# k-means algorithm implementation on Hadoop

Efstratios Gounidellis
stratos.gounidellis [at] gmail.com


Lamprini Koutsokera
lkoutsokera [at] gmail.com

Course: "Big Data Management Systems"
Professor: Damianos Chatziantoniou


Department of Management Science & Technology


School of Business
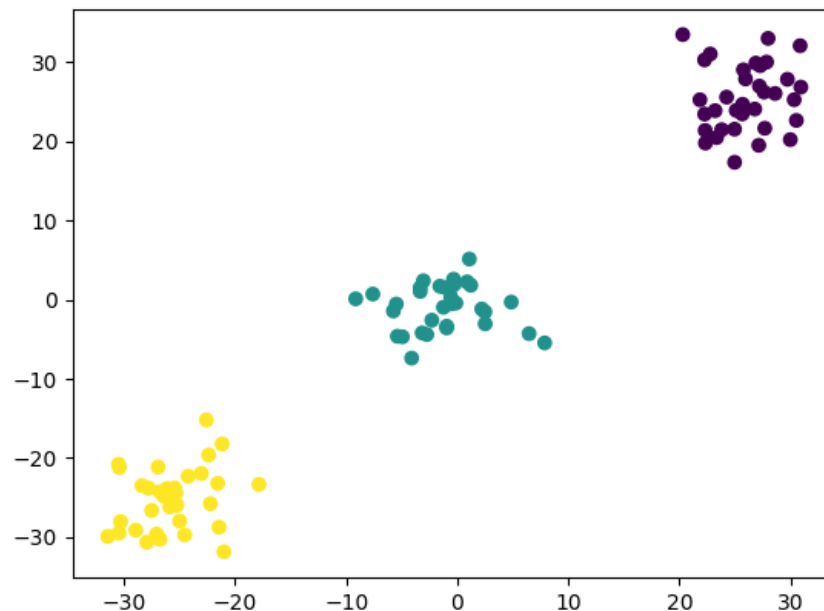Athens University of Economics & Business




