# Clustering

Standard problem:

Regression and classification:

We want to find the function that satisfies the training data. In Case of clustering the task is clustering.

Task: It is to group the similar data points.
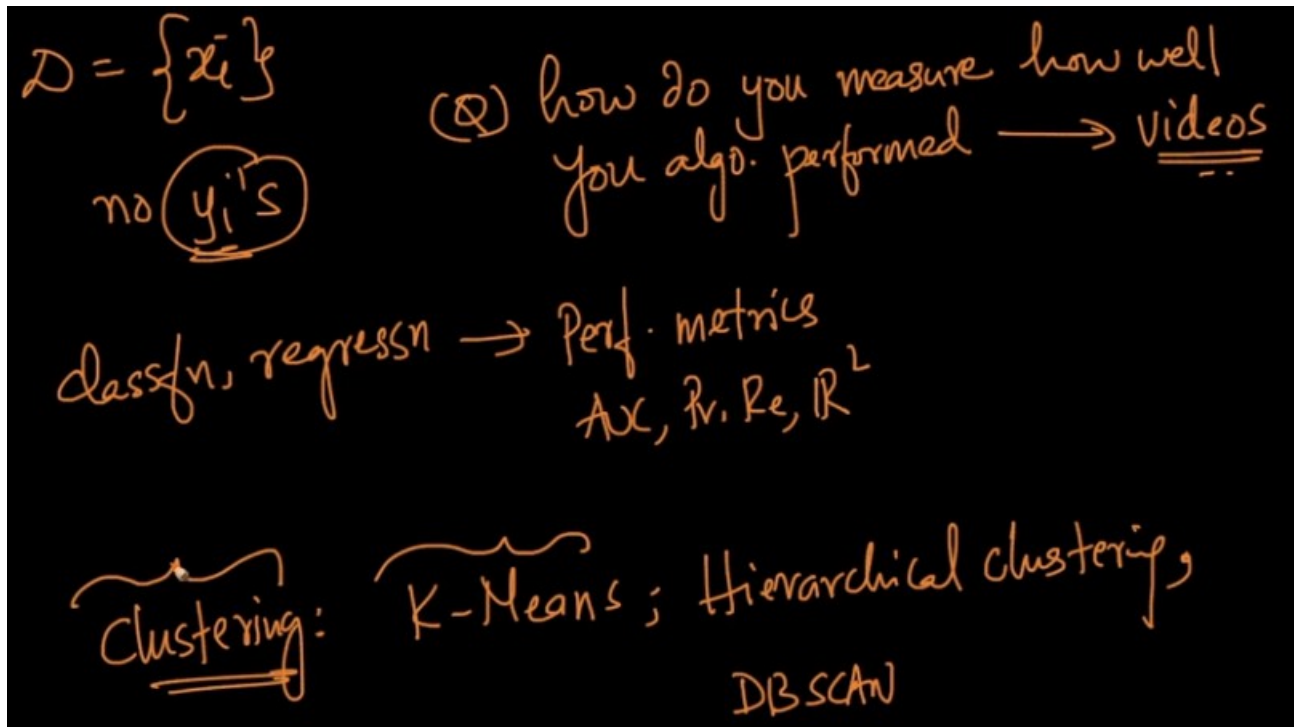


Consider the data points as follows:

The task of clustering is to group the similar data points. All these points are grouped together. The points are grouped together. The similar is very much task specific.

Here all the data given are the data points in the mathematical form.
Measuring the performance of the clustering algorithm. There are various metrics that can measure the performance of clustering.
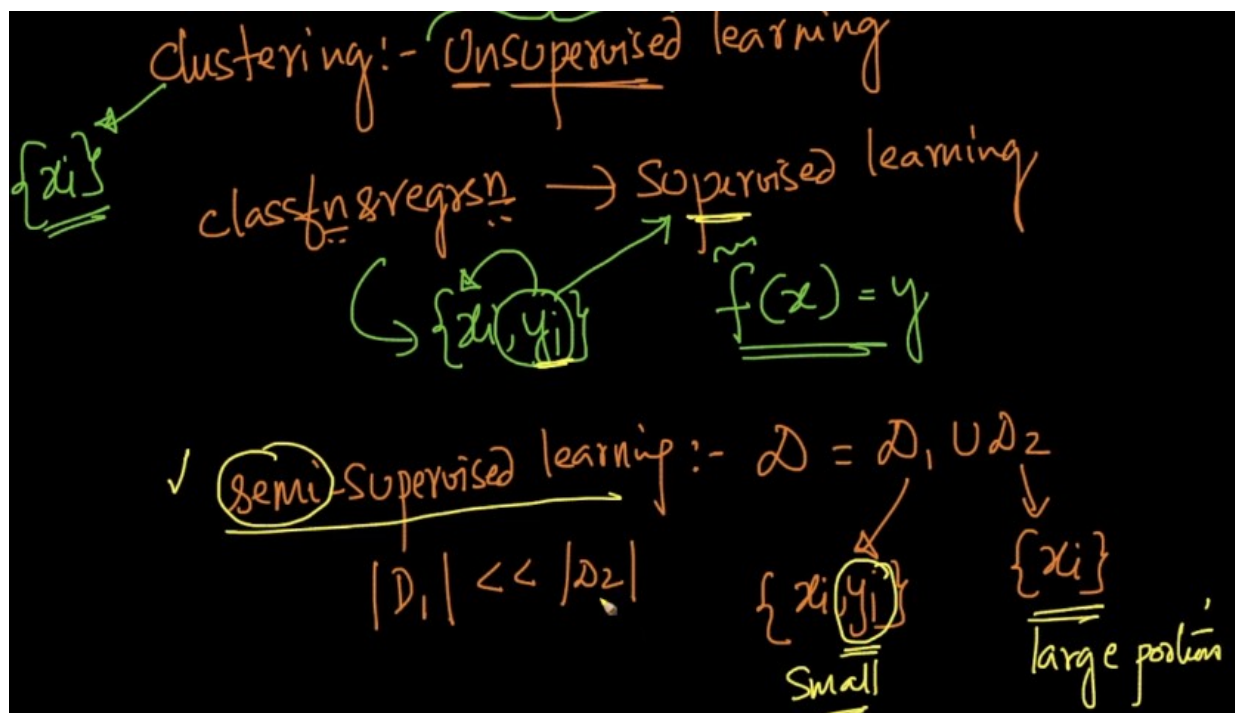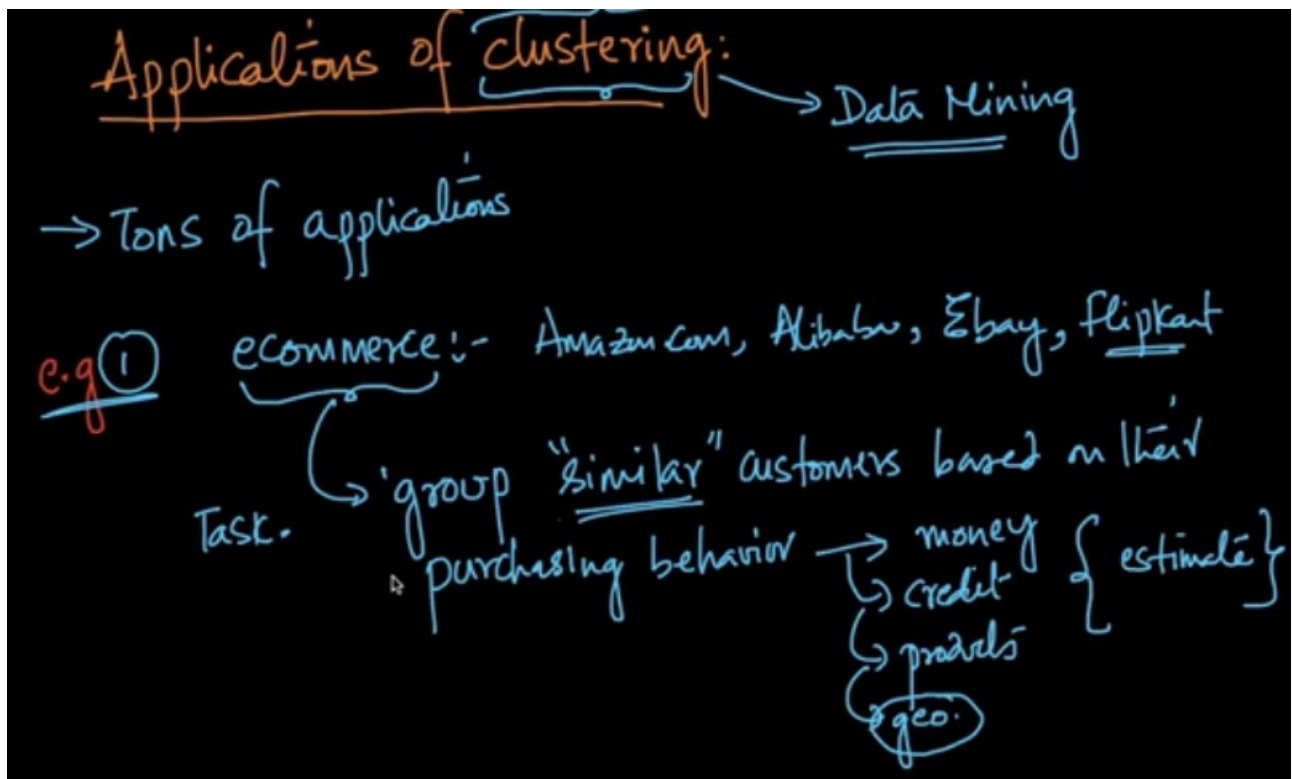Clustering algos: K- Means, Hierarchical, DBSCAN.



