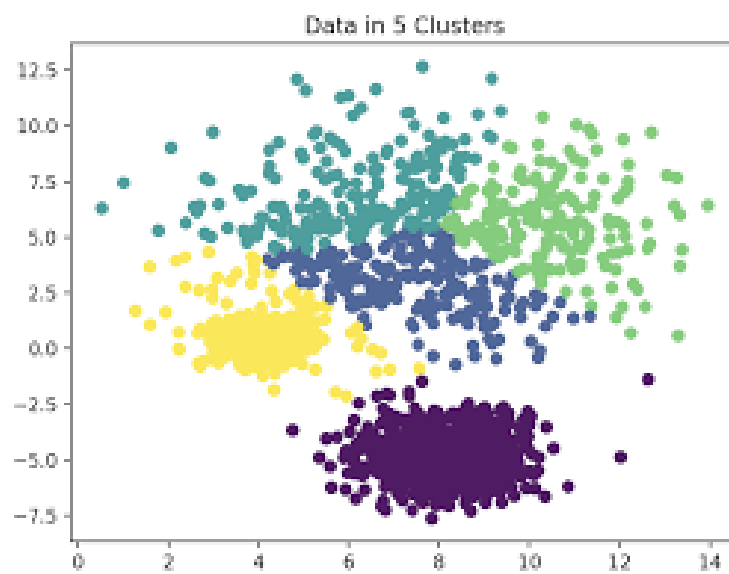


machine learning : K-means

Ben Amira Rawia

23 novembre 2022



Université de Toulon, UFR Sciences et Techniques

M1 DID

2022-2023

Table des matières

1	Introduction	3
2	EXERCICE1	3
2.1	calcul distance	3
2.2	initialisation des centroides	3
2.3	assign cluster	3
2.4	compute centroids	3
2.5	display	3
2.6	Resultats et analyse	3
3	EXERCICE2	6
4	EXERCICE3	6
5	Conclusion	11

1 Introduction

Ce rendu est élaboré dans le cadre du module I143 pour le TP K-means.

2 EXERCICE1

2.1 calcul distance

Pour calculer la distance entre un point et un centroïde, plusieurs méthodes existent. Dans mon code j'ai implémenté deux méthodes classiques la distance euclidienne ainsi que la distance de manhattan. Je n'ai pas remarqué de différence dans les résultats obtenues avec les deux méthodes même si que dans la littérature apparemment la distance euclidienne est la mieux appropriée pour des données sous forme de coordonnées spatiales.

2.2 initialisation des centroïdes

Ma première démarche pour calculer les centroïdes était de générer des nombres aléatoires entre le maximum et le minimum des données, mais vite j'ai changé de méthode car une exception est levée suite à une division par zéro si un cluster est vide. La deuxième démarche était de générer un entier x aléatoirement entre 0 et la longueur de la vx , et de considérer que $vx[x]$ est un centroïde tiré aléatoirement.

2.3 assign cluster

Cette fonction va renvoyer une matrice contenant les points de chaque cluster. Par exemple pour $k=3$ la matrice contiendra dans `matrice[0]` les points les plus proches du premier centroïde etc.

2.4 compute centroids

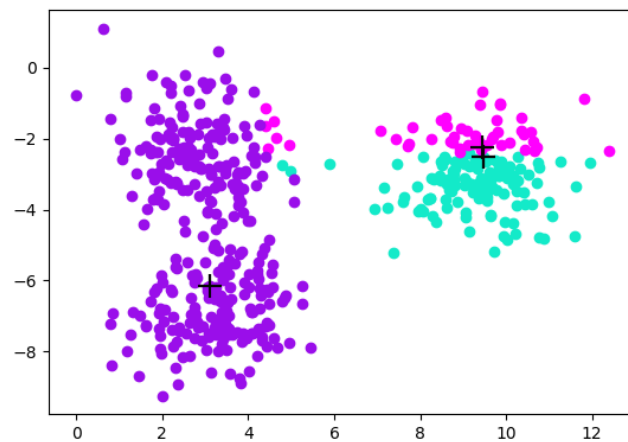
Cette fonction va contenir comme paramètres k le nombre de clusters, vx les données, vc les centroïdes tirés aléatoirement et enfin un nouveau paramètre a été ajouté qui est le nombre d'itérations. En effet ce dernier permet d'arrêter le calcul des nouveaux centroïdes.

