

UNIVERSITÉ PARIS DAUPHINE

DATA ANALYTICS

---

# Implementation de KMeans avec Spark

---

*Etudiants*

Elie ABI HANNA DAHER  
Bilal EL CHAMI  
Badr ERRAJI

*Professeur*

M. Benjamin  
NEGREVERGNE

29 avril 2018



## Table des matières

<b>1</b>	<b>Travail effectué</b>	<b>2</b>
1.1	Implementation . . . . .	2
1.1.1	Kmeans . . . . .	2
1.1.2	Kmeans++ . . . . .	3
1.1.3	Generateur . . . . .	4
1.1.4	Plot . . . . .	7
1.2	Simulation . . . . .	7
1.2.1	Resultat plot - Sepal . . . . .	9
1.2.2	Resultat plot - Petal . . . . .	13
<b>2</b>	<b>Analyse</b>	<b>15</b>
2.1	Distance intra-cluster . . . . .	15
2.2	Convergence . . . . .	15
2.3	Tableau de performance . . . . .	15
<b>3</b>	<b>Application</b>	<b>15</b>

La méthode K-means est une des méthodes de clustering les plus utilisées lors de l'implémentation d'algorithmes cherchant à regrouper un ensemble de données disparates. Cette méthode repose principalement sur le *unsupervised learning*. L'objectif étant alors de grouper ces dernières en implémentant l'algorithme KMean avec Spark (python) et en évaluant la performance de notre implémentation basée sur des données que nous générons.

## 1 Travail effectué

### 1.1 Implementation

#### 1.1.1 Kmeans

Toutes les fonctions demandées étaient implémentées en respectant le format donné. Ce format nous a énormément aidé pour comprendre et développer l'algorithme en suivant le paradigme MapReduce en Spark.

Les méthodes de l'algorithme sont les suivantes :

- *customSplit(...)*
- *loadData(...)*
- *initCentroids(...)*
- *assignToCluster(...)*
- *calculateDistance(...)*
- *minDist(...)*
- *computeCentroids(...)*
- *reCalculating(...)*
- *hasConverged(...)*
- *computeIntraClusterDistance(...)*

Les méthodes *customSplit* et *loadData* préparent les données passées par l'utilisateur et les formatent en affectant un identifiant pour chaque enregistrement.

*initCentroid* est la fonction critique de l'algorithme qui sert à initialiser les centroïdes. Dans un premier temps, nous avons choisis ces centres aléatoirement parmi les données existantes en utilisant la fonction *takeSample(...)* de la classe RDD. Nous avons implémenté aussi l'algorithme K-means++ qui propose une procédure d'initialisation plus pratique (voire la partie suivante).

Ensuite, la méthode *assignToCluster()* prends les données et calcule la distance relative de chaque centroïdes. Puis elle choisit le centroïde ayant la distance minimale, en utilisant les méthodes *calculateDistance()* et *minDist()*, et l'affecte comme cluster de l'enregistrement. Le format utilisé nous a permis de profiter de l'utilisation de la méthode *groupByKey* de la classe RDD ce qui améliore les performances en passant à l'échelle.

Dans la fonction *computeCentroids()*, nous avons regroupé les résultats d'affectation des clusters, en utilisant autant que possible les fonctions de la classe

RDD ( *join()* et *groupByKey()*). Après avoir grouper les données par cluster, nous avons recalculer les nouvelles valeurs des centroïdes à l'aide de la méthode *reCalculating()*.

La méthode *computeIntraClusterDistance()* permet de calculer la distance intra-cluster. Cette méthode calcule la somme des moyennes des distances des points aux clusters dont qu'il appartient.

Finalement, le programme recalcule à chaque fois les nouvelles valeurs des centroïdes et vérifie si l'algorithme a convergé. La condition de convergence est renvoyé par la méthode *hasConverged()*.

#### Observation importante

Une amélioration considérable des performances a été remarqué lors de la modification de l'instruction suivante qui casse le lineage du RDD centroids et crée un nouveau RDD pour l'itération suivante : *centroids = sc.parallelize(newCentroids.collect())*. Alors que l'initialisation de centroïdes pouvait se faire de manière aléatoire, ce qui restait fonctionnel.

