

---

# **Optimisation de l'algorithme de K-means en Spark**

---

Olivier RANDAVEL  
Adèle THORNER  
Valentin VU VAN

28 mars 2020

# Table des matières

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Analyse de la solution existante</b>	<b>3</b>
2.1	L'algorithme K-means . . . . .	3
2.2	Description de la solution existante . . . . .	3
<b>3</b>	<b>Exploration d'améliorations</b>	<b>6</b>
3.1	Optimisations liées au code . . . . .	6
3.1.1	Join des trois centroïdes par point . . . . .	6
3.1.2	Retirer le calcul de la racine dans la distance . . . . .	7
3.1.3	Éviter le groupByKey pour calculer le minimum . . . . .	7
3.2	Optimisation de l'algorithme des k-means . . . . .	8
3.2.1	Initialisation des centroïdes et k-means++ . . . . .	8
3.2.2	Utilisation de la distance de Manhattan . . . . .	9
3.2.3	Exploration aléatoire des features . . . . .	10
3.2.4	Condition d'arrêt . . . . .	11
3.2.5	Création d'une métrique pour évaluer la précision du résultat . . . . .	12
<b>4</b>	<b>Passage à l'échelle et analyse des résultats</b>	<b>13</b>
4.1	Jeu de données - Classer des voitures par type . . . . .	13
4.2	Protocole d'expérimentation . . . . .	13
4.3	Résultats . . . . .	14
<b>5</b>	<b>Implémentation avec des dataframes</b>	<b>14</b>
5.1	Fonction principale . . . . .	15
5.2	df_get_nb_switch() . . . . .	16
5.3	Résultats . . . . .	16
<b>6</b>	<b>Conclusion</b>	<b>17</b>
<b>7</b>	<b>Annexes</b>	<b>18</b>
7.1	RDD . . . . .	18
7.1.1	Fonctions . . . . .	18
7.1.2	Algorithme initial . . . . .	20
7.1.3	Algorithme optimisé . . . . .	22
7.1.4	Main . . . . .	24
7.2	Dataframe . . . . .	26
7.2.1	Fonctions . . . . .	26
7.2.2	Algorithme . . . . .	28
7.2.3	Main . . . . .	29
7.3	Jeu de test pour les dataframes . . . . .	30
7.4	Tableau des résultats . . . . .	31

# 1 Introduction

Les algorithmes de machine learning nous proposent des solutions d'apprentissage efficaces. En utilisant un grand volume de données, les algorithmes deviennent d'autant plus performants. Cependant, si les résultats théoriques s'améliorent avec la quantité de données, la rapidité de calcul se voit ralentie. En effet, on ne peut plus imaginer réaliser nos calculs de la même manière, les temps d'exécution seraient trop longs. C'est pour cela que l'on fait appel aux technologies de Big Data et du calcul distribué. Hadoop nous propose des solutions pour mettre en place des solutions stockées et exécutées de manière distribuée, partagée sur plusieurs machines. Spark est un framework permettant de créer et gérer des applications distribuées inspiré par l'algorithme MapReduce d'Hadoop. Ce framework est plus performant que les outils proposés par Hadoop mais nécessite d'être connecté à un système de gestion de fichiers distribués (celui d'Hadoop par exemple).

L'objectif de ce projet est de réaliser une implémentation efficace de l'algorithme de clustering des k-means, en utilisant la technologie Spark. Un code de base nous a été fourni et il nous faut l'analyser afin de proposer une solution améliorée qui met en avant les avantages du calcul distribué. Afin de mener à bien ce projet, nous avons analysé ensemble le code fourni, ligne par ligne et mis en oeuvre plusieurs optimisations, tout en gardant un historique de nos modifications, pour comparer l'évolution de nos fonctionnalités.

Au cours de ce rapport, nous analyserons dans un premier temps la solution existante. Ensuite, nous présenterons l'intégralité de nos optimisations en distinguant les améliorations du code et les optimisations de l'algorithme des k-means. Ensuite, nous décrirons l'ensemble des résultats que nous avons obtenu, en termes de temps d'exécution, avec le jeu de données fourni et un jeu plus volumineux. Depuis les résultats précédents, nous présenterons le passage des Resilient Distributed Datasets aux Dataframes et les résultats de performance obtenus.

## 2 Analyse de la solution existante

### 2.1 L'algorithme K-means

L'algorithme des K-moyennes (ou K-means en anglais) est un algorithme d'apprentissage non supervisé qui regroupe  $n$  observations  $x_i$  en  $K$  clusters. L'algorithme consiste à minimiser la distance de chaque point au centre (ou centroïde) de son cluster, de manière itérative. La fonction à minimiser est donc :

$$J(\mu, z) = \sum_{i=1}^n \sum_{j=1}^K z_j^i \|x_i - \mu_j\|^2$$

avec  $\mu_j$  le centre du cluster  $j$ , et  $z^i$  le vecteur d'affectation de l'observation  $x_i$  :  $z_j^i = 1$  si l'observation  $x_i$  appartient au cluster  $j$ , 0 sinon.

L'algorithme se déroule de la manière suivante. A l'initialisation, on choisit  $K$  centroïdes de manière aléatoire parmi les  $n$  observations. A chaque itération, on calcule la distance de chaque point aux  $K$  centroïdes et on affecte chacun de ces points au cluster dont il est le plus proche. On met ensuite à jour les centroïdes en calculant les moyennes de chaque cluster. L'algorithme s'arrête lorsque les clusters n'ont pas changé d'une itération à une autre.

### 2.2 Description de la solution existante

L'algorithme proposé initialement est défini par la fonction *simpleKmeans*. Cette fonction prend en entrée la donnée préalablement indexée et le nombre de labels. Cette solution a pour but de séparer la donnée en un nombre de groupes défini par le paramètre d'entrée *nb\_clusters*.

L'algorithme s'exécute ainsi :

Dans un premier temps, les centroïdes sont définis de façon aléatoire. Les centroïdes sont sélectionnés parmi les points donnés et sont tous différents. Enfin, ils sont numérotés pour pouvoir associer chaque observation à un centroïde. Nous obtenons donc un RDD constitué d'un id et des coordonnées du centroïde.

```
1 def simpleKmeans(data, nb_clusters):
2     centroides = sc.parallelize(data.takeSample('withoutReplacment', nb_clusters)) \
3         .zipWithIndex() \
4         .map(lambda x: (x[1], x[0][1][: -1]))
5     # (0, [4.4, 3.0, 1.3, 0.2])
```

Une fois les centroïdes initialisés, la boucle principale de l'algorithme commence. Voici les différentes étapes. Tout d'abord, la donnée est associée à chaque centroïde à l'aide d'un produit cartésien. La donnée est donc répétée autant de fois que le nombre de centroïdes. Ensuite, un calcul de distance est effectué

dans le but de fournir la distance entre chaque donnée et un centroïde. Puis, une opération *groupBy* permet de relier l'id d'une donnée à la distance la séparant des trois centroïdes. Enfin, il ne reste plus qu'à choisir la distance minimale et à assigner cette distance à la donnée.

```

1  def simpleKmeans(data, nb_clusters):
2      while not clusteringDone:
3          joined = data.cartesian(centroïdes)
4          # [(0, [5.1, 3.5, 1.4, 0.2, 'Iris-setosa']), (0, [6.0, 3.0, 4.8, 1.8]),
5          # (1, [4.9, 3.0, 1.4, 0.2, 'Iris-setosa']), (0, [6.0, 3.0, 4.8, 1.8]), ...]
6
7          dist = joined.map(lambda x: (x[0][0], (x[1][0],
8          computeDistance(x[0][1][: -1], x[1][1])))
9          # (0, (0, 0.866025403784438))
10
11         dist_list = dist.groupByKey().mapValues(list)
12         # (0, [(0, 0.866025403784438), (1, 3.7), (2, 0.5385164807134504)])
13
14         min_dist = dist_list.mapValues(closestCluster)
15         # (0, (2, 0.5385164807134504))
16
17         assignment = min_dist.join(data)
18         # (0, ((2, 0.5385164807134504), [5.1, 3.5, 1.4, 0.2, 'Iris-setosa']))

```

Le code au-dessus utilise deux fonctions. *computeDistance* permet de calculer la distance euclidienne entre deux coordonnées. Dans notre cas, cette fonction est utilisée pour calculer la distance entre les coordonnées d'un point donné et un centroïde.

$$\sqrt{\sum_{i,j=1}^n (x_i - y_j)^2}$$

La seconde fonction est *closestcluster*. Elle prend en entrée l'id de la donnée en question et les distances calculées aux trois centroïdes. Cette fonction retourne la distance minimale entre le point et les trois centroïdes.

```

1  def computeDistance(x,y):
2      return sqrt(sum([(a - b)**2 for a,b in zip(x,y)]))
3
4  def closestCluster(dist_list):
5      cluster = dist_list[0][0]
6      min_dist = dist_list[0][1]
7      for elem in dist_list:
8          if elem[1] < min_dist:
9              cluster = elem[0]
10             min_dist = elem[1]
11     return (cluster, min_dist)