March 31, 2017

# Contents

# 1 Data points generation

## 1.1 createDataPoints.py

The initial task of the project is to generate a set of more than one million data points to be used later as input for the k-means clustering algorithm. Using this python script three isotropic Gaussian blobs for clustering are generated. More specifically, the centers are the following data points [25, 25], [-1, -1], [-25, -25]. Additionally, the data points are presented visually with the use of a scatter plot.



```python
1  """createDataPoints.py: Generate data points for clustering."""
2
3  import argparse
4  import matplotlib.pyplot as plt
5  import os
6  import pandas as pd
7  from sklearn.datasets.samples_generator import make_blobs
8
9  __author__ = "Stratos Gounidellis, Lamprini Koutsokera"
10 __copyright__ = "Copyright 2017, BDSMasters"
11
12
13 class DataGenerator():
14
15     def generateData(self, points, dataFile):
16         """Generate the input data points.
17
18         :param self: An instance of the class DataGenerator.
19         :param points: The number of data points to be generated.
```

```
20        :param dataFile: The file to save the data points.
21        """
22        centers = [[25, 25], [-1, -1], [-25, -25]]
23        X, labels_true = make_blobs(n_samples=long(points),
24                                    centers=centers, cluster_std=3.5,
25                                    n_features=2)
26
27        df = pd.DataFrame(X)
28        df.to_csv(dataFile, header=False, index=False, sep=" ")
29
30        plt.scatter(X[:, 0], X[:, 1], c=labels_true)
31        directory = "../images"
32        if not os.path.isdir(directory):
33            os.makedirs(directory)
34        plt.savefig("../images/data_points.png")
35
36
37 if __name__ == "__main__":
38     parser = argparse
39     parser = argparse.ArgumentParser()
40     parser.add_argument("dataFile", type=str,
41                         help="File to save the generated data points.")
42
43     parser.add_argument("points", type=int,
44                         help="Number of data points to create.")
45     args = parser.parse_args()
46     instanceDataGenerator = DataGenerator()
47     instanceDataGenerator.generateData(args.points, args.dataFile)
```
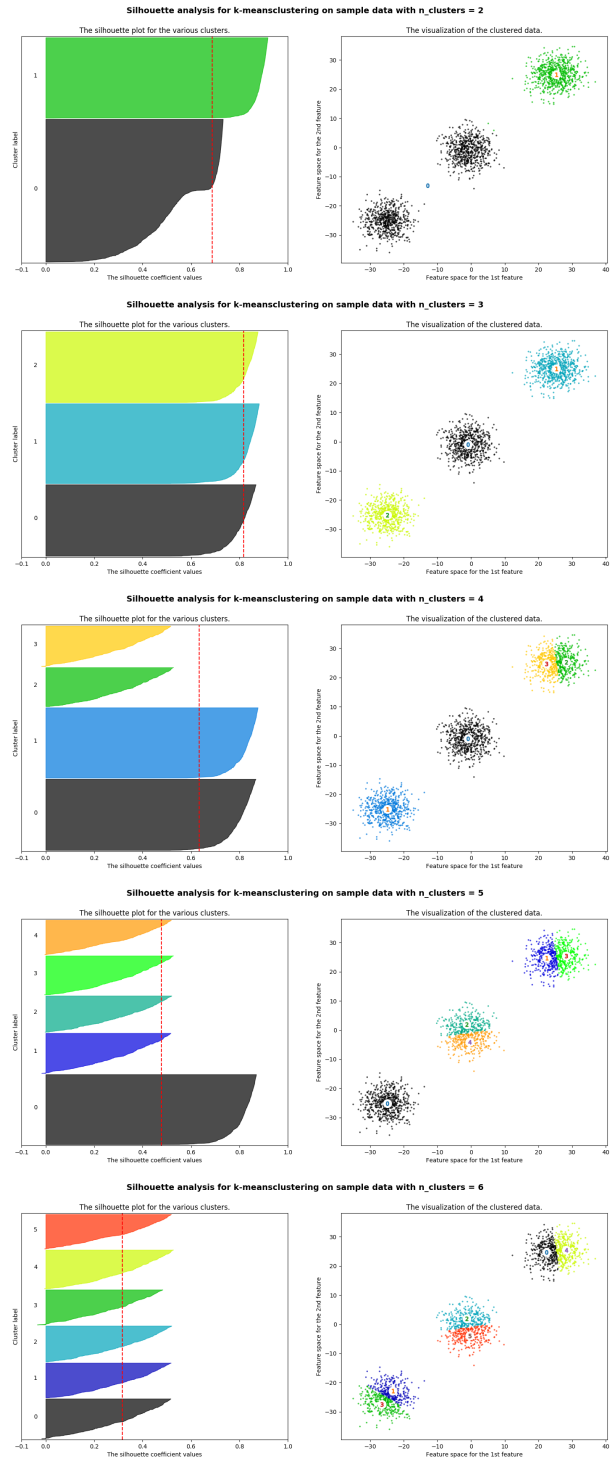
## 2  Number of clusters

### 2.1  plotSilhouetteScore.py

The silhouette score constitutes a useful criterion for determining the proper number of clusters and it was firstly suggested by Peter J. Rousseeuw. The silhouette shows which objects lie well within their cluster, and which ones are merely somewhere in between clusters. A silhouette close to 1 implies the datum is in an appropriate cluster, while a silhouette close to $-1$ implies the datum is in the wrong cluster.

The following python script calculates the silhouette score for different numbers of clusters ranging from 2 to 6. With this script not only the average silhouette score of each cluster is visualized but also the thickness (i.e. the number of data points) of each cluster. The number of clusters which leads to clusters of more or less similar thickness and silhouette score above the average could be the optimal number of clusters for the k-means algorithm.

