

Clustering

Unsupervised machine learning:

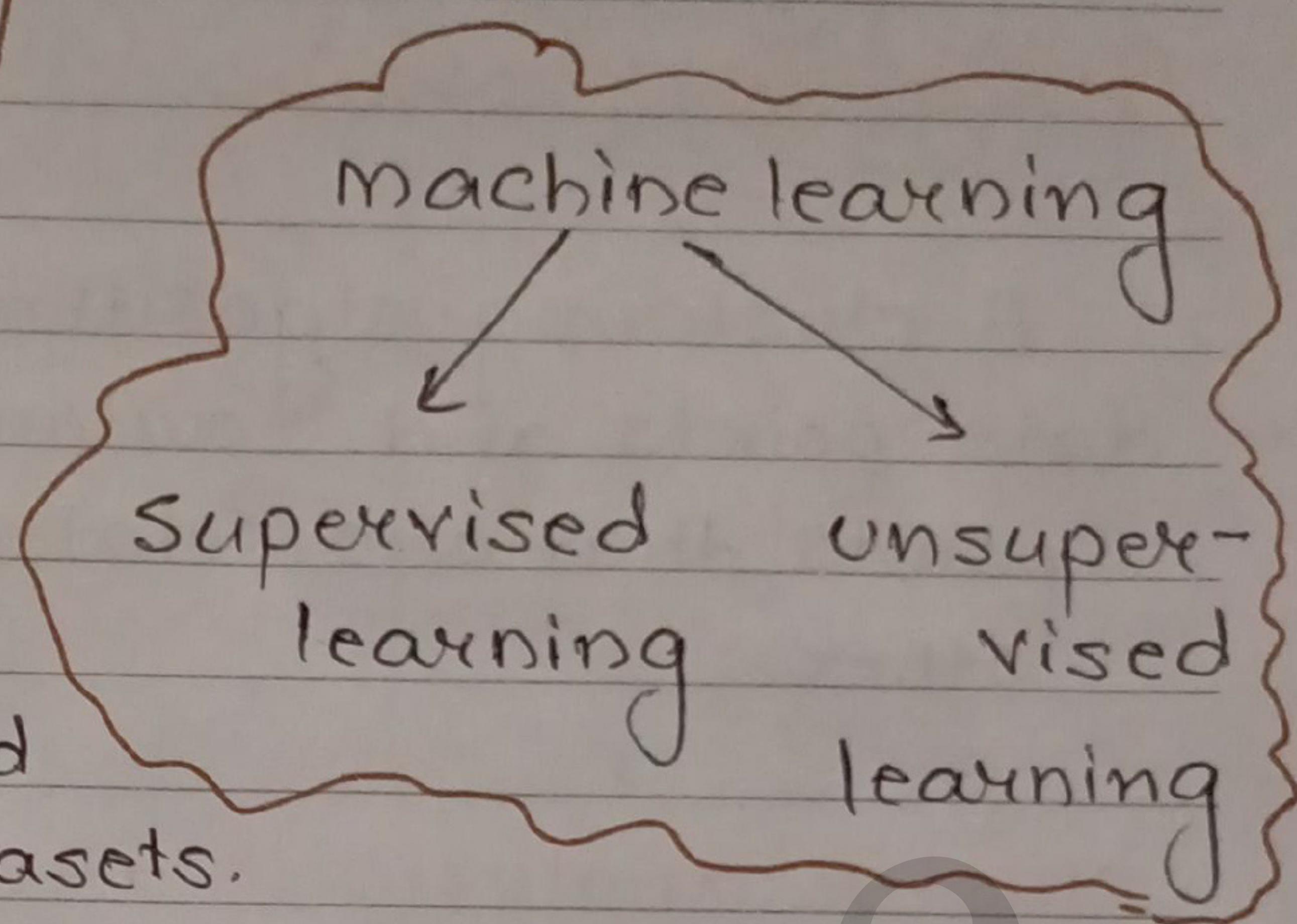
UNSUPERVISED ML uses ML algorithm to analyze and cluster unlabelled datasets.

These algorithms discovers hidden patterns as data groups without the need of human interaction.

Its ability to discover similarities and differences in information make it the ideal solution for EDA, cross-selling strategies, customer segmentation, and image recognition.

Unsupervised learning models are used for these main tasks

- clustering
- association
- dimensionality reduction



What is clustering?

A clustering algorithm looks at a number of data points and automatically finds data points that are related or similar to each other.

Cluster analysis is a technique used in data mining and ML to group similar objects into clusters.

K-means clustering → the aim is to partition a set of objects into k-clusters in such a way that the sum of the squared distances between the objects into k-clusters in such a way and their assigned cluster mean is minimized.

" Clustering is the process of dividing the entire data into groups (also known as clusters) based on the patterns in the data.