```

Cette première étape terminée, vient alors le calcul des nouveaux centroïdes. Dans un premier temps, les coordonnées de chaque point sont rattachées à un numéro de centroïde. Puis, un comptage est effectué dans le but de comparer la population liée à chaque centroïde. Enfin, les nouveaux centroïdes sont calculés en prenant la moyenne de chaque variable de toutes les données liées au centroïde en question.

```

1 clusters = assignment.map(lambda x: (x[1][0][0], x[1][1][: -1]))
2 # (2, [5.1, 3.5, 1.4, 0.2])
3
4 count = clusters.map(lambda x: (x[0],1)).reduceByKey(lambda x,y: x+y)
5 # count = (1, 100), (2, 200), (3, 50) # (numero de cluster, nb d'element)
6
7 somme = clusters.reduceByKey(sumList)
8 # (1,2000), (2, 6000), (3, 1500) # (numero de cluster, somme des elements)
9
10 centroidesCluster = somme.join(count).map(lambda x: (x[0], moyenneList(x[1][0], x
11 [1][1])))
# (0, 4.1, 2.5, 2.4, 0.8) # (numero de cluster, average)

```

On retrouve deux fonctions utilisées dans le code ci-dessus. La première est *sumList*. Elle permet de sommer deux à deux les éléments de deux listes.

$$x = [x_1, x_2], y = [y_1, y_2]$$

$$z = [x_1 + y_1, x_2 + y_2]$$

La seconde fonction *moyenneList* retourne la moyenne des éléments d'une liste.

```

1 def sumList(x,y):
2     return [x[i]+y[i] for i in range(len(x))]
3
4 def moyenneList(x,n):
5     return [x[i]/n for i in range(len(x))]

```

La dernière étape permet d'évaluer le changement d'état et vérifie si l'algorithme doit s'arrêter. A partir de la deuxième boucle, il est possible de comparer les répartitions précédentes et actuelles. Si on constate un nombre de changements nul ou que le nombre de boucles est de 100, alors l'algorithme s'arrête et renvoie l'erreur de classification ainsi que la donnée classée. Dans le cas contraire, l'algorithme effectue une nouvelle boucle en se basant sur les nouveaux centroïdes. L'erreur ainsi calculée est retournée par l'algorithme. La formule de l'erreur suit cette formulation : prendre la racine carrée de la somme des distances entre les centroïdes et les points, divisée par le nombre de points.

$$erreur = \sqrt{\frac{\sum_i^n (dist_i)}{n}}$$

```

1  if number_of_steps > 0:
2      switch = prev_assignment.join(min_dist)\
3          .filter(lambda x: x[1][0][0] != x[1][1][0])\
4          .count()
5  else:
6      switch = 150
7  if switch == 0 or number_of_steps == 100:
8      clusteringDone = True
9      error = sqrt(min_dist.map(lambda x: x[1][1])
10         .reduce(lambda x,y: x + y))/nb_elem.value
11  else:
12      centroides = centroidesCluster
13      prev_assignment = min_dist

```

### 3 Exploration d'améliorations

En explorant différentes pistes d'amélioration de la solution proposée, nous avons mis en place un certain nombre d'optimisations, de deux types différents. Les premières améliorations sont liées à de l'optimisation du code, notamment afin d'optimiser les requêtes Spark. Les secondes améliorations mettent en avant l'optimisation de l'algorithme des K-moyennes pour obtenir de meilleurs temps d'exécution. Nous avons réalisé l'ensemble de notre implémentation sur la plateforme Databricks, afin d'utiliser des clusters distants.

#### 3.1 Optimisations liées au code

##### 3.1.1 Join des trois centroïdes par point

**Solution existante** L'algorithme des k-means nécessite le calcul, à chaque itération, des distances de chaque point aux différents centroïdes. Dans la solution proposée, on crée un RDD contenant un point et un cluster à l'aide d'un produit cartésien. De cette manière, le RDD obtenu est redondant dans le sens où chaque point est présent sur 3 lignes.

```

1  joined = data.cartesian(centroides)

```

**Notre amélioration** Nous avons donc choisi de regrouper ces trois lignes en une seule afin de proposer un RDD contenant un point et une liste des trois centroïdes. Le RDD *joined* devient alors :

```

1  joined = join_centroide_data(data,3)

```

avec

```

1 def join_centroide_data(data_indexed, nb_clusters):
2     centroides = centroides_input(data_indexed)
3     centroides = centroides.collect()
4     joined = data_indexed.map(lambda x : (x, centroides))
5     return joined

```

### 3.1.2 Retirer le calcul de la racine dans la distance

Nous avons fait le choix de nous séparer du calcul de la racine carrée de la formule la distance euclidienne, car notre utilisation de la distance est simplement comparative. La fonction racine étant croissante, si  $\sqrt{a} > \sqrt{b}$ , alors  $a > b$ . Cela permet d'effectuer un calcul de moins lors du calcul de chaque distance, qui est une opération répétée de nombreuses fois lors du déroulé de l'algorithme. Nous avons ainsi réécrit la fonction de distance de la manière suivante :

```

1 def computeDistance_improve(x,y):
2     return sum([(a - b)**2 for a,b in zip(x,y)])

```

### 3.1.3 Éviter le groupByKey pour calculer le minimum

**Solution existante** A chaque itération de l'algorithme, il est nécessaire de calculer la distance de chaque point aux centroïdes et de comparer ces trois distances afin de choisir la distance minimale. La solution proposée était la suivante :

```

1 joined = data.cartesian(centroides)
2 dist = joined.map(lambda x: (x[0][0],(x[1][0], computeDistance(x[0][1][: -1], x[1][1])))
3 dist_list = dist.groupByKey().mapValues(list)
4 min_dist = dist_list.mapValues(closestCluster)
5 assignment = min_dist.join(data)

```

Le code ci-dessus crée d'abord un RDD avec un point et un cluster comme nous l'avons vu précédemment. Ensuite, pour chaque couple de points il calcule la distance à l'aide de la fonction `compute_distance` et réalise un `groupByKey` pour associer à chaque point les trois distances aux centroïdes. Enfin, on associe à chaque point le centroïde le plus proche et on effectue un join avec les données de départ pour récupérer les coordonnées du point.

**Notre amélioration** L'idée de notre amélioration était de se débarrasser du `groupByKey` que l'on sait si coûteux. Nous utilisons l'amélioration vue dans la section 3.1.1 qui consistait à regrouper les trois centroïdes pour un point. Ainsi, pour calculer les distances, on parcourt une fois le RDD et pour chaque ligne on effectue 3 calculs de distance. Nous voyons une vraie optimisation au moment du calcul du minimum. En effet, au lieu de réaliser un `groupByKey`, il nous suffit de comparer les trois distances associées au point



et de choisir la plus petite. Notre fonction `min_` ci-dessous nous permet de combiner les calculs des distances et leur comparaison en une fonction. Ainsi, on limite les temps de lecture du RDD, il nous suffit d'accéder une fois à l'élément pour obtenir son cluster.

```
1 joined = join_centroide_data(data,3)
2 min_dist = min_(joined)
3 assignment = min_dist.join(data)
```

avec

```
1 def min_(joined) :
2     dist_list = joined.map(lambda x : (x[0][0],
3         ((x[1][0][0], computeDistance_improve(x[0][1][: -1], x[1][0][1])),
4         (x[1][1][0], computeDistance_improve(x[0][1][: -1], x[1][1][1])),
5         (x[1][2][0], computeDistance_improve(x[0][1][: -1], x[1][2][1]))
6         )))
7     min_dist = dist_list.mapValues(closestCluster)
8     return min_dist
```

## 3.2 Optimisation de l'algorithme des k-means

### 3.2.1 Initialisation des centroïdes et k-means++

L'idée de ne pas choisir aléatoirement les centroïdes initiaux nous est venue après avoir fait tourner la version "basique" du code. En effet, après différents essais, nous avons remarqué que les temps d'exécution pouvaient s'avérer beaucoup plus longs selon les choix des premiers centroïdes, et parfois mener à des erreurs de type "Executor heartbeat timed out after 146532 ms". L'une des explications à cette erreur est le choix initial des centroïdes, qui sont choisis aléatoirement (cf partie 2.2) et il arrive donc que ceux-ci soient trop proches, entraînant un grand nombre d'itérations supplémentaires. Nous avons alors pensé à une autre méthode pour éviter ce genre de "départs ratés", dont l'idée est très simple : prendre le point A le plus grand possible, le point B le plus éloigné de ce point A et enfin le point C étant au milieu de A et B. Cette méthode d'initialisation existe en fait dans un algorithme nommé k-means ++, qui est, comme son nom l'indique, une optimisation de k-means.

