

Clustering

Question 1: We can cluster in one dimension as well as in many dimensions. In this problem, we are going to cluster numbers on the real line. The particular numbers (data points) are 1, 4, 9, 16, 25, 36, 49, 64, 81, and 100, i.e., the squares of 1 through 10. We shall use a k-means algorithm, with two clusters. You can verify easily that no matter which two points we choose as the initial centroids, some prefix of the sequence of squares will go into the cluster of the smaller and the remaining suffix goes into the other cluster. As a result, there are only nine different clusterings that can be achieved, ranging from $\{1\}\{4,9,\dots,100\}$ through $\{1,4,\dots,81\}\{100\}$.

We then go through a reclustering phase, where the centroids of the two clusters are recalculated and all points are reassigned to the nearer of the two new centroids. For each of the nine possible clusterings, calculate how many points are reclassified during the reclustering phase. List out pair of initial centroids that results in *exactly one* point being reclassified.

Sol:

Given data points are 1, 4, 9, 16, 25, 36, 49, 64, 81, and 100.

There can be 9 different clusters as follows:

1. $\{1\}, \{4, 9, 16, 25, 36, 49, 64, 81, 100\}$,
2. $\{1, 4\}, \{9, 16, 25, 36, 49, 64, 81, 100\}$,
3. $\{1, 4, 9\}, \{16, 25, 36, 49, 64, 81, 100\}$,
4. $\{1, 4, 9, 16\}, \{25, 36, 49, 64, 81, 100\}$,
5. $\{1, 4, 9, 16, 25\}, \{36, 49, 64, 81, 100\}$,
6. $\{1, 4, 9, 16, 25, 36\}, \{49, 64, 81, 100\}$,
7. $\{1, 4, 9, 16, 25, 36, 49\}, \{64, 81, 100\}$,
8. $\{1, 4, 9, 16, 25, 36, 49, 64\}, \{81, 100\}$,
9. $\{1, 4, 9, 16, 25, 36, 49, 64, 81\}, \{100\}$.

Only one point has to be shifted between clusters if we change the centroid. Let the initial values be 36 and 100. Mean of 36 and 100 = $(36+100)/2 = 68$.

So, the clusters will be $\{1, 4, 9, 16, 25, 36, 49, 64\}, \{81, 100\}$.

Centroids of these clusters are 25.5 and 90.5. Mean = $(25.5 + 90.5)/2 = 58$.

Now the clusters will be $\{1, 4, 9, 16, 25, 36, 49\}, \{64, 81, 100\}$. Here the only one element is shifted between clusters.

Question 2: Suppose we want to assign points to one of two cluster centroids, either (0,0) or (100,40). Depending on whether we use the L_1 or L_2 norm, a point (x,y) could be clustered with a different one of these two centroids. For this problem, you should work out the conditions under which a point will be clustered with the centroid (0,0) when the L_1 norm is used, but clustered with the centroid (100,40) when the L_2 norm is used. List out those points.

Sol:

Given centroids are (0,0), (100, 40).

Given a point (x, y) which could be clustered with a different one of these two centroids.

L1 norm is the Manhattan Distance and L2 norm is the Euclidean Distance.

After L1 norm and L2 norm are calculated the values of x and y are 55, 5.

When L1 norm is applied on point (55, 5), the point is clustered with centroid (0, 0).

When L2 norm is applied on point (55, 5), the point is clustered with centroid (100, 40)

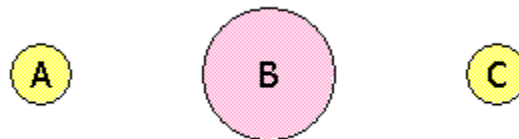
Question 3: Suppose our data set consists of the perfect squares 1, 4, 9, 16, 25, 36, 49, and 64, which are points in one dimension. Perform a hierarchical clustering on these points, as follows. Initially, each point is in a cluster by itself. At each step, merge the two clusters with the closest centroids, and continue until only two clusters remain. Which centroid of a cluster that exists at some time during this process? Positions are represented to the nearest 0.1.