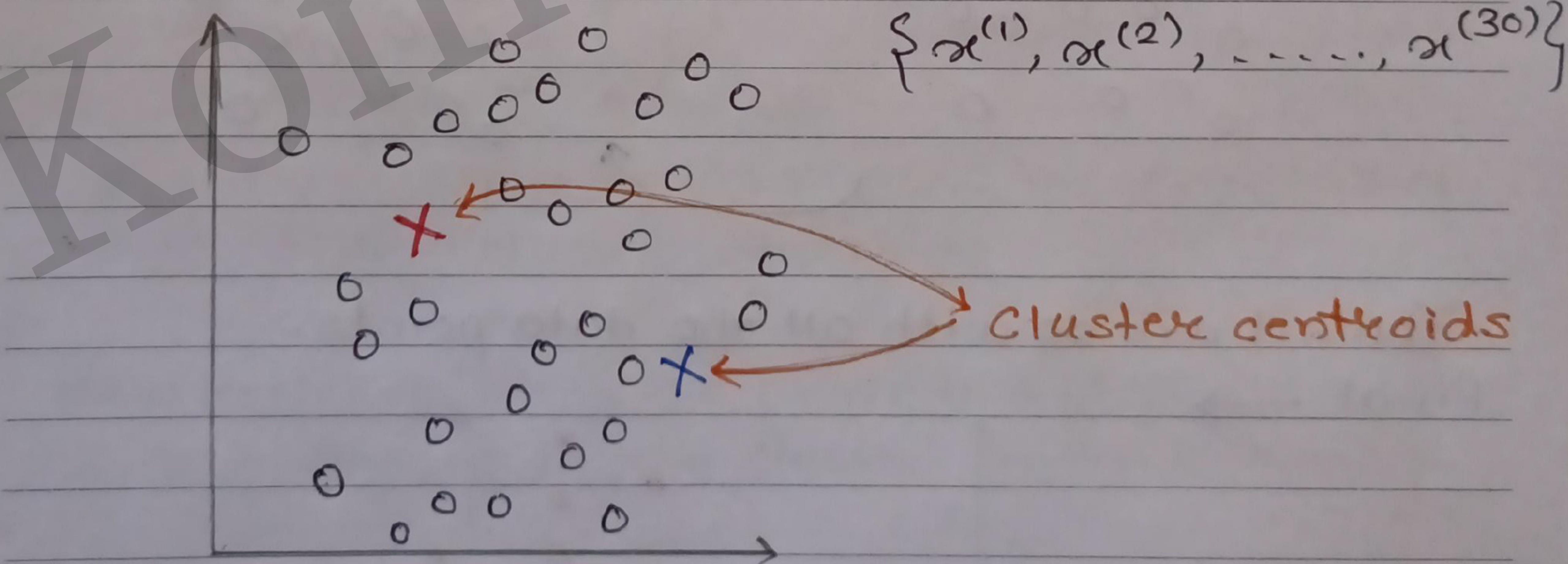
K-means Intuition:

plotted a data set with
30 unlabeled training
examples →

K-means will repeatedly do
two different things.

1. assign points to cluster centroids
2. move cluster centroids

- ① Now, randomly pick two points, at where
might be the centers of two different
clusters.



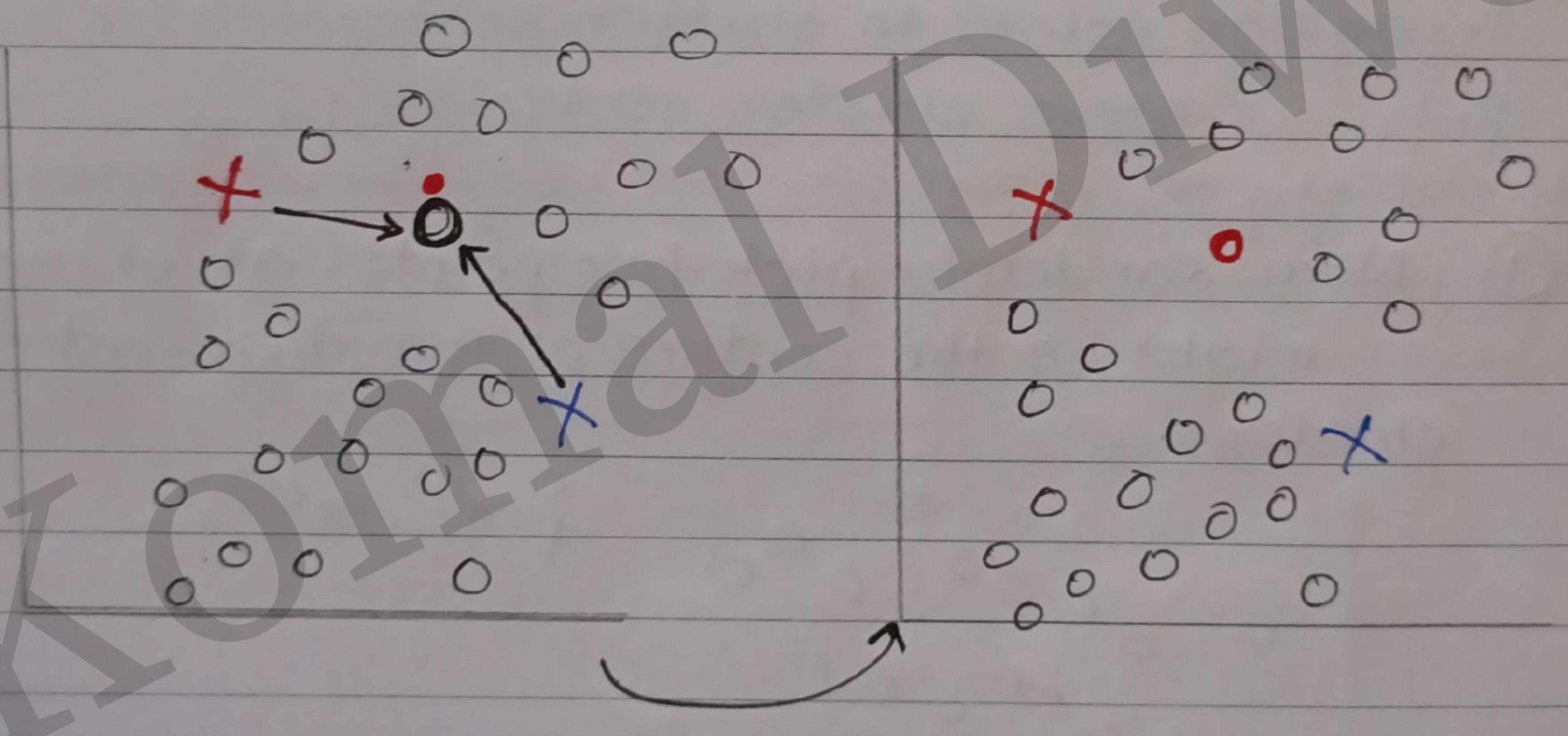
* this is just a random
guess.