### 1.1.2 Kmeans++

L'initialisation de centroïdes pouvait se faire de manière aléatoire, ce qui restait fonctionnel.

En effet, le fait de générer des centroïdes de manière aléatoire pouvait parfois mener à des situations bloquantes tel que le choix de deux centroids avec exactement les même coordonnées. Par ailleurs, le clustering peut donner des résultats différents selon les centroids initiaux choisis qui peuvent parfois être très loin de l'optimum rallongeant les calculs donc avant convergence.

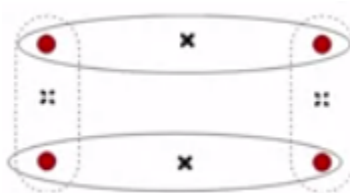


FIGURE 1 – Initialize centroids

Comme nous pouvons le voir sur l'exemple ci-dessus, le clustering peut se faire dans ce cas verticalement ou horizontalement selon points initiaux.

Nous avons décidé de suivre une autre technique certes quelque peu plus complexe mais qui réduisait le nombre d'itérations considérablement. Nous nous sommes basés sur l'algorithme K-means++ pour implémenter la méthode init centroïdes :

Kmeans ++ a été proposé en 2007 par David Arthur et Sergei Vassilvitskii. Le premier centre du cluster est choisi de manière aléatoire uniforme parmi les points de données. Ensuite, chaque centre de cluster suivant est choisi parmi les

points de données restants avec une probabilité proportionnelle à sa distance au carré du centre de cluster existant le plus proche du point.

Malheureusement, nous n'avons pas pu implémenter ce dernier en MapReduce à cause de quelques limitations de calcul qui utilisent la librairie numpy. Donc une diminution de performance peut affecter la partie d'initialisation des centroïdes.

### 1.1.3 Générateur

Afin d'évaluer l'algorithme KMeans implémenté, un générateur de données est mis en place.

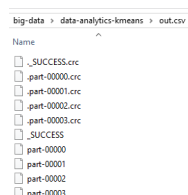
On a mis en place 2 fichiers de python : *generator.py* et *generator-noise.py*. Le fichier *generator.py* permet de générer des points et de les sauvegarder dans un fichier.

Le deuxième fichier, *generator-noise.py*, permet de générer les points ainsi que des points bruits qui sont générés aléatoirement.

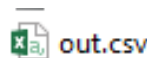
Afin de sauvegarder les données dans un fichier externe, on a implémenté 2 façons différentes :

- *saveAsTextFile()* : qui est une méthode déjà présente qui sauvegarde un rdd
- *write-into-csv()* qui est une méthode que nous avons créé qui permet d'écrire directement dans un fichier les valeurs du RDD

La méthode que nous avons créé *write-into-csv*, permet de générer un seul fichier mais n'est pas efficiente et ne profite pas des caractéristiques du RDD. C'est pour cela que la méthode *saveAsTextFile* est meilleur en terme d'efficacité. Cette dernière crée plusieurs fichiers.



(a) Méthode *saveAsTextFile*



(b) Méthode *write-into-csv*

FIGURE 2 – Fichier output de chaque méthode du générateur de données

En exécutant, le *generator.py* avec la commande suivantes :

```
$ spark-submit generator.py out 9 3 2 10
```

Les points sont générés, et en affichant les résultats dans le graphique on obtient :