These are the most used algos. There are algos for various data.
Unsupervised learning: Clustering is referred to as unsupervised learning. Both classification and regression are the supervised learning.
There is also an area called semi-supervised learning, the data sets are union of D1 and D2. This happens when the cost of labeling the data is expensive.

Applications of clustering: These are more studied in Data mining, than machine learning.
There are many applications of clustering.
They want to group the similar customers based on the purchasing behavior.
**They can decide the income level using the purchasing behavior.**



Assume the clusters are grouped to various classes. By grouping customers we can give different offers.
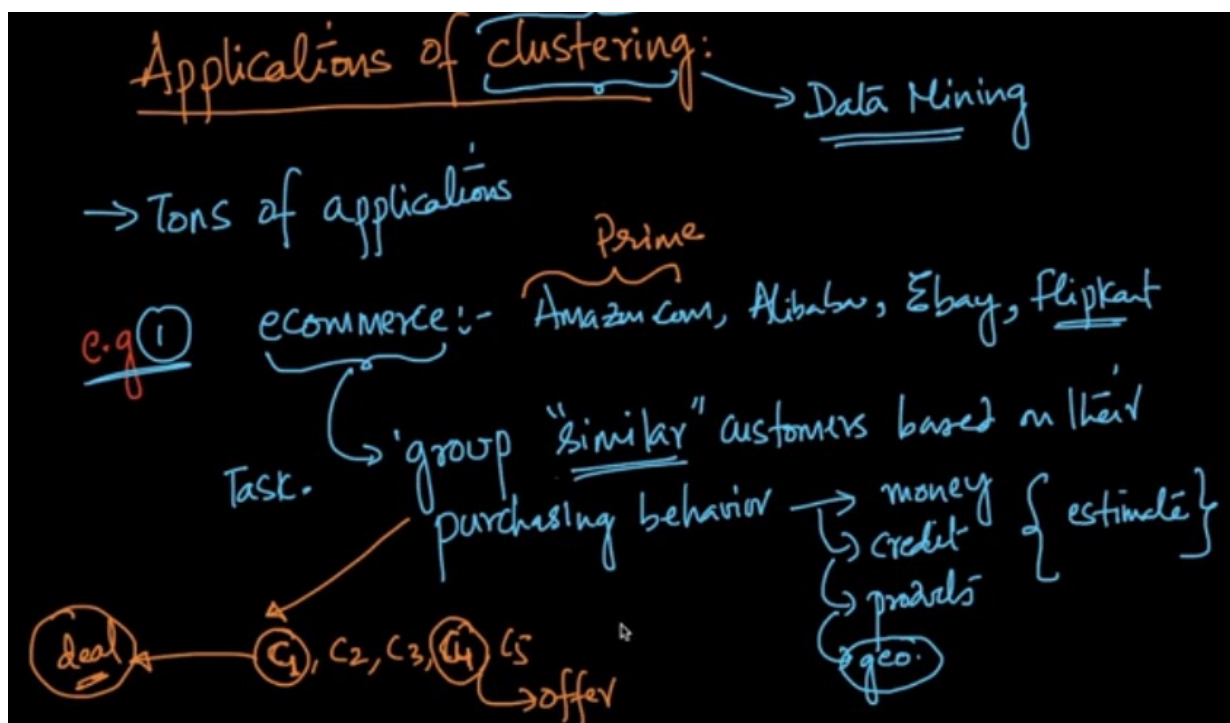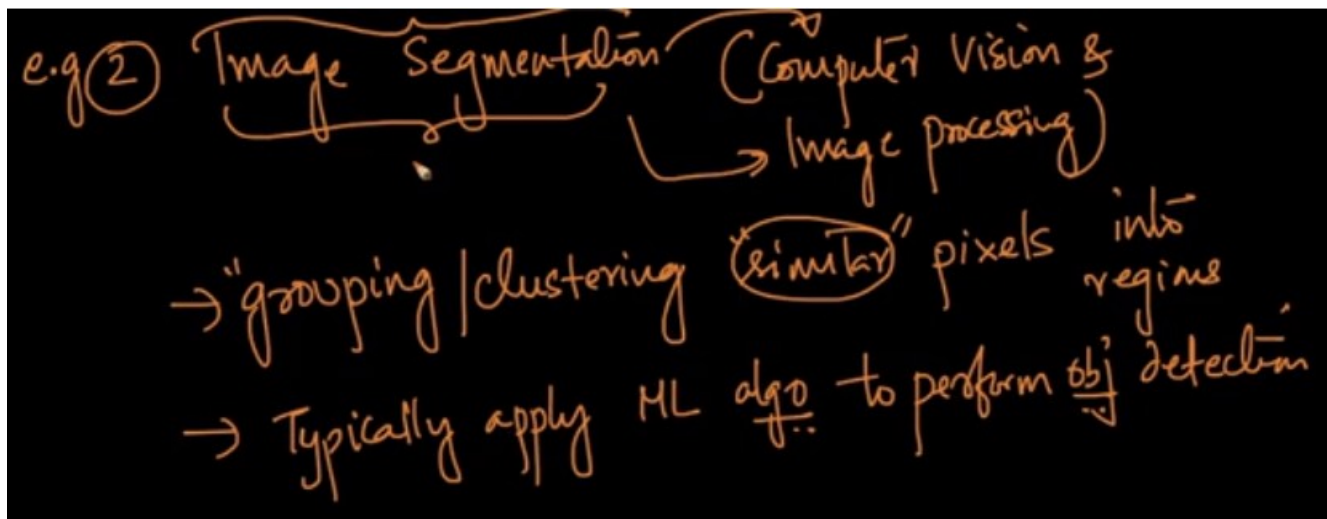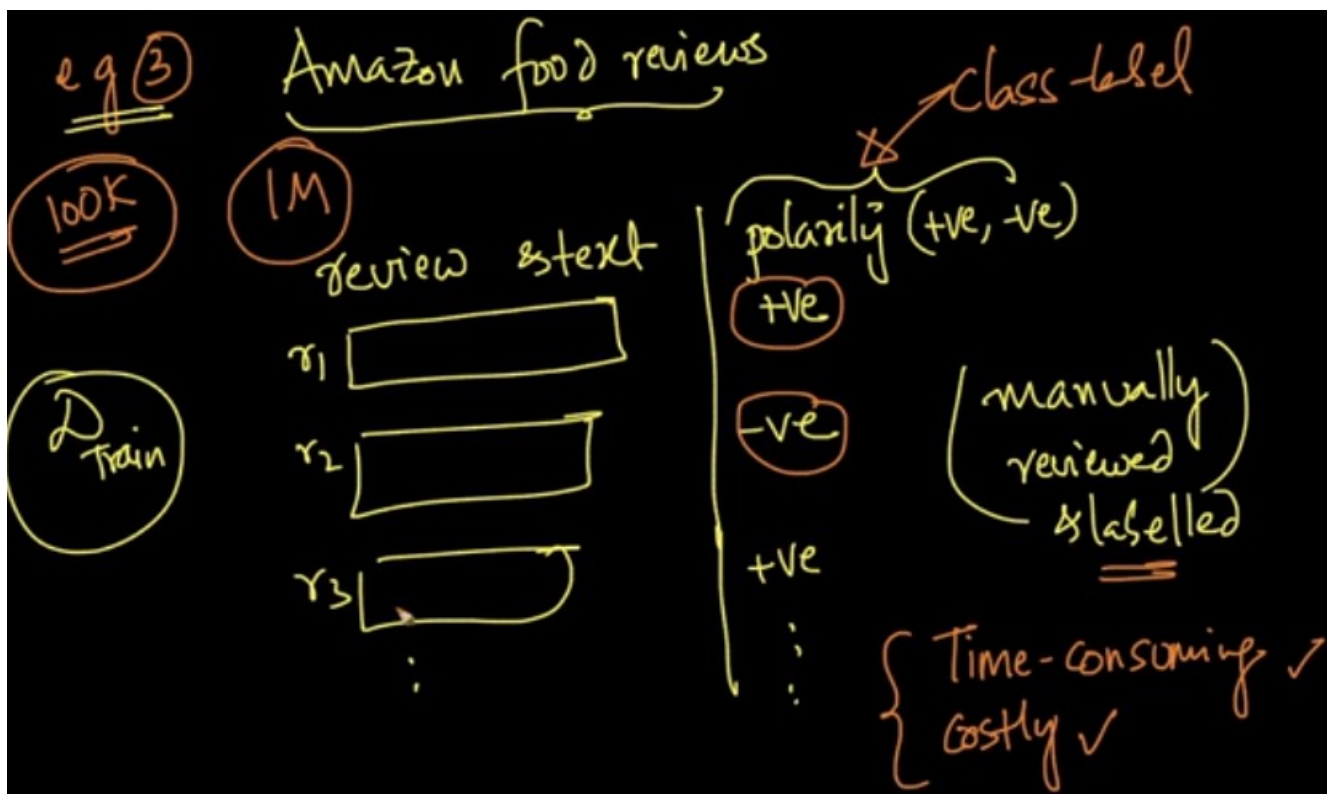
Image segmentation: This is a problem in computer vision and image processing.
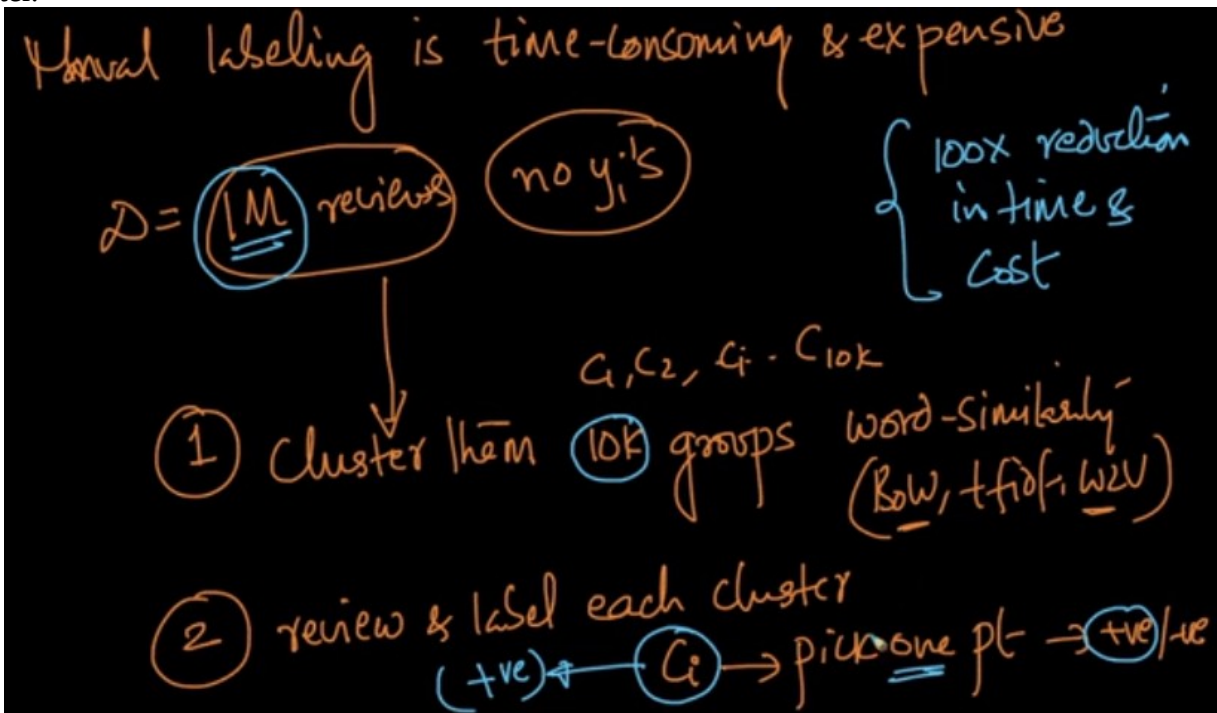
This is all about segmenting the pixels of the image. The clustering can help in seg the image.
We can apply the ML algos to apply the object detection.