- After it's made an initial guess at where the

cluster centroids is, it will go through all of these examples, $x^{(1)}, \dots, x^{(30)}$ datapoints.

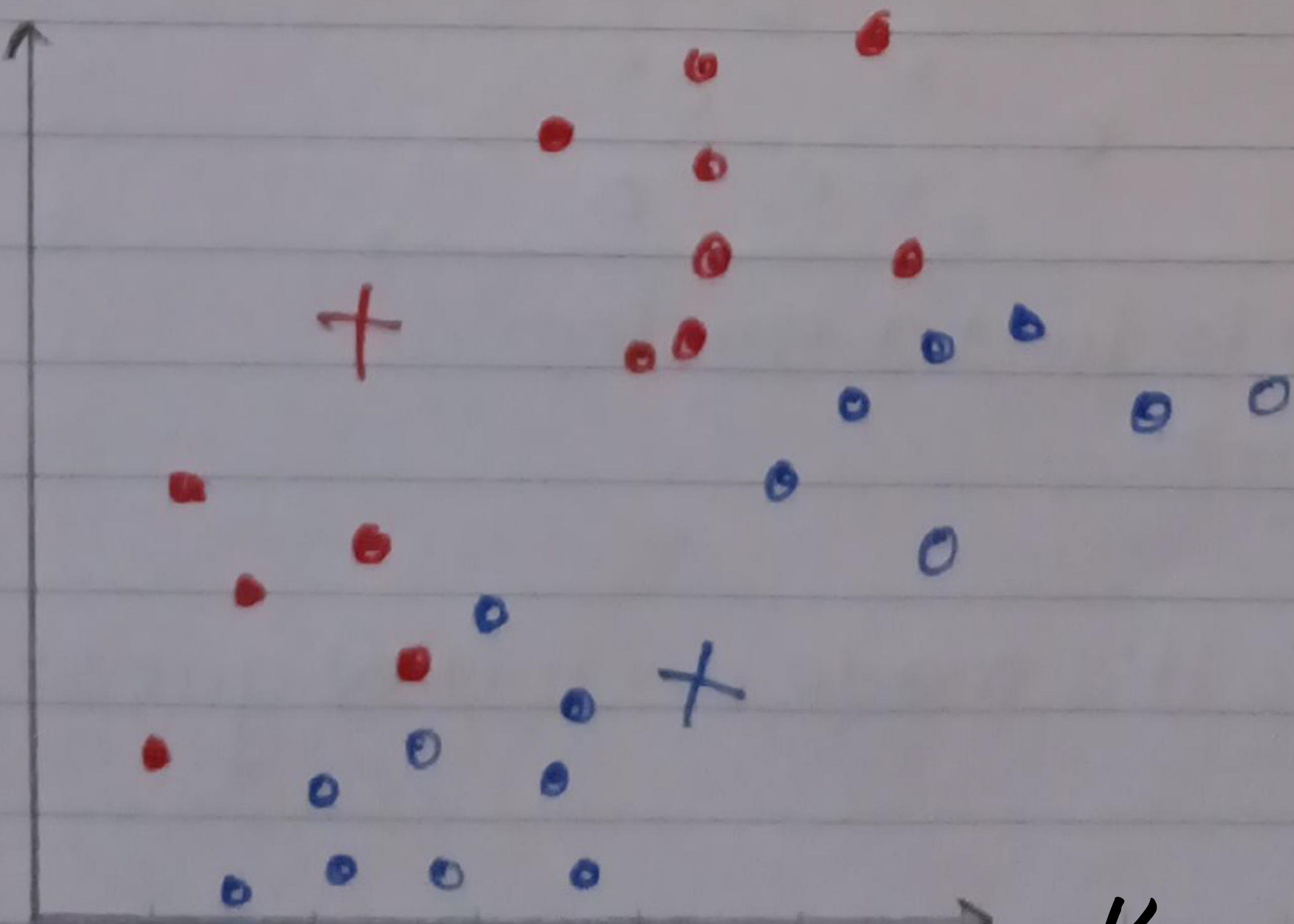
And for each of them it will check if it is closer to the fed cluster centroid or blue cluster centroid.

Step 01: assign each point to its closest centroid.



Thus, it will do with all the data points.

