



DEPARTMENT OF
COMPUTER SCIENCE AND ENGINEERING

Title: K-Means Clustering Algorithm Implementation

ARTIFICIAL INTELLIGENCE LAB
CSE 404



GREEN UNIVERSITY OF BANGLADESH

1 Objective(s)

- To understand K-Means algorithm to solve clustering of points
- To understand how to assign clusters and steps to finalize the cluster using Euclidean distance

2 Problem analysis

K-Means Clustering is an Unsupervised Learning algorithm, which groups the unlabeled dataset into different clusters. Here K defines the number of pre-defined clusters that need to be created in the process, as if $K=2$, there will be two clusters, and for $K=3$, there will be three clusters, and so on. It allows us to cluster the data into different groups and a convenient way to discover the categories of groups in the unlabeled dataset on its own without the need for any training.

It is a centroid-based algorithm, where each cluster is associated with a centroid. The main aim of this algorithm is to minimize the sum of distances between the data point and their corresponding clusters.

The algorithm takes the unlabeled dataset as input, divides the dataset into k-number of clusters, and repeats the process until it does not find the best clusters. The value of k should be predetermined in this algorithm.

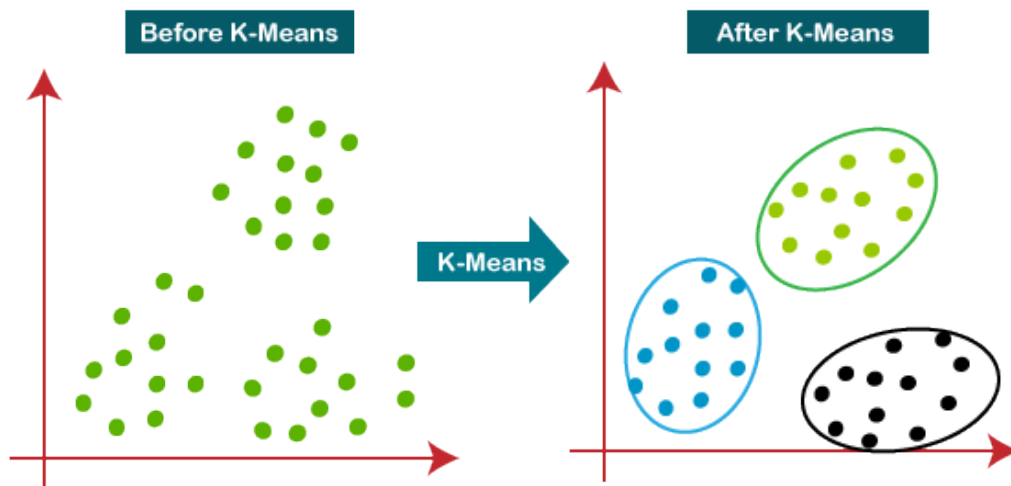


Figure 1: Clustering using K-Means

3 Algorithm

Step 1: Select the number K to decide the number of clusters.

Step 2: Select random K points or centroids. (It can be other from the input dataset).

Step 3: Assign each data point to their closest centroid, which will form the predefined K clusters.

Step 4: Calculate the variance and place a new centroid of each cluster.

Step 5: Repeat the third steps, which means reassign each datapoint to the new closest centroid of each cluster.

Step 6: If any reassignment occurs, then go to step-4 else go to FINISH.

3.1 Working example of K-Means clustering

Let's take number k of clusters, i.e., $K=2$, to identify the dataset and to put them into different clusters. It means here we will try to group these datasets into two different clusters.

We need to choose some random k points or centroid to form the cluster like in figure 2a. These points can be either the points from the dataset or any other point. So, here we are selecting the below two points as k points, which are not the part of our dataset.

Now we will assign each data point of the scatter plot to its closest K-point or centroid. We will compute it by applying some mathematics that we have studied to calculate the distance between two points. So, we will draw a median between both the centroids (figure 2b).

