
TP N° 5 : Clustering

Pour ce travail vous devez déposer un **unique** fichier au format `nom_prenom.ipynb` sur le site moodle du cours (si moodle ne l'accepte pas, zipper cet unique fichier en fichier `.zip`)

Vous devez charger votre fichier sur Moodle, avant le vendredi 23/10/2020, 23h59.

La note totale est sur **20** points répartis comme suit :

- qualité des réponses aux questions : **15** pts,
- qualité de rédaction, de présentation et d'orthographe : **2** pts,
- indentation, style PEP8, commentaires adaptés, etc. : **2** pts,
- absence de bug : **1** pt.

Les personnes qui n'auront pas rendu leur devoir avant la limite obtiendront **zéro**.

Rappel : aucun travail par mail ne sera accepté !

- INTRODUCTION -

Le but de cette séance est la mise en œuvre d'algorithmes de *clustering*. Le *clustering* est un problème non-supervisé. Étant donné un ensemble d'apprentissage $\mathbf{x}_1, \dots, \mathbf{x}_n$ de points dans \mathbb{R}^d , le but est de regrouper les points par groupes (si possible homogènes), ou *clusters*. Notez qu'à la différence d'un problème supervisé il n'y a pas d'étiquettes associées dans ce contexte.

- MISE EN ŒUVRE -

Les algorithmes de *clustering* de `scikit-learn` sont disponibles dans le module `sklearn.cluster`. Les modèles de mélanges Gaussiens (*Gaussian Mixture Models* ou GMM) peuvent être vus comme des estimateurs de densité et sont disponibles dans le module `sklearn.mixture`, mais ils peuvent aussi servir pour faire du *clustering*.

K-Means

Dans cette section on s'intéressera aux algorithmes où le nombre de *clusters* (noté K dans la suite) est fixé au préalable par l'utilisateur.

L'algorithme de K -Means partitionne les points en K groupes disjoints $\{\mathcal{C}_1, \dots, \mathcal{C}_K\}$ en minimisant la variance intra-cluster. Le critère minimisé G est appelé l'*inertie* :

$$G_X(\mathcal{C}_1, \dots, \mathcal{C}_K) = \sum_{k=1}^K \sum_{i \in \mathcal{C}_k} \|\mathbf{x}_i - \mu_k\|_2^2$$

où les $\mu_k \in \mathbb{R}^d$ sont les centroïdes des K classes ($|\mathcal{C}_k|$ le cardinal de chaque classe) :

$$\mu_k = \frac{1}{|\mathcal{C}_k|} \sum_{i \in \mathcal{C}_k} \mathbf{x}_i, \quad \forall k \in [K],$$

et où la i° observation a pour classe celle du centroïde le plus proche, *i.e.*, $i \in \mathcal{C}_j$ si et seulement si :

$$j = \arg \min_{k \in [K]} \|\mathbf{x}_i - \mu_k\|_2 = \arg \min_{k \in [K]} \|\mathbf{x}_i - \mu_k\|_2^2.$$

Rem : on départage les ex-æquo aléatoirement si besoin.

L'apprentissage des *clusters* se fait en alternant itérativement les deux étapes suivantes :

- I) une étape d'assignation où, connaissant les $(\mu_k)_{k=1,\dots,K}$, on détermine les labels de chaque point.
 - II) une étape de mise à jour des centroïdes où, connaissant les labels, on détermine les centres de classe.
- On arrête l'algorithme quand l'inertie ne décroît plus beaucoup (il est à noter que sans critère d'arrêt l'algorithme K -Means termine de toute manière en un nombre fini d'étape cf. [BMDG05]).

L'inertie est un critère non-convexe, la solution trouvée dépend donc de l'initialisation, comme souvent pour de tels problèmes. En conséquence, l'algorithme est souvent lancé plusieurs fois avec des initialisations différentes, et l'on ne garde alors que la solution dont l'inertie est la plus faible.

Pour plus d'information sur l'algorithme K -means et ses généralisations on pourra consulter [BMDG05].

- 1) En vous basant sur le squelette de code dans le fichier `kmeans.py`, implémentez l'algorithme K -Means. Le fichier termine avec un exemple de code qui vous servira de validation. Il utilise la fonction `make_blobs` qui permet de simuler un dataset.
- 2) L'inertie décroît-elle bien au cours des itérations? Vous afficherez la décroissance de l'inertie en fonction des itérations.
- 3) En faisant varier l'initialisation du générateur de nombres aléatoires, observez que la solution trouvée n'est pas toujours la même. Vous pourrez augmenter le nombre de clusters simulés pour mettre d'avantage en évidence le phénomène.
- 4) Comparez votre implémentation avec celle de `scikit-learn`, en terme de résultat et temps de calcul. Vous pourrez utiliser la commande `%timeit` d'IPython pour évaluer le temps de calcul d'une cellule du notebook. Voici un exemple d'utilisation de `scikit-learn` :

```
from sklearn import cluster
kmeans = cluster.KMeans(n_clusters=3, n_init=10)
kmeans.fit(X)
labels = kmeans.labels_
```

Calcul optimal du nombre de cluster

Dans la section précédente, nous avons supposé que le nombre de clusters K était fixé par l'utilisateur. Néanmoins, il peut être intéressant en pratique de déterminer le nombre optimal de clusters. Nous proposons de définir une méthode pour le calculer.