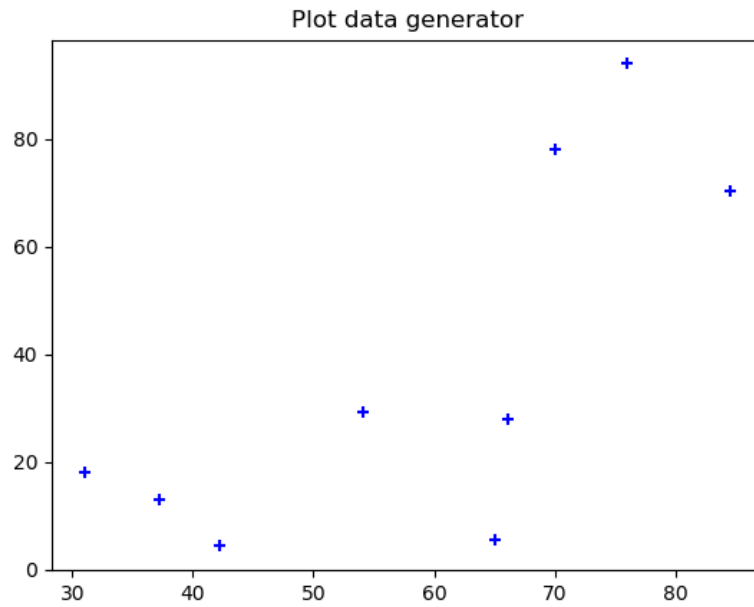


FIGURE 3 – Générateur des points

En executant, le *generator-noise.py* avec la commande suivante :

```
$ spark-submit generator_noise.py out 9 3 2 10
```

Le nombre de bruits est le double du nombre des points à générer. Les points sont générés, et en affichant les résultats dans le graphique on obtient :

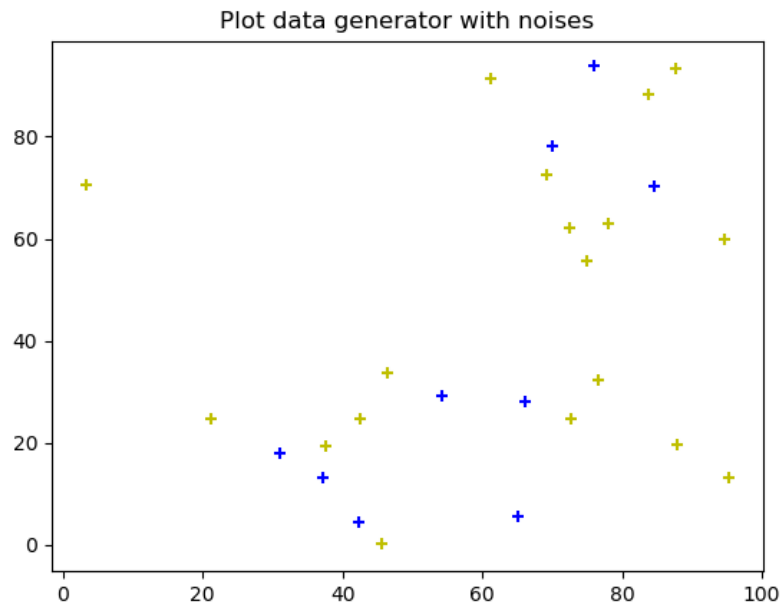


FIGURE 4 – Générateur des points avec des bruits

### 1.1.4 Plot

Afin de générer les graphiques, nous avons utilisé la librairie *matplotlib.pyplot*. Pour les fichiers KmeansPlusPlus, Kmeans nous avons créé 2 fichiers pour chacun, 1 sans plot et 1 avec plot.

Pour chaque itération, une image est sauvegardée dans le répertoire afin d'analyser l'évolution du résultat.

Les graphes générés sont de 2 dimensions, ce qui indique que si les points ont plus de 2 coordonnées ça ne sera pas très bien représenté graphiquement. Afin de voir des exemples d'exécution, veuillez voir la section Simulation.

## 1.2 Simulation

En exécutant, le *kmeans-plot.py* avec la commande suivantes :

```
$ spark-submit kmeans-plot.py data/iris-small.dat 3 10
```

Le graphe suivant représente les résultats de la première itération et en affichant les 2 premiers coordonnées : sepal width et sepal height

