# Problem Statement

We have given a collection of 8 points.

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 = Cluster#1=C1 and m2=P8=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  class KMeans:
 2
 3
 4    def __init__(self, dataset, centers):
 5      k = len(centers)
 6      centers, center2datapoints = self.action(dataset, k, centers)
 7      pass
 8
 9
10    def action(self, dataset, k, centers):
11      print('\nK Means Clustering\n===========')
12
13      from math import sqrt
14      import matplotlib.pyplot as plt
15
16      centers_old = centers.copy()
17      #iteration_count = 1
18
19      for iteration_count in range(0,5):
20
21        print('\nIteration Count = {iter}'.format(iter=iteration_count))
22
23        print('Centers = {centers}'.format(centers=centers_old))
24
25        # Get datapoints closest to the cluster center
26        center2datapoints = {}
```

```
27        cluster = []
28
29        # For each datapoint
30        for datapoint in dataset:
31
32          # Get distance of datapoint from each cluster center
33          center2distance = {}
34
35          # For each cluster center
36          for center_index, center in enumerate(centers_old):
37
38            # Calculate distance
39            distance = 0
40            for datapoint_dim, center_dim in zip(datapoint,center):
41              distance = distance + (datapoint_dim-center_dim)**2
42            distance = distance**0.5
43
44            # Save the distance
45            center2distance[center_index] = distance
46
47          # Find closest center to the datapoint
48          closest_center_index = 0
49          closest_center_distance = center2distance[closest_center_index]
50          for center_index in center2distance:
51            if center2distance[center_index] < closest_center_distance:
52              closest_center_index = center_index
53              closest_center_distance = center2distance[center_index]
54          cluster.append(closest_center_index)
55
56          print()
57          print('\tDatapoint = {datapoint}'.format(datapoint=datapoint))
58          print('\tDistance from each cluster centre  = {center2distance}'.format(ce
59          print('\tClosest Center = {closest_center_index}'.format(closest_center_in
60
61
62          # Save datapoint to nearest center in center2datapoints
63          if closest_center_index not in center2datapoints:
64            center2datapoints[closest_center_index] = [datapoint]
65          else:
66            center2datapoints[closest_center_index].append(datapoint)
67
68        # Compute new centers by taking center of each set of datapoints in center2d
69        centers_new = []
70        for center_index in center2datapoints:
71          nearest_datapoints = center2datapoints[center_index]
72          x_center, y_center = 0, 0
73          for x,y in nearest_datapoints:
74            x_center, y_center = x_center+x, y_center+y
75          x_center, y_center = x_center/len(nearest_datapoints), y_center/len(neares
76          centers_new.append((x_center,y_center))
77
78        '''
```

```
 78
 79        plt.scatter([x for x,y in dataset],[y for x,y in dataset],c=cluster)
 80        plt.scatter([x for x,y in centers_old],[y for x,y in centers_old],c='red')
 81        plt.scatter([x for x,y in centers_new],[y for x,y in centers_new],c='orange'
 82        plt.show()
 83        '''
 84
 85        print()
 86        print('Old Centers = {centers}'.format(centers=centers_old))
 87        print('New Centers = {centers}'.format(centers=centers_new))
 88
 89
 90        # Compare the old and the new centers, break the loop if no change
 91        if centers_old == centers_new:
 92          return centers_old, center2datapoints
 93        else:
 94          centers_old = centers_new
 95
 96        centers_old = centers_new
 97
 98      return centers_old, center2datapoints
 99
100
101 ##############################################################################
102 dataset = [
103   (0.1,0.6),
104   (0.15,0.71),
105   (0.08,0.9),
106   (0.16, 0.85),
107   (0.2,0.3),
108   (0.25,0.5),
109   (0.24,0.1),
110   (0.3,0.2)
111 ]
112 centers = [
113   (0.1,0.6),
114   (0.3,0.2)
115 ]
116 kmeans = KMeans(dataset, centers)
```

```
    K Means Clustering
    ============

    Iteration Count = 0
    Centers = [(0.1, 0.6), (0.3, 0.2)]

            Datapoint = (0.1, 0.6)
            Distance from each cluster centre  = {0: 0.0, 1: 0.44721359549995787}
            Closest Center = 0

            Datapoint = (0.15, 0.71)
            Distance from each cluster centre  = {0: 0.12083045973594571, 1: 0.531
```

```
        Closest Center = 0

        Datapoint = (0.08, 0.9)
        Distance from each cluster centre  = {0: 0.3006659275674582, 1: 0.7337
        Closest Center = 0

        Datapoint = (0.16, 0.85)
        Distance from each cluster centre  = {0: 0.2570992026436488, 1: 0.6649
        Closest Center = 0

        Datapoint = (0.2, 0.3)
        Distance from each cluster centre  = {0: 0.31622776601683794, 1: 0.141
        Closest Center = 1

        Datapoint = (0.25, 0.5)
        Distance from each cluster centre  = {0: 0.18027756377319945, 1: 0.304
        Closest Center = 0

        Datapoint = (0.24, 0.1)
        Distance from each cluster centre  = {0: 0.5192301994298868, 1: 0.1166
        Closest Center = 1

        Datapoint = (0.3, 0.2)
        Distance from each cluster centre  = {0: 0.44721359549995787, 1: 0.0}
        Closest Center = 1

Old Centers = [(0.1, 0.6), (0.3, 0.2)]
New Centers = [(0.148, 0.712), (0.24666666666666667, 0.20000000000000004)]

Iteration Count = 1
Centers = [(0.148, 0.712), (0.24666666666666667, 0.20000000000000004)]

        Datapoint = (0.1, 0.6)
        Distance from each cluster centre  = {0: 0.12185236969382252, 1: 0.426
        Closest Center = 0

        Datapoint = (0.15, 0.71)
        Distance from each cluster centre  = {0: 0.0028284271247461927, 1: 0.5
        Closest Center = 0

        Datapoint = (0.08, 0.9)
        Distance from each cluster centre  = {0: 0.19991998399359684, 1: 0.719
        Closest Center = 0

        Datapoint = (0.16, 0.85)
```