codes_by_town(density_threshold,filename='city_information.tsv')`: renvoie deux ensembles de zips codes des villes dont la densité est supérieur et inférieur au `density_threshold`
- **`city_village_classifier`**(`density_threshold,data`): fait la classification des entrees de la base de données en fonction de sa `density_threshold`. La convention est 1 pour les grand villes et -1 pour le villages.
- `word_count_by_question(data,yes_no_questions)`: count tous les mots des questions textuels par question et les met dans un array.
- `word_count_total(dictionary_array)::` combine tous les elements de plusieurs dictionnaires stockes dans un array.
- **`get_most_used_words`**(`data, yes_no_array, word_count_array, number_of_words`): elle prend en arguments un array de entrees, un array qui indique les reponses aux yes/no questions et aux QCM et le nombre de mots a prendre par question. Elle renvoie l'array d'origine dont les dictionnaires des questions ont été transformés en une liste de taille `number_of_words` avec le nombre d'ocurrences des mots qui coincident avec `word_count_array`
- **`get_set_features`**(`data,column`): renvoie une colonne de l'array de data.