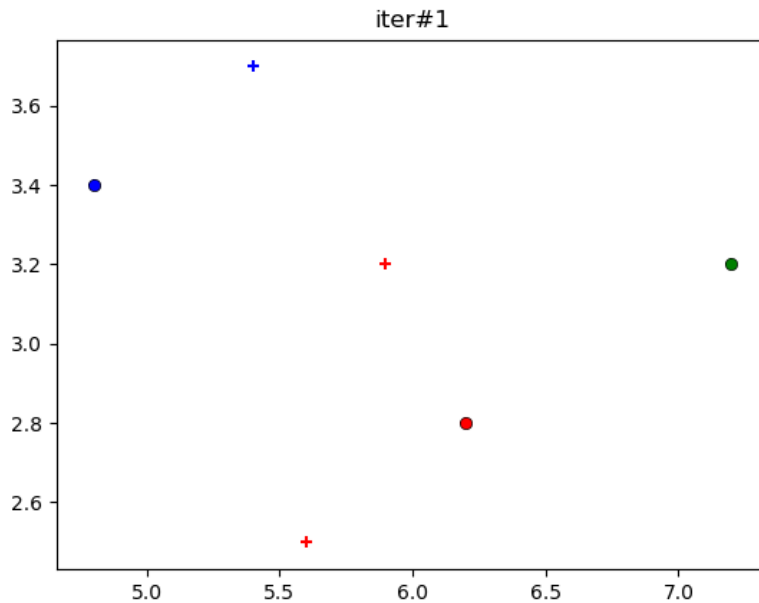


FIGURE 5 – Iteration 1

Pour trouver le résultat final on a eu besoin de 2 itérations avec une distance finale de 0.9711 d'où le résultat final est le suivant :



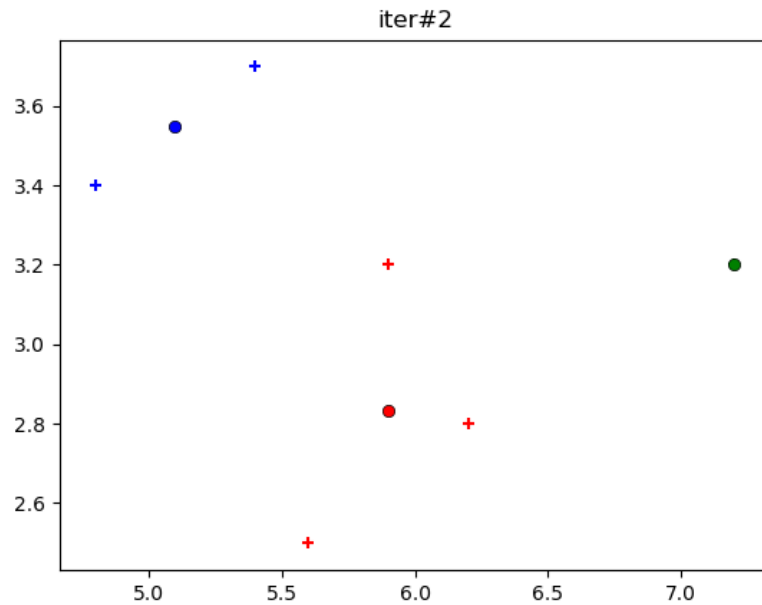


FIGURE 6 – Iteration finale

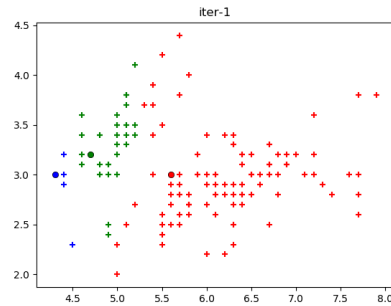
Et sur la console on obtient le résultat suivant :

```
c:\workspaces\big-data\data-analytics-kmeans>spark-submit kmeans.py data/iris_small.dat 3 10
18/04/29 17:40:30 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform
asses where applicable
iter #1: 1.3175589008869322
iter #2: 1.3795241880176898
Elapsed time: 0:00:24.838000
Number of iterations: 2
Final distance: 1.3795241880176898
```