*Remarque :* la méthode, telle que nous l'avons implémentée, n'est applicable que dans notre cas car le nombre de clusters à former est inférieur ou égal à 3. Pour ce faire, nous avons défini les fonctions suivantes :

#### Création des centroïdes

```
1 def centroides_input(data) :
2     centroide1 = data.map(lambda x : (0, (sum(x[1][: -1]), x[1][: -1])))
3     .map(lambda x : (x[0], x[1][1])).max(lambda x: x[1][0])
4     y = centroide1[1]
5     centroide2 = data.map(lambda x : (1, (computeDistance_improve(y, x[1][: -1]), x[1][: -1])))
6     .map(lambda x : (x[0], x[1][1])).min(lambda x: x[1])
```

```

6  y2 = centroide2[1]
7  point = [x + y for x, y in zip(y, y2)]
8  point = [x / 2 for x in point]
9  centroide = sc.parallelize([centroide1, centroide2, (2, point)])
10 return centroide

```

On commence par sommer la valeur de chaque feature pour chaque point de notre jeu de données : on obtient une valeur (la somme de chacun des 4 features pour chaque point). On en extrait alors le maximum, que l'on stocke dans la variable *y*.

Ensuite, grâce à la fonction *computeDistance\_improve* présentée plus haut, on identifie *y2*, le point le plus éloigné du point *y*. Pour finir, on calcule le point situé au milieu de *y* et *y2*. Ces trois points seront les centroïdes initiaux de notre algorithme.

*Remarque* : le troisième centroïde n'est pas forcément un point appartenant à notre jeu de données. Ce n'est pas un problème car au fur et à mesure des itérations, le centroïde de chaque cluster évolue et n'est généralement plus équivalent à un point du jeu de données.

*Limites* : Si le jeu de données n'est pas propice à une bonne initialisation des centroïdes malgré cette méthode, le fait de relancer l'algorithme ne pourra pas permettre d'avoir un meilleur résultat, car les clusters initiaux seront toujours les mêmes. Pour cela, nous avons inséré dans notre fonction principale un paramètre de valeur *True* ou *False*, qui permet d'activer ou non cette initialisation truquée.

### 3.2.2 Utilisation de la distance de Manhattan

Nous avons vu dans la partie précédente que le calcul de la distance euclidienne pouvait être optimisé en supprimant le passage à la racine, ce dernier n'étant pas nécessaire pour obtenir la distance minimale. Il existe un grand nombre de distances, plus ou moins complexes à implémenter. Nous cherchons à améliorer le temps d'exécution de l'algorithme. Nous avons choisi de tester une distance différente, pour observer son impact sur la vitesse de convergence de l'algorithme.

#### La distance de Manhattan

$$D(x, y) = \sum_{i=0}^n |x_i - y_i|$$

Pour pouvoir utiliser cette distance, en la combinant avec les différentes améliorations (centroïdes regroupés, random features, etc), nous avons implémenté les fonctions suivantes, qui sont utilisées en substitutions des fonctions de départ qui utilisent la distance euclidienne.

```

1  # if the centroids are not joined in 1 row
2  def computeDistance_manhattan(x, y):
3      return sum([(a - b) for a, b in zip(x, y)])
4
5  # use manhattan distance with random features and joined centroids
6  def compute_distance_manhattan_features(x, y, features):
7      sum = 0
8      for f in features:
9          sum += (x[f] - y[f])

```

```

10     return sum
11 #use manhattan distance with all features and joined centroides
12 def compute_distance_manhattan(x,y):
13     return compute_distance_manhattan(x,y,[0,1,2,3])
14
15 # get min of manhattan distances with all features
16 def min_manhattan(joined) :
17     # [(0, ((0, 0.30000000000000016), (1, 5.31224999411737), (2, 3.057776970284131)) ) )
18     dist_list = joined.map(lambda x : (x[0][0],
19         ((x[1][0][0], compute_distance_manhattan(x[0][1][: -1], x[1][0][1])),
20         (x[1][1][0], compute_distance_manhattan(x[0][1][: -1], x[1][1][1])),
21         (x[1][2][0], compute_distance_manhattan(x[0][1][: -1], x[1][2][1]))
22         )))
23     min_dist = dist_list.mapValues(closestCluster)
24     return min_dist
25
26 # get min of manhattan distances with random features
27 def min_manhattan_selectedfeatures(joined, features) :
28     dist_list = joined.map(lambda x : (x[0][0],
29         ((x[1][0][0], compute_distance_manhattan_features(x[0][1][: -1], x[1][0][1], features)),
30         (x[1][1][0], compute_distance_manhattan_features(x[0][1][: -1], x[1][1][1], features)),
31         (x[1][2][0], compute_distance_manhattan_features(x[0][1][: -1], x[1][2][1], features))
32         )))
33     min_dist = dist_list.mapValues(closestCluster)
34     return min_dist

```

### 3.2.3 Exploration aléatoire des features

A chaque itération, l'algorithme évalue la distance de chaque point aux différents centroïdes en calculant la distance euclidienne. Cette dernière dépend de l'ensemble des features des observations. Nous avons voulu mettre en place une solution permettant de réduire le temps d'exécution en évaluant seulement une distance partielle à chaque itération.

L'idée était la suivante. Pendant les premières itérations de l'algorithme, choisir aléatoirement 2 features parmi les 4 et calculer les distances de l'itération par rapport à ces deux caractéristiques uniquement. Ainsi, on effectue deux fois moins de calculs que si l'on prenait en compte toutes les itérations.

De cette idée ont découlé plusieurs problématiques. Tout d'abord, nous nous sommes rendu compte que l'algorithme ne convergerait pas toujours au sens des K-moyennes (aucun changement dans les affectations au cluster) en regardant seulement un échantillon des features. Il fallait donc déterminer une solution pour calculer la distance "complète", c'est-à-dire la distance qui prend en compte les 4 features, à la fin de l'algorithme. Nous avons souhaité prendre en paramètre un seuil, nommé dans le code ci-dessous "epsilon\_change\_nbfeatures", qui détermine lorsque l'on doit changer le nombre de features à prendre en compte (2, puis 3, puis 4). Ce seuil permet d'évaluer lorsque l'algorithme est en train de stagner. Nous avons établi que l'algorithme stagnait lorsque le nombre d'observations qui changeaient de cluster étaient en-dessous de ce seuil à deux itérations de suite. En effet, on considère alors que nous n'avons plus assez d'informations (de features) pour améliorer l'affectation.

Une autre problématique qui nous est venue à l'esprit était que l'algorithme risquait de partir dans la mauvaise direction en ne prenant en compte que certaines features. Pour remédier à cela, nous avons ajouté une option qui permet d'effectuer une itération de "contrôle" régulièrement. On peut alors préciser un entier qui détermine l'intervalle entre deux "contrôles". On appelle une itération de "contrôle" une itération qui prend en compte l'ensemble des features pour le calcul de la distance.

```
1      # start by checking if we need to change the number of features
2      if (number_of_steps > 0) and switch < epsilon_change_nbfeatures and switch_prev <
        epsilon_change_nbfeatures and nb_features_selected < 4:
3          nb_features_selected += 1
4          random.seed(number_of_steps)
5          features = random.sample([0,1,2,3], k=nb_features_selected) #choose randomly the
        features to observe
6
7          min_dist = min_selectedfeatures(joined, features) # specify the features selected
```

*features* va ainsi prendre un nombre de valeurs aléatoires entre 0 et 3 selon la valeur de *k*. Il est important de changer les features sur lesquels sont faits les calculs à chaque itération pour ne pas partir dans une mauvaise direction. *k* commence avec une valeur de 2 et augmente lorsque l'algorithme commence à stagner. Cette méthode par paliers permet de converger plus vite vers une solution acceptable, et finit par effectuer ses calculs sur tous les features (lorsque  $k = 4$ ) pour s'assurer de finir sur un résultat correct, c'est-à-dire arriver à convergence au sens de k-means.

**Utilisation de `random.seed()`** Afin de pouvoir comparer de manière pertinente les résultats obtenus sur plusieurs exécutions de l'algorithme, nous utilisons la fonction `random.seed()` qui retournera le même nombre si on l'appelle avec le même entier. Ici on choisit en paramètre le numéro de l'itération, les features choisis différeront d'une itération à une autre mais d'un run à un autre, on choisira les mêmes features à la même itération.

```
1  def min_selectedfeatures(joined, features) :
2      dist_list = joined.map(lambda x : (x[0][0],
3          ((x[1][0][0], compute_distance_selectedfeatures_nosqrt(x[0][1][: -1], x[1][0][1],
4              features)),
5          (x[1][1][0], compute_distance_selectedfeatures_nosqrt(x[0][1][: -1], x[1][1][1],
6              features)),
7          (x[1][2][0], compute_distance_selectedfeatures_nosqrt(x[0][1][: -1], x[1][2][1],
8              features))
9          )))
10     min_dist = dist_list.mapValues(closestCluster)
11     return min_dist
```

### 3.2.4 Condition d'arrêt

La condition d'arrêt initialement implémentée se base sur deux aspects. Le premier concerne une limite sur le nombre de fois où les centroïdes sont calculés, cette limite est fixée à 100. La seconde permet

de stopper l'algorithme si aucune donnée n'a été labélisée différemment (l'état présent est identique au précédent).