The polarity is the class label. This is very time consuming. We can apply the clustering algorithm to get the labels of the reviews.
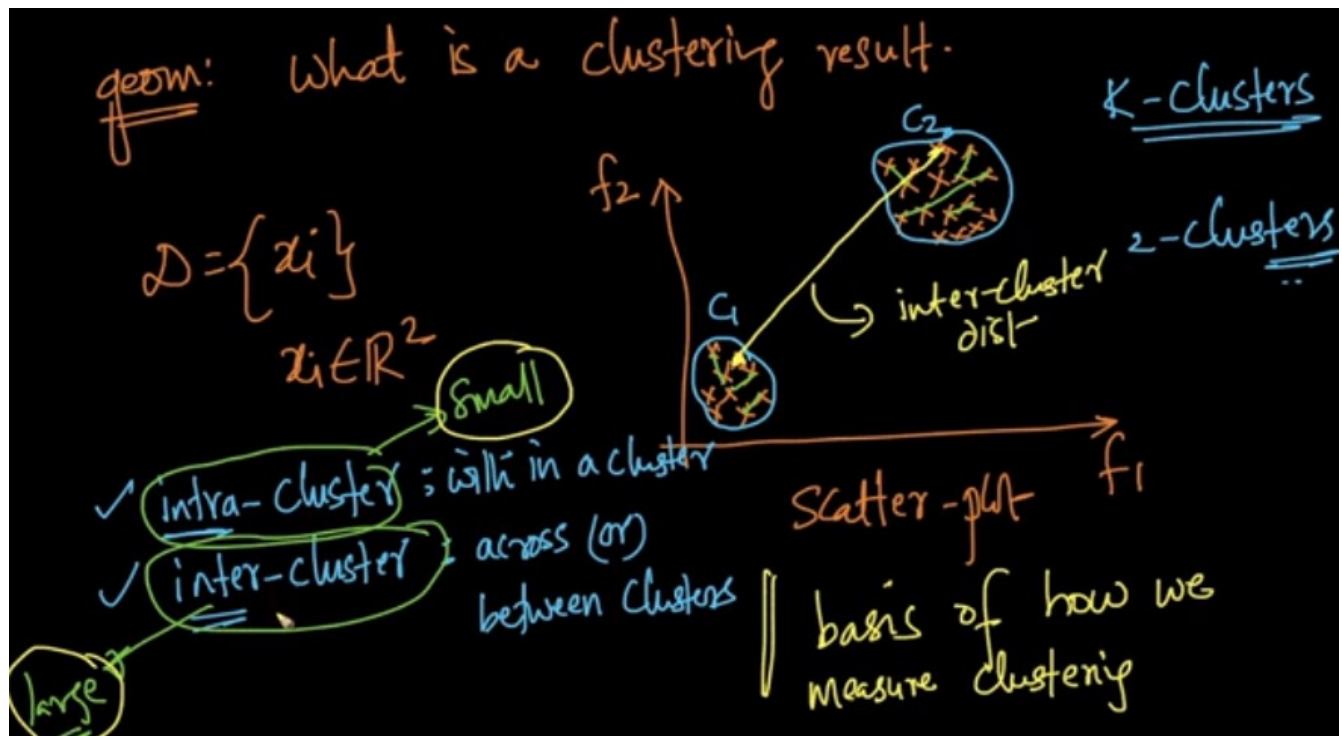
The manual labeling is time consuming and expensive. We can pick a cluster it is review and label each cluster.