As expected creating three clusters is the optimal solution in this case.

Silhouette analysis for k-meansclustering on sample data with n_clusters = 2

Silhouette analysis for k-meansclustering on sample data with n_clusters = 3

Silhouette analysis for k-meansclustering on sample data with n_clusters = 4

Silhouette analysis for k-meansclustering on sample data with n_clusters = 5

Silhouette analysis for k-meansclustering on sample data with n_clusters = 6

```
1  """
2  plotSilhouetteScore.py: Selecting the number of clusters with
3                  silhouette analysis on k-means clustering.
4
5  Silhouette analysis can be used to study the separation distance between the
6  resulting clusters. The silhouette plot displays a measure of how close each
7  point in one cluster is to points in the neighboring clusters and thus
```

```
       provides
 8 a way to assess parameters like number of clusters visually. This measure has
       a
 9 range of [-1, 1].
10
11 Silhouette coefficients (as these values are referred to as) near +1 indicate
12 that the sample is far away from the neighboring clusters. A value of 0
13 indicates that the sample is on or very close to the decision boundary between
14 two neighboring clusters and negative values indicate that those samples might
15 have been assigned to the wrong cluster.
16
17 Source:
18 http://scikit-learn.org/stable/auto_examples/cluster/
       plot_kmeans_silhouette_analysis.html
19
20 """
21
22 import argparse
23 from kmeans import KmeansRunner
24 import matplotlib.cm as cm
25 import matplotlib.pyplot as plt
26 import numpy as np
27 import os
28 import PIL
29 from sklearn.cluster import KMeans
30 from sklearn.metrics import silhouette_samples, silhouette_score
31
32 __author__ = "Scikit-Learn"
33
34
35 class SilhouetteScore():
36
37     def calculateSilhouetteScore(self, dataFile):
38         """Calculate the silhouette score for different numbers of clusters.
39
40         :param self: An instance of the class SilhouetteScore.
41         :param dataFile: An array with the input data points.
42         :return: A list with the names of the image files created.
43         """
44         instanceKmeans = KmeansRunner()
45         X = instanceKmeans.retrieveData(dataFile)
46         if (X.shape[0] > 10000):
47             size = round(X.shape[0] * 0.001)
48             idx = np.random.randint(X.shape[0], size=size)
49             subset = X[idx, :]
50             X = subset
51         range_n_clusters = [2, 3, 4, 5, 6]
52         list_images = []
53
54         for n_clusters in range_n_clusters:
55
56             fig, (ax1, ax2) = plt.subplots(1, 2)
57             fig.set_size_inches(18, 7)
58
```

```python
59            ax1.set_xlim([-0.1, 1])

60

61            ax1.set_ylim([0, len(X) + (n_clusters + 1) * 10])

62

63            clusterer = KMeans(n_clusters=n_clusters, random_state=10)
64            cluster_labels = clusterer.fit_predict(np.array(X))

65

66            silhouette_avg = silhouette_score(X, cluster_labels)
67            print("For n_clusters =", n_clusters,
68                  "The average silhouette_score is :", silhouette_avg)

69

70            sample_silhouette_values = silhouette_samples(X, cluster_labels)

71

72            y_lower = 10
73            for i in range(n_clusters):

74

75                ith_cluster_silhouette_values = \
76                    sample_silhouette_values[cluster_labels == i]

77

78                ith_cluster_silhouette_values.sort()

79

80                size_cluster_i = ith_cluster_silhouette_values.shape[0]