Dans notre cas pour chaque itération, comme présenté précédemment, nous n'utilisons pas toujours toutes les caractéristiques pour constituer nos groupes. Donc, une nouvelle condition s'ajoute. Nous vérifions que toutes les variables explicatives disponibles ont été utilisées si aucune observation n'a changé d'affectation. Ces points sont présentés dans le code ci-dessous.

```
1     if (switch == 0 and nb_features_selected == 4) or number_of_steps > 50:
2         clusteringDone = True
3         error_improve = erreur_improve(assignment)
4         error = sqrt(min_dist.map(lambda x: x[1][1]).reduce(lambda x,y: x + y))/nb_elem
           .value
```

### 3.2.5 Création d'une métrique pour évaluer la précision du résultat

L'algorithme du k-means permet de retourner les données labélisées. Chaque donnée est alors associée à un chiffre entre 0 et le nombre de labels distincts. Nous avons identifié deux solutions pour évaluer la précision de l'algorithme. Le but consiste à rapprocher un chiffre d'un label.

La première solution se déroule ainsi. Pour chaque label, nous calculons la moyenne des caractéristiques, puis nous évaluons la distance euclidienne entre ces moyennes et chaque centroïde résultant de l'algorithme des k-means. Enfin, le label est associé au centroïde le plus proche. Cette première approche s'est avérée peu concluante. En effet, deux labels pouvaient être rapprochés d'un même centroïde. De plus, le choix de la moyenne comme moyen de comparaison est discutable, la médiane ou encore un calcul de quartile auraient aussi pu être envisagés.

Au vu de ces problèmes, une deuxième méthode a été implémentée. Celle-ci consiste à retourner l'erreur de classification optimale. Dans le cadre d'un jeu de données labélisées avec trois labels distincts, il suffit de faire 3! (3 factoriel) opérations (6 opérations) et de prendre la valeur pour laquelle le moins d'erreurs ont été effectuées. Cette méthode est présentée ci-dessous.

```
1     def erreur_improve(assignment):
2         list1 = [1,1,2,2,3,3]
3         list2 = [2,3,3,1,2,1]
4         list3 = [3,2,1,3,1,2]
5         label = assignment.map(lambda x: x[1][1][4]).distinct().collect()
6         percentage_error_list = []
7         for a,b,c in zip(list1, list2, list3) :
8             combinaison = [(label[0],a),(label[1],b),(label[2],c)]
9             pourcentage_error = 0
10            combinaison_assignment = assignment.map(lambda x: (x[1][1][4], x[1][0][0])).
collect()
11            for i in range(len(combinaison_assignment)) :
12                if combinaison_assignment[i] not in (combinaison) :
13                    pourcentage_error += 1
14            pourcentage_error /= len(combinaison_assignment)
```

```

15     percentage_error_list.append(pourcentage_error)
16     return min_list(pourcentage_error_list)
17
18     def min_list(array) :
19         min = array[0]
20         for i in range(1, len(array)) :
21             if min > array[i] :
22                 min = array[i]
23     return min

```

## 4 Passage à l'échelle et analyse des résultats

Nous avons souhaité implémenter les différentes optimisations décrites précédemment et comparer les résultats obtenus avec les résultats donnés par la solution initiale. Pour ce faire, nous avons créé une fonction globale `custom_kmeans` qui prend en paramètres des booléens permettant de préciser les différentes optimisations à mettre en oeuvre. Cette fonction retourne les résultats de l'algorithme, c'est-à-dire les différentes erreurs calculées, le nombre d'itération effectuées et les centres des clusters.

A chaque appel du *main*, on précise combien de fois on veut exécuter les algorithmes, et on indique si on veut faire tourner l'algorithme initial, l'algorithme amélioré ou les deux. Le *main* retourne le temps d'exécution moyen de chacun de ces algorithmes ainsi que chacun de ces temps d'exécution. Cette fonction globale nous a permis de garder une trace des différentes optimisations et éventuellement de combiner des optimisations compatibles. Nous avons choisi de comparer les résultats en termes de temps d'exécution.

### 4.1 Jeu de données - Classer des voitures par type

Afin d'obtenir des résultats significatifs, nous avons choisi d'effectuer nos expériences sur un jeu de données plus conséquent que le jeu iris. Étant donné que l'algorithme initial et, par conséquent, l'algorithme amélioré, étaient tous deux adaptés à un problème à 4 features et 3 clusters, nous avons cherché un jeu de données du même format.

Nous avons trouvé un jeu de données de plus de 6000 lignes portant sur la classification de voitures en trois gammes (inférieure, luxe et supérieure), selon 4 caractéristiques (puissance maximale, consommation, émission de CO2 et poids du véhicule à vide). Ces données proviennent du site gouvernemental <https://www.data.gouv.fr/fr/datasets/emissions-de-co2-et-de-polluants-des-vehicules-commercialises-en-france/>.

### 4.2 Protocole d'expérimentation

Nous avons implémenté le protocole suivant pour obtenir des résultats comparables de nos différents algorithmes.

**Définition des différentes améliorations** Nous avons commencé par définir une liste de fonctions que l'on souhaitait comparer, c'est-à-dire une liste de combinaisons d'améliorations. Nous avons réalisé un tableau permettant de mettre en évidence les valeurs à donner aux paramètres de la fonction `custom_kmeans`, pour chaque expérience. Ce tableau est donné en annexe.

Pour que nos résultats soient significatifs, nous avons choisi de lancer chaque expérience 5 fois de suite, sur les datasets `iris` et `voitures_types`. Un tableau regroupe l'ensemble de ces résultats, pour chaque expérience et chaque dataset. Ce tableau est également ajouté en annexe de ce rapport.

## 4.3 Résultats

L'exécution de nos différentes expériences nous a permis de mettre en valeur plusieurs résultats.

**Suppression du `groupBy`** En comparant les temps d'exécution moyens, nous avons pu constater qu'une de nos fonctionnalités améliorait considérablement la performance de l'algorithme. Il s'agit de notre option `join_centroids_in_rdd` qui nous permet de représenter les trois centres pour une observation plutôt que d'exploser ces informations en trois lignes. Comme nous l'avons précisé dans la section 3.1.1, cette fonctionnalité nous permet d'éviter de réaliser une opération `groupBy`, que l'on sait très coûteuse. Les résultats nous prouvent cette affirmation. En effet, on constate que lorsque cette option n'est pas activée, les algorithmes n'ont pas réussi à converger. On obtient, sans cette amélioration, une erreur d'exécution due à un temps trop long.

**Choix des centroïdes initiaux** Cette deuxième fonctionnalité permet d'améliorer les performances du modèle, lorsqu'elle est combinée à la précédente. En effet, le fait de choisir les premiers centroïdes biaise le modèle et permet de converger plus rapidement.

**Random features** Contrairement à ce que l'on pensait lors de l'implémentation de cette fonctionnalité, `random_features` n'améliore pas la performance de l'algorithme. Au contraire, les temps d'exécution moyens observés sont même un peu ralentis. Ceci est dû au fait que la sélection d'un sous ensemble des features entraîne une convergence en plus d'itérations. De plus, cette fonctionnalité ne paraît pas très pertinente pour un jeu de données à seulement 4 features. La différence de temps de calcul de la distance avec 2 ou 4 features est négligeable. En revanche, il serait intéressant de tester cette option sur des jeux de données comportant un grand nombre de features.

**Distances** Les résultats nous montrent que l'utilisation d'une distance plutôt qu'une autre ne changeait pas les temps d'exécution.

## 5 Implémentation avec des dataframes

Jusqu'ici, nous avons uniquement manipulé nos données sous la forme de Resilient Distributed Datasets (RDD). Dans cette partie, nous présenterons notre migration vers l'utilisation de dataframes. Une dataframe est une structure de données qui s'apparente à un tableau. Cette structure de données a été pensée sur Spark pour pouvoir être manipulée de manière distribuée.

Comme décrit dans la partie 3 de ce rapport, nous avons mis en oeuvre un grand nombre d'optimisations du code initialement fourni. L'analyse des performances de ces différentes optimisations nous a permis de mettre en avant la meilleure combinaison de fonctionnalités, en termes de performance. Les différentes optimisations liées à la gestion propre des RDD ne peut se retranscrire avec des dataframes. En revanche, les modifications de l'algorithme en lui même peuvent être maintenues lors du passage en dataframes. Nous avons choisi de ne mettre en oeuvre que la meilleure combinaison d'optimisations liées à l'algorithme des k-means, lors du passage aux dataframes.