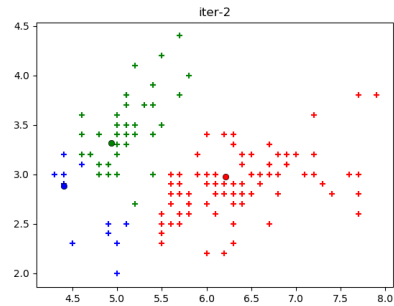
FIGURE 7 – Console

### 1.2.1 Resultat plot - Sepal

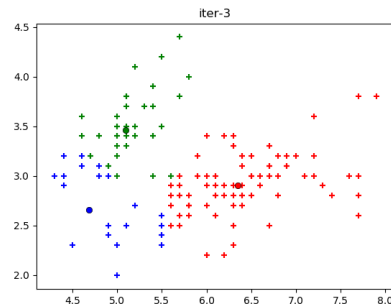
Ci-dessus vous pouvez voir la progression et les résultats de chaque itération sur l'exécution de l'algorithme KMeans sur les données iris clustering avec seulement les 2 coordonnées de sepal : sepal width en cm et sepal height en



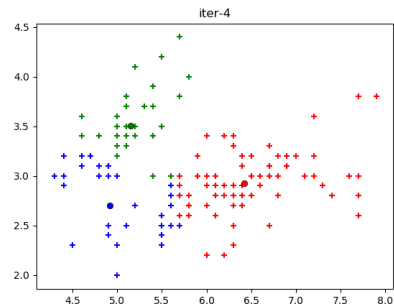
(a) Iteration 1



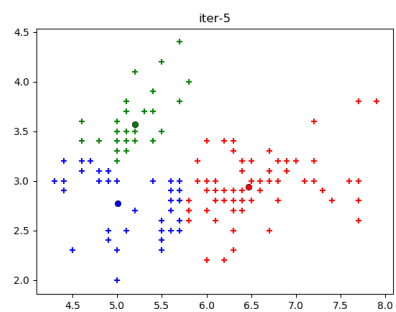
(b) Iteration 2



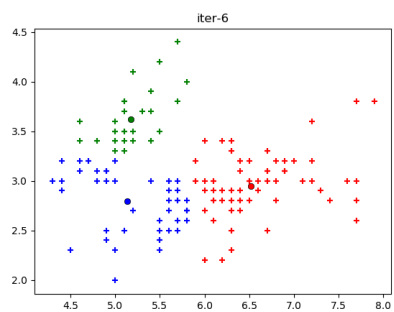
(a) Iteration 3



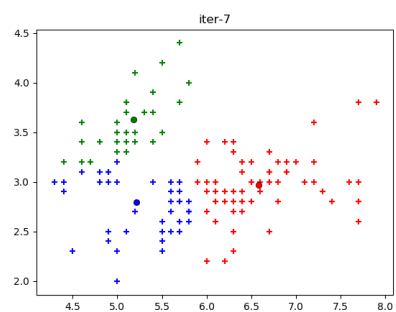
(b) Iteration 4



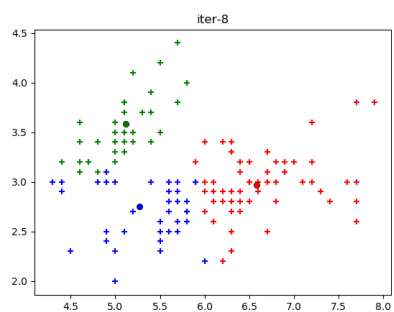
(a) Iteration 5



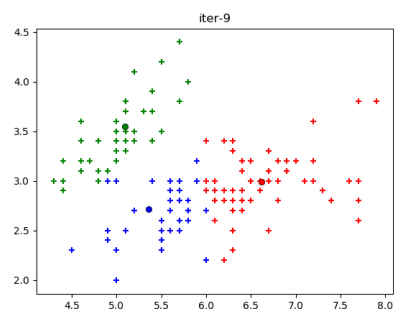
(b) Iteration 6



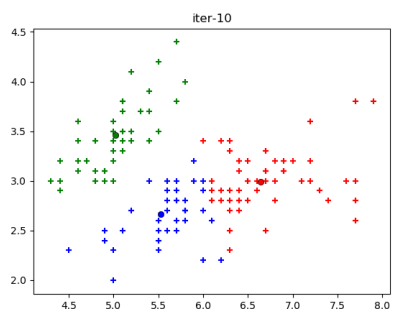
(a) Iteration 7



