

**Datamining : TP6**  
M2 Informatique, Université Paris 7  
Anne-Claire Haury

**Exercice 1**

1. Reprendre les données de spams texto. Appliquer leur la transformation TF-IDF.
2. Ecrire une fonction qui prend en entrée les données TF-IDF, les labels et un entier  $k$  et renvoie les  $k$  mots avec les valeurs les plus significativement différentes entre les spams et les non-spams.
3. Réduire la matrice de données à cet ensemble pour  $k = 100$ .
4. En utilisant la classe `pca` de la library `sklearn.decomposition.PCA`, appliquer à ces données une transformation en composantes principales.
5. Représenter graphiquement les données sur les deux premières composantes principales.

**Exercice 2 : travail en groupe sur vos projets**