Sol: 1, 4, 9, 16, 25, 36, 49 and 64

1 4 9 16 25 36 49 64

2.5 4.66 20.5 42.5 49.6

Question 4: Suppose that the true data consists of three clusters, as suggested by the diagram below:



There is a large cluster B centered around the origin (0,0), with 8000 points uniformly distributed in a circle of radius 2. There are two small clusters, A and C, each with 1000 points uniformly distributed in a circle of radius 1. The center of A is at (-10,0) and the center of C is at (10,0).

Suppose we choose three initial centroids x, y, and z, and cluster the points according to which of x, y, or z they are closest. The result will be three *apparent* clusters, which may or may not coincide with the *true* clusters A, B, and C. Say that one of the true clusters is *correct* if there is an apparent cluster that consists of all and only the points in that true cluster. Assuming initial centroids x, y, and z are chosen independently and at random, what is the probability that A is correct? What is the probability that C is correct? What is the probability that both are correct?

Sol:

Given centroids are x, y, z

We can assign each of x, y, z to A, B, C in 27 possible ways.

Chance of being in A is $1000/10000 = 0.1$

Chance of being in B is $8000/10000 = 0.8$

Chance of being in C is $1000/10000 = 0.1$

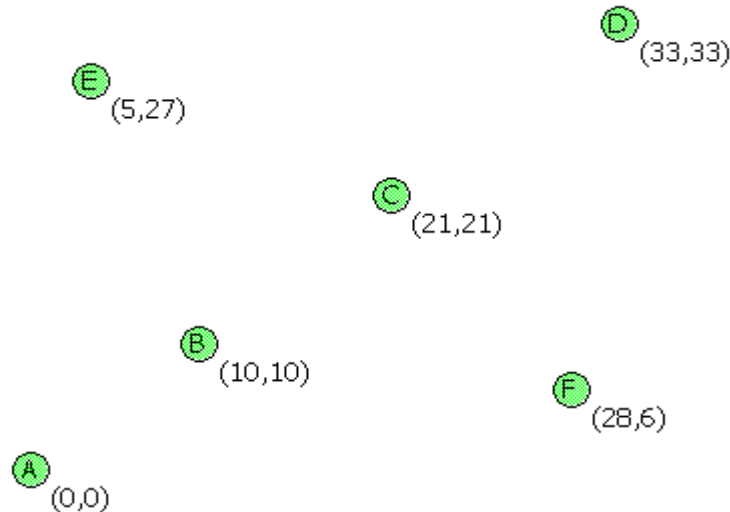
There are 6 different cases to interchange x, y, z in A, B, C which will be total 27.

The probability that A is correct is 24%

The probability that C is correct is 24%

The probability that A & C are correct is 4.8%

Question 5: Perform a hierarchical clustering of the following six points:



using the *complete-link* proximity measure (the distance between two clusters is the largest distance between any two points, one from each cluster). Find out a cluster at some stage of the agglomeration?

Sol:

	A	B	C	D	E	F
A	0	14.1	29.6	46.6	27.4	28.6
B		0	15.5	39.5	17.7	18.4
C			0	16.9	17	16.5
D				0	28.5	27.4
E					0	31.1
F						0

As A and B are low, clustering will be done as follows:

A and B will be clustered with C $\rightarrow AC - 29.6$ and $BC - 15.5$

A and B will be clustered with D $\rightarrow AD - 46.6$ and $BD - 32.5$

A and B will be clustered with E $\rightarrow AE - 27.4$ and $BE - 17.7$

A and B will be clustered with F \rightarrow AF – 28.6 and BF – 18.4

C and D will be clustered \rightarrow CD – 16.9

D and E will be clustered \rightarrow DE – 28.6

D and F will be clustered \rightarrow DF – 27.4

In CD, DE, DF as CD is low, clustering as follows:

C and D will be clustered with E \rightarrow CE – 17 and DE – 28.6

C and D will be clustered with F \rightarrow CF – 16.5 and DE – 27.4