2.5 display

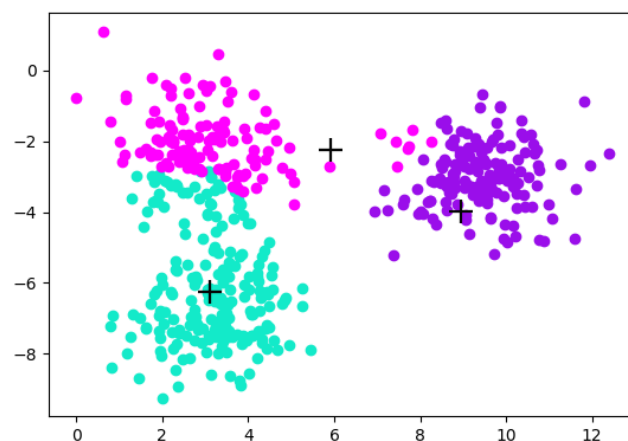
Cette fonction fait tous les affichages nécessaires. La fonction `display 2 data` a été supprimée et l'affichage se fait une fois pour toute dans cette fonction.

2.6 Resultats et analyse

Les premiers centroïdes générés sont :

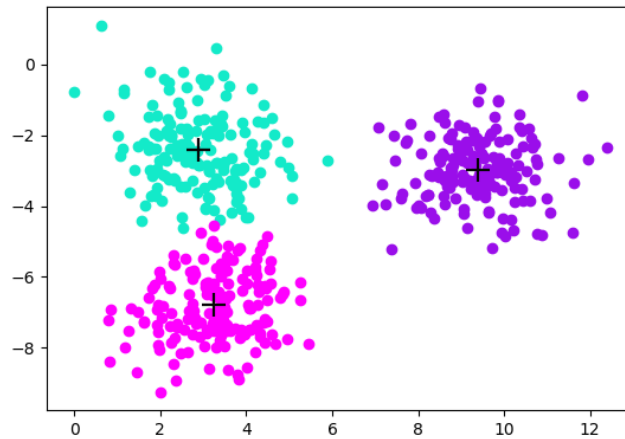


les croix representent les centroides tirés aleatoirement au debut de l'algorithme ainsi que leurs clusters correspondant. Ensuite je decide de recalculer les centroides en deux autres iterations et j'obtiens ces resultats :



les donnees sont visiblement mieux reparties que dans le cas precedent. en effet j'ai fait le choix de programmer avec un nombre d'iterations qu'on peut fixer pour voir à partir de combien d'iterations on converge. Dans une ancienne version j'avais la version qui recalculer les centroides jusqu'à egalité des centroides entre deux iterations et puisque on est sur des flottants ceci prenait beaucoup de temps.

En testant avec nombre d'iterations egale à 5 j'obtient ce resultat final qui est bien parlant et qui represente le resultat attendu :



En reflechissant encore plus sur la convergence de l'algorithme j'ai ajouté une fonction qui qui arrete le recalcul des clusters jusqu'a obtention de difference entre les centroides egales a une valeur prise en parametres et pour pousser mes idées j'ai aussi ajouté une variable type compteur qui s'incrémente à chaque nouveau calcul pour savoir apres combien de tour on obtient bien une convergence.