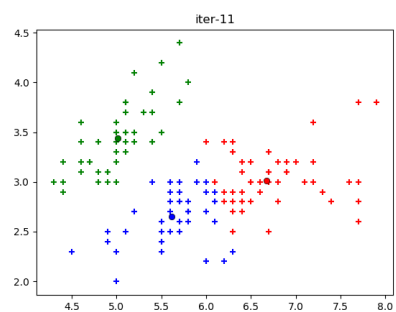
(b) Iteration 8



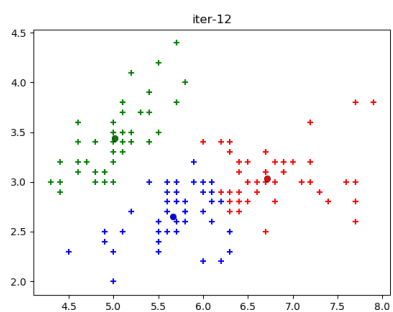
(a) Iteration 9



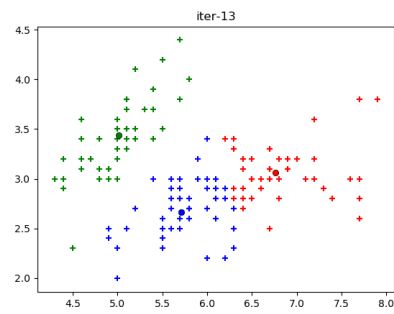
(b) Iteration 10



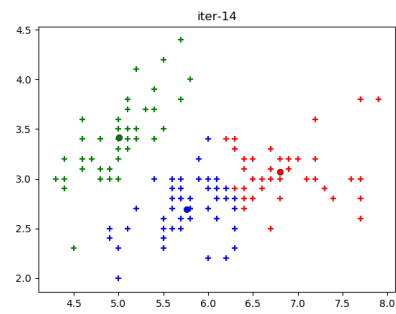
(a) Iteration 11



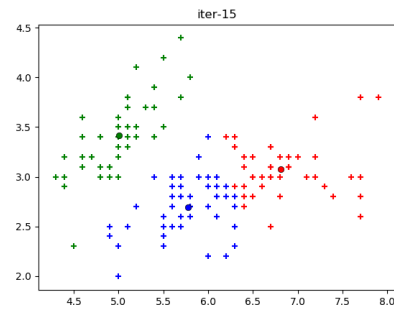
(b) Iteration 12



(a) Iteration 13



(b) Iteration 14

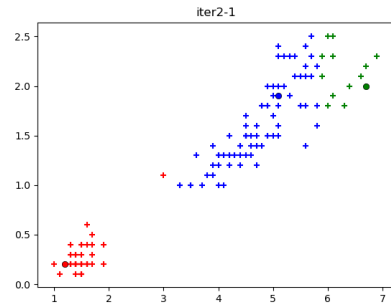


(a) Iteration 15

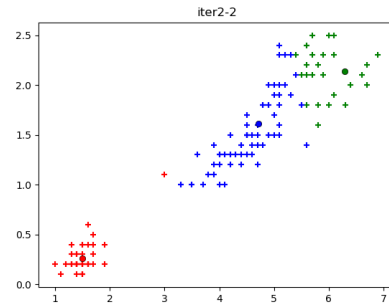
- Temps d'exécution : 4 min 20 sec
- Itérations : 15
- Distance : 1.259

### 1.2.2 Resultat plot - Petal

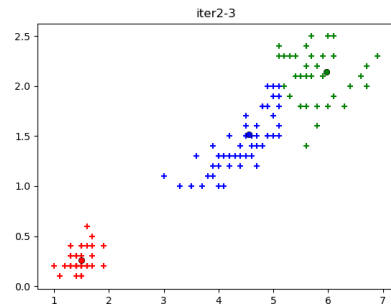
Ci-dessus vous pouvez voir la progression et les résultats de chaque itération sur l'exécution de l'algorithme KMeans sur les données iris clustering avec seulement les 2 coordonnées de sepal : petal width en cm et petal height en