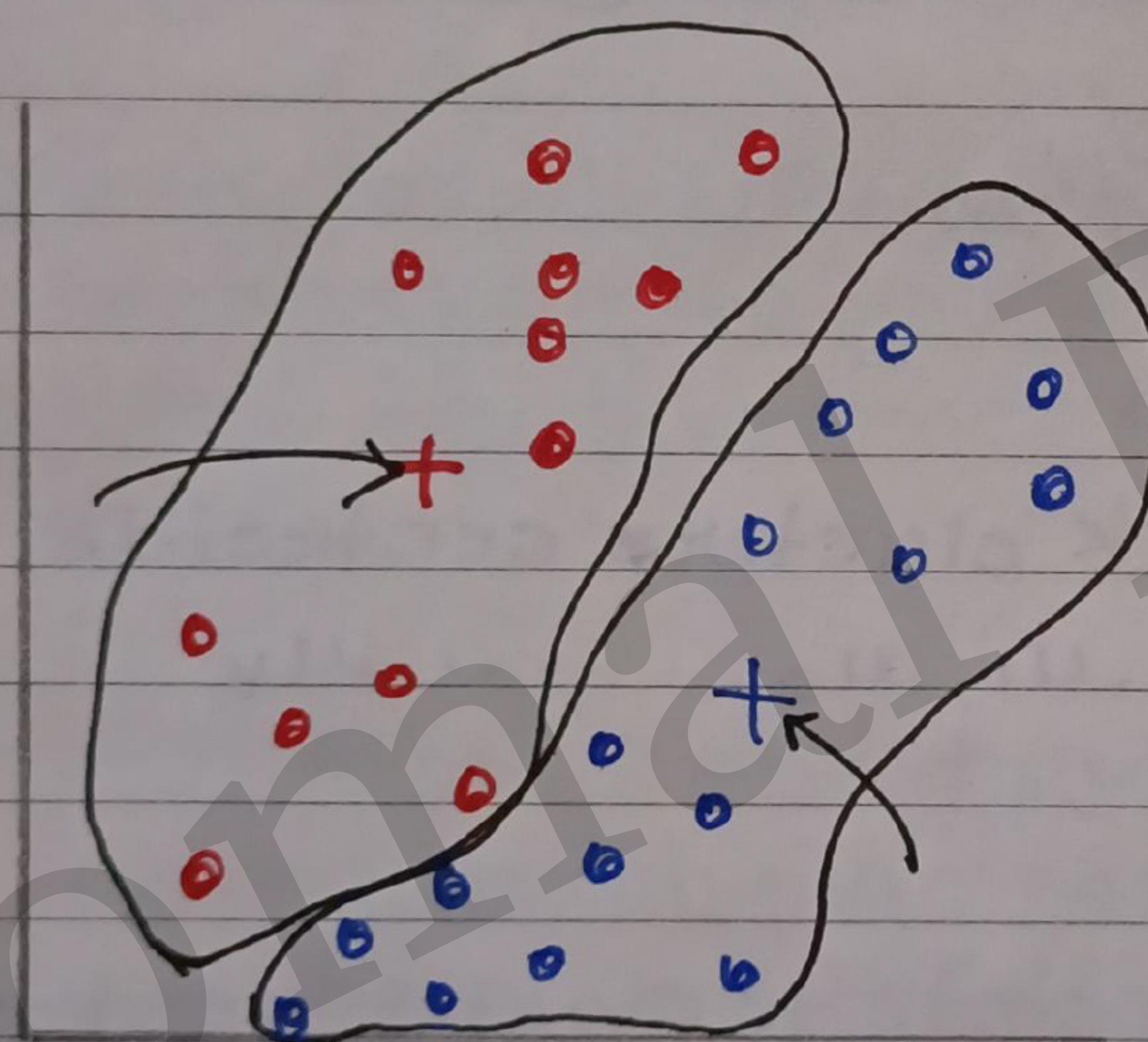
Final \rightarrow



Step 02: Recompute the centroids

It'll look at all the red points and take an average of them. and it will move the red cross to whatever is the average location of the red dots.

Same thing for blue points/dots.



Slightly improved guesses for the locations of the cluster centroids.

Now keep repeating the above two steps till there is no change in the cluster centroids location.

If we keep on doing those two steps, we find that there are no more changes to the colors of the points or to the locations of the clusters centroids.

And so this means that at this point the k-means clustering algorithm has converged.

Because applying those two steps over and over, results in no further changes to either the assignment of the point to the centroids or the location of the cluster centroids.

K-means Algorithm

Randomly initialize K clusters' centroids
 $\mu_1, \mu_2, \dots, \mu_K$

Repeat {

assign points to cluster centroids
for $i=1$ to m

$c^{(i)} := \text{index (from } 1 \text{ to } k \text{) of cluster centroid closest to } x^{(i)}$

$$\min_k \|x^{(i)} - \mu_k\|$$

move cluster centroids

for $k=1$ to k

$\mu_k := \text{average (mean) of points assigned to cluster } k$

}

1. Randomly initialize k cluster centroids, $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k$, where k is the number of clusters.
2. Assign points to clusters by computing the distance between each point and each centroid and assigning the point to the closest centroid.
3. Update the cluster centroids by computing the mean of all the points assigned to each cluster, which becomes the new centroid location.
4. Repeat steps 2 and 3 until the cluster assignments no longer change. or a maximum number of iterations is reached.

