

Problem Statement

We have given a collection of 8 point

$P1=[0.1,0.6]$, $P2=[0.15,0.71]$, $P3=[0.08,0.9]$, $P4=[0.16, 0.85]$, $P5=[0.2,0.3]$, $P6=[0.25,0.5]$, $P7=[0.24,0.1]$, $P8=[0.3,0.2]$.

Perform the k-mean clustering with initial centroids as $m1=P1 = \text{Cluster\#1}=C1$ and $m2=P8=\text{cluster\#2}=C2$.

Answer the following

1. Which cluster does P6 belong to?
2. What is the population of cluster around $m2$?
3. What is updated value of $m1$ and $m2$?

K Means Clustering Key Points

- Unsupervised Machine Learning Algorithm
- Partition Clustering
- k = number of classes, $k > 2$
- A cluster is a collection of data points aggregated together because of certain similarities.
- in the beginning the centroids are taken randomly
- ends when
 - centroids have stabilized
 - iteration has exceeded the limit
- k can be determined using elbow method
 - Plot a graph of k vs Distortion
 - Identify the Elbow of the curve
- Advantages
 - Relatively simple to implement.
 - Scales to large data sets.
 - Guarantees convergence.
- Disadvantages
 - Choosing k manually
 - Clustering data of varying sizes and density.

- Clustering outliers.
- Curse of Dimensionality
- Dependant on initial cluster center values

▼ K Means Clustering Algorithm

Algorithm

1. for each datapoint in dataset
 - 1.1 for each center in centers
 - 1.1.1. get distance between datapoint and center
 - 1.2 select center closest to datapoint
 - 1.3 assign cluster based on closest center
2. for each cluster
 - 2.1 assign new center as centroid of datapoints in cluster
3. if new centers = old centers then return new centers else go to step 1

```

1 #importing libraries
2 import numpy as np
3 import numpy.matlib
4 import matplotlib.pyplot as plt
5 import pandas as pd

```

```

1 data = [
2     [0.1, 0.6],
3     [0.15, 0.71],
4     [0.08, 0.9],
5     [0.16, 0.85],
6     [0.2, 0.3],
7     [0.25, 0.5],
8     [0.24, 0.1],
9     [0.3, 0.2]
10 ]
11
12 data = pd.DataFrame(data, columns = ['x', 'y'])
13
14 centroids = [
15     [0.1, 0.6],
16     [0.3, 0.2]
17 ]
18
19 k = 2

```

```

1 #calculate distance between 2 points
2 def calc_distance(x1, x2):
3     return (sum((x1 - x2)**2))**0.5

1 #check the point is closer to which centroid
2 def assign_clusters(centroids, data):
3     clusters = []
4     for i in range(data.shape[0]):
5         distances = []
6         for centroid in centroids:
7             distances.append(calc_distance(centroid, data.iloc[i]))
8         print(distances)
9         cluster = [z for z, val in enumerate(distances) if val==min(distances)]
10        clusters.append(cluster[0])
11
12    return clusters

1 #new centroid = mean of all the points belonging to that cluster
2 def calc_centroids(clusters, data):
3     new_centroids = []
4     cluster_df = pd.concat([pd.DataFrame(data), pd.DataFrame(clusters, columns=['c
5     for c in set(cluster_df['cluster']):
6         current_cluster = cluster_df[cluster_df['cluster']==c][cluster_df.columns[
7         cluster_mean = current_cluster.mean(axis=0)
8         new_centroids.append(cluster_mean)
9     return new_centroids

1 clusters = assign_clusters(centroids, data)
2 print(clusters)

[0.0, 0.44721359549995787]
[0.12083045973594571, 0.5316013544000805]
[0.3006659275674582, 0.7337574531137656]
[0.2570992026436488, 0.6649060083951716]
[0.31622776601683794, 0.14142135623730948]
[0.18027756377319945, 0.30413812651491096]
[0.5192301994298868, 0.11661903789690602]
[0.44721359549995787, 0.0]
[0, 0, 0, 0, 1, 0, 1, 1]

1 centroids = calc_centroids(clusters, data)
2 print(centroids)

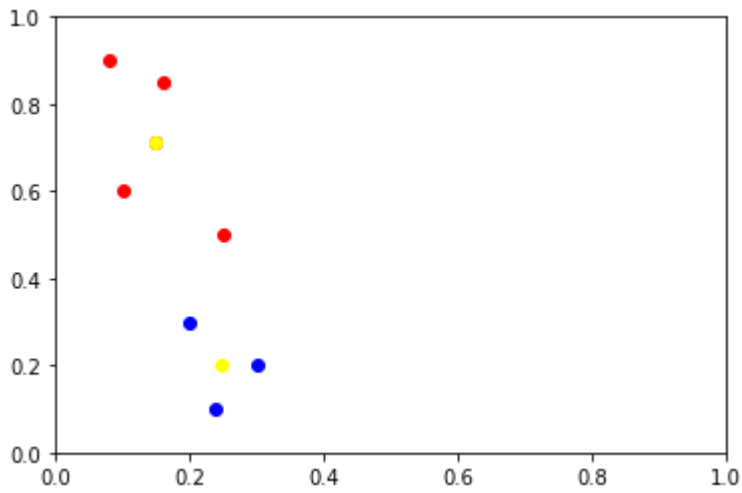
[x      0.148
 y      0.712
dtype: float64, x      0.246667
 y      0.200000
dtype: float64]

```

```

1 plt.plot()
2 colors = ['red', 'blue']
3
4 for i in range(data.shape[0]):
5     plt.scatter(data.iloc[i]['x'], data.iloc[i]['y'], c = colors[clusters[i]])
6
7 for i in centroids:
8     plt.scatter(i[0], i[1], c = 'yellow')
9
10 plt.axis([0, 1, 0, 1])
11 plt.show()

```



```

1 print(f'P6 belongs to cluster C{clusters[5]+1} coloured in {colors[clusters[5]]}')
    P6 belongs to cluster C1 coloured in red

```

```

1 print(f'Population of cluster around P8 is {clusters.count(clusters[7])} (Cluster
    Population of cluster around P8 is 3 (Cluster is shown in blue colour)

```

```

1 print(f'Updated values of centroids are ({centroids[0][0]}, {centroids[0][1]}) and
    Updated values of centroids are (0.148, 0.712) and (0.24666666666666667, 0.20000

```



<https://towardsdatascience.com/understanding-k-means-clustering-in-machine-learning-6a6e67336aa1>

```

1 from sklearn.cluster import KMeans
2 Kmean = KMeans(n_clusters=2)
3 Kmean.fit(X)

```