## 5.1 Fonction principale

```

1  def df_kmeans(data, nb_steps_max):
2      clusteringDone = False
3      number_of_steps = 0
4      current_error = float("inf")
5      nb_elem = sc.broadcast(data.count())
6      switch = -1
7      # compute initial centroids
8      centroide1, centroide2, centroide3 = centroide_input_df(data)
9      data = data.drop("_c0").drop("_c1").drop("_c2").drop("_c3").drop('total')
10     while(not clusteringDone and number_of_steps < nb_steps_max):
11         min_dist_df = df_compute_min_distances(data, centroide1, centroide2, centroide3)
12         #index, coordonnees, dist1, dist2, dist3, cluster
13         new_centroids_df = df_update_centroids(min_dist_df) # cluster, avg_coords
14         centroide1, centroide2, centroide3 = df_get_new_centroids(new_centroids_df)
15         if (number_of_steps > 0):
16             switch = df_get_nb_switch(min_dist_df_prev, min_dist_df)
17             min_dist_df_prev = min_dist_df
18             number_of_steps += 1
19
20     if (switch == 0):
21         clusteringDone = True
22         labels = data.select("label").distinct().collect()
23         error = df_compute_error(min_dist_df, labels)
24         if number_of_steps % 10 == 0:
25             print("number_of_steps :", number_of_steps)
26         assignment = min_dist_df.drop("dist1").drop("dist2").drop("dist3")
27         centroides = [centroide1, centroide2, centroide3]
28         return (assignment, error, number_of_steps, centroides)

```

Le développement de cette fonction est évidemment fortement similaire à celui réalisé pour l'algorithme présenté dans la partie 3. Nous avons simplement pu profiter de certaines spécificités des dataframes. Il est par exemple inutile de copier la totalité des coordonnées de chaque cluster pour chaque élément de la base de données (ce qui rajouterait en effet 12 colonnes, 4 colonnes pour chaque cluster). Ces informations sont stockées dans les variables *centroidei*.

Concernant la valeur de *switch*, elle est ici obtenue en effectuant une jointure entre les clusters de l'étape précédente et de l'étape suivante. On ne garde que les éléments ayant subi un changement dans le dataframe et on compte ce nombre d'éléments.



## 5.2 df\_get\_nb\_switch()

```
1 def df_get_nb_switch(df_prev, df_suiv):
2     join = df_prev.withColumnRenamed("cluster", "cluster_prev")
3     .join(df_suiv, "index")
4     .filter(F.col("cluster_prev") != F.col("cluster"))
5     return join.count()
```

## 5.3 Résultats

Après avoir lancé la fonction principale exécutant le code avec les dataframes, les résultats obtenus ne se sont pas avérés satisfaisants. Nous obtenons une convergence systématique de l'algorithme mais les temps obtenus sont supérieurs à ceux obtenus avec l'utilisation des RDD. Deux explications logiques nous viennent en tête :

1. La présence d'un join pour le calcul du switch, qui est en effet lourde en terme de temps de calcul;
2. Les centroïdes stockés dans des variables, nous ne savons pas comment le système gère cela de manière distribuée.

## 6 Conclusion

De nombreuses solutions efficaces existent déjà en Python pour l'algorithme des k-means. Malheureusement, lorsque notre jeu de données est trop conséquent, ces approches "classiques" ne sont plus adaptées. Il est alors nécessaire de faire appel au calcul distribué, et c'est dans ce cadre que se situe ce projet. Ces méthodes de calcul distribué impliquent une réflexion en amont plus poussée. Si certaines solutions nous paraissent plus intuitives, elles ne seront pas nécessairement performantes. Cette notion est parfaitement illustrée par le code fourni pour commencer le projet : logique mais peu efficace. Ce rapport résume les différentes idées que nous avons eues et implémentées pour répondre à ce problème.

Les résultats obtenus nous permettent de mettre en avant deux optimisations du code. La première consiste à éviter l'utilisation d'un `groupBy`. La deuxième permet de choisir les centroïdes initiaux pour amener plus rapidement à convergence (cette approche est inspirée de l'algorithme K-means++). De manière à mieux observer nos améliorations, il serait intéressant d'adapter puis d'appliquer notre code à des jeux de données comportant plus de features, et répartis dans plus de clusters. En particulier, notre exploration aléatoire des features (cf 3.2.3) apporterait probablement une meilleure performance en grande dimension.

## 7 Annexes

Le code peut être trouvé au lien suivant :

<https://databricks-prod-cloudfront.cloud.databricks.com/public/4027ec902e239c93eaaa8714f173bcfc/3631309044632400/4360712146874161/8024126452304013/latest.html>

### 7.1 RDD

#### 7.1.1 Fonctions

```
1 from pyspark import SparkContext, SparkConf
2 from math import sqrt
3 import pyspark.sql.functions as f
4 import time
5 from time import gmtime, strftime
6 from datetime import datetime
7 import pandas as pd
8 import random
9
10 def seconds_between(d1, d2):
11     d1 = datetime.strptime(d1, "%H:%M:%S")
12     d2 = datetime.strptime(d2, "%H:%M:%S")
13     return abs((d2 - d1).seconds)
14
15 def centroides_input(data) :
16     centroide1 = data.map(lambda x : (0, (sum(x[1][: -1]), x[1][: -1]))).map(lambda x : (x[0],
17     x[1][1])).max(lambda x: x[1][0])
18     y = centroide1[1]
19     centroide2 = data.map(lambda x : (1, (computeDistance_improve(y, x[1][: -1]), x[1][: -1] )))
20     .map(lambda x : (x[0], x[1][1])).min(lambda x: x[1])
21     y2 = centroide2[1]
22     point = [x + y for x, y in zip(y, y2)]
23     point = [x / 2 for x in point]
24     centroide = sc.parallelize([centroide1, centroide2, (2, point)])
25     return centroide
26
27 def join_centroide_data(data_indexed, nb_clusters):
28     centroides = centroides_input(data_indexed)
29     centroides = centroides.collect()
30     joined = data_indexed.map(lambda x : (x, centroides))
31     return joined
32
33 def computeDistance_improve(x, y):
34     return sum([(a - b)**2 for a, b in zip(x, y)])
35
36 def min_(joined) :
37     dist_list = joined.map(lambda x : (x[0][0],
38     ((x[1][0][0], computeDistance_improve(x[0][1][: -1], x[1][0][1])),
39     (x[1][1][0], computeDistance_improve(x[0][1][: -1], x[1][1][1])),
40     (x[1][2][0], computeDistance_improve(x[0][1][: -1], x[1][2][1]))
41     )))
```

```

41 min_dist = dist_list.mapValues(closestCluster)
42 return min_dist
43
44 def min_list(array) :
45     min = array[0]
46     for i in range(1, len(array)) :
47         if min > array[i] :
48             min = array[i]
49     return min
50
51 def graphe_display(nb_runs, custom_algo_times) :
52     d = {'times': pd.Series(range(1, nb_runs+1)), 'time': custom_algo_times}
53     df = pd.DataFrame(data=d)
54     display(df)
55
56 def erreur_improve(assignment):
57     list1 = [1,1,2,2,3,3]
58     list2 = [2,3,3,1,2,1]
59     list3 = [3,2,1,3,1,2]
60     print("assignment :", assignment)
61     label = assignment.map(lambda x: x[1][1][4]).distinct().collect()
62     print("label :", label)
63     percentage_error_list = []
64     for a,b,c in zip(list1, list2, list3) :
65         combinaison = [(label[0],a),(label[1],b),(label[2],c)]
66         pourcentage_error = 0
67         combinaison_assignment = assignment.map(lambda x: (x[1][1][4], x[1][0][0])).collect()
68         print("combinaison_assignment :", combinaison_assignment)
69         for i in range(len(combinaison_assignment)) :
70             if combinaison_assignment[i] not in (combinaison) :
71                 pourcentage_error += 1
72         pourcentage_error /= len(combinaison_assignment)
73         percentage_error_list.append(pourcentage_error)
74     return min_list(percentage_error_list)
75
76 def compute_distance_selectedfeatures_nosqrt(x,y, features):
77     sum = 0
78     for f in features:
79         sum += (x[f]-y[f])**2
80     return sum
81
82 def min_selectedfeatures(joined, features) :
83     dist_list = joined.map(lambda x : (x[0][0],
84         ((x[1][0][0], compute_distance_selectedfeatures_nosqrt(x[0][1][: -1], x[1][0][1],
85             features)),
86         (x[1][1][0], compute_distance_selectedfeatures_nosqrt(x[0][1][: -1], x[1][1][1],
87             features)),
88         (x[1][2][0], compute_distance_selectedfeatures_nosqrt(x[0][1][: -1], x[1][2][1],
89             features))
90     )))
91     min_dist = dist_list.mapValues(closestCluster)
92     return min_dist
93
94 # if the centroides are not joined in 1 row
95 def computeDistance_manhattan(x,y):
96     return sum([(a - b) for a,b in zip(x,y)])
97
98 # use manhattan distance with random features and joined centroides