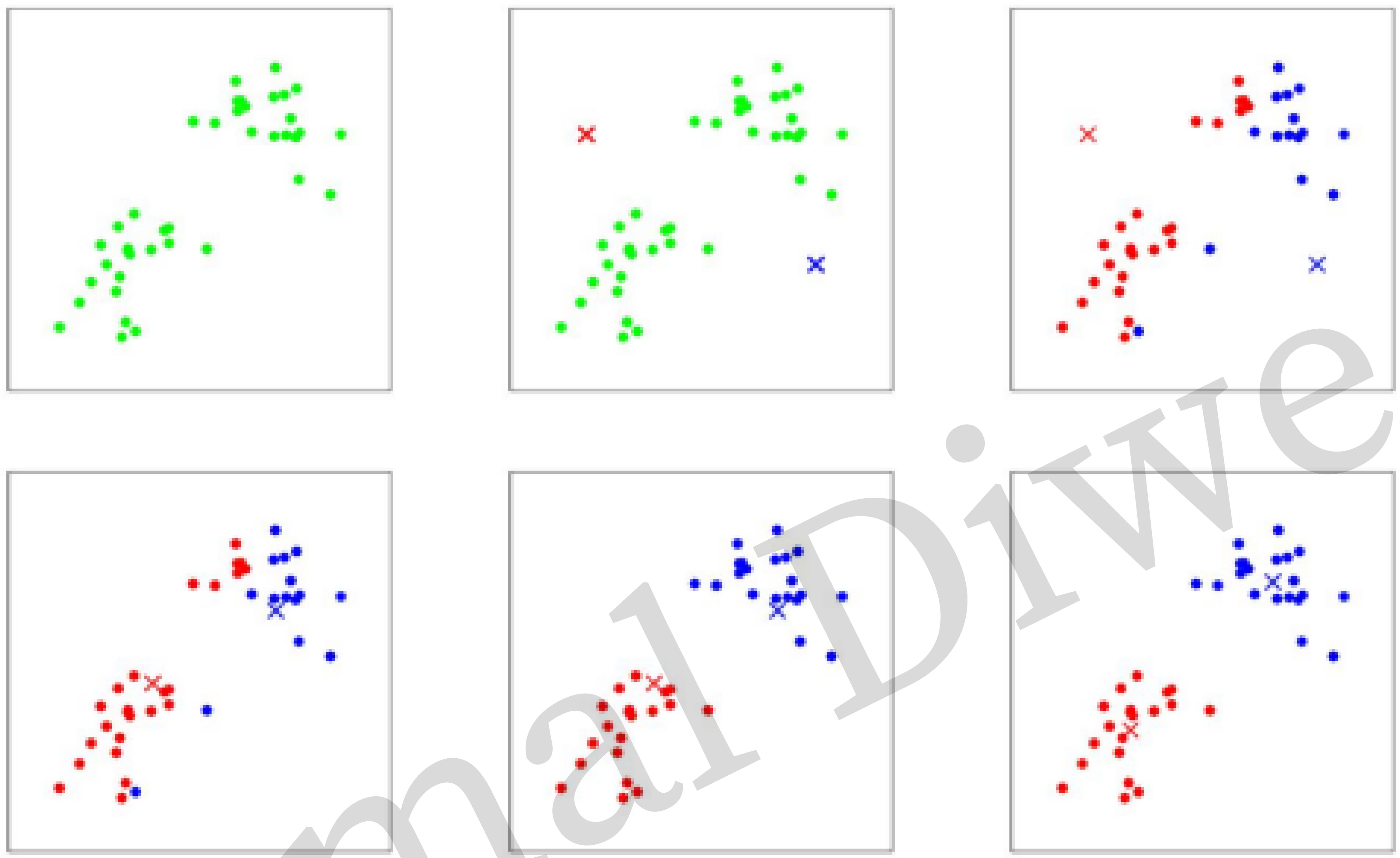


Figure 10.1: K-means algorithm. Training examples are shown as dots, and cluster centroids are shown as crosses. (a) Original dataset. (b) Random initial cluster centroids (in this instance, not chosen to be equal to two training examples). (c-f) Illustration of running two iterations of k -means. In each iteration, we assign each training example to the closest cluster centroid (shown by “painting” the training examples the same color as the cluster centroid to which is assigned); then we move each cluster centroid to the mean of the points assigned to it. (Best viewed in color.) Images courtesy Michael Jordan.

K-means Optimization objective:

$c^{(i)}$ = index of cluster ($1, 2, \dots, k$) to which example $\alpha^{(i)}$ is currently assigned

u_k = cluster centroid k

$u_{c(i)}$ = cluster centroid of cluster to which example $\alpha^{(i)}$ has been assigned

cost function

$$J(c^{(1)}, \dots, c^{(m)}, u_1, \dots, u_k) = \frac{1}{m} \sum_{i=1}^m \| \alpha^{(i)} - u_{c(i)} \|^2$$

$$\min_{\substack{c^{(1)}, \dots, c^{(m)} \\ u_1, \dots, u_k}} J(c^{(1)}, \dots, c^{(m)}, u_1, \dots, u_k)$$

This function is also known as distortion function.

