

## Finding Optimal Number of clusters in k-Means

Even before we start the algorithm of k-Means we need to furnish 'k'. (number of clusters)

Now, how many clusters are there is not known in advance.

So k value is not known in advance  
# clusters

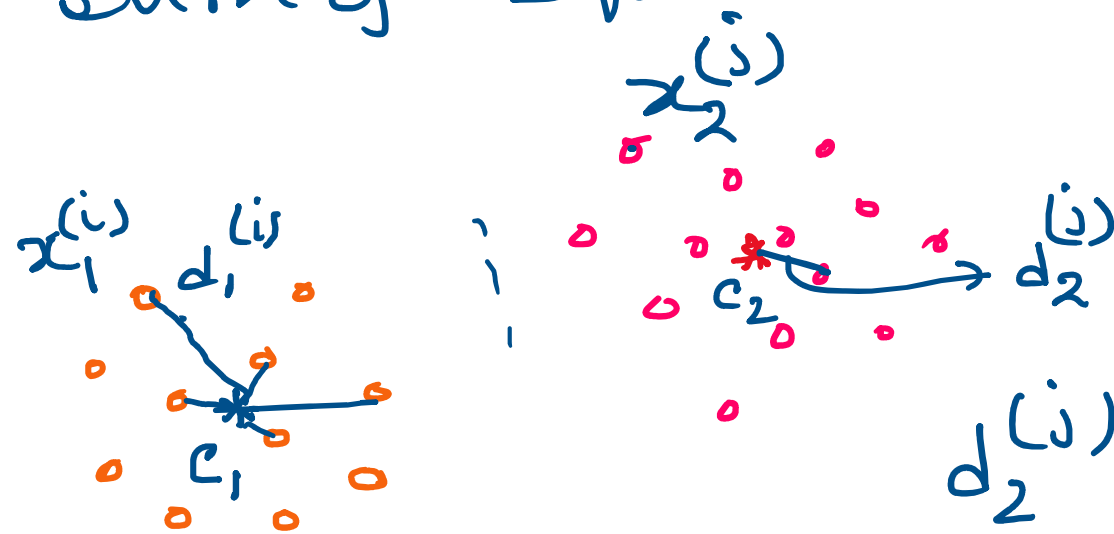
$f_1$	$f_2$	$f_3$	...	$f_n$	
..		.	-	-	
-	-	-	-	-	
-	-	-	-	-	

How can we select optimum number of cluster?

Elbow method

## Elbow method:-

1. First choose a range of values of  $k$ . (say,  $k = 2, 3, 4, \dots, 10$ )
2. For each value of  $k$  obtain the clustering.
3. Calculate sum of squared distance (SSD) from the cluster.



$$d_1^{(i)} = \text{dist}(x_1^{(i)}, c_1)$$

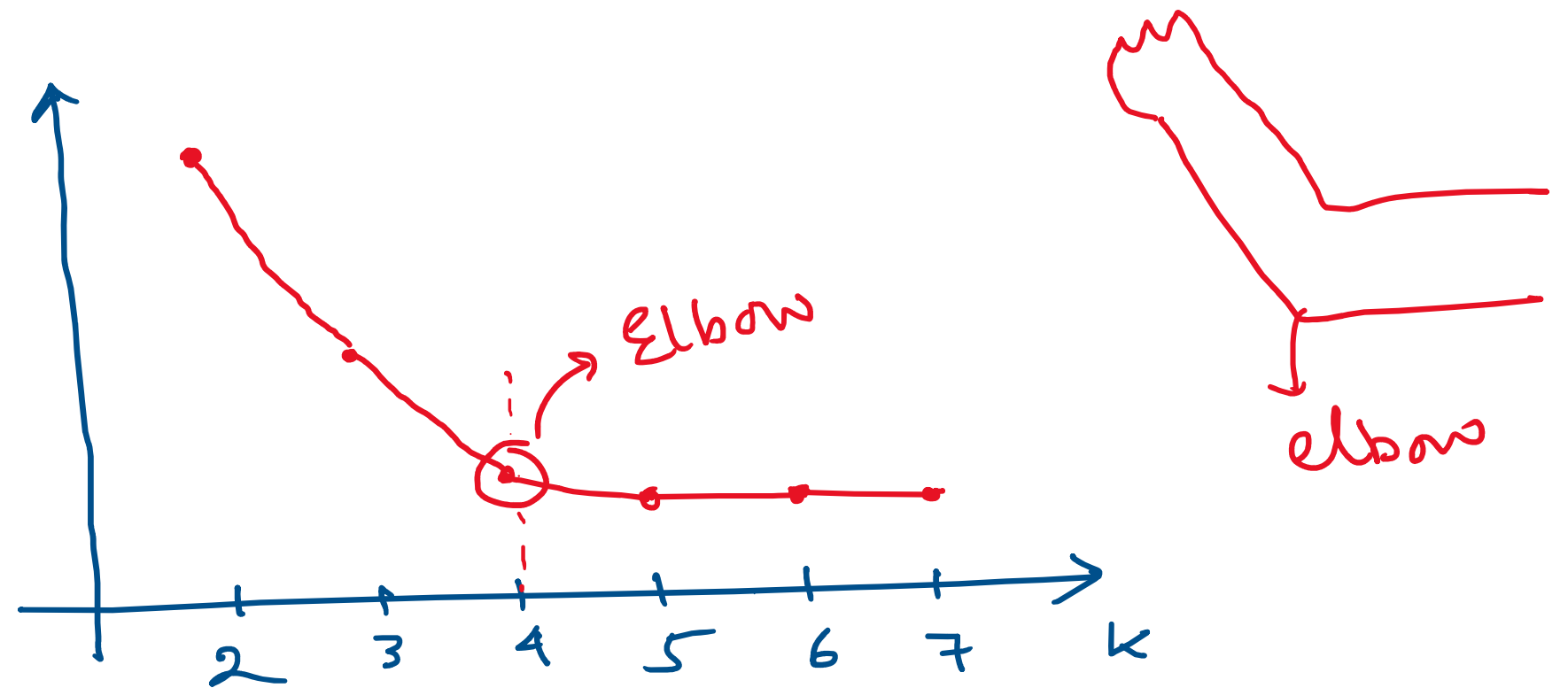
$$d_2^{(i)} = \text{dist}(x_2^{(i)}, c_2)$$