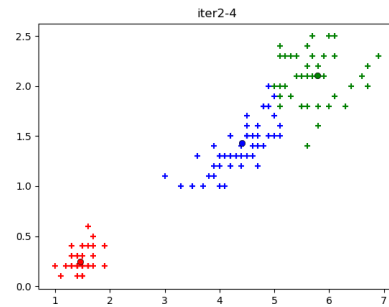
(a) Iteration 1



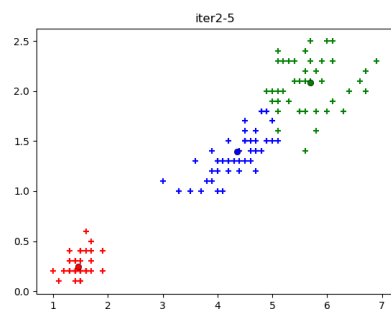
(b) Iteration 2



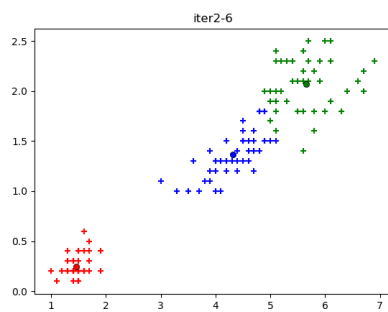
(a) Iteration 3



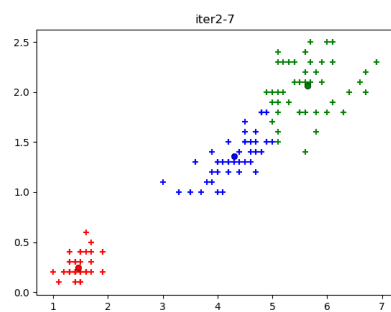
(b) Iteration 4



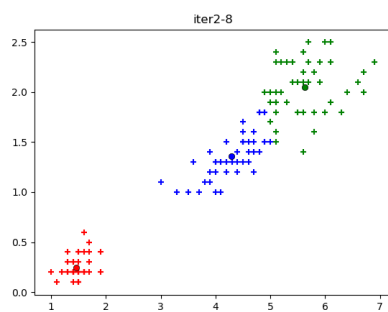
(a) Iteration 5



(b) Iteration 6



(a) Iteration 7



(b) Iteration 8

- Temps d'exécution : 2 min 20 sec
- Itérations : 8
- Distance : 1.1143

## 2 Analyse

### 2.1 Distance intra-cluster

En exécutant le programme, nous avons bien remarqué que la distance intra-cluster diminue progressivement jusqu'à atteindre un minimum durant la dernière iteration. On aperçoit parfois une légère augmentation suivi d'une diminution notable.

### 2.2 Convergence

L'algorithme s'arrête dès que les centroïdes gardent leurs valeurs durant deux itérations consécutives.

K-means avec une initialisation aléatoire des centres converge en 7 iterations, tandis que k-means++ réussit à converger en 6 iterations.

### 2.3 Tableau de performance

TABLE 1 – 100 executions

Algorithme	Temps d'exécution	Iterations	Distance finale
KMeans - implémenté	10.25 s	6.75	2.41
KMeans ++- implémenté	9.5 s	6.35	2.41

## 3 Application

Afin de tester sur plus de données, nous avons créé notre propre données. Nous avons collecté la taille ainsi que le poids de 92 joueurs de Basket de la NBA et leur position représente leurs clusters. (voir dans le répertoire data : nba.csv).

Donc notre but c'est de voir si c'est 2 coordonnées : taille (m) et poids (kg) ont une forte relation avec la position des joueurs.

