

## ***Job Title Classification***

L'algorithme développé permet la classification d'intitulé de poste afin de déterminer s'il s'agit d'un Data Scientist ou pas. Il s'agit là d'obtenir un résultat binaire.

Dans un premier temps les données ont été chargées via pandas library afin de les manipuler comme dataframe. Une fois les deux fichiers JSON chargés ("**list\_job\_tech\_dataScience\_NO**" et "**list\_job\_tech\_dataScience\_YES**") on concatène nos deux dataframes.

Après la visualisation des données chargées, j'ai remarqué que les données sont déséquilibrées avec une majorité pour les classes de Data Scientist. Donc j'ai opté pour l'oversampling aléatoire afin d'équilibrer mes données avant d'entraîner les différents modèles de classification. Une fois cela effectué j'ai mis en place un pipeline. Tout d'abord il s'agit d'appliquer le CountVectoriser où les occurrences de chaque mot sont comptabilisées, cela permet de représenter les données textuelles sous forme de vecteurs numériques. Ensuite, on applique le TF-IDF, il s'agit là de calculer la fréquence de mots en divisant le nombre d'occurrence de chaque mot par le nombre total de mots.

Enfin, j'ai testé deux modèles de classification qui ont donné le même résultat en commençant par le RandomForestClassifier puis la Régression Linéaire.