```

```

96 def compute_distance_manhattan_features(x,y, features):
97     sum = 0
98     for f in features:
99         sum += abs(x[f]-y[f])
100     return sum
101
102 #use manhattan distance with all features and joined centroides
103 def compute_distance_manhattan(x,y):
104     return compute_distance_manhattan(x,y,[0,1,2,3])
105
106 # get min of manhattan distances with all features
107 def min_manhattan(joined) :
108     dist_list = joined.map(lambda x : (x[0][0],
109         ((x[1][0][0], compute_distance_manhattan(x[0][1][: -1], x[1][0][1])),
110         (x[1][1][0], compute_distance_manhattan(x[0][1][: -1], x[1][1][1])),
111         (x[1][2][0], compute_distance_manhattan(x[0][1][: -1], x[1][2][1]))
112         )))
113     min_dist = dist_list.mapValues(closestCluster)
114     return min_dist
115
116 # get min of manhattan distances with random features
117 def min_manhattan_selectedfeatures(joined, features) :
118     dist_list = joined.map(lambda x : (x[0][0],
119         ((x[1][0][0], compute_distance_manhattan_features(x[0][1][: -1], x[1][0][1], features)
120         ),
121         (x[1][1][0], compute_distance_manhattan_features(x[0][1][: -1], x[1][1][1], features))
122         ),
123         (x[1][2][0], compute_distance_manhattan_features(x[0][1][: -1], x[1][2][1], features))
124         )))
125     min_dist = dist_list.mapValues(closestCluster)
126     return min_dist
127
128 def computeDistance(x,y):
129     return sqrt(sum([(a - b)**2 for a,b in zip(x,y)]))
130
131 def closestCluster(dist_list):
132     cluster = dist_list[0][0]
133     min_dist = dist_list[0][1]
134     for elem in dist_list:
135         if elem[1] < min_dist:
136             cluster = elem[0]
137             min_dist = elem[1]
138     return (cluster, min_dist)
139
140 def sumList(x,y):
141     return [x[i]+y[i] for i in range(len(x))]
142
143 def moyenneList(x,n):
144     return [x[i]/n for i in range(len(x))]

```

### 7.1.2 Algorithmme initial

```

1 def simpleKmeans(data, nb_clusters):

```

```

2 clusteringDone = False
3 number_of_steps = 0
4 current_error = float("inf")
5 nb_elem = sc.broadcast(data.count())
6
7 #####
8 # Select initial centroids #
9 #####
10
11 centroids = sc.parallelize(data.takeSample('withoutReplacment', nb_clusters))\
12     .zipWithIndex()\
13     .map(lambda x: (x[1], x[0][1][: -1]))
14 print("initial_centroide", centroids.collect())
15 # In the same manner, zipWithIndex gives an id to each cluster
16
17 while not clusteringDone:
18     joined = data.cartesian(centroids)
19
20     # We compute the distance between the points and each cluster
21     dist = joined.map(lambda x: (x[0][0], (x[1][0], computeDistance(x[0][1][: -1], x[1][1])
22     )))
23
24     dist_list = dist.groupByKey().mapValues(list)
25
26     # We keep only the closest cluster to each point.
27     min_dist = dist_list.mapValues(closestCluster)
28
29     # assignment will be our return value : It contains the datapoint,
30     # the id of the closest cluster and the distance of the point to the centroid
31     assignment = min_dist.join(data)
32
33     #####
34     # Compute the new centroid of each cluster #
35     #####
36
37     clusters = assignment.map(lambda x: (x[1][0][0], x[1][1][: -1]))
38
39     count = clusters.map(lambda x: (x[0], 1)).reduceByKey(lambda x, y: x+y)
40     somme = clusters.reduceByKey(sumList)
41     centroidesCluster = somme.join(count).map(lambda x : (x[0], moyenneList(x[1][0], x
42     [1][1])))
43
44     #####
45     # Is the clustering over ? #
46     #####
47
48     # Let's see how many points have switched clusters.
49     if number_of_steps > 0:
50         switch = prev_assignment.join(min_dist)\
51             .filter(lambda x: x[1][0][0] != x[1][1][0])\
52             .count()
53     else:
54         switch = 150
55     if switch == 0 or number_of_steps == 100:
56         clusteringDone = True
57         error_improve = erreur_improve(assignment) #adding to compare
58         error = sqrt(min_dist.map(lambda x: x[1][1]).reduce(lambda x, y: x + y))/nb_elem.
59         value
60     else:

```

```

57     centroides = centroidesCluster
58     prev_assignment = min_dist
59     if number_of_steps % 10 == 0 :
60         print("number_of_steps :", number_of_steps)
61         number_of_steps += 1
62
63     return (assignment, error, number_of_steps, centroides, error_improve)

```

### 7.1.3 Algorithme optimisé

```

1  def custom_kmeans(data, nb_clusters=3, join_centroids_inrdd=False,
2      choose_initial_centroid=False, random_features=False, epsilon_change_nbfeatures=0,
3      with_regular_check=False, interval_regular_check=10, distance="euclidian"):
4      """
5      This function executes the kmeans algorithm with combinations of the different
6      implemented improvements.
7
8      Parameters :
9      * data
10     * nb_clusters
11     * join_centroids_inrdd (bool) : if set to True, at each iteration, one rdd line
12     contains a point and the 3 centroids (and no longer 3 lines for 3 centroids)
13     * choose_initial_centroid (bool) : if set to True, (initialize according to kmeans+
14     algo) the initial centroids are chosen with our improvement (1 : max, 2 : furthest
15     from 1, 3 : middle between 1 & 2)
16     * random_features (bool) : if set to True, the algorithm will choose a subset of the
17     features to compute the distance (2 to start, then 3, then 4)
18     * epsilon_change_nbfeatures (int) : if random_features is set to True, this parameter
19     tells when to increase the number of features to select
20     * with_regular_check (bool) : if set to True and with random_features, performs a
21     step with all features regularly
22     * interval_regular_check (int) : sets the number of iterations between two full
23     feature checks
24     * distance (string) : the name of the distance to use (supported yet are "euclidian"
25     and "manhattan")
26
27     Returns :
28     (assignment, error, number_of_steps)
29
30     """
31     clusteringDone = False
32     number_of_steps = 0
33     current_error = float("inf")
34     nb_elem = sc.broadcast(data.count())
35     nb_features_selected = 2
36
37     #####
38     # Select initial centroides #
39     #####
40     if choose_initial_centroid == True :
41         centroides = centroides_input(data)
42         print("initial_centroide", centroides.collect())
43     else :

```

```

33 centroids = sc.parallelize(data.takeSample('withoutReplacment',nb_clusters))\
34     .zipWithIndex()\
35     .map(lambda x: (x[1],x[0][1][: -1]))
36
37 while not clusteringDone:
38
39     #####
40     # Assign points to clusters #
41     #####
42     if join_centroids_inrdd == True :
43         joined = join_centroide_data(data,3)
44         if random_features == True :
45             features = []
46             if with_regular_check == True and (number_of_steps%interval_regular_check == 0)
47 :
48             features = [0,1,2,3]
49         else:
50             # start by checking if we need to change the number of features
51             if (number_of_steps > 0) and switch < epsilon_change_nbfeatures and
52 switch_prev < epsilon_change_nbfeatures and nb_features_selected < 4:
53                 nb_features_selected += 1
54                 random.seed(number_of_steps)
55                 features = random.sample([0,1,2,3], k=nb_features_selected) #choose randomly
56 the features to observe
57
58         if (distance == "manhattan"):
59             min_dist = min_manhattan_selectedfeatures(joined , features)
60         else :
61             min_dist = min_selectedfeatures(joined , features) # specify the features
62 selected
63
64         else:
65             if (distance == "manhattan"):
66                 min_dist = min_manhattan(joined)
67             else :
68                 min_dist = min_(joined)
69
70         assignment = min_dist.join(data)
71
72     else :
73         joined = data.cartesian(centroids)
74
75         # We compute the distance between the points and each cluster
76         if (distance == "manhattan"):
77             dist = joined.map(lambda x: (x[0][0],(x[1][0], computeDistance_manhattan(x
78 [0][1][: -1], x[1][1]))))
79         else:
80             dist = joined.map(lambda x: (x[0][0],(x[1][0], computeDistance(x[0][1][: -1], x
81 [1][1]))))
82
83         dist_list = dist.groupByKey().mapValues(list)
84
85         # We keep only the closest cluster to each point.
86         min_dist = dist_list.mapValues(closestCluster)
87
88         # assignment will be our return value : It contains the datapoint,
89         # the id of the closest cluster and the distance of the point to the centroid

```



```

85     assignment = min_dist.join(data)
86
87     #####
88     # Compute the new centroid of each cluster #
89     #####
90
91     clusters = assignment.map(lambda x: (x[1][0][0], x[1][1][: -1]))
92
93     count = clusters.map(lambda x: (x[0],1)).reduceByKey(lambda x,y: x+y)
94     #count = (1 , 100),(2 , 200), (3, 50)
95     #      (numero de cluster, nb d'element)
96     somme = clusters.reduceByKey(sumList)
97     # somme = (1,sum of rows du clust 1), (2, 6000) , (3, 1500)
98     #      (numero de cluster, sommedes elements )
99     centroidesCluster = somme.join(count).map(lambda x : (x[0],moyenneList(x[1][0],x
100 [1][1])))
101     # (numero de cluster, average)
102     #####
103     # Is the clustering over ? #
104     #####
105
106     # Let's see how many points have switched clusters.
107     if number_of_steps > 0:
108         switch_prev = switch
109         switch = prev_assignment.join(min_dist)\
110             .filter(lambda x: x[1][0][0] != x[1][1][0])\
111             .count()
112
113     else:
114         switch = 150
115         switch_prev = switch
116     if (switch == 0 and nb_features_selected == 4) or number_of_steps > 3:
117         clusteringDone = True
118         error_improve = erreur_improve(assignment)
119         error = sqrt(min_dist.map(lambda x: x[1][1]).reduce(lambda x,y: x + y)/nb_elem
120 .value
121     else:
122         if nb_features_selected != 4 and switch == 0:
123             nb_features_selected += 1
124             centroides = centroidesCluster
125             prev_assignment = min_dist
126             if number_of_steps % 10 == 0 :
127                 print("number_of_steps :", number_of_steps)
128                 number_of_steps += 1
129
130     return (assignment, error, number_of_steps, centroides, error_improve)

```

#### 7.1.4 Main

```

1 def main_function(path=' /FileStore/tables/iris_data-04fff.txt',nb_runs=3, teacher_runs =
2     True, student_runs=True):
3     """
4     Global function that includes all the improvements.

