# K-Means clustering algorithm in MapReduce

Assignment 2

Panourgia Evangelia (t8190130)
Papadatos Iwannis (t8190314)
Professor: Damianos Chatziantwniou

Latest Update: April 9, 2022

School of Management Science and Technology,
Athens University of Economics and Business

# Contents

# 1   Hadoop Installation

As far as the installation of Hadoop, we decided to download Hadoop packaged by Bitnami [1] and Oracle VirtualBox[2]. When the package for Hadoop was downloaded, we import it to Oracle VirtualBox and run it. After, these steps you can see in your pc the image shows in fig. 1. We note that as default username, password is the word bitname(look red letters). Now, you should set the scene by installing git[2] (command : sudo apt-get install git), pip3(sudo apt install python3-pip)[3] and mrjob(pip3 install mrjob)[4]. So, now you have the necessary tools - software so as to run the code.



Figure 1: Install Bitnami

# 2  Dataset Creation

The file dataGenerate.py was written for generating data - points in the form (x, y). Specifically, it reads data from file centroids.csv which contains three centroids in the form (x, y) and then generate 8000 data points using python library skewnorm[5]. Last but not least, user can see the the generated data graphically by entering python3 generateData.py -d, ehere -d parameters call suitable function to draw[6] data points(see fig 2.).
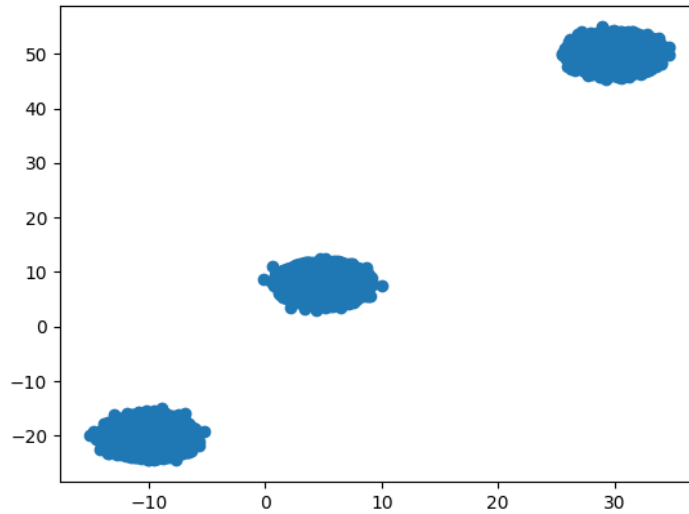


Figure 2: Dataset Creation

Now, we can see the code for data generation.

# 3  K-Means Clustering Algorithm

The file kmeans.py was written to implement the algorithm kmeans[6] using MapReduce. For this aim, we used the python library mrjob. Firstly, we defined configuration and steps. Then, we write the function for mapper and reducer. For mapper, we calculate euclidean distance and for each point (x, y) we calculate all distances from all centroids and we hold the minimum distance, so we assign point to suitable centroid. Last but not least, the only work of reducer is to revise cluster centers as mean of assigned observations.

Now, we can see the code for kmeans implementation.

```python
1   from math import sqrt
2   from mrjob.job import MRJob
3   from mrjob.step import MRStep
4
5   class KMeans(MRJob):
6
7       def configure_args(self):
8           super(KMeans, self).configure_args()
9
10          self.add_file_arg(
11              '--centroids-file',
12              dest='centroids_file',
13              help='path to the file containing the centroids.'
14          )
15
16      def steps(self):
17          return[
18              MRStep(mapper_init=self.load_centroids,
19                     mapper=self.mapper,
20                     reducer=self.reducer)
21          ]
22
23      # This method is executed before mappers process any input.
24      def load_centroids(self):
25          self.__centroids = []
26
27          with open(self.options.centroids_file) as centroids_file:
28              for line in centroids_file:
29                  x, y = line.strip().split(',')
30                  centroid = (float(x), float(y))
31                  self.__centroids.append(centroid)
```

Figure 3: kmeans implementation

```
33      def __calculate_euclidean_dist(self, point, centroid):
34          x1, y1 = point
35          x2, y2 = centroid
36          return sqrt((x2 - x1)**2 + (y2 - y1)**2)
37
38      # The line will be a raw line of the input file, with newline (\n) stripped.
39      def mapper(self, _, line):
40          x, y = line.split(',')
41          point = (float(x), float(y))
42
43          min_euclidean_dist = float('inf')
44          closest_centroid = None
45          for centroid in self.__centroids:
46              euclidean_dist = self.__calculate_euclidean_dist(point, centroid)
47              if euclidean_dist < min_euclidean_dist:
48                  min_euclidean_dist = euclidean_dist
49                  closest_centroid = centroid
50
51          yield closest_centroid, point
52
53      def reducer(self, centroid, points):
54          n, sum_x, sum_y = 0, 0, 0
55          for x, y in points:
56              sum_x += x
57              sum_y += y
58              n += 1
59          mean_x = sum_x / n
60          mean_y = sum_y / n
61
62          yield centroid, (mean_x, mean_y)
63
64  if __name__ == "__main__":
65      KMeans.run()
```

Figure 4: kmeans implementation

## 3.1  HDFS

## 3.2  Run

So, git clone + commands

# References

[1] Hadoop packaged by Bitnami, *https://bitnami.com/stack/hadoop/virtual-machine*

[2] Oracle VirtualBox, *https://www.virtualbox.org/*

[3] Install pip3, *https://linuxize.com/post/how-to-install-pip-on-ubuntu-18.04/*

[4] Install mrjob, *https://pypi.org/project/mrjob/*

[5] skewnorm, *https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.skewnorm.html*

[6] matplotlib, *https://matplotlib.org/*