81                y_upper = y_lower + size_cluster_i

82

83                color = cm.spectral(float(i) / n_clusters)
84                ax1.fill_betweenx(np.arange(y_lower, y_upper),
85                                  0, ith_cluster_silhouette_values,
86                                  facecolor=color, edgecolor=color, alpha=0.7)

87

88                ax1.text(-0.05, y_lower + 0.5 * size_cluster_i, str(i))

89

90                y_lower = y_upper + 10

91

92            ax1.set_title("The silhouette plot for the various clusters.")
93            ax1.set_xlabel("The silhouette coefficient values")
94            ax1.set_ylabel("Cluster label")

95

96            ax1.axvline(x=silhouette_avg, color="red", linestyle="--")

97

98            ax1.set_yticks([])
99            ax1.set_xticks([-0.1, 0, 0.2, 0.4, 0.6, 0.8, 1])

100

101           colors = cm.spectral(cluster_labels.astype(float) / n_clusters)
102           ax2.scatter(X[:, 0], X[:, 1], marker=".", s=30, lw=0, alpha=0.7,
103                       c=colors)

104

105           centers = clusterer.cluster_centers_
106           ax2.scatter(centers[:, 0], centers[:, 1],
107                       marker="o", c="white", alpha=1, s=200)

108

109           for i, c in enumerate(centers):
110               ax2.scatter(c[0], c[1], marker="$%d$" % i, alpha=1, s=50)

111

112           ax2.set_title("The visualization of the clustered data.")
```

```
113            ax2.set_xlabel("Feature space for the 1st feature")
114            ax2.set_ylabel("Feature space for the 2nd feature")
115
116            plt.suptitle(("Silhouette analysis for k-means"
117                          "clustering on sample data "
118                          "with n_clusters = %d" % n_clusters),
119                         fontsize=14, fontweight="bold")
120            fig.savefig("cluster_" + str(n_clusters) + ".png")
121            list_images.append("cluster_" + str(n_clusters) + ".png")
122        return list_images
123
124    def silhouetteScoretoPNG(self, list_images):
125        """Save the results of the plots in asingle image file.
126
127        :param self: An instance of the class SilhouetteScore.
128        :param list_images: A list with the name of the image files created.
129        """
130        clusterImages = [PIL.Image.open(i) for i in list_images]
131        minSize = sorted([(np.sum(i.size), i.size)
132                          for i in clusterImages])[0][1]
133
134        imagesCombination = np.vstack((np.asarray(i.resize(minSize))
135                                       for i in clusterImages))
136        imagesCombination = PIL.Image.fromarray(imagesCombination)
137        directory = "../images"
138        if not os.path.isdir(directory):
139            os.makedirs(directory)
140        imagesCombination.save("../images/clustersScore.png")
141        for image in list_images:
142            os.remove(image)
143        print ("The silhouette score for the number of"
144               " clusters ranging from 2 "
145               "to 6 has been saved in the file clustersScore.png!")
146
147
148 if __name__ == "__main__":
149
150    parser = argparse
151    parser = argparse.ArgumentParser()
152    parser.add_argument("dataFile", type=str,
153                        help="File to retrieve the generated data points.")
154    args = parser.parse_args()
155    instanceSilhouetteScore = SilhouetteScore()
156    images = instanceSilhouetteScore.calculateSilhouetteScore(args.dataFile)
157    instanceSilhouetteScore.silhouetteScoretoPNG(images)
```