```

```

4
5 Parameters :
6     path (string) : path to the data file
7     nb_runs (int) : number of runs to compute the mean for the time spent
8 Returns :
9     avg_custom (float) : the average of execution time for the custom algo
10    avg_basic (float) : the average of execution time for the original algo
11    """
12
13    # include parameters in case we use different data : nb of clusters ?, nb of features ?
14
15    startTimeQuery = time.clock()
16    lines = sc.textFile(path)
17    data = lines.map(lambda x: x.split(',')\
18                    .map(lambda x: [float(i) for i in x[:4]]+[x[4]])\
19                    .zipWithIndex()\
20                    .map(lambda x: (x[1],x[0])))
21    custom_algo_times = []
22    basic_algo_times = []
23    if student_runs == True :
24        for i in range(nb_runs) :
25            print("")
26            start =strftime("%H:%M:%S", gmtime())
27            clustering = custom_kmeans(data, nb_clusters=3, join_centroids_inrdd=True,
28            choose_initial_centroid=True, random_features=True, epsilon_change_nbfeatures=5,
29            with_regular_check=True, interval_regular_check=10,distance="manhattan") # specify
30            here the improvements to do in the parameters
31            print("CUSTOM :", "pourcentage_miss_labelised :",clustering[4], "error :",
32            clustering[1], " clustering_nbr_loop :", clustering[2], " clusters :", clustering
33            [3].collect())
34            custom_algo_times.append(seconds_between(start , strftime("%H:%M:%S", gmtime()))))
35
36    if teacher_runs == True :
37        for i in range(int(nb_runs/10)) :
38            #compare to the teacher's function
39            start =strftime("%H:%M:%S", gmtime())
40            clustering = simpleKmeans(data, 3)
41            print("")
42            print("TEACHER :", "pourcentage_miss_labelised :",clustering[4], "error :",
43            clustering[1], " clustering_nbr_loop :", clustering[2], " clusters :", clustering
44            [3].collect())
45            basic_algo_times.append(seconds_between(start , strftime("%H:%M:%S", gmtime()))))
46
47    avg_custom = sum(custom_algo_times)/nb_runs
48    avg_basic = 0
49    print("Nb of runs = ",nb_runs)
50    if(teacher_runs == True):
51        avg_basic = sum(basic_algo_times)/int(nb_runs/10)
52        print("Average of execution time for the original algo :",avg_basic)
53        print("Execution times for the initial algo :",basic_algo_times)
54        print("Average of execution time for the custom algo :",avg_custom)
55        print("Execution times for the custom algo :",custom_algo_times)
56        return avg_custom ,avg_basic ,custom_algo_times ,basic_algo_times
57
58    # _____
59    # MAIN
60    # _____
61
62

```

```

55 #IRIS DATASET STUDENT ALGO
56
57 if __name__ == "__main__":
58     path = '/FileStore/tables/iris_data-04fff.txt'
59     nb_runs = 5
60     data = sc.textFile(path).map(lambda x: x.split(',')).map(lambda x: [float(i) for i in x[:4]]+[x[4]])
61     avg_custom, avg_basic, custom_algo_times, basic_algo_times = main_function(path=path,
62                                     nb_runs=nb_runs, student_runs = True, teacher_runs = False)
63
64 # -----
65 #             MAIN — VOITURES
66 # -----
67 #VOITURES DATASET STUDENT ALGO
68
69 if __name__ == "__main__":
70     path = '/FileStore/tables/voiture_type.txt'
71     nb_runs = 5
72     data = sc.textFile(path).map(lambda x: x.split(',')).map(lambda x: [float(i) for i in x[:4]]+[x[4]])
73     avg_custom, avg_basic, custom_algo_times, basic_algo_times = main_function(path=path,
74                                     nb_runs=nb_runs, student_runs = True, teacher_runs = False)

```

## 7.2 Dataframe

### 7.2.1 Fonctions

```

1 from pyspark.sql.functions import udf
2 from pyspark.sql import SQLContext
3 from pyspark.sql.functions import row_number, monotonically_increasing_id
4 from pyspark.sql import Window
5 import pyspark.sql.functions as F
6 import numpy as np
7
8 def centroide_input_df(df):
9     df = df.withColumn('total', sum(df[col] for col in ["_c0", "_c1", "_c2", "_c3"]))
10    row1 = df.agg({"total": "max"}).collect()[0]
11    max_tot = row1["max(total)"]
12    max_point = df.where(F.col('total') == max_tot).select('_c0', '_c1', '_c2', '_c3')
13    list_max_point_coordonees = list(max_point.collect()[0][0:4])
14    row2 = df.agg({"total": "min"}).collect()[0]
15    min_tot = row2["min(total)"]
16    min_point = df.where(F.col('total') == min_tot).select('_c0', '_c1', '_c2', '_c3')
17    list_min_point_coordonees = list(min_point.collect()[0][0:4])
18    list_moy_point_coordonees = [(x+y)/2 for x,y in zip(list_min_point_coordonees,
19                                                         list_max_point_coordonees)]
19    return list_min_point_coordonees, list_moy_point_coordonees, list_max_point_coordonees
20
21
22 """
23 Création de la première dataframe

```

```

24 """
25
26 def get_dataframe(path):
27     sqlContext = SQLContext(sc)
28     df = sqlContext.read.format('com.databricks.spark.csv').options(header='false',
29                               inferSchema='true')\
30         .load(path).withColumnRenamed('_c4','label')
31
32     df = df.withColumn("index", row_number().over(Window.orderBy(
33         monotonically_increasing_id()))-1)
34     df = df.withColumn("coordonnees", F.array(F.col("_c0"), F.col("_c1"), F.col("_c2"), F.col(
35         "_c3")))
36     return df
37
38 @udf("float")
39 def euclidian_distance(x,y):
40     d = 0
41     for i in range(len(x)):
42         d += (x[i] - y[i])**2
43     return d
44
45 @udf("integer")
46 def min_index(d1,d2,d3):
47     if(d1 < d2 and d1 < d3):
48         return 1
49     if(d2 < d1 and d2 < d3):
50         return 2
51     return 3
52
53 def df_compute_min_distances(df,c1,c2,c3): # df contains id of the observation + features
54     of the observation + features of each centroid as a list
55     df = df.withColumn("c1",F.array([F.lit(i) for i in c1]))\
56         .withColumn("c2",F.array([F.lit(i) for i in c2]))\
57         .withColumn("c3",F.array([F.lit(i) for i in c3]))
58     distances = df.withColumn("dist1",euclidian_distance(F.col("coordonnees"),F.col("c1")))\
59         \
60         .withColumn("dist2",euclidian_distance(F.col("coordonnees"),F.col("c2")))\
61         \
62         .withColumn("dist3",euclidian_distance(F.col("coordonnees"),F.col("c3")))\
63         \
64         .drop("c1").drop("c2").drop("c3")
65     min_df = distances.withColumn("cluster",min_index(F.col("dist1"),F.col("dist2"),F.col("
66         dist3")))
67     return min_df
68
69 @udf("array<float>")
70 def average_list(list_of_list):
71     center = [0,0,0,0]
72     nb = len(list_of_list)
73     for i in range(nb):
74         for j in range(0,4):
75             center[j] += list_of_list[i][j]/nb
76     return center
77
78 def df_update_centroids(df):
79     df_gp = df.groupBy("cluster").agg(F.collect_list(F.col("coordonnees")).alias("
80         list_of_coords"))
81     return df_gp.withColumn("avg_coords",average_list(F.col("list_of_coords"))).drop("