En lançant le programme à plusieurs reprises, j'obtient ces nombres d'iterations :

```

ra@ra-VirtualBox: ~/Desktop/baseApprenti
l.centroide = new_centroid(clusters, k)#nouveaux centroid
File ~/home/ra/Desktop/baseApprenti/exercice1.py, line 120, in new_centroid
    centr.append(s_x/len(data[i]))
ZeroDivisionError: division by zero
ra@ra-VirtualBox:~/Desktop/baseApprenti$ python3 exercice1.py
le nombre d'iterations est egale à 1
ra@ra-VirtualBox:~/Desktop/baseApprenti$ python3 exercice1.py
le nombre d'iterations est egale à 1
ra@ra-VirtualBox:~/Desktop/baseApprenti$ python3 exercice1.py
le nombre d'iterations est egale à 1
ra@ra-VirtualBox:~/Desktop/baseApprenti$ python3 exercice1.py
le nombre d'iterations est egale à 5
ra@ra-VirtualBox:~/Desktop/baseApprenti$ python3 exercice1.py
le nombre d'iterations est egale à 1
ra@ra-VirtualBox:~/Desktop/baseApprenti$ python3 exercice1.py
le nombre d'iterations est egale à 3
ra@ra-VirtualBox:~/Desktop/baseApprenti$ python3 exercice1.py
le nombre d'iterations est egale à 2
ra@ra-VirtualBox:~/Desktop/baseApprenti$ python3 exercice1.py
le nombre d'iterations est egale à 1
ra@ra-VirtualBox:~/Desktop/baseApprenti$ python3 exercice1.py
le nombre d'iterations est egale à 1
ra@ra-VirtualBox:~/Desktop/baseApprenti$ python3 exercice1.py
le nombre d'iterations est egale à 2
ra@ra-VirtualBox:~/Desktop/baseApprenti$ python3 exercice1.py
le nombre d'iterations est egale à 0
ra@ra-VirtualBox:~/Desktop/baseApprenti$ python3 exercice1.py
le nombre d'iterations est egale à 1
ra@ra-VirtualBox:~/Desktop/baseApprenti$

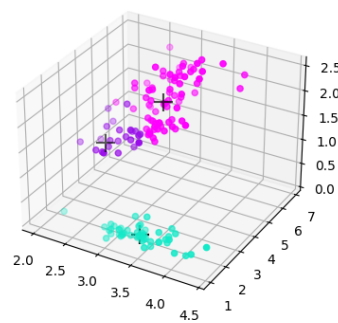
```

ensuite dans une dernière version j'ai ajouté le fait que cette fonction peut à la fois arrêter le calcul jusqu'à convergence et en même temps elle nous affiche le nombre d'itérations pour arriver à convergence.

3 EXERCICE2

Dans l'exercice 2 on reprend les mêmes concepts juste on ajoute une coordonnée z à x et y. Certaines fonctions implémentées dans l'exercice 1 ont été appliquées et d'autres ont été adaptées pour traiter des données à 3 features.

`kmeans3d(k, conv)` reprend le même concept que `kmeans(k, conv)` juste les données sont à 3 features et l'affichage est ainsi :



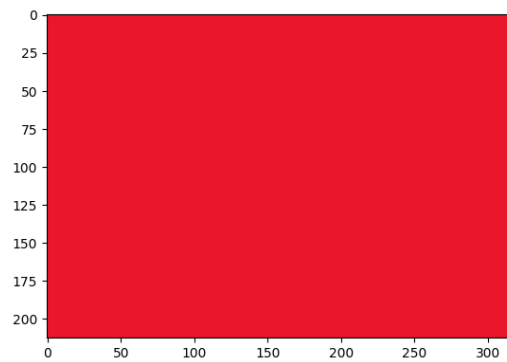
4 EXERCICE3

Dans cet exercice on modifie des images png et ceci en agissant sur les codes RGB. Des problèmes de division par zéro se sont posés lors du traitement par exemple j'ai tenté de changer l'image du drapeau de la Tunisie (drapeau rouge et blanc) avec un clustering à 4 et j'obtiens erreur de division par zéro ce qui est totalement logique en fait tous les points seront forcément attribués au cluster 1 (blanc) ou cluster 2 (rouge) et les deux autres seront à vide.