From the figure 2b, it is clear that points left side of the line is near to the K1 or blue centroid, and points to the right of the line are close to the yellow centroid. Let's color them as blue and yellow for clear visualization which is given in figure 2c

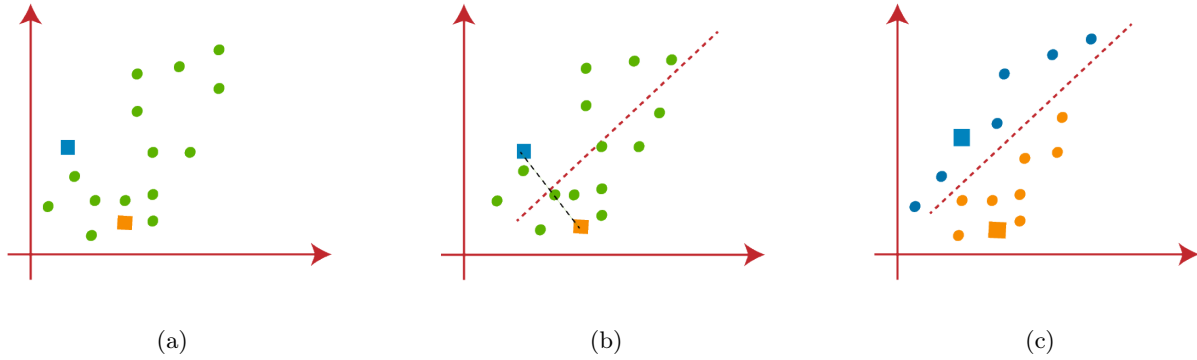


Figure 2

As we need to find the closest cluster, so we will repeat the process by choosing a new centroid. To choose the new centroids, we will compute the center of gravity of these centroids, and will find new centroids as in figure 3a

Next, we will reassign each datapoint to the new centroid. For this, we will repeat the same process of finding a median line. The median will be like as in figure 3b.

From the figure 3b, we can see, one yellow point is on the left side of the line, and two blue points are right to the line. So in figure 3c, these three points will be assigned to new centroids.

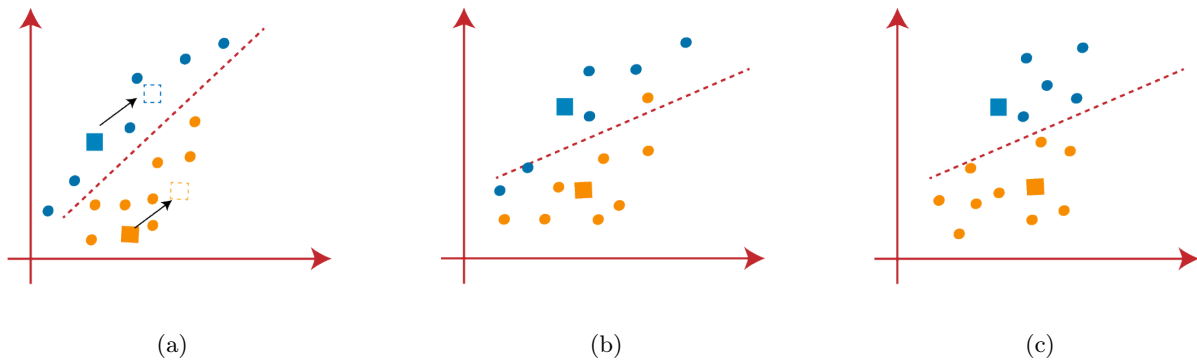


Figure 3

As reassignment has taken place, so we will again go to the step-4, which is finding new centroids or K-points (figure 4a)

We will repeat the process by finding the center of gravity of centroids, so the new centroids will be as shown in the figure 4b:

As we got the new centroids so again will draw the median line and reassign the data points. So, from figure 4c, we can see there are no dissimilar data points on either side of the line, which means our model is formed.