# 3   k-means clustering algorithm

## 3.1   kmeans.py

This python script calls the k-means algorithm implemented on hadoop. However, before implementing k-means the initial centroids are computed using the k-means++ algorithm proposed in

2007 by Arthur and Vassilvitskii.

---

**Algorithm 1** k-means++ algorithm

---

1. Take one center $c_1$, chosen uniformly at random from $X$.

2. Take a new center $c_1$, choosing $x \in X$ with probability $\dfrac{D(x')^2}{\sum_{x \in X} D(x')^2}$.

3. Repeat Step 2. until we have taken $k$ centers altogether.

4. Proceed as with the standard k-means algorithm.

---

After determining the initial centroids, k-means algorithm is called in order to detetermine the new centroids of the clusters and the results are saved as an image file.

```python
"""kmeans.py: Run the k-means algorithm."""

import argparse
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import os
import random
import re
import sys
sys.tracebacklimit = 0

__author__ = "Stratos Gounidellis, Lamprini Koutsokera"
__copyright__ = "Copyright 2017, BDSMasters"


class KmeansRunner():

    def retrieveData(self, file):
        """Retrieve the data points from the input file.

        :param self: An instance of the class KmeansRunner.
        :param file: A file with the input data.
        :return: An array with the input data points.
        """
        df_points = pd.read_csv(file, header=None, names=["x", "y"], sep=" ")
        if (len(df_points.index) < 1):
            raise Exception("The input file is empty!")
        data = [tuple(row) for row in df_points.values]
        points = np.array([data_point for data_point in data])
        return points

    def initialCentroids(self, file, nclusters):
        """Calculate the initial centroids to be used by the k-means
            clustering algorithm.
```

```python
        :param self: An instance of the class KmeansRunner.
        :param file: A file with the input data.
        :param nclusters: The number of clusters.
        :return: A list with the initial centroids.
        """
        points = self.retrieveData(file)
        initial_centroids = [list(random.choice(points))]
        dist = []
        if nclusters < 2:
            raise Exception("Error the number of clusters should be" +
                            " greater than or equal to 2!")
        for i in range(2, nclusters + 1):
            dist.append([np.linalg.norm(np.array(point) -
                        initial_centroids[i - 2])**2 for point in points])
            min_dist = dist[0]
            if (len(dist) > 1):
                min_dist = np.minimum(
                    min_dist, (dist[index] for index in range(1, len(dist))))

            sumValues = sum(min_dist)
            probabilities = [float(value) / sumValues for value in min_dist]
            cumulative = np.cumsum(probabilities)

            random_index = random.random()
            index = np.where(cumulative >= random_index)[0][0]
            initial_centroids.append(list(points[index]))

        return initial_centroids

    def retrieveCentroids(self, file):
        """Retrieve the centroids coordinated from the centroids file.

        :param self: An instance of the class KmeansRunner.
        :param file: A file with the centroids.
        :return: A list with the centroids.
        """
        with open(file, "r") as inputFile:
            output_data = inputFile.readlines()

        centroids = []
        for point in output_data:
            p = re.search("\[(.*?)\]", point).group()
            p = p.replace("[", "").replace("]", "")
            p.strip()
            axisx, axisy = p.split(",")
            axisx = float(axisx)
            axisy = float(axisy)
            point_list = [axisx, axisy]
            centroids.append(point_list)
        return centroids

    def retrieveLabels(self, dataFile, centroidsFile):
        """Retrieve the labels of the imput data points.
```

```python