On s'intéresse toujours à un jeu de données X avec n points. Soit $X_t, t \in [T]$, T tirages indépendants de n points aléatoires tirés selon une distribution uniforme dans la boîte englobante de X . À partir des inerties G_X et G_{X_t} des K -means appliqués à X et à X_t , nous définissons la différence :

$$\delta(t) = \text{Gap}(t) - (\text{Gap}(t+1) - \sigma(t+1))$$

où $\text{Gap}(t) \in \mathbb{R}$ correspond à la différence entre l'espérance du log des inerties de G_{X_t} et le log de l'inertie G_X :

$$\text{Gap}(t) = \mathbb{E}[\log(G_{X_t})] - \log(G_X)$$

et où $\sigma(t)$ est défini par :

$$\sigma(t) = \sqrt{\frac{T+1}{T}} \text{std}(\log(X_t))$$

avec `std` qui représente l'écart-type.

En vous basant sur le squelette de code fourni dans le fichier `gap.py` :

- 5) Complétez la fonction `compute_log_inertia_rand` de sorte qu'elle estime la moyenne et l'écart type des $\log(G_{X_t})$.
- 6) Complétez la fonction `compute_gap` qui calcule, pour un jeu de donnée X fixé, les valeurs de $\text{Gap}(t)$ et de $\delta(t)$ pour différentes valeurs de t .
- 7) Analyser les valeurs de $\text{Gap}(t)$ et de $\delta(t)$ pour différentes valeurs de k et proposez une méthode pour permettre de définir automatiquement le nombre K des K -means. Vous utiliserez la fonction `sklearn.datasets.make_blobs` qui permet de générer un ensemble de K clusters et la fonction `plot_gap` pour visualiser correctement vos différentes données.
- 8) Complétez la fonction `optimal_n_clusters_search` avec la méthode que vous avez proposée et vérifiez son fonctionnement sur différents jeux de données.

Application à la compression d'images

Les algorithmes de *clustering* peuvent aussi servir à la quantification vectorielle.

Une image est représentée par un tenseur de dimension 3. L'image a $n_1 \times n_2$ pixels, et pour chaque pixel (x_1, x_2) la valeur associée est un triple $(I_R, I_G, I_B) \in \mathbb{R}^3$ qui indique l'intensité sur les trois couleurs avec l'encodage suivant :

- I_R qui code pour l'intensité du rouge (🇬🇧 : *red*)
- I_G qui code pour l'intensité du vert (🇬🇧 : *red*)
- I_B qui code pour l'intensité du bleu (🇬🇧 : *blue*)

L'intensité de couleur pour chaque pixel est codé par un octet : un entier entre 0 et 255 ($= 2^8 - 1$).

Rem: il arrive qu'on utilise aussi le RGBA, qui utilise un terme supplémentaire, I_A , qui code le niveau de transparence (élément souvent nommé **alpha** et qui influe comme ici¹). Le α est alors lui codé par un **float** sur $[0, 1]$.

L'idée de la quantification vectorielle est de coder chaque pixel sur seulement K couleurs, avec K bien plus faible que les 256^3 valeurs possible¹. L'image peut ainsi être codée de façon plus compacte.

Nous proposons de comparer les résultats obtenus avec différentes techniques de compression :

- 9) Appliquer l'algorithme K -Means pour réduire le nombre de niveaux de couleurs de l'image à $K=8$ valeurs possibles. Visualiser le résultat sous forme d'une image. À partir de quelle valeur de K ne voit-on plus de différence avec l'image d'origine ?
- 10) Appliquer l'algorithme de clustering spectral² pour réduire le nombre de niveaux de couleurs de l'image à 8 valeurs possibles. Attention, pour des raisons de performances et d'utilisation mémoire, il faut absolument utiliser la version du clustering spectral utilisant les matrices creuses. Visualisez et comparez les résultats obtenus avec K -Means.

Indication : Pour lire l'image 'china.jpg' fournie, on pourra utiliser la commande :

```
import matplotlib.pyplot as plt
img = plt.imread('china.jpg')
```

Il faudra aussi effectuer un **reshape** de l'image afin de travailler avec un **array** à 2 dimensions de taille " $n = n_1 \times n_2 \times 3$ " (en effet, il y a trois couleurs pour un codage RGB, et on négligera le niveau de **alpha**).

Références

- [BMDG05] A. Banerjee, S. Merugu, I. S. Dhillon, and J. Ghosh. Clustering with Bregman divergences. *J. Mach. Learn. Res.*, 6 :1705–1749, 2005. 2

1. on ne modifiera pas le niveau de alpha ici

2. <https://scikit-learn.org/stable/modules/generated/sklearn.cluster.SpectralClustering.html>