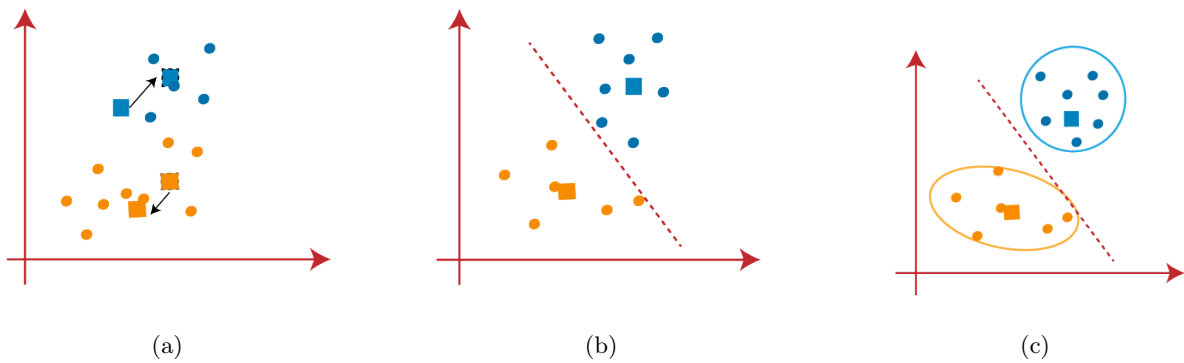


Figure 4

4 Implementation in Java

```

1 package k_means;
2
3 import java.util.ArrayList;
4 import java.util.Random;
5 import java.util.Scanner;
6
7 /**
8  *
9  * @author Jargis Ahmed
10 */
11 class points {
12
13     int x, y, cl;
14
15     points(int x, int y) { // points on the grid (x,y)
16         this.x = x;
17         this.y = y;
18     }
19
20     void setClass(int c) { // cluster number in which this (x,y) point reside
21         cl = c;
22     }
23 }
24
25 class St {
26
27     int pt, ks; // pt is the number of points and ks is the number of
28                 // clusters
29     int mat[][];
30     points[] p, k;
31     ArrayList<points[]> fkc;
32     Random rand = new Random();
33
34     St(int points, int clusters) {
35         pt = points;
36         ks = clusters;
37         mat = new int[pt][pt]; // mat is the 2d cartesian grid or matrix
38         fkc = new ArrayList<points[]>();
39         Start();
40     }
41
42     void Start() {
43         // ... (code for starting the algorithm) ...
44     }
45 }

```

```

40
41 void Start() {
42     p = new points[pt]; // create random points coordinate
43     for (int i = 0; i < pt; i++) {
44         int x = rand.nextInt(pt);
45         int y = rand.nextInt(pt);
46         p[i] = new points(x, y);
47     }
48     k = new points[ks]; // create random cluster points coordinate
49     for (int i = 0; i < ks; i++) {
50         int x = rand.nextInt(pt);
51         int y = rand.nextInt(pt);
52         k[i] = new points(x, y);
53     }
54     int i, j;
55     double min;
56     int count = 0;
57     while (true) {
58         // starting cluster allocation of point p(x,y) based on minimum
59         distance from each cluster points
60         for (i = 0; i < pt; i++) {
61             for (j = 0, min = 10000; j < ks; j++) {
62                 double dl = Math.sqrt(Math.pow((double) (k[j].x - p[i].x),
63                     2) + Math.pow((double) (k[j].y - p[i].y), 2));
64                 if (dl < min) {
65                     p[i].setClass(j);
66                     min = dl;
67                 }
68             }
69             points[] kdup = new points[ks]; // creating duplicate set of
70             cluster points to store starting cluster coordinates
71             for (i = 0; i < ks; i++) {
72                 kdup[i] = new points(k[i].x, k[i].y);
73             }
74             // calculating mean for each points in different clusters
75             for (j = 0; j < ks; j++) {
76                 int x = 0, y = 0, ci = 0;
77                 for (i = 0; i < pt; i++) {
78                     if (p[i].cl == j) {
79                         x += p[i].x;
80                         y += p[i].y;
81                         ci++;
82                     }
83                 }
84                 if (ci != 0) { // allocating the mean as new cluster point
85                     coordinate
86                     k[j].x = x / ci;
87                     k[j].y = y / ci;
88                 }
89             }
90             int err = 0;
91             // calculating error between previous and present cluster points
92             coordinates
93             for (i = 0; i < ks; i++) {
94                 err += Math.abs(k[i].x - kdup[i].x) + Math.abs(k[i].y - kdup[i].
95                     y);
96             }