Donc en exécutant l'algorithme on obtient les résultats suivantes :

- Temps d'exécution : 48 sec
- Itérations : 2
- Distance : 14.9



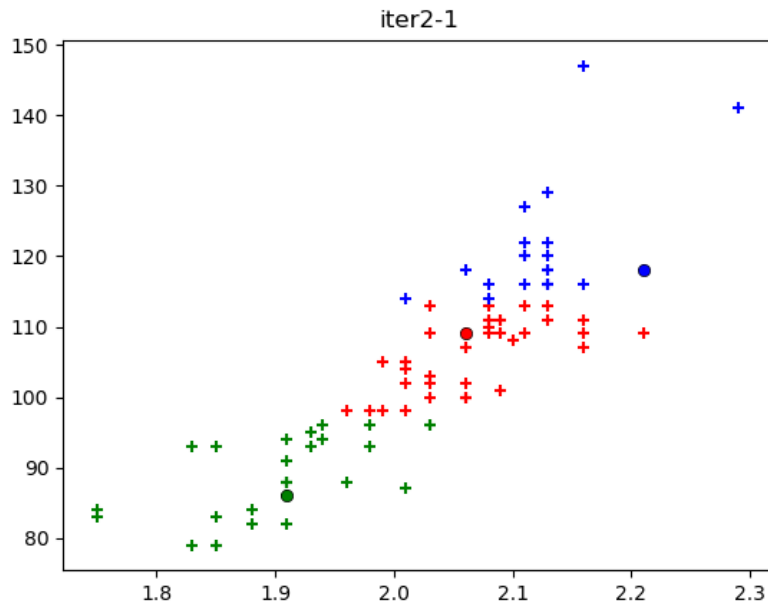


FIGURE 20 – Iteration 1

On remarque ici que les données sont assez réparties :

- les joueurs qui ont un petit poid et une petit taille sont regroupé (Guards)
- les joueurs qui ont un moyen poid et une moyenne taille sont regroupé (Forwards)
- les joueurs qui ont un grand poin et une grande taille sont regroupé (Center)

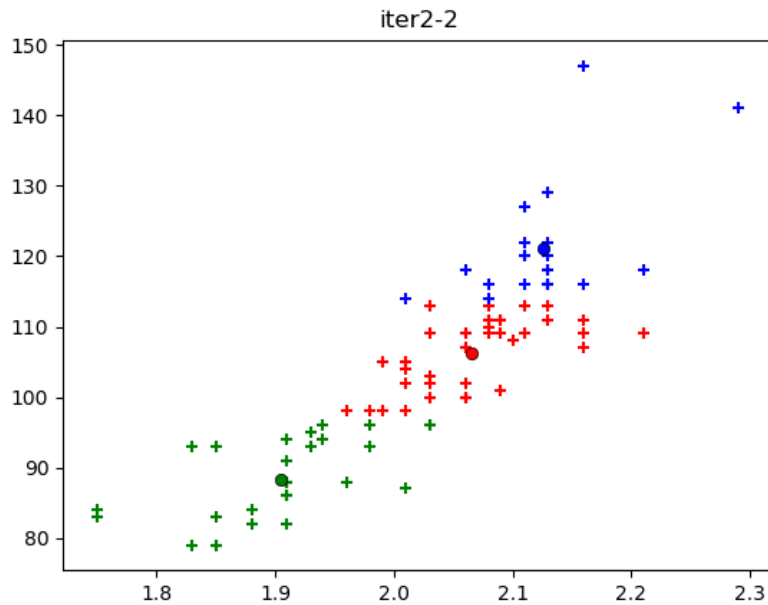


FIGURE 21 – Iteration 2

## Références

- [1] KMeans ++ in python [Online].  
Disponible : <https://datasciencelab.wordpress.com/2014/01/15/improved-seeding-for-clustering-with-k-means/>
- [2] Donnée de test sur Archive ICS UCI. [Online].  
Disponible : <https://archive.ics.uci.edu/ml/machine-learning-databases/soybean/soybean-small.data>
- [3] Initialization of k means clustering [Online].  
Disponible : <https://www.coursera.org/learn/cluster-analysis/lecture/bPyB1/3-3-initialization-of-k-means-clustering>
- [4] K-means clustering : how it works [Online].  
Disponible : [https://www.youtube.com/watch?v=\\_aWzGGNrcic](https://www.youtube.com/watch?v=_aWzGGNrcic)
- [5] Scalable K Means++ . [Online].  
Disponible : <https://www.youtube.com/watch?v=KHu8uNy4pV8&t=924s>