$$\underline{SSD_1} = \sum_i (d_1^{(i)})^2 = \sum_{x_i \in C_1} [\text{dist}(x_i, c_1)]^2$$

$$\underline{SSD_2} = \sum_j (d_2^{(j)})^2$$

$$SSD_1 + SSD_2 = SSD$$

$k$	$SSD$
2	:
3	.
4	.
5	.
:	.
10	.



4. Plot  $SSD$  vs  $k$

5. We will choose the value of  $k$  where we get the elbow.

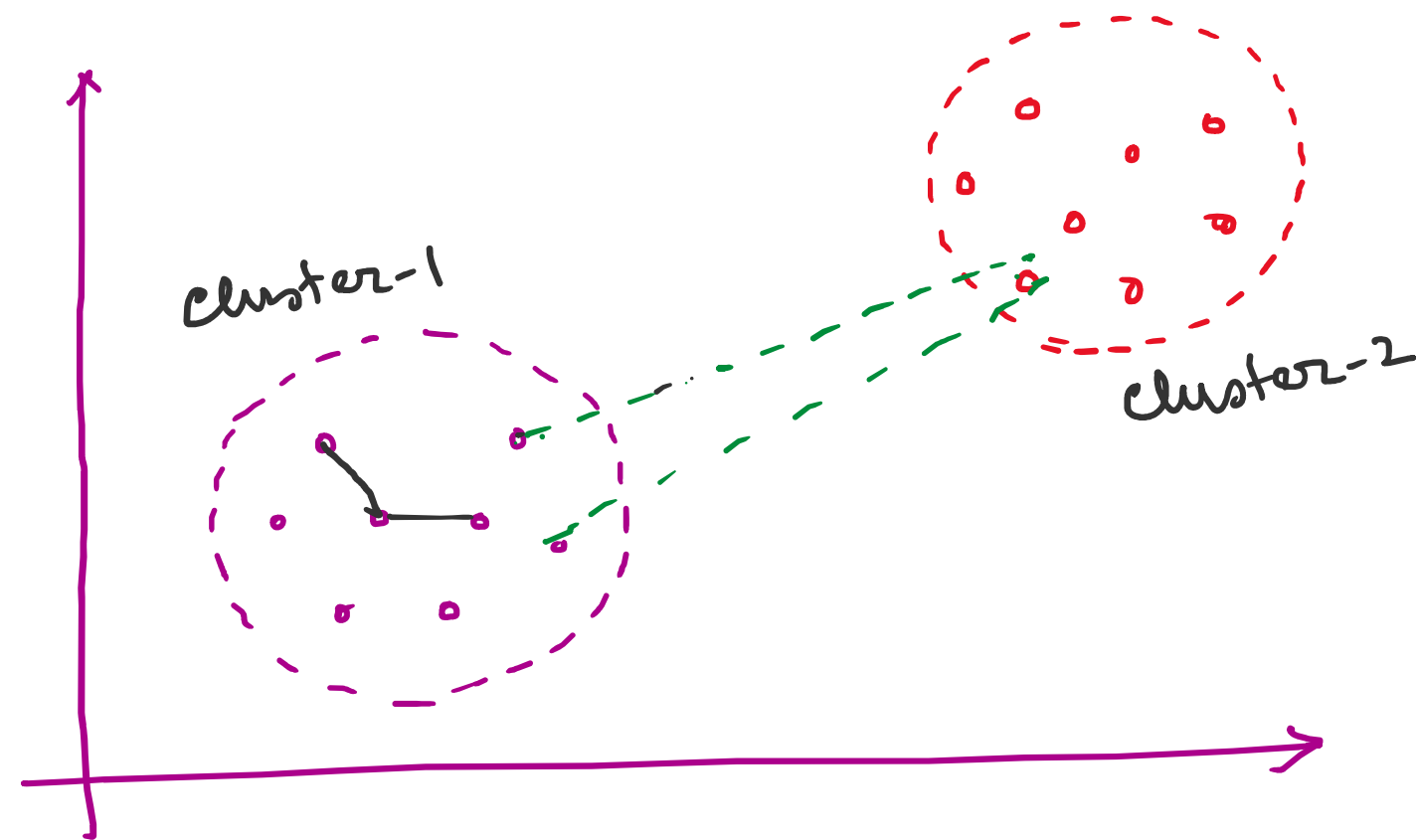
## How to measure goodness of clustering:-

### Silhouette Coefficient

It is a metric used to calculate the goodness of a clustering technique.

Its value ranges from  $-1$  to  $+1$

- $1$ : Clusters are well apart from each other & are clearly distinguished.
- $0$ : Clusters are overlapping.
- $-1$ : Clusters are assigned in the wrong way.



$$\text{Silhouette Score} = \frac{b - a}{\max(a, b)}$$

$a$  = average intra-cluster distance, i.e. average of the distances between each point within a cluster.

$b$  = average inter-cluster distance, i.e. average of the distances between the points of clusters.