```

```

    list_of_coords")
73
74 def df_get_new_centroids(df):
75     array = []
76     for i in range(1,4):
77         array.append(df.filter(F.col("cluster")==i).collect()[0]["avg_coords"])
78     return array[0],array[1],array[2]
79
80 def df_get_nb_switch(df_prev,df_suiv):
81     join = df_prev.withColumnRenamed("cluster","cluster_prev").join(df_suiv,"index").filter
82     (F.col("cluster_prev")!=F.col("cluster"))
83     return join.count()
84
85 def df_compute_error(df, labels):
86     list1 = [1,1,2,2,3,3]
87     list2 = [2,3,3,1,2,1]
88     list3 = [3,2,1,3,1,2]
89     percentage_error_list = []
90     for a,b,c in zip(list1,list2,list3):
91         combinaison = [(labels[0][0],a),(labels[1][0],b),(labels[2][0],c)]
92         pourcentage_error = 0
93         combinaison_assignment = df.select("label","cluster").collect()
94         for i in combinaison_assignment:
95             i = (i[0], i[1])
96             if i not in combinaison:
97                 pourcentage_error += 1
98             pourcentage_error /= len(combinaison_assignment)
99     percentage_error_list.append(pourcentage_error)
100     return min_list(percentage_error_list)

```

## 7.2.2 Algorithmme

```

1 def df_kmeans(data,nb_steps_max):
2     clusteringDone = False
3     number_of_steps = 0
4     current_error = float("inf")
5     nb_elem = sc.broadcast(data.count())
6     switch = -1
7     centroide1,centroide2,centroide3 = centroide_input_df(data)
8     data = data.drop("_c0").drop("_c1").drop("_c2").drop("_c3").drop('total')
9     while(not clusteringDone and number_of_steps < nb_steps_max):
10         min_dist_df = df_compute_min_distances(data,centroide1,centroide2,centroide3)
11         new_centroids_df = df_update_centroids(min_dist_df)
12         centroide1,centroide2,centroide3 = df_get_new_centroids(new_centroids_df)
13         if(number_of_steps > 0):
14             switch = df_get_nb_switch(min_dist_df_prev,min_dist_df)
15             min_dist_df_prev = min_dist_df
16             number_of_steps += 1
17
18         if(switch == 0):
19             clusteringDone = True
20             labels = data.select("label").distinct().collect()
21             error = df_compute_error(min_dist_df, labels)

```

```

22     if number_of_steps %10 ==0 :
23         print("number_of_steps :", number_of_steps)
24     assignment = min_dist_df.drop("dist1").drop("dist2").drop("dist3")
25     centroides = [centroide1,centroide2,centroide3]
26     return (assignment, error, number_of_steps, centroides)

```

### 7.2.3 Main

```

1  def df_main(path='/FileStore/tables/iris_data-04fff.txt',nb_runs=3):
2      """
3      Global function that includes all the improvements.
4
5      Parameters :
6          path (string) : path to the data file
7          nb_runs (int) : number of runs to compute the mean for the time spent
8      Returns :
9          avg_df (float) : the average of execution time for the kmeans dataframe algo
10     """
11
12     startTimeQuery = time.clock()
13     path = '/FileStore/tables/voiture_type.txt'
14
15     data = get_dataframe(path)
16     df_algo_times = []
17     for i in range(nb_runs) :
18         print("")
19         start =strftime ("%H:%M:%S", gmtime())
20
21         clustering = df_kmeans(data,150)
22         print("DATAFRAME :", "error :", clustering[1], " clustering_nbr_loop :",
23             clustering[2], " clusters :", clustering[3])
24         df_algo_times.append(seconds_between(start, strftime ("%H:%M:%S", gmtime())))
25
26     avg_df = sum(df_algo_times)/nb_runs
27     print("Nb of runs = ",nb_runs)
28
29     print("Average of execution time for the dataframe algo :",avg_df)
30     print("Execution times for the dataframe algo :",df_algo_times)
31     return avg_df,df_algo_times
32
33 # -----
34 #             MAIN
35 # -----
36
37 #IRIS DATASET STUDENT ALGO
38
39 if __name__ == "__main__":
40     path = '/FileStore/tables/iris_data-04fff.txt'
41     nb_runs = 5
42     df_main(path,nb_runs)
43
44 # -----
45 #             MAIN — VOITURES

```

```

45 #
46
47 #VOITURES DATASET STUDENT ALGO
48
49 if __name__ == "__main__":
50     path = '/FileStore/tables/voiture_type.txt'
51     nb_runs = 5
52     df_main(path, nb_runs)

```

### 7.3 Jeu de test pour les dataframes

```

1  clust1 = [0,0,0,0]
2  clust2 = [1,1,1,1]
3  clust3 = [2,2,2,2]
4
5  from pyspark.sql.types import *
6
7  cSchema = StructType([ StructField("index", IntegerType())\
8                          ,StructField("x1",FloatType())\
9                          ,StructField("x2",FloatType())\
10                         ,StructField("x3",FloatType())\
11                         ,StructField("x4",FloatType())])
12
13  test_list = [[1,0.0,0.0,0.0,0.0],[2,0.0,0.0,0.0,0.0],[3,1.0,1.0,1.0,1.0],
14               [4,2.0,2.0,2.0,2.0],[5,1.2,1.2,1.2,1.2]]
15
16  df = spark.createDataFrame(test_list,schema=cSchema)
17  df = df.select(F.col("index"),F.array(F.col("x1"),F.col("x2"),F.col("x3"),F.col("x4")).
18              alias("coordonnees"))
19
20  df.show()
21  dis = df_compute_min_distances(df,clust1,clust2,clust3)
22  dis.show()
23  gp_df = df_update_centroids(dis)
24  gp_df.show()
25  print("NEW CENTROIDS = ",df_get_new_centroids(gp_df))
26  print("NB switch = ",df_get_nb_switch(dis,dis))

```

#### **7.4 Tableau des résultats**



Expérimentation	Numéro	Requête	Conditions					
			NB CLUSTERS	join_centroids_in_rdd	choose_initial_centroid	random_features	epsilon_change_nb_features	with_regular_check
Algorithme initial	0	simplekmeans(data, nb_clusters)	3	X		X	X	X
	1	custom_kmeans(data,nb_clusters=3,join_centroids_in_rdd=False,choose_initial_centroids=False,random_features=False,epsilon_change_nb_features=5,with_regular_check=False,interval_regular_check=10,distance=euclidian)		False				
	2	custom_kmeans(data,nb_clusters=3,join_centroids_in_rdd=True,choose_initial_centroids=False,random_features=False,epsilon_change_nb_features=5,with_regular_check=False,interval_regular_check=10,distance=euclidian)		True				
	3	custom_kmeans(data,nb_clusters=3,join_centroids_in_rdd=False,choose_initial_centroids=True,random_features=False,epsilon_change_nb_features=5,with_regular_check=False,interval_regular_check=10,distance=euclidian)		False		False		False
	4	custom_kmeans(data,nb_clusters=3,join_centroids_in_rdd=True,choose_initial_centroids=True,random_features=False,epsilon_change_nb_features=5,with_regular_check=False,interval_regular_check=10,distance=euclidian)			True			
	5	custom_kmeans(data,nb_clusters=3,join_centroids_in_rdd=True,choose_initial_centroids=False,random_features=True,epsilon_change_nb_features=5,with_regular_check=False,interval_regular_check=10,distance=euclidian)	3				5	10
	6	custom_kmeans(data,nb_clusters=3,join_centroids_in_rdd=True,choose_initial_centroids=False,random_features=True,epsilon_change_nb_features=5,with_regular_check=True,interval_regular_check=10,distance=euclidian)			False			True
	7	custom_kmeans(data,nb_clusters=3,join_centroids_in_rdd=True,choose_initial_centroids=True,random_features=True,epsilon_change_nb_features=5,with_regular_check=False,interval_regular_check=10,distance=euclidian)		True		True		False
	8	custom_kmeans(data,nb_clusters=3,join_centroids_in_rdd=True,choose_initial_centroids=True,random_features=True,epsilon_change_nb_features=5,with_regular_check=True,interval_regular_check=10,distance=euclidian)			True			True
Algorithme personnalisé	9	custom_kmeans(data,nb_clusters=3,join_centroids_in_rdd=True,choose_initial_centroids=True,random_features=True,epsilon_change_nb_features=5,with_regular_check=True,interval_regular_check=10,distance=manhattan)	3	True	True	False	X	False
								X

FIGURE 1: Description des expériences

		Temps moyen sur 5 essais						Temps moyen sur 5 essais					
Numéro	Format	1	2	3	4	5	Moyenne	1	2	3	4	5	Moyenne
0	RDD	Iris						Vehicule_types					
1	RDD	Executor heartbeat timed out after 136447 ms					400	x	x	x	x	x	Aucun résultat après 2h ..
2	RDD	5	6	6	6	5	5,6	9	10	10	9	9	9,4
3	RDD	170	172	170	170	171	170,6	x	x	x	x	x	Aucun résultat après 2h ..
4	RDD	5	4	5	4	5	4,6	6	5	5	5	5	5,2
5	RDD	23	22	25	27	22	23,8	21	22	25	21	23	22,4
6	RDD	27	27	28	27	27	27,2	34	34	33	35	37	34,6
7	RDD	18	18	18	17	18	17,8	20	20	20	19	20	19,8
8	RDD	26	25	26	25	25	25,4	33	35	33	33	34	33,6
9	RDD	29	30	30	31	30	30	33	33	34	33	35	33,6
	Dafaframe							81	85	87			84,33333333

FIGURE 2: Résultats obtenus en secondes