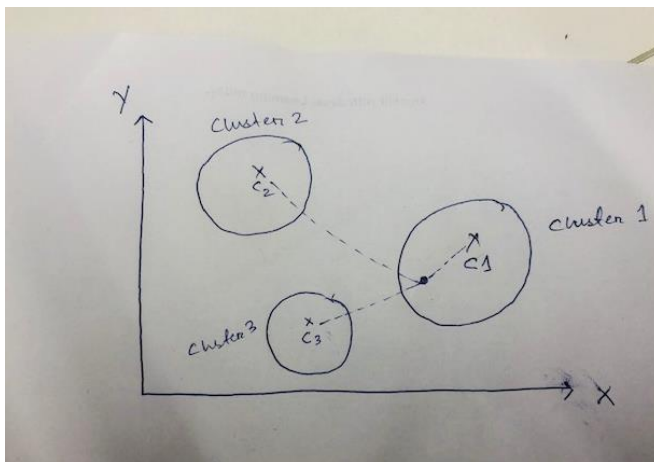


Assume we have completed a k-means algorithm for two columns of data say X and Y, and we obtained two clusters as shown in figure. Blue cluster is cluster 1 with centroid C1 and cluster 2 is the red cluster with centroid C2.



Let us take one single observation as shown in the picture. Now the distance between this observation and its own centroid that is C1 is computed, and the distance between this observation and the nearest neighbour that is basically cluster 2 with c2 centroid is computed.

If there are 3 clusters, say c1, c2 and c3 then the distance between this observation 1 and its own cluster i.e c1 is measured. And out of this C2 and C3 whichever is closer to this observation only that distance will be taken into consideration. The observation belongs to that cluster where the distance is less according to the law of distance. Here it is evident for us visually that the distance to c1 cluster centroid is lesser than c2 cluster centroid. Hence mapping of observation 1 to cluster 1 was correctly done by the k-means algorithm.

There is a simple way to compute a standardized metric to arrive at this conclusion. The following formula that is:

$$\text{Sil-Width} = \frac{b-a}{\max(a,b)}$$

b = distance between observation and the neighbour cluster centroid (c2)

a = distance between observation and its own cluster centroid(c1)

Assume b = 5 and a = 3

$$\text{Sil - width} = (5-3) / 5 = 2/5$$

$$= 0.4$$

If sil-width is a positive value, then we say the mapping of the observation is correct to its current centroid. And we will get a negative value, if distance b is less than distance a.

So for example if we take another observation 2 where distance between observation and the neighbour cluster centroid is lesser than distance between observation and its own cluster centroid means if $b < a$, then sil-width will return a negative value.

If $b=2$ and $a=5$,

$$\text{Sil-width} = (2-5) / 5 = -0.6$$

And sil-width can have minimum of -1 and maximum of +1

Now, if we take another observation 3 where it is place right on top of C1, a will be zero and b will be the distance from c2 to observation 3. In this case,

$$\text{Sil-width} = b / \max(a,b) = 1$$

And if we take another observation 4 where it is place right on top of C2, b will be zero and a will be the distance from c1 to observation 4. In this case,

$$\text{Sil-width} = -a / \max(a,b) = -1$$

Now if we have 10 observations, we can calculate sil-width for each observation, and when we take the average of sil-widths that is called as silhouette score for a dataset.

If the silhouette score is close to +1 then we can say the clusters are well separated from each other on an average.

If the silhouette score is close to 0, then we can say the clusters are not separated from each other.

If the silhouette score is close to -1 then we can say the model has done a blunder in terms of clustering the data