Now, we can train the machine learning algorithms.
Metrics for measuring Clustering: What is a good clustering result?
All the points inside the cluster is called the intra clusters. All the outside the cluster is called the inter clusters. The intra cluster distance is kept small and the inter cluster distance kept large. This leads to the good clustering. **This is the basis of clustering effectiveness.**

In an ideal world, we want the inter clustering distance to be very high and intra to be lvery less.

Dunn – index: The numerator is the max inter cluster distance. The denominator is the intra cluster distance.



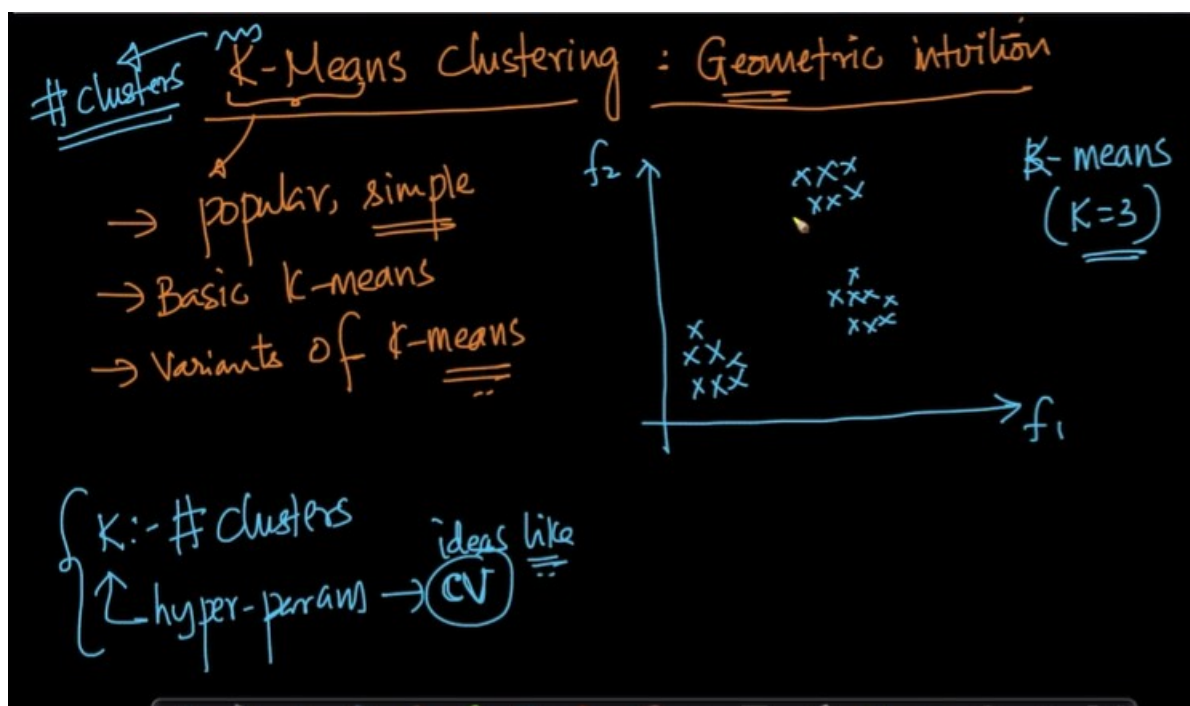The distances are calculated as follows and the farthest points are chosen.

Denominator: The farthest distance of the cluster is computed, then we take the max over the distances d'1 and d'2.
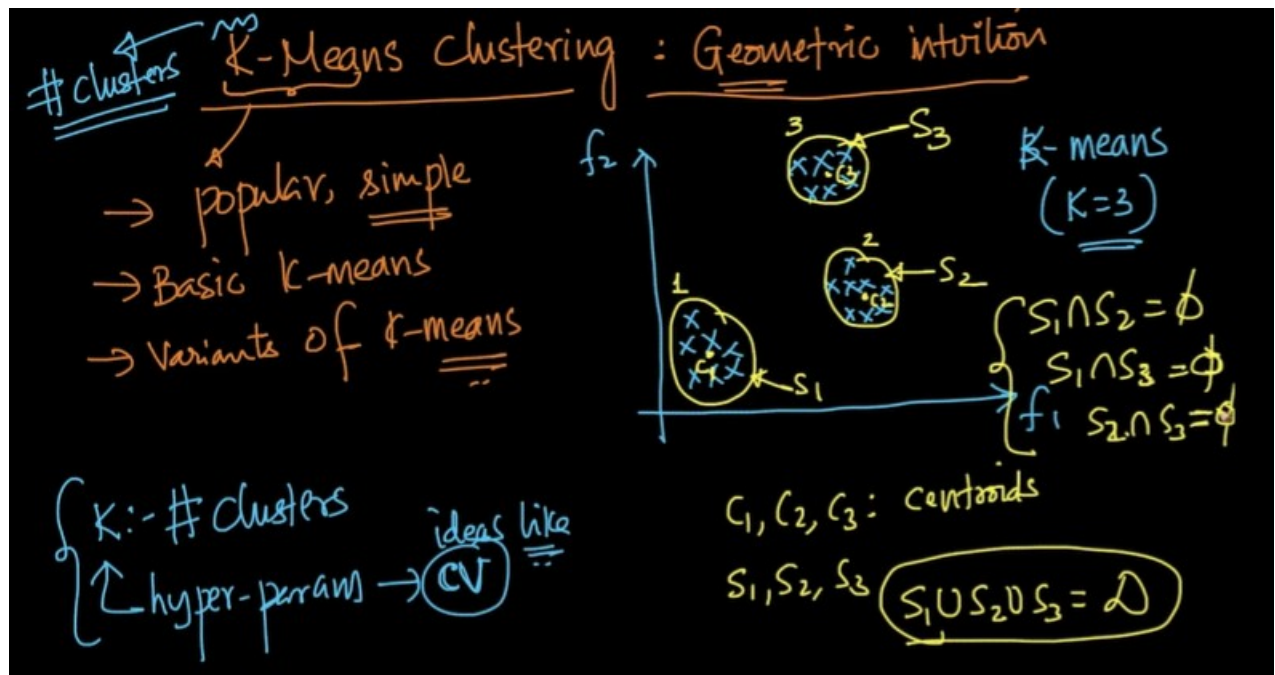


K – Means: Geometric intuition, Centroids
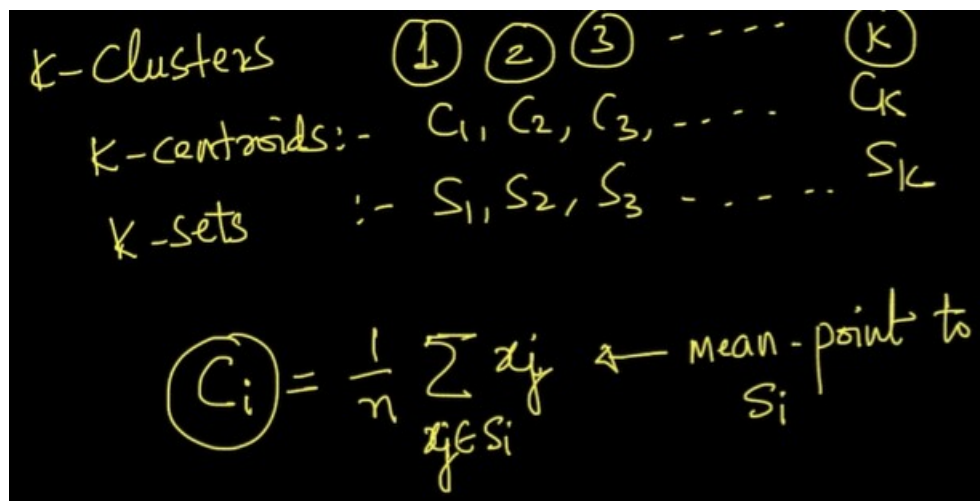It is the popular clustering algorithm, variants of K- Means.
The K is the number of clusters which is the hyper parameter int K – Means.

