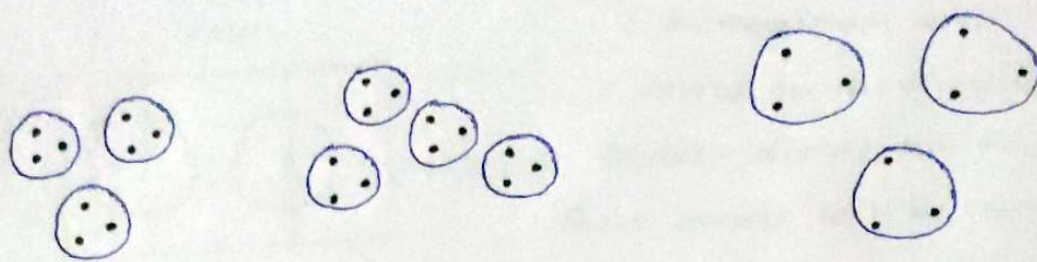


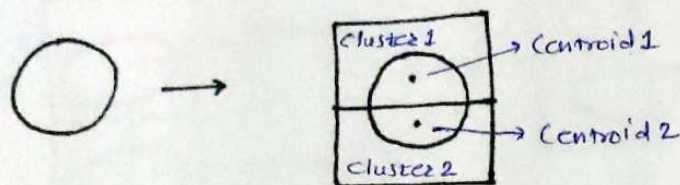
Q2 Find all well-separated clusters in the set of points shown below

Ans:



Q6. For following set of two dimensional points .....

a)  $k=2$  how many possible ways are there

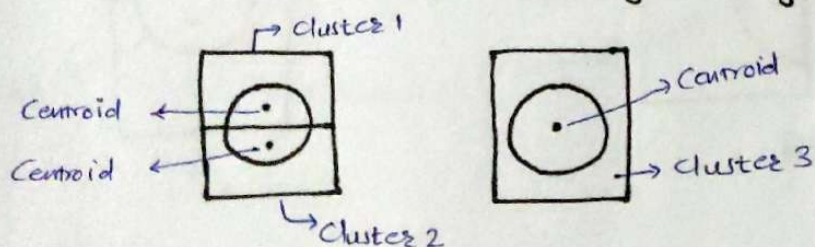


We can partition given circle in infinite way into two clusters.

Simple way is to draw diameter and partition circle into two parts will give us two clusters. As we can draw infinite diameters so there are infinite ways to partition the circle. And all the partitions will have almost same global minimum error.

If we draw a perpendicular to a diameter and then midpoint between the centre and the circumference on both the halves of the circle will be the position of Centroids of the cluster.

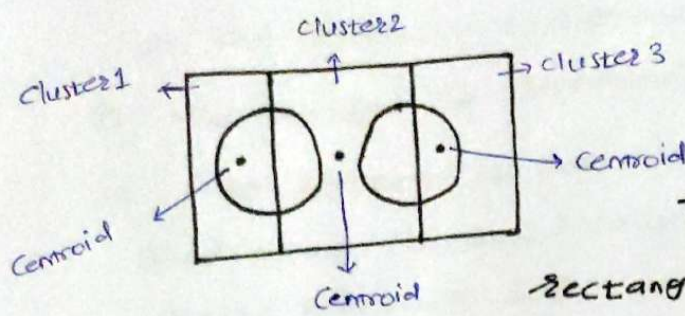
b)  $k=3$  distance between edges is greater than radii of circles



In this case we can make one cluster as circle as a cluster and divide another one into two parts. In first case center will become a centroid and in second case midpoint between the center and circumference on both halves of circle will be the position of 2<sup>nd</sup> & 3<sup>rd</sup> Centroids so that all the partitions will have same global minimum error.



c)  $K=3$ . The distance between the edges of circle is much less than radius of circle.