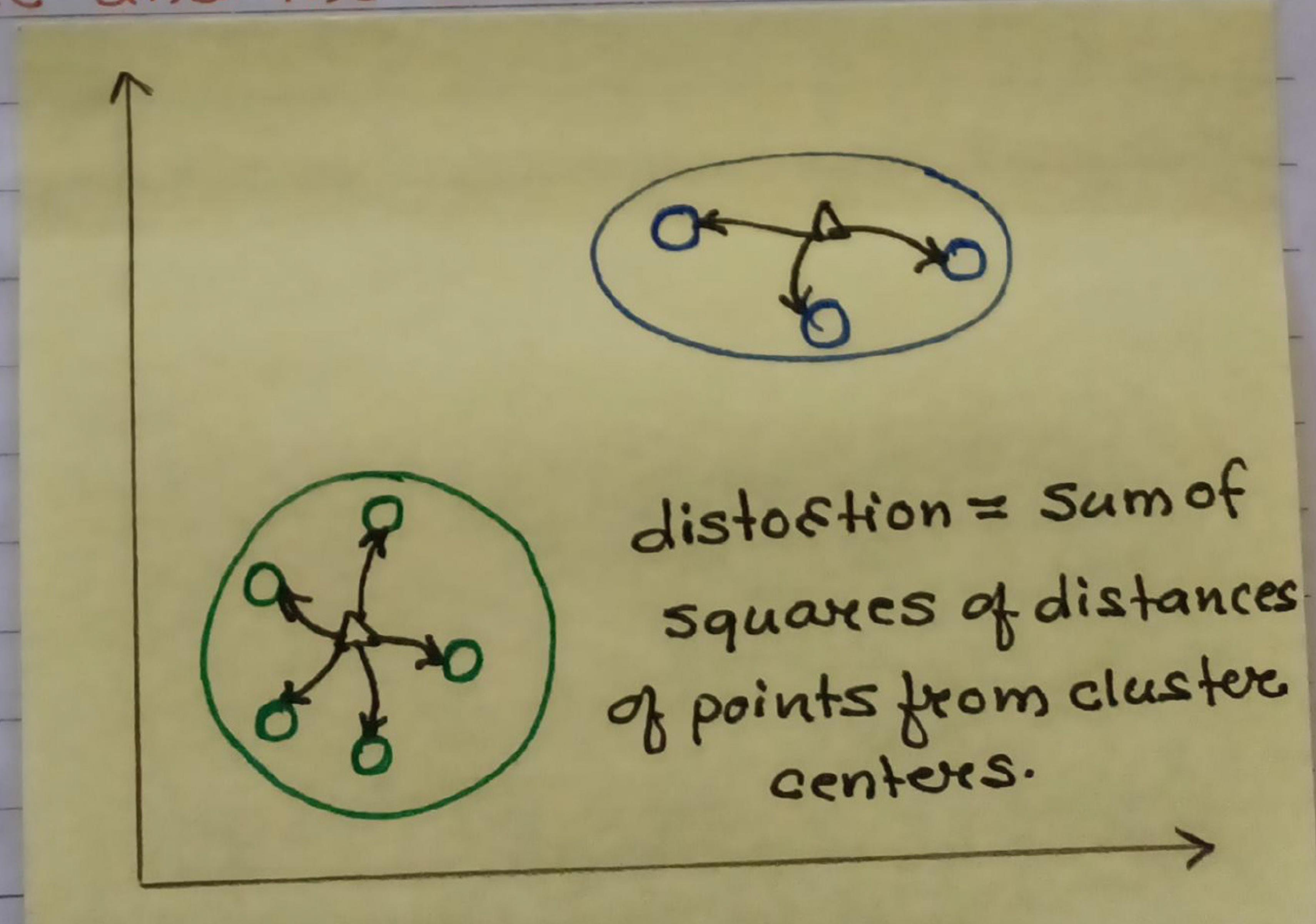
J measures the sum of square distances between each training example $\alpha^{(i)}$ and the cluster centroid $u_{c(i)}$ to which it has been assigned.

The k-means algorithm is optimizing a specific cost function called the **distortion function**.

This function measures the average square distance between every training example and the location of the cluster centroid to which it has been assigned.

The algorithm first updates the cluster assignments to minimize the cost function while holding the centroid locations constant, and then updates the centroid locations to minimize the cost function while holding the cluster assignments constant.

During the first step, the algorithm assigns a training example to the closest centroid to minimize the distance or square distance betn the example and the centroid.



Initializing K-means

Random initialization:

choose $k < m$ \rightarrow no. of cluster center's.
 $\frac{1}{m}$ \rightarrow training examples

For $i = 1$ to $100 \{$

Randomly initialize k-means.

Run k-means. Get $c^{(1)}, \dots, c^{(m)}, u_1, \dots, u_K$
compute cost function (distortion)

$J(c^{(1)}, \dots, c^{(m)}, u_1, \dots, u_K)$

}

Pick set of clusters that gave lowest cost J .

- The first step in k-means clustering is to choose random locations as initial guess.

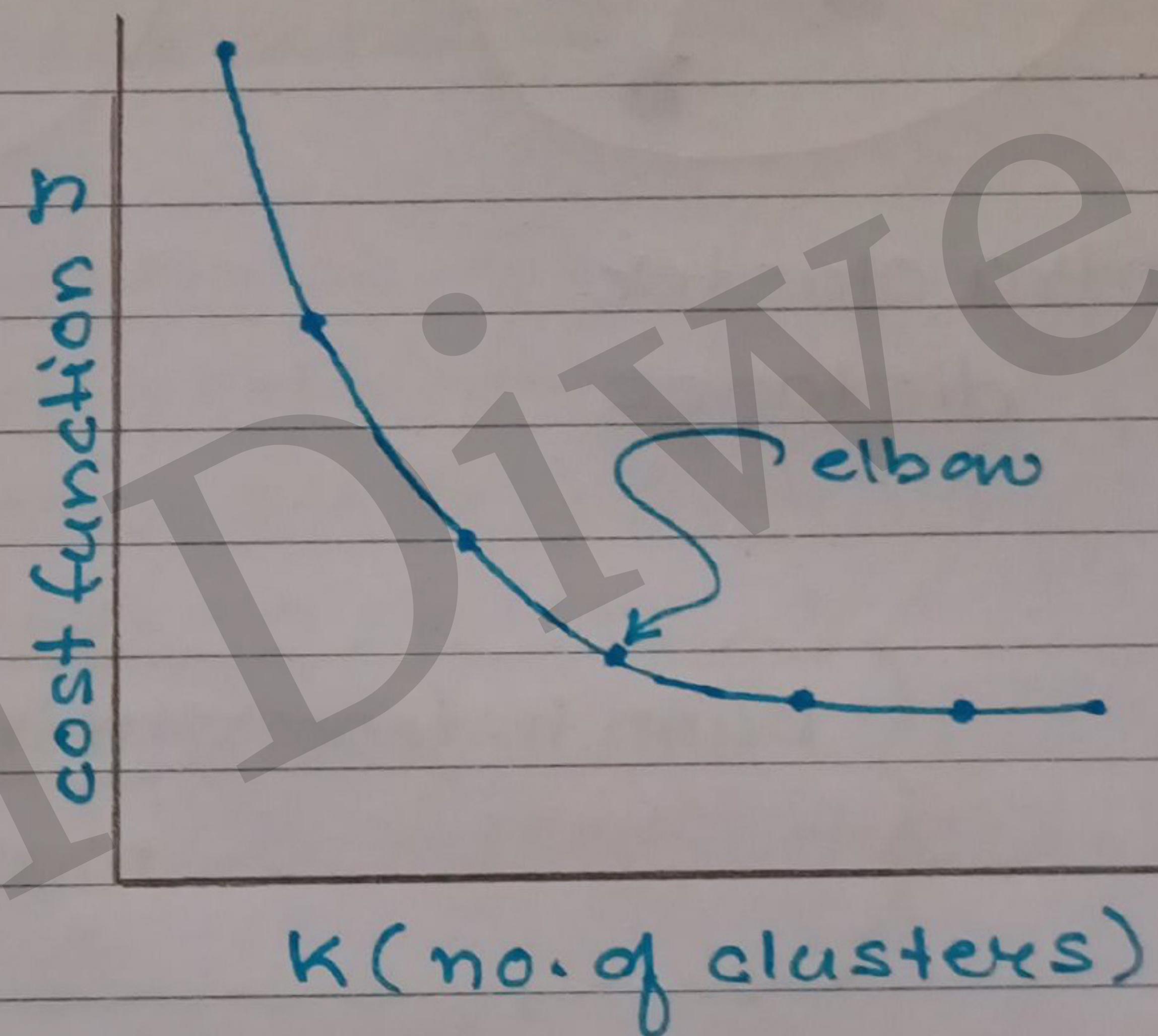
To avoid poor results, k-means can be run multiple times with different random initializations.