K-Means groups every cluster, In this case we have three clusters. The intersection of the clusters is the null set.
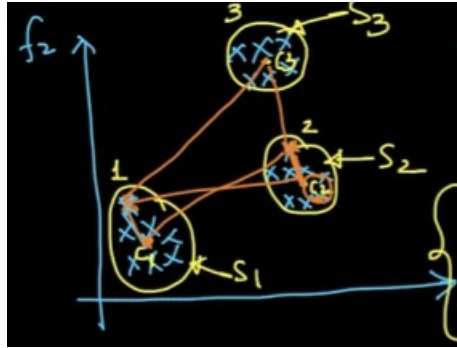


In K-Clusters, we have K centroids. K- clusters, K-sets of points. It is the geometric mean of the each cluster.
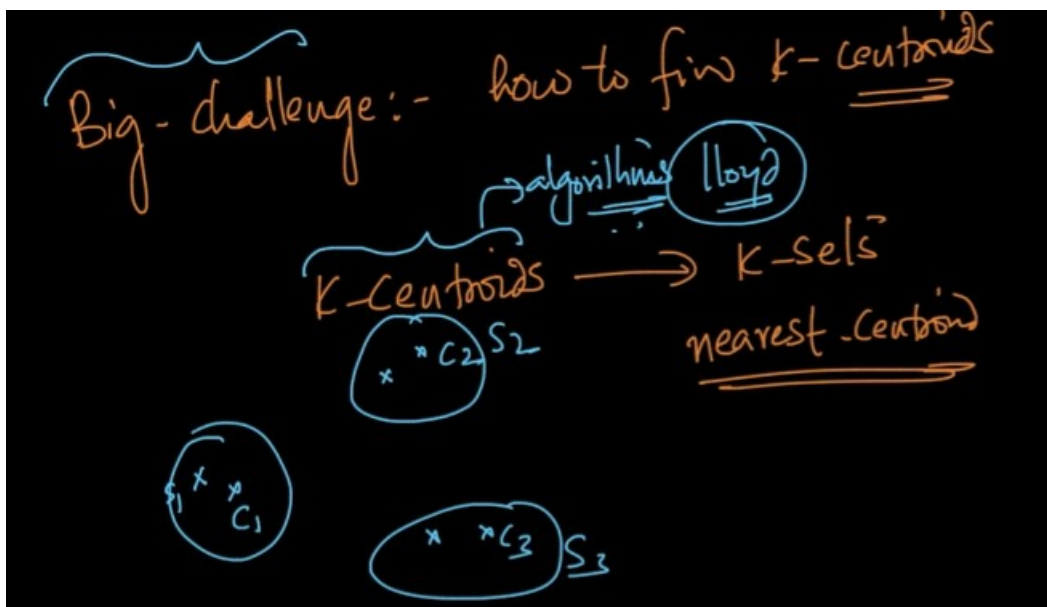


K – Means clustering is the centroid based scheme.
There are other clustering called

Heriarical based clsutersing
DBSCAN.

The challenge is how to find the K – Centroids. Once we get the K-centroids, we can compute the K-set nearest centroid. There are algos to find the K centroids.



To find the K central points and assign those points to the cluster.

K – Means: Mathematical formulation: Objective function

K-means : Mathematical formulation

$$\mathcal{D} = \{x_1, x_2, \ldots x_n\}$$

Task:
K-centroids :- $C_1, C_2, \ldots C_k$

Sets :- $S_1, S_2, \ldots S_k$

$\begin{cases} \forall i \quad \underline{x_i \in S_j} \\ \forall i,j \quad \underline{S_i \cap S_j = \phi} \end{cases}$

This is the optimization problem -

$$\underset{C_1, C_2, C_i \cdots C_k}{arg\ min} \sum_{i=1}^{k} \sum_{x \in S_i} ||x - C_i||^2 \longleftarrow \text{intra-clust. dist is minimized}$$

Sq. dist of $x$ fw $C_i$

s.t $\quad x \in S_i$
$\quad\quad S_i \cap S_j \neq \phi$ $\Big\} \longrightarrow$ constraints

proximity $C_1, C_2 \ldots C_k$
$\downarrow$
$S_1, S_2, \ldots S_k$

inter-clust. dist

This is sum of squared distances in the cluster I,

$$\underset{C_1, C_2, \ldots C_k}{\arg\min} \sum_{i=1}^{k} \sum_{x \in S_i} \|x - c_i\|^2 \qquad \text{v.v. hard to solve}$$

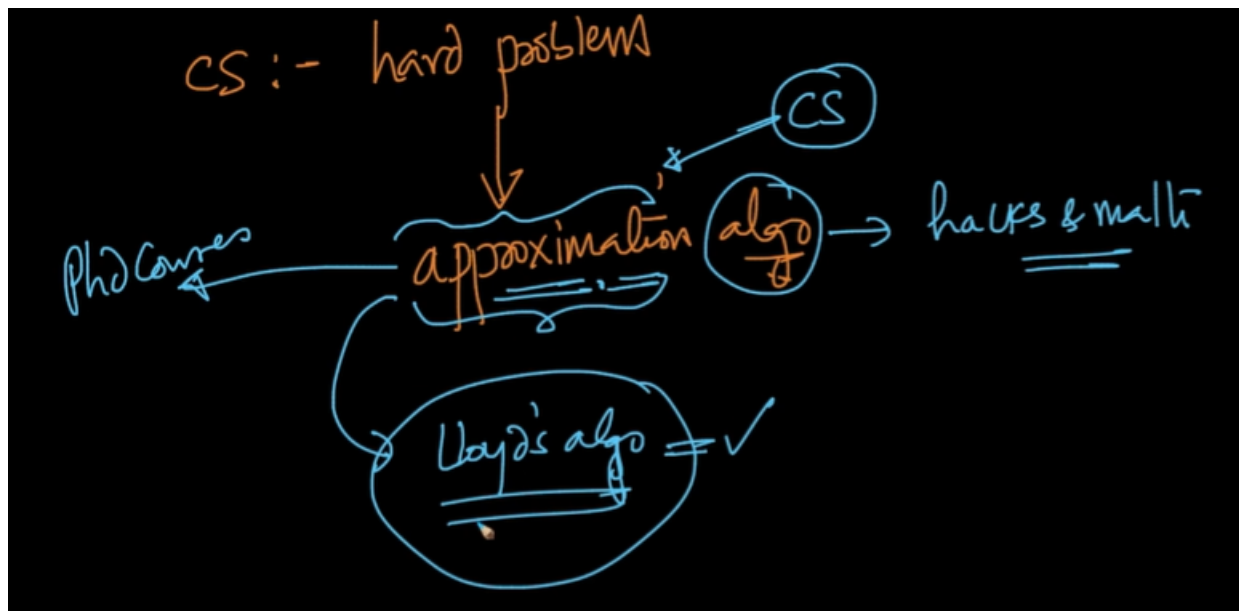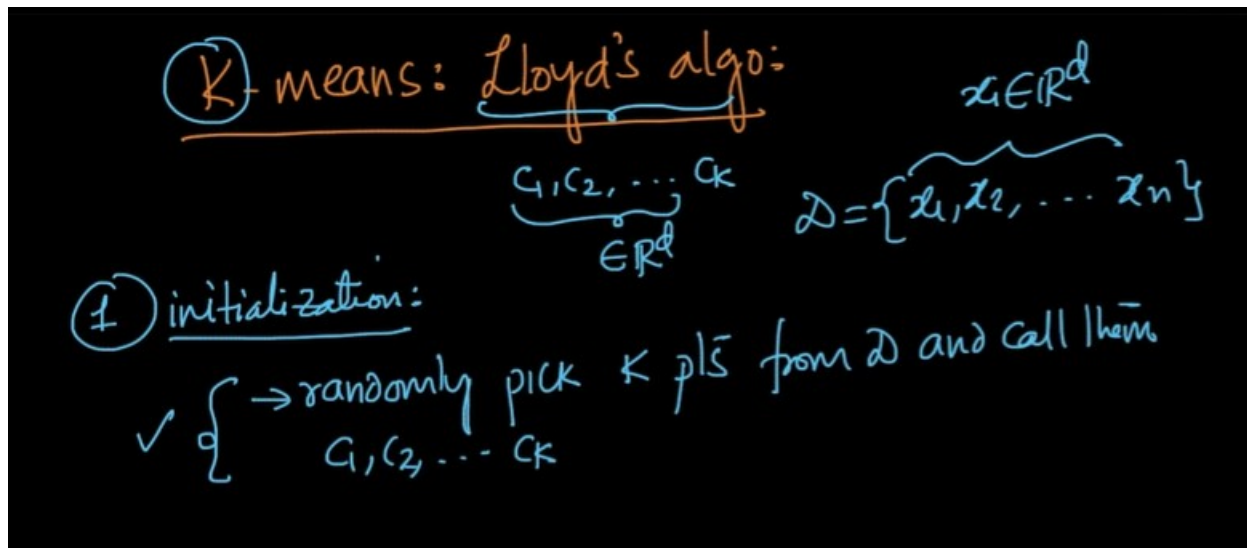↑ all clusters — Sum of Sq. dist fr centroid in cluster i