```

```

92         count++;
93         // 0 error between previous and present cluster points coordinates
          indicates clustering is finalized
94         if (err == 0) {
95             break;
96         }
97     }
98
99     double IntraDis;
100    for (i = 0; i < ks; i++) {
101        for (j = 0, IntraDis = 0; j < pt; j++) {
102            if (p[j].cl == i) {
103                IntraDis += Math.sqrt(Math.pow((double) (k[i].x - p[j].x),
104                    2) + Math.pow((double) (k[i].y - p[j].y), 2));
105            }
106            System.out.println("Cluster " + (i + 1) + " Intra-distance = " +
107                IntraDis);
108        }
109        for (i = 0; i < pt; i++) {
110            mat[p[i].x][p[i].y] = p[i].cl + 1;
111        }
112        for (i = 0; i < ks; i++) {
113            for (j = 0; j < p.length; j++) {
114                if (p[j].cl == i) {
115                    System.out.println("Point (" + p[j].x + ", " + p[j].y + ") "
116                        + "Cluster - " + (p[j].cl + 1));
117                }
118            }
119        }
120        for (i = 0; i < k.length; i++) {
121            System.out.println("Cluster " + (i + 1) + " - (" + k[i].x + ", " + k
122                [i].y + ")");
123        }
124        //printGrapgh(); // print a 2D graph alike matrix showing the points
          and
125    }
126 }
127
128 public class K_Means {
129
130     public static void main(String[] args) {
131         // TODO code application logic here
132         Scanner sc = new Scanner(System.in);
133         System.out.println("Insert number of points");
134         int points = sc.nextInt();
135         System.out.println("Insert number of clusters");
136         int clusters = sc.nextInt();
137         St s = new St(points, clusters);
138     }
139 }

```

5 Implementation in Python

```

1 import math
2 import random

```

```

3
4 class Point:
5     def __init__(self, x, y):
6         self.x = x
7         self.y = y
8         self.cluster = None
9
10 class KMeans:
11     def __init__(self, points, clusters):
12         self.points = points
13         self.clusters = clusters
14         self.mat = [[0 for _ in range(points)] for _ in range(points)]
15         self.fkc = []
16         self.start()
17
18     def start(self):
19         self.p = [Point(random.randint(0, self.points - 1), random.randint(0,
20             self.points - 1)) for _ in range(self.points)]
21         self.k = [Point(random.randint(0, self.points - 1), random.randint(0,
22             self.points - 1)) for _ in range(self.clusters)]
23
24         count = 0
25         while True:
26             for i in range(self.points):
27                 min_dist = float('inf')
28                 for j in range(self.clusters):
29                     dl = math.sqrt((self.k[j].x - self.p[i].x) ** 2 + (self.k[j]
30                         ].y - self.p[i].y) ** 2)
31                     if dl < min_dist:
32                         self.p[i].cluster = j
33                         min_dist = dl
34
35             kdup = [Point(self.k[i].x, self.k[i].y) for i in range(self.clusters
36                 )]
37             for j in range(self.clusters):
38                 x, y, ci = 0, 0, 0
39                 for i in range(self.points):
40                     if self.p[i].cluster == j:
41                         x += self.p[i].x
42                         y += self.p[i].y
43                         ci += 1
44                     if ci != 0:
45                         self.k[j].x = x // ci
46                         self.k[j].y = y // ci
47
48             err = 0
49             for i in range(self.clusters):
50                 for j in range(2): # Only 2 dimensions (x and y)
51                     err += abs(self.k[i].x - kdup[i].x) + abs(self.k[i].y - kdup
52                         [i].y)
53
54             count += 1
55             if err == 0:
56                 break
57
58         for i in range(self.clusters):
59             intra_distance = sum(math.sqrt((self.k[i].x - p.x) ** 2 + (self.k[i]
60                 ].y - p.y) ** 2) for p in self.p if p.cluster == i)