After running k-means multiple times, the set of cluster with the lowest distortion or cost function can be selected.

How to find the right k?

- No absolute method to find right number of clusters (k) in k-means clustering.
- Elbow method

The elbow method, which involves plotting the distortion function J as a function of k , is one technique used to determine k .

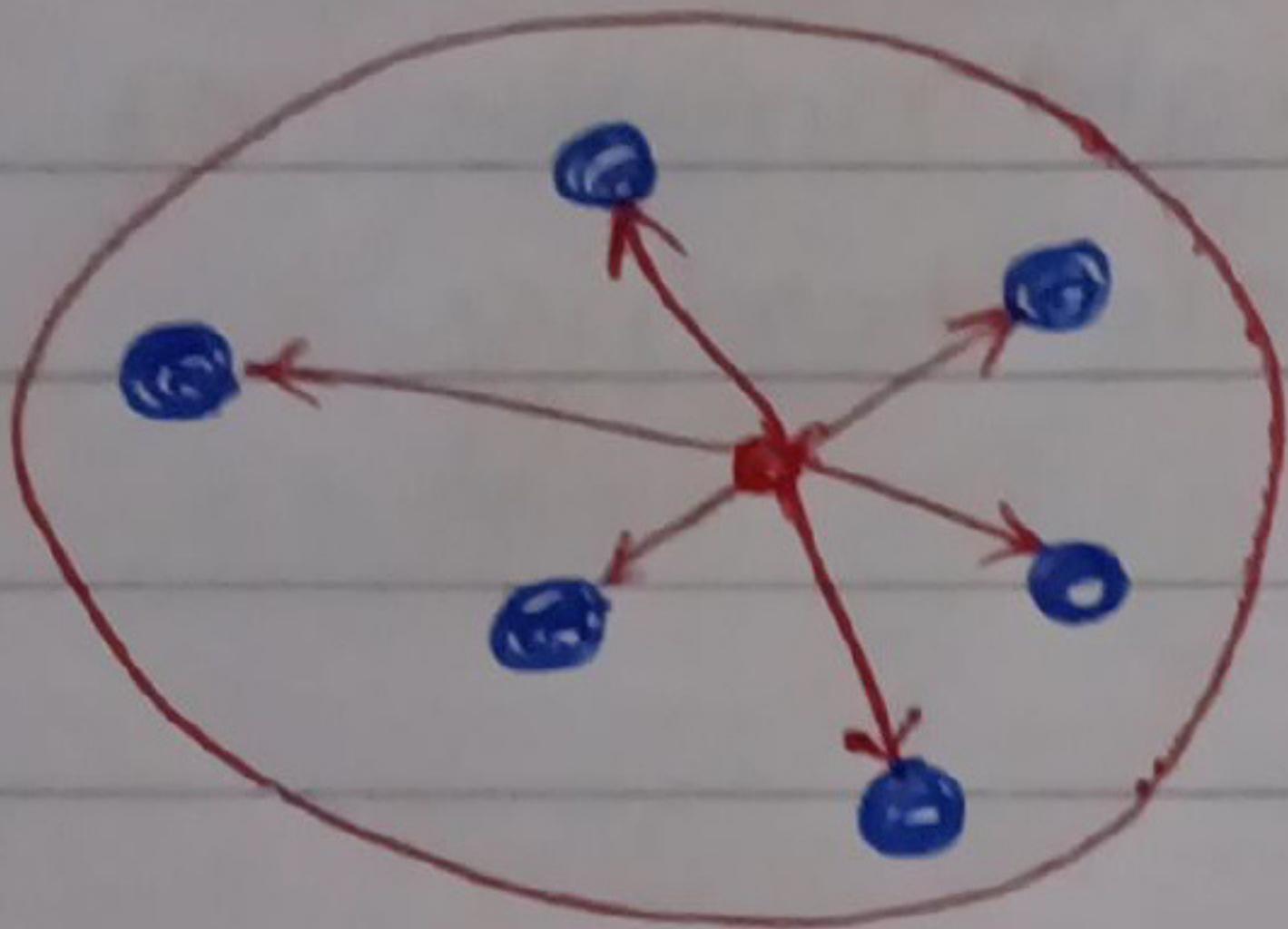


However, for many applications, the appropriate number of clusters is ambiguous, so the downstream purpose should be considered when evaluating k-means.