$$\text{st} \quad x_i \in S_i$$
$$S_i \cap S_j = \phi$$

This problem is very hard to solve. If there is a hard problem then we use the approximation of algo.
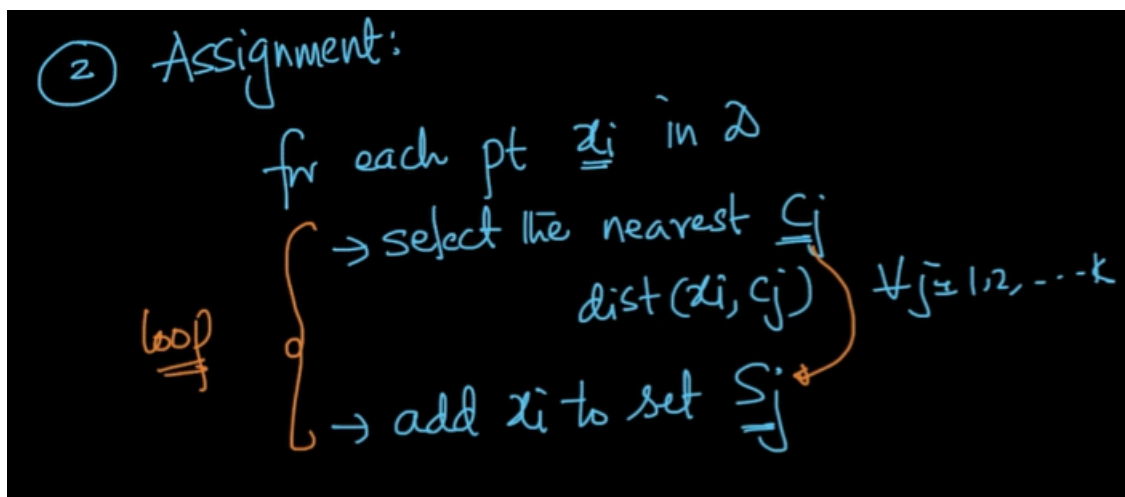Using some hacks and math.
Lloyds algorithm: This solves the problem.

CS :- hard problem

PhD Course → approximation algo → hacks & math

Lloyd's algo = ✓

Lloyd's algorithm: K-means (we find the K centroids randomly)

$\textcircled{K}$ means: $\underline{Lloyd's\ algo:}$     $x_i \in \mathbb{R}^d$

$\underbrace{c_1, c_2, \dots, c_k}_{\in \mathbb{R}^d}$     $D = \{x_1, x_2, \dots x_n\}$

$\textcircled{1}$ $\underline{initialization:}$

$\checkmark$ $\left\{ \begin{array}{l} \rightarrow randomly\ pick\ K\ pts\ from\ D\ and\ call\ them \\ c_1, c_2, \dots c_k \end{array} \right.$

The second step is assignment -

$\textcircled{2}$ Assignment:

for each pt $\underline{x_i}$ in $D$

loop $\left\{ \begin{array}{l} \rightarrow select\ the\ nearest\ c_j \\ \qquad dist(x_i, c_j)\ \forall\ j = 1, 2, \dots k \\ \rightarrow add\ x_i\ to\ set\ \underline{S_j} \end{array} \right.$

The third stage is called recompute centroid state: This is also called the update stage of centroids.
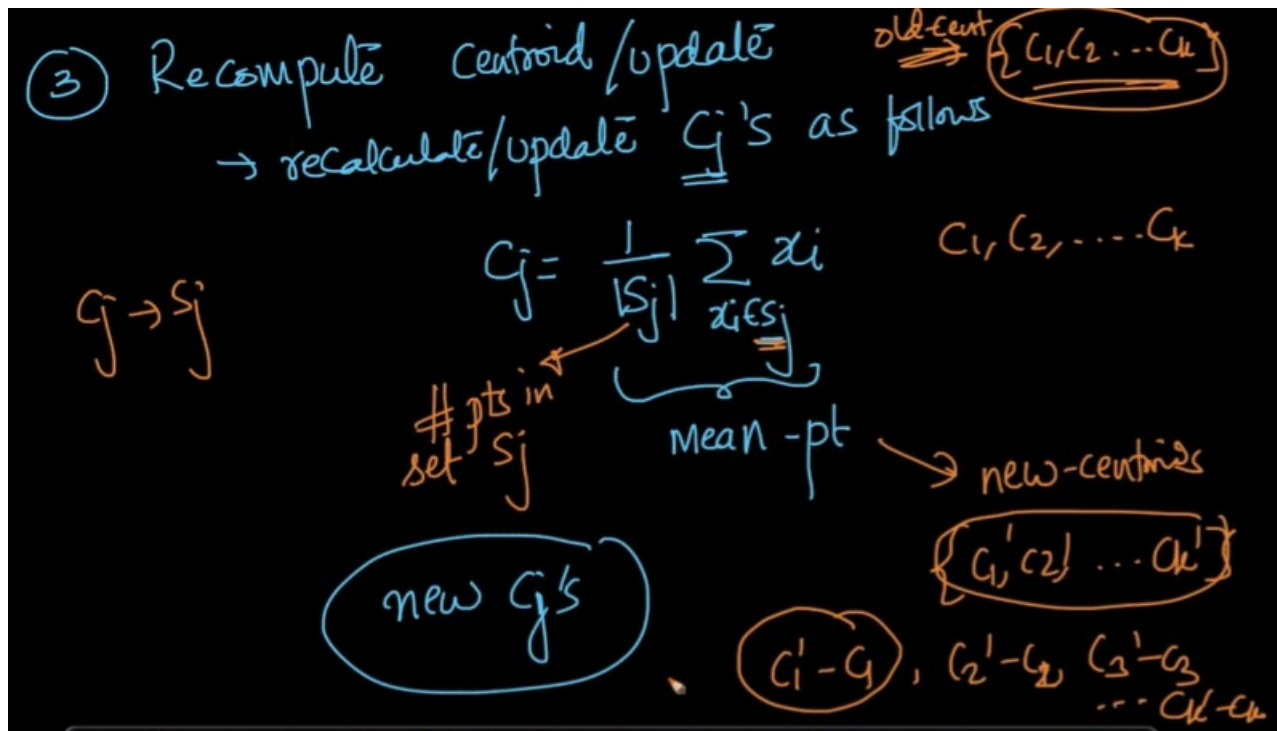


Step – 4:

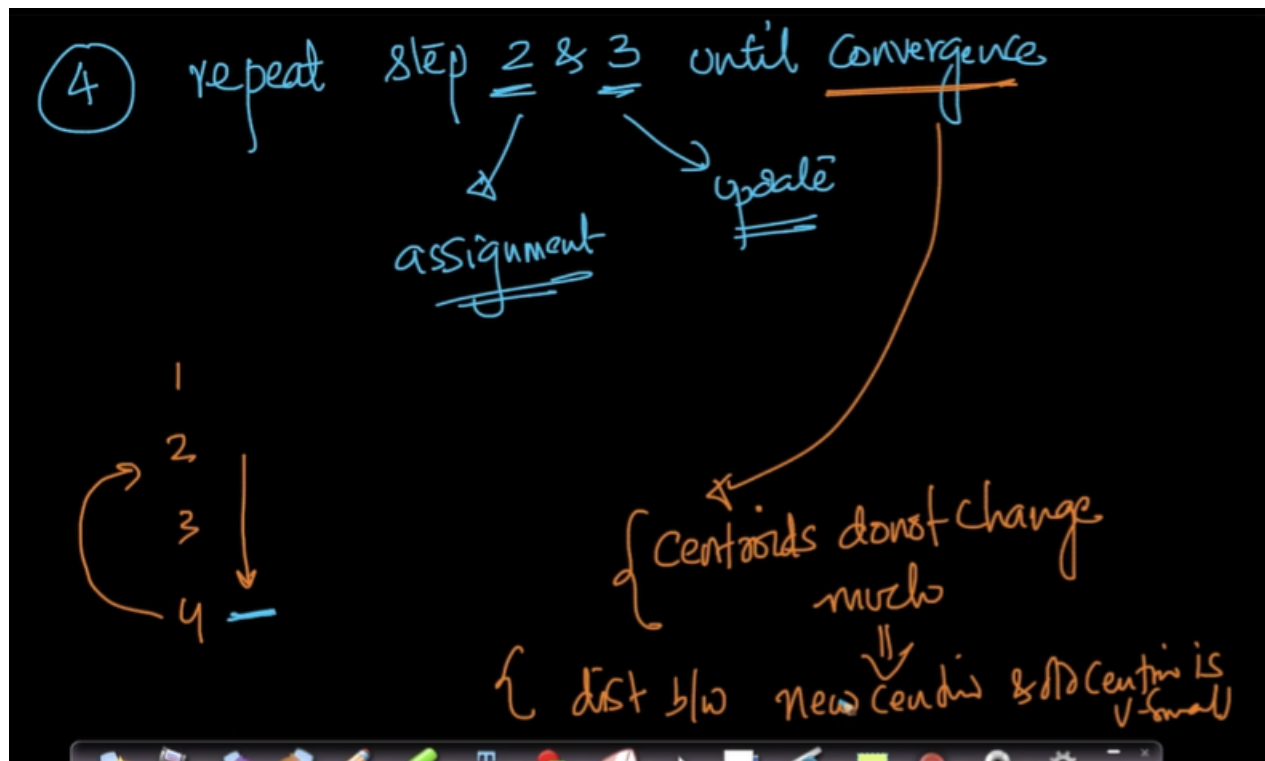We can repeat the step -2 and step -3, until convergence.