            :param self: An instance of the class KmeansRunner.
            :param dataFile: A file with the input data points.
            :param centroidsFile: A file with the centroids.
            :return: A list with the labels.
            """
            data_points = self.retrieveData(dataFile)
            centroids = self.retrieveCentroids(centroidsFile)
            labels = []
            for data_point in data_points:
                distances = [np.linalg.norm(data_point - centroid)
                             for centroid in centroids]
                cluster = np.argmin(distances)
                labels.append(int(cluster))
            return labels

    def writeCentroids(self, centroids, file):
        """Write centroids to a file.

            :param self: An instance of the class KmeansRunner.
            :param centroids: A list with the centroids.
            :param file: A file to write the centroids.
            """
            f = open(CENTROIDS_FILE, "w+")
            for item in centroids:
                f.write("%s\n" % str(item))
            f.close()

    def plotClusters(self, data_points, centroids, labels):
        """Plot the clusters with the centroids and save the plot as an image.

            :param self: An instance of the class KmeansRunner.
            :param data_points: An array with the input data points.
            :param centroids: A list with the centroids.
            :param labels: The labels of the input data points.
            """
            plt.scatter(data_points[:, 0], data_points[:, 1], c=labels)
            for i in range(len(centroids)):
                label = "Centroid " + str(i)
                colors = ["red", "green", "blue"]
                plt.scatter(centroids[i][0], centroids[i][1], s=50,
                            c=colors[i], label=label)
            plt.legend(loc="best", fancybox=True)
            fig = plt.gcf()
            plt.show()
            directory = "../images"
            if not os.path.isdir(directory):
                os.makedirs(directory)
            fig.savefig("../images/clusters.png")


CENTROIDS_FILE = "centroids.txt"
OUTPUT_FILE = "output.txt"

if __name__ == "__main__":
```

```python
    parser = argparse
    parser = argparse.ArgumentParser(description="k-means algorithm"
                                     " implementation on Hadoop",
                                     epilog="Go ahead and try it!")
    parser.add_argument("inputFile", type=str,
                        help="Input data points for the clustering algorithm."
    )
    parser.add_argument("centroids", type=int,
                        help="Number of clusters.")
    args = parser.parse_args()

    data = args.inputFile
    k = args.centroids
    instanceKmeans = KmeansRunner()
    centroids = instanceKmeans.initialCentroids(data, int(k))
    instanceKmeans.writeCentroids(centroids, CENTROIDS_FILE)

    outputFile = open(OUTPUT_FILE, "w+")
    outputFile.close()

    i = 1
    while True:
        print "k-means iteration #%i" % i

        command = "python kmeansAlgorithm.py < " \
                  + data + " --k=" \
                  + str(k) + " --centroids=" \
                  + CENTROIDS_FILE + " > " + OUTPUT_FILE \
                  + " -r hadoop"
        os.popen(command)

        new_centroids = instanceKmeans.retrieveCentroids(OUTPUT_FILE)

        if sorted(centroids) != sorted(new_centroids):
            centroids = new_centroids
            instanceKmeans.writeCentroids(centroids, CENTROIDS_FILE)
        else:
            break
        i += 1

    os.remove(OUTPUT_FILE)
    labels = instanceKmeans.retrieveLabels(data, CENTROIDS_FILE)
    labelsFile = open("labels.txt", "w+")
    for label in labels:
        labelsFile.write("%s\n" % str(label))
    labelsFile.close()
    data_points = instanceKmeans.retrieveData(data)
    instanceKmeans.plotClusters(data_points, centroids, labels)
```
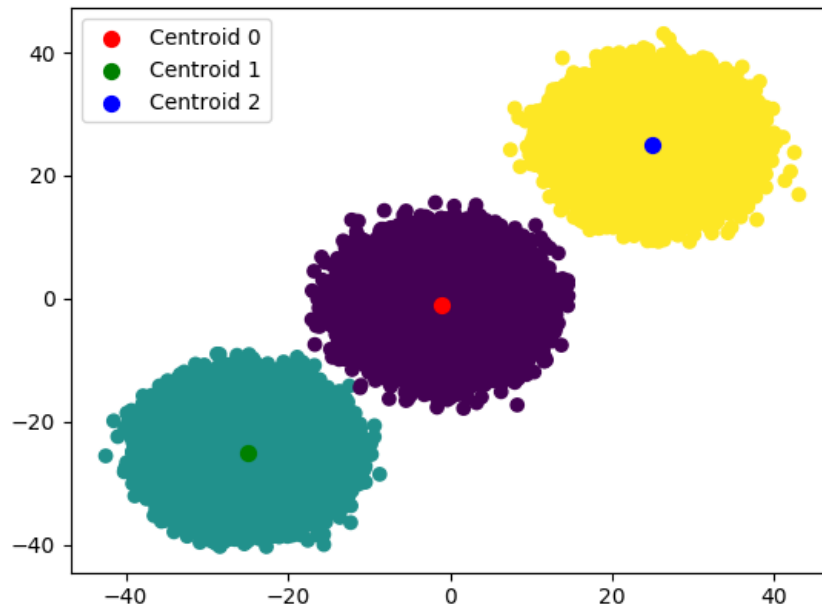
## 3.2 kmeansAlgorithm.py