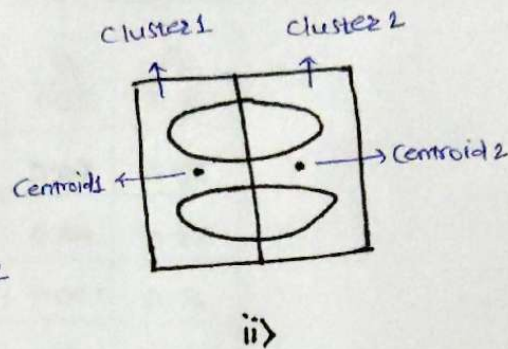
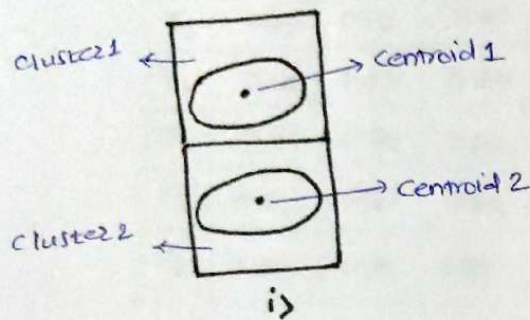
Take maximum distance between the 2 points on the circumference of the circles

Equally divide the line into 3 parts.

These points will be centroids and 3

rectangles that have these points as the center will be our three clusters.

d)  $K=2$

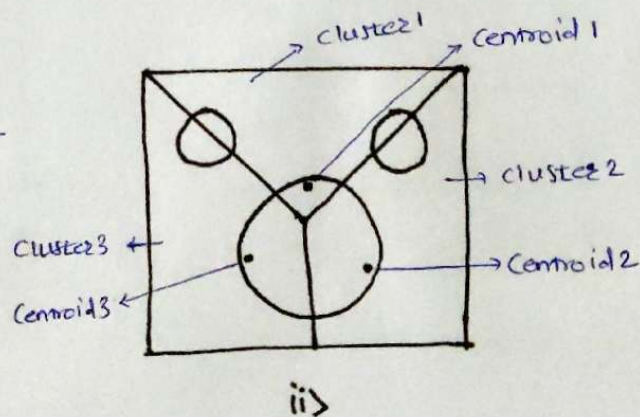
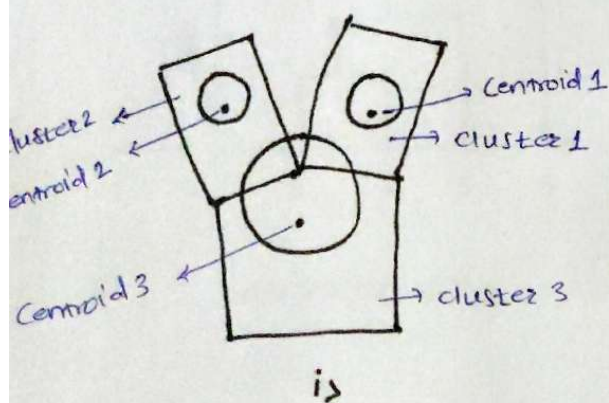


Above are two solutions for the problem where  $K=2$ .

i) produces local minimum error (left figure)

ii) right figure will produce global minimum error.

e)  $K=3$  Hint: Use the symmetry of the situation



There are two solutions to the problem i) produces global minimum error and ii) produces local minimum error.



11> Total SSE is the sum of the SSE . . . . .

- i> If a given attribute is constant, its SSE is constant for all the clusters. And it has minimal impact on clustering result.
- ii> If a given attribute variable has low SSE for just one cluster then that variable dominates in defining the cluster.
- iii> If a given attribute variable has high SSE for all the clusters then there is a high probability that this attribute is noise.
- iv> If a given attribute variable has high SSE for only one cluster then this variable doesn't help in defining the cluster, and the attribute that has low SSE for that cluster dominated this cluster.
- v> Per variable SSE information helps us eliminate attribute that has little impact in defining the clusters. Attributes that have low SSE for all the clusters are effectively constant and have a low impact on defining the cluster. Attributes that have high SSE for all the clusters are essentially noise, and they impact overall SSE.