Here, we have the set of centroids, of new and old.
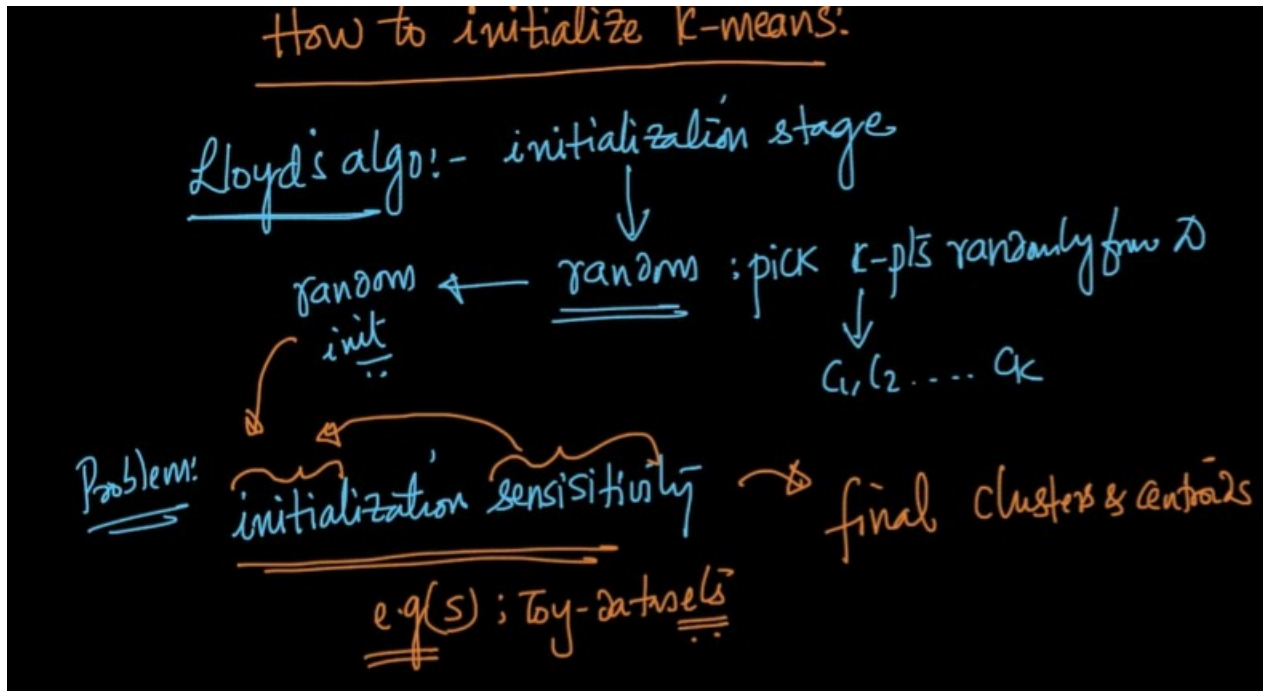If there are no much change in the distance of the old and new points distances.



Distance between new centroids and old has no change.

The actual mathematical formulation is very hard. This is the lloyd algorithm.
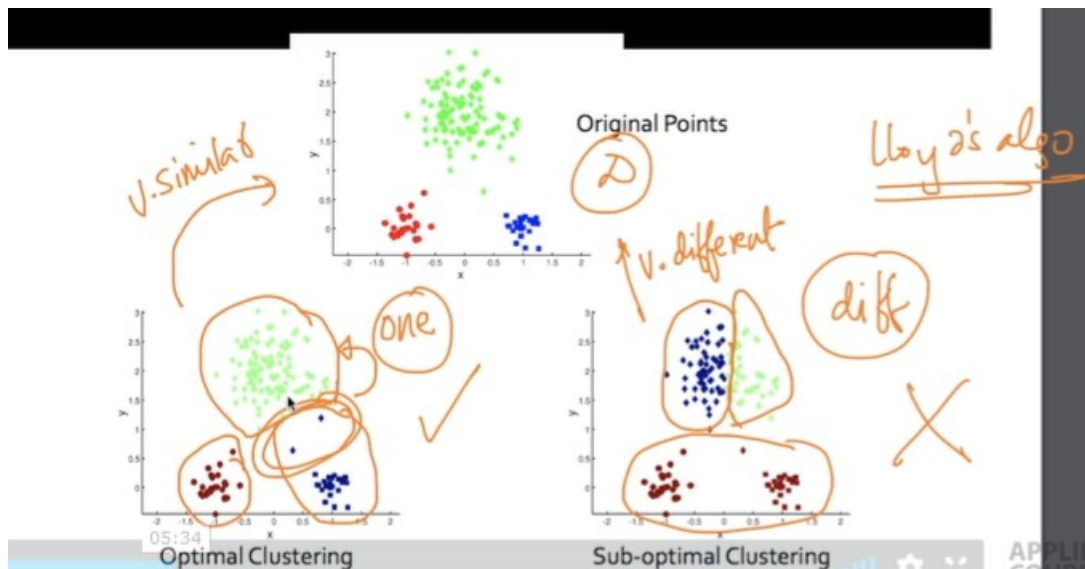
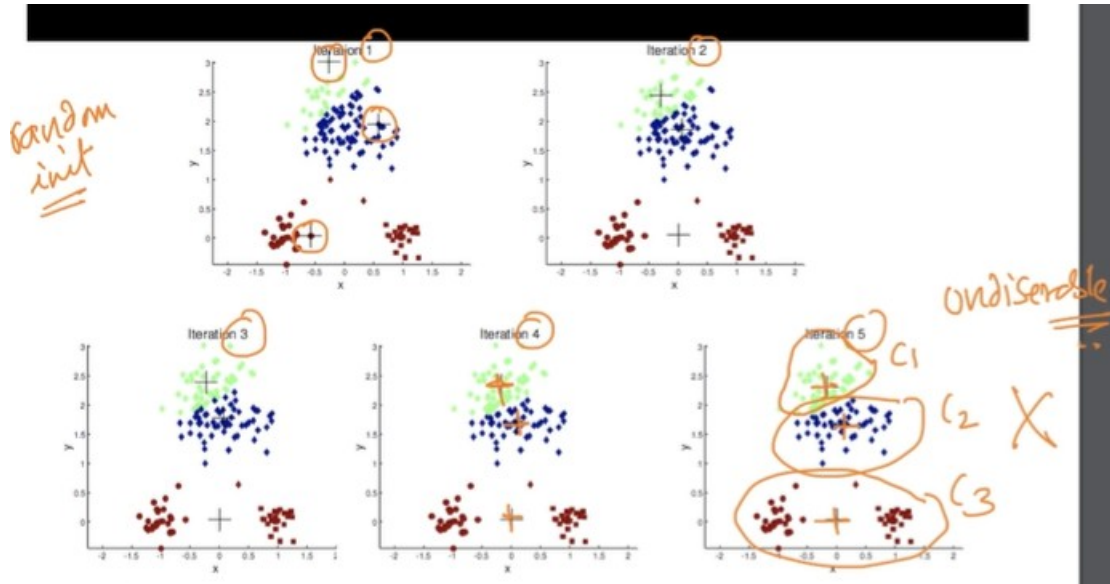How to initialize: K-Means++
We can do the random initialization from the dataset and make the K centroids. There is one problem, initialization sensitivity.