In order to implement k-means algorithm on hadoop mrjob is used. Mrjob is a python package, which allows to write multi-step MapReduce jobs in pure Python and run them on a hadoop cluster. In our case mrjob run on a single-node cluster. (The script can also be run locally by commenting the argument "-r hadoop".)

---

**Algorithm 2** k-means algorithm

---

1. Define the number of clusters, k.

2. Select k data points as initial centroids.

3. Assign each data object to the closest cluster centroid.

4. Recalculate the clusters' centroids.

5. If the centroids remain unchanged the algorithm terminates. Otherwise, the steps are repeated from Step 2.

---

**Algorithm 3** k-means algorithm - MapReduce

---

1. The mapper function returns each data point and the cluster, to which it belongs.

2. The combiner function returns partial sums of batches of data points belonging to the same cluster.

3. The reducer returns the new centroids of each cluster.

4. If the centroids remain unchanged the algorithm terminates. Otherwise, the steps are repeated from the beginning.

---

```python
1  """kmeansAlgorithm.py: Implement the k-means clustering
2      algorithm on the input data."""
3
4  from mrjob.job import MRJob
5  from mrjob.job import MRStep
6  import numpy as np
7  import re
8
9  __author__ = "Stratos Gounidellis, Lamprini Koutsokera"
10 __copyright__ = "Copyright 2017, BDSMasters"
11
12
13 class KmeansAlgorithm(MRJob):
14     def configure_options(self):
15         """Set the arguments for the class KmeansAlgorithm.
16
17         :param self: A instance of the class KmeansAlgorithm.
18         """
19         super(KmeansAlgorithm, self).configure_options()
20         self.add_passthrough_option(
21             "--k", type="int", help="Number of clusters.")
22         self.add_file_option("--centroids")
23
24     def retrieveCentroids(self, file):
25         """Retrieve the centroids coordinated from the centroids file.
26
27         :param self: An instance of the class KmeansAlgorithm.
28         :param file: A file with the centroids.
29         :return: A list with the centroids.
30         """
```

```python
31          with open(file, "r") as inputFile:
32              output_data = inputFile.readlines()
33
34          centroids = []
35          for point in output_data:
36              p = re.search("\[(.*?)\]", point).group()
37              p = p.replace("[", "").replace("]", "")
38              p.strip()
39              axisx, axisy = p.split(",")
40              axisx = float(axisx)
41              axisy = float(axisy)
42              point_list = [axisx, axisy]
43              centroids.append(point_list)
44          return centroids
45
46      def assignPointtoCluster(self, _, line):
47          """Assign each point to its closest cluster – Mapper Function.
48
49          :param self: An instance of the class KmeansAlgorithm.
50          :param line: A line from the input data, with data points in
51              the form [axisx axisy]
52          :yield: The identifier of a cluster and a point belonging to it.
53          """
54          axisx, axisy = line.split()
55          data_point = np.array([float(axisx), float(axisy)])
56          centroids = self.retrieveCentroids(self.options.centroids)
57          distances = [np.linalg.norm(data_point – centroid)
58                       for centroid in centroids]
59          cluster = np.argmin(distances)
60          yield int(cluster), data_point.tolist()
61
62      def calculatePartialSum(self, cluster, data_points):
63          """Calculate the partial sum of the data points belonging to
64              each cluster – Combiner Function.
65
66          :param self: An instance of the class KmeansAlgorithm.
67          :param cluster: An identifier for each cluster.
68          :param data_points: A list of points belonging to each cluster.
69          :yield: The identifier of a cluster, the partial sum of its
70              data points and their number.
71          """
72          sum_points = np.array(data_points.next())
73          counter = 1
74          for data_point in data_points:
75              sum_points += data_point
76              counter += 1
77          yield cluster, (sum_points.tolist(), counter)
78
79      def calculateNewCentroids(self, cluster, partial_sums):
80          """Calculate the new centroids of the clusters – Reduce Function.
81
82          :param self: An instance of the class KmeansAlgorithm.
83          :param cluster: An identifier for each cluster.
84          :param partial_sums: A list with the partial sum of the
```

```
85          data points of a cluster and their number.
86      :yield: The identifier of a cluster and its new centroid.
87      """
88      total_sum, total_counter = partial_sums.next()
89      total_sum = np.array(total_sum)
90      for partial_sum, counter in partial_sums:
91          total_sum += partial_sum
92          total_counter += counter
93      yield cluster, (total_sum / total_counter).tolist()
94
95  def steps(self):
96      """Set the steps of the MRJob.
97
98      :param self: An instance of the class KmeansAlgorithm.
99
100     :return: a list of steps constructed with MRStep().
101     """
102     return [MRStep(mapper=self.assignPointtoCluster,
103             combiner=self.calculatePartialSum,
104             reducer=self.calculateNewCentroids)]
105
106
107 if __name__ == "__main__":
108     KmeansAlgorithm.run()
```