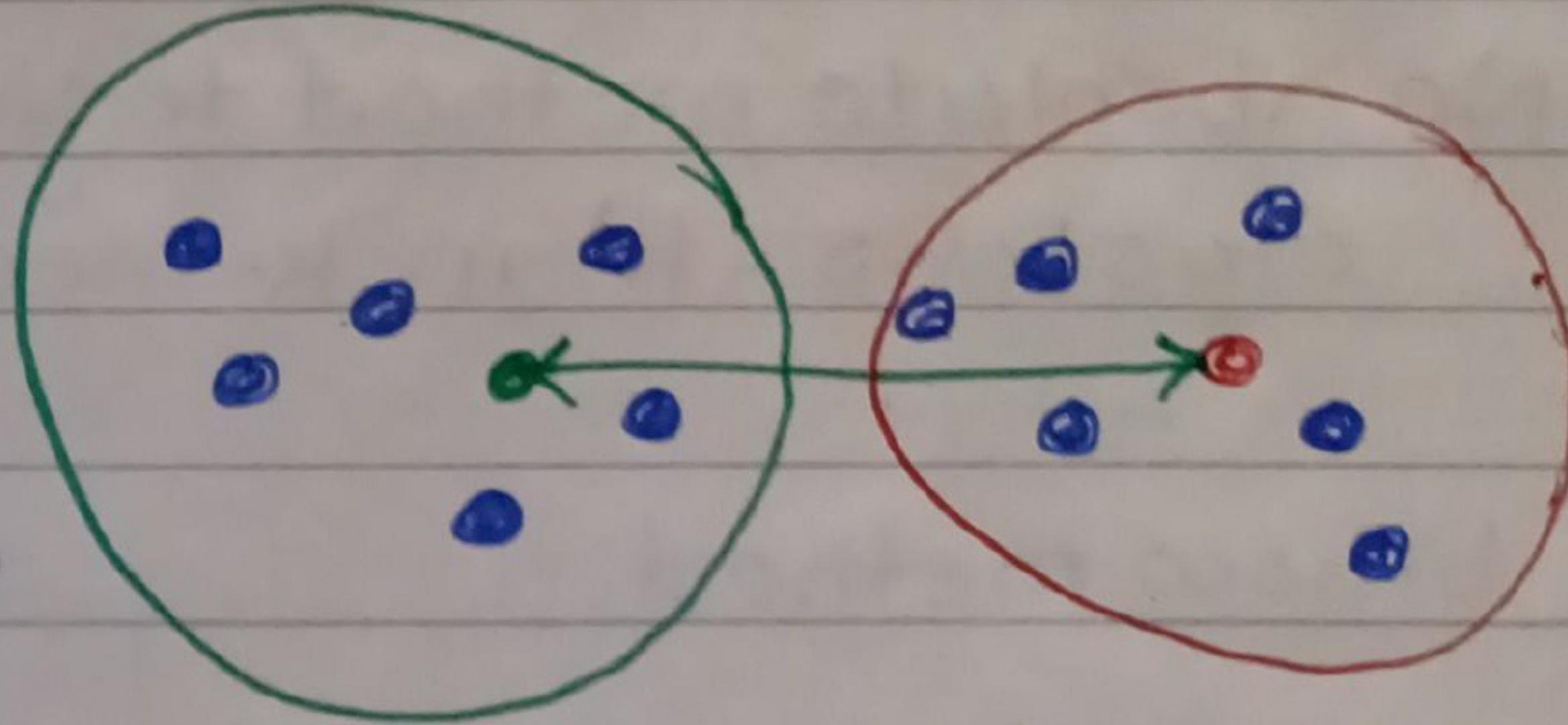
How to validate cluster?

1. Dunn Index (Dunn Index)
2. silhouette score

Dunn Index



Intra cluster
distance



Inter cluster distance

$$\text{Dunn Index} = \frac{\min(\text{Inter cluster distance})}{\max(\text{Intra cluster distance})}$$

Dunn Index is the ratio of the minimum of inter-cluster distances and maximum of intra-cluster distances.

We want to maximize the Dunn Index. The more the value of the Dunn Index, the better the clusters will be.

Let's Intuition behind the dunn Index:

$$\text{Dunn Index} = \frac{\min(\text{Inter cluster distance})}{\max(\text{Intra cluster distance})}$$

clusters are far apart

In order to maximize the value of the Dunn Index,
the numerator should be maximum.

Here, we are taking the minimum of the inter cluster distances. So, the distance between even the closest clusters should be more. which will eventually make sure that the clusters are far away from each other.

$$\text{Dunn Index} = \frac{\min(\text{Inter cluster distance})}{\max(\text{Intra cluster distance})}$$

clusters are compact

also, the denominator should be minimum to maximize the dunn Index.

Here, we are taking the maximum of all intracluster distances.

Again, the intuition is same here. The maximum distance between the cluster centroids and the points should be minimum, eventually ensuring that the clusters are compact.

Silhouette score

The Silhouette score measures the similarity of each point to its own cluster compared to other clusters,

and the silhouette plot visualizes these scores for each sample.

$$\text{Silhouette score} = \frac{b_i - a_i}{\max(a_i, b_i)}$$

where, $a_i \rightarrow$ Intra-cluster distance
 $b_i \rightarrow$ Nearest-cluster distance

Silhouette score ranges from -1 to 1.

A high silhouette score indicates that the clusters are well separated, and each sample is more similar to the samples in its own cluster than to sample in other clusters.

A silhouette score close to 0 suggests overlapping clusters, and a -ve score suggests poor clustering solutions.