12> Advantages of Leader Algorithm:

- i> In this algorithm as each object is compared to the final set of centroids so it is computationally efficient and its complexity is  $O(n)$  where  $n$  is number of elements.
- ii> It will always return the same result if the order of the input elements are same.

Disadvantages:

- i> We cannot have cluster number in advance as we have it in k-means. So even if we know the value of  $k$  we can't control the algo to make exact number of clusters.
- 2> This simplistic algo does not take SSE into consideration so the sum of error is high. k-means will almost always have a better result.

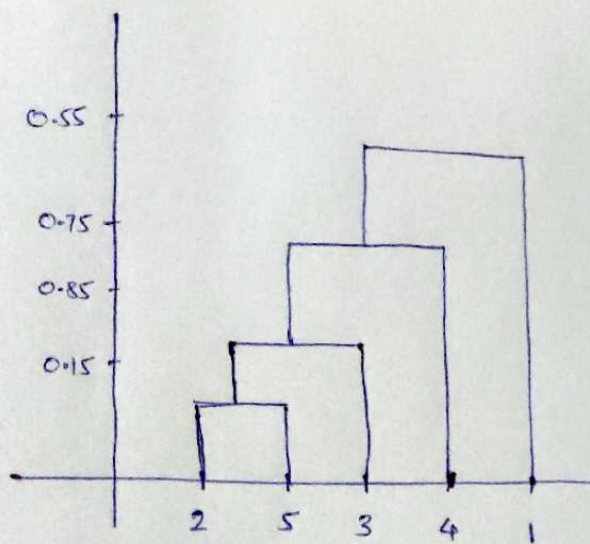


ways to improve leadership algo:

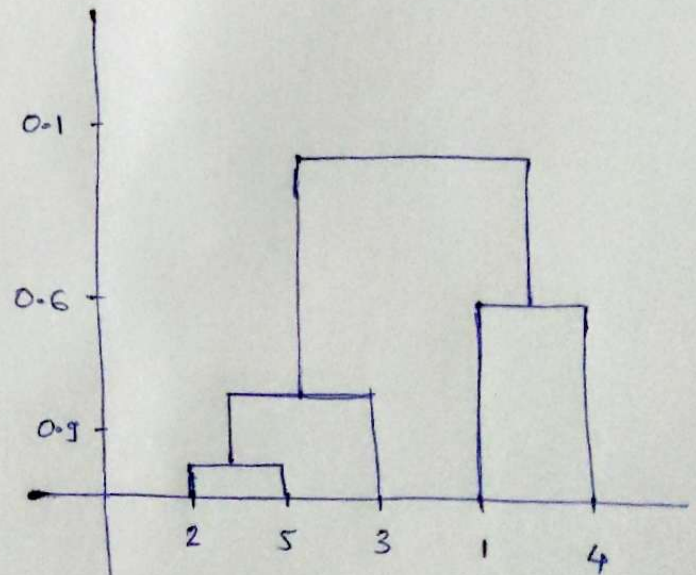
- 1) Run the algorithm in 1st Pass with given threshold as mean of all the consecutive distances
- 2) Now change the threshold to get the predefined k clusters
- 3) Also after each pass we can calculate total SSE and by modifying threshold values we can check which configuration giving minimum SSE.

16) Single and complete link hierarchical clustering.

	$P_1$	$P_2$	$P_3$	$P_4$	$P_5$
$P_1$	1.00	0.10	0.41	0.55	0.35
$P_2$	0.10	1.00	0.64	0.47	0.98
$P_3$	0.41	0.64	1.00	0.44	0.85
$P_4$	0.55	0.47	0.44	1.00	0.76
$P_5$	0.35	0.98	0.85	0.76	1.00



a) Single Link



b) Complete Link