# 4    Testing Functionality

## 4.1    test.py

```
1  import unittest
2  from createDataPoints import DataGenerator
3  from kmeans import KmeansRunner
4
5  __author__ = "Stratos Gounidellis, Lamprini Koutsokera"
6  __copyright__ = "Copyright 2017, BDSMasters"
7
8
9  class TestStringMethods(unittest.TestCase):
10
11     def test_dataPoints(self):
12         instanceData = DataGenerator()
13         fname = "test.txt"
14         instanceData.generateData(100, fname)
15         with open(fname) as f:
16             for i, l in enumerate(f):
17                 pass
18         i + 1
19         self.assertEqual(100, i+1)
20
21     def test_exceptionClustersNumber(self):
22         fname = "test.txt"
23         instanceKmeans = KmeansRunner()
```

```
24          with self.assertRaises(Exception) as context:
25              instanceKmeans.initialCentroids(fname, 1)
26          self.assertIn("Error the number of clusters should be greater" +
27                          "than or equal to 2!", "".join(context.exception))
28
29      def test_fileLength(self):
30          fname = "test.txt"
31          instanceKmeans = KmeansRunner()
32          testFile = open(fname, "w+")
33          testFile.close()
34          with self.assertRaises(Exception) as context:
35              instanceKmeans.retrieveData(fname)
36          self.assertIn("The input file is empty!", "".join(context.exception))
37
38
39 if __name__ == "__main__":
40      unittest.main()
```

# References

[1] Michael G. Noll. *Running Hadoop On Ubuntu Linux (Single-Node Cluster)*. http://www.michael-noll.com/tutorials/running-hadoop-on-ubuntu-linux-single-node-cluster/ [*Accessed 29 Mar. 2017*].

[2] Stackoverflow.com. *Error launching job using mrjob on Hadoop*. http://stackoverflow.com/questions/25358793/error-launching-job-using-mrjob-on-hadoop [*Accessed 29 Mar. 2017*].

[3] David Arthur, and Sergei Vassilvitskii (2007). *k-means++: the advantages of careful seeding* Proceedings of the eighteenth annual ACM-SIAM Symposium on Discrete Algorithms, New Orleans, LA, January 7-9, 2007. 1st ed. New York: ACM, pp.1027–1035.

[4] Nlp.stanford.edu. *Kmeans*. http://nlp.stanford.edu/IR-book/html/htmledition/k-means-1.html [*Accessed 15 Mar. 2017*].

[5] Home.deib.polimi.it. *Clustering - K-means*. http://home.deib.polimi.it/matteucc/Clustering/tutorial_html/kmeans.html [*Accessed 8 Mar. 2017*].

[6] Kyuseok Shim. *MapReduce Algorithms for Big Data Analysis*. VLDB Conference, 2012.