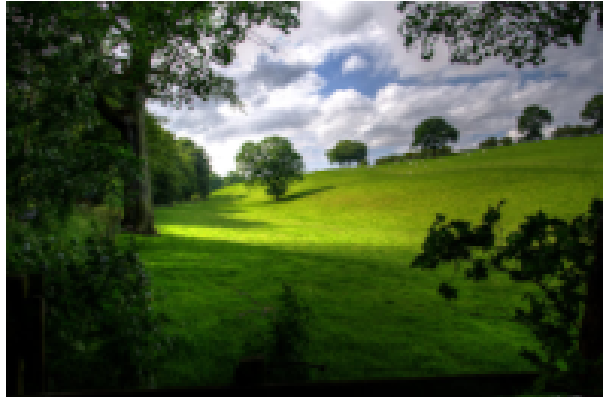
Drapeau de la tunisie :



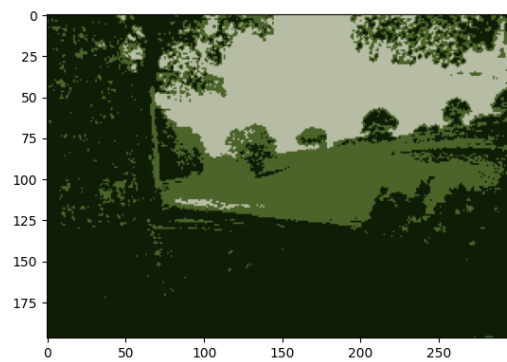
Pour $k=1$ on obtient cet affichage tout logique :

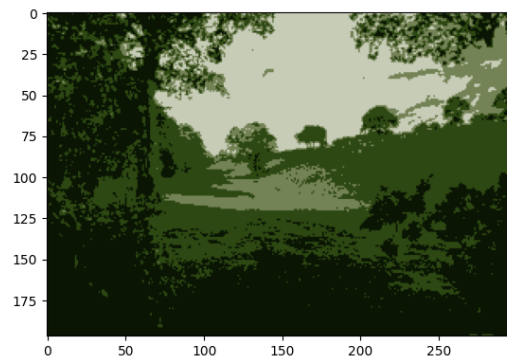
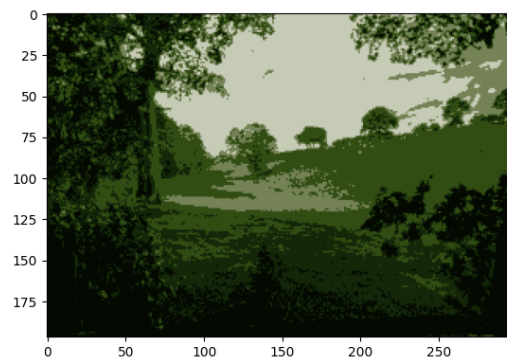


les images fournies avec le tp :
L'original :

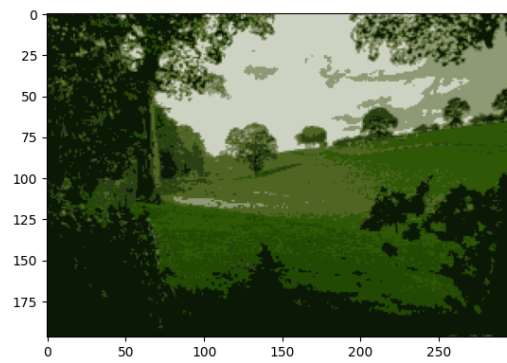


K=3

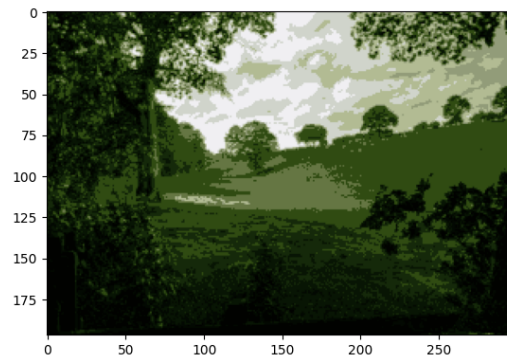


$K=4$  $K=5$ 

K=8



K=10



autre exemple



k=5



k=2



5 Conclusion

k-means est un algorithme de clustering très répandu, son implémentation sur python ne pose pas vraiment de soucis et ceci grâce à numpy. Avoir les affichages dans le 3ème exercice aide beaucoup à la compréhension de l'algorithme en détail.