```

```

55         print(f"Cluster {i + 1} Intra-distance = {intra_distance}")
56
57     for p in self.p:
58         self.mat[p.x][p.y] = p.cluster + 1
59
60     for i in range(self.clusters):
61         for p in self.p:
62             if p.cluster == i:
63                 print(f"Point ({p.x}, {p.y}) Cluster - {p.cluster + 1}")
64
65     for i, k in enumerate(self.k):
66         print(f"Cluster {i + 1} - ({k.x}, {k.y})")
67
68 def main():
69     points = int(input("Insert number of points: "))
70     clusters = int(input("Insert number of clusters: "))
71     kmeans = KMeans(points, clusters)
72
73 if __name__ == "__main__":
74     main()

```

6 Sample Input/Output (Compilation, Debugging & Testing)

Input:

```

1 Insert number of points
2 20
3 Insert number of clusters
4 4

```

Output:

```

1 Cluster 1 Intra-distance = 42.004763409818885
2 Cluster 2 Intra-distance = 16.82842712474619
3 Cluster 3 Intra-distance = 10.537319187990756
4 Cluster 4 Intra-distance = 9.930106595800748
5 Point (8, 6) Cluster - 1
6 Point (12, 10) Cluster - 1
7 Point (13, 10) Cluster - 1
8 Point (17, 4) Cluster - 1
9 Point (6, 9) Cluster - 1
10 Point (16, 1) Cluster - 1
11 Point (10, 11) Cluster - 1
12 Point (16, 6) Cluster - 1
13 Point (12, 4) Cluster - 1
14 Point (8, 0) Cluster - 2
15 Point (0, 0) Cluster - 2
16 Point (2, 2) Cluster - 2
17 Point (4, 0) Cluster - 2
18 Point (10, 0) Cluster - 2
19 Point (15, 14) Cluster - 3
20 Point (9, 15) Cluster - 3
21 Point (18, 16) Cluster - 3
22 Point (5, 16) Cluster - 4
23 Point (1, 12) Cluster - 4
24 Point (0, 18) Cluster - 4
25 Cluster 1 - (12, 6)

```



```
26 Cluster 2 - (4, 0)
27 Cluster 3 - (14, 15)
28 Cluster 4 - (2, 15)
```

7 Discussion & Conclusion

Based on the focused objective(s) to understand the K-Means clustering problem and its steps to assign different points under some cluster points or centroids, the additional lab exercise will increase confidence towards the fulfilment of the objectives(s).

8 Lab Exercise (Submit as a report)

- Rewrite K-Means clustering algorithm where -
 - Construct a data file containing 100 cartesian points $P(x, y)$ and 10 cluster points $C(x, y)$
 - Calculate distance using "Manhattan distance" method
 - Make a 2D visualization of points and its cluster points by only using print functionality. (Hint: Use a 2D Matrix which will hold the points information and clusters. Print the matrix content accordingly)

9 Policy

Copying from internet, classmate, seniors, or from any other source is strongly prohibited. 100% marks will be *deducted* if any such copying is detected for lab exercise.