Given this ideal dataset, when we initialize the centroids. We get the optimal and sub-optmal clustering states.
Applying the same lloyd algorithm, we get the complete different solution.

random
init

Iteration 1    Iteration 2

Iteration 3    Iteration 4    Iteration 5

undiserable

$C_1$
$C_2$ $\times$
$C_3$

① repeat (k-means) multiple times with different initializations

↳ pick the best clustering based on
✓ { smaller intra-clust
✓ { larger inter clust }

The second way is K-Means ++, Instead of using the random initialization we use the smart initialization.
Step – 1:
Pick K centroids, we will pick the first cenroid randomly C1.



Step – 2:

The chances of the points near to the centroid is very low.
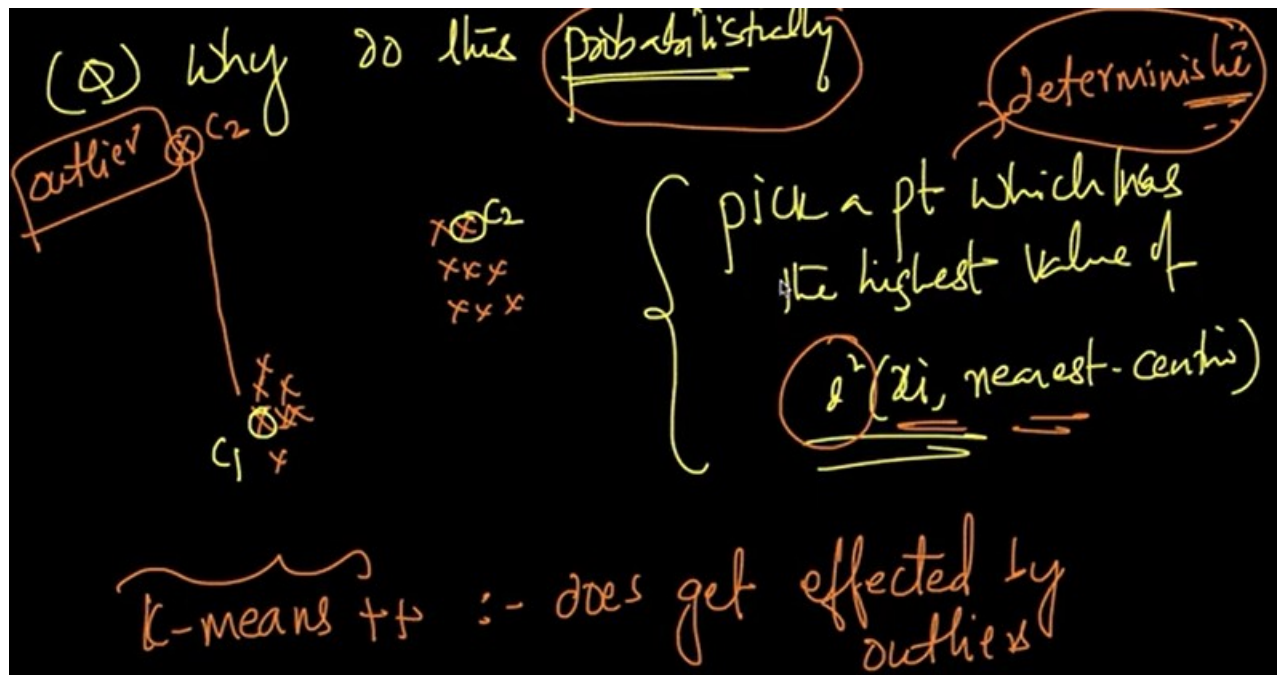In the initialization we want to pick the points that are far as possible from other centroids picked up.



The points far to the chosen center has more probability to be the centroid.
We are trying to pick the points as far as possible that are already picked up.
Why cant I just pick the point that has highest value from the nearest center. We can have outliers.

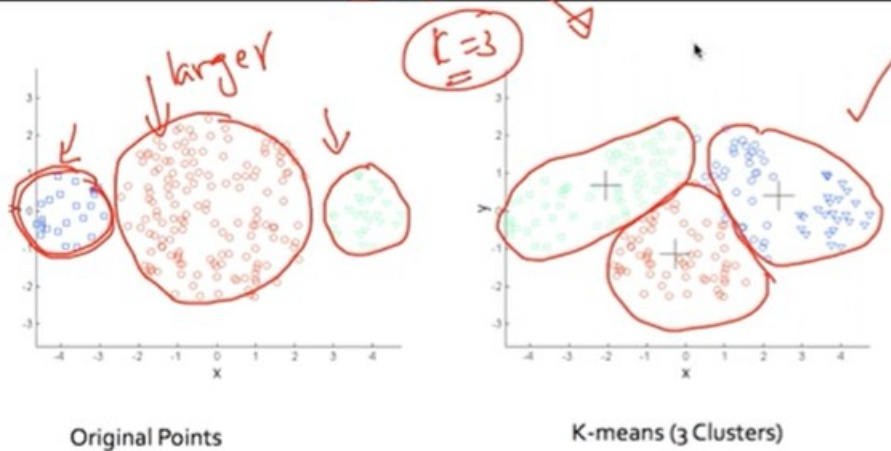We will pick the outliers as the centroid. K – Means ++ does get affected by outliers.

That is the whole reason we do it probabilistic-ally.  To reduce the mitigation of being an outlier.

Failure cases/Limitations:
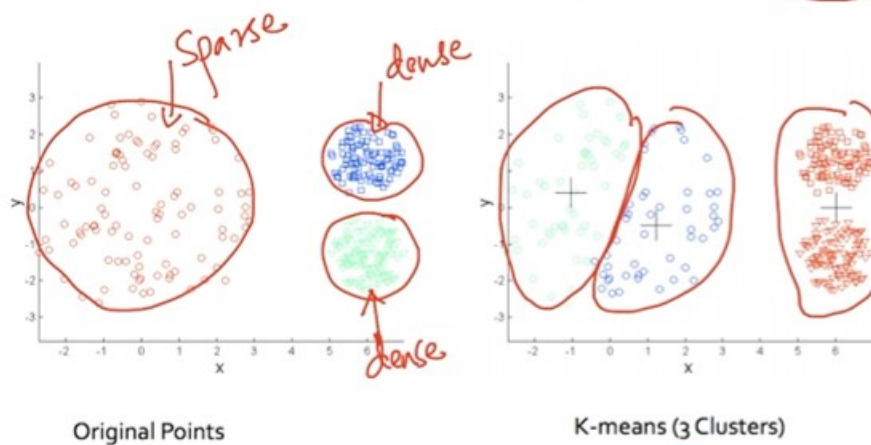When we have clusters of different sizes, different Densities and Non – gobular shapes.



Different densities:

K means tend to fail in case of different densities.

Non – gobular shapes:

If we give this data to the K – Means, we cannot work with gobular data.



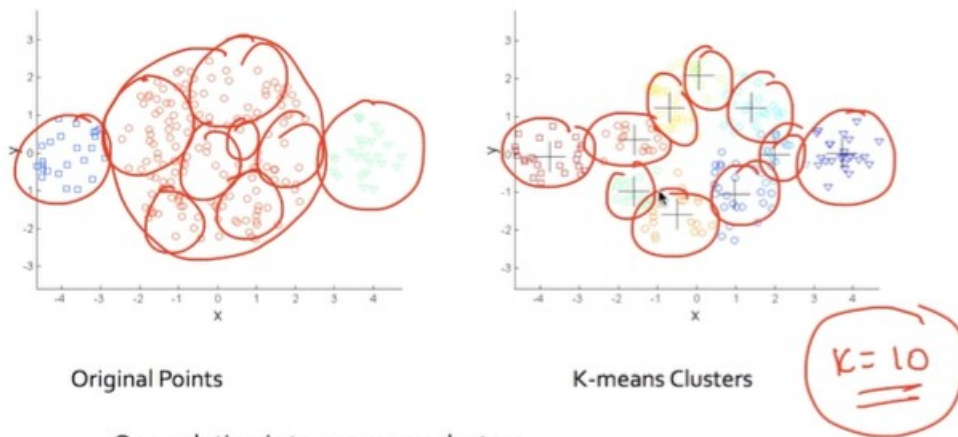If we are continuing this K means then, if we keep k = 10, w get clusters like as follows:



TO find the parts of cluster, but need to put together.

Clusters of different densities:



Original Points            K-means Clusters

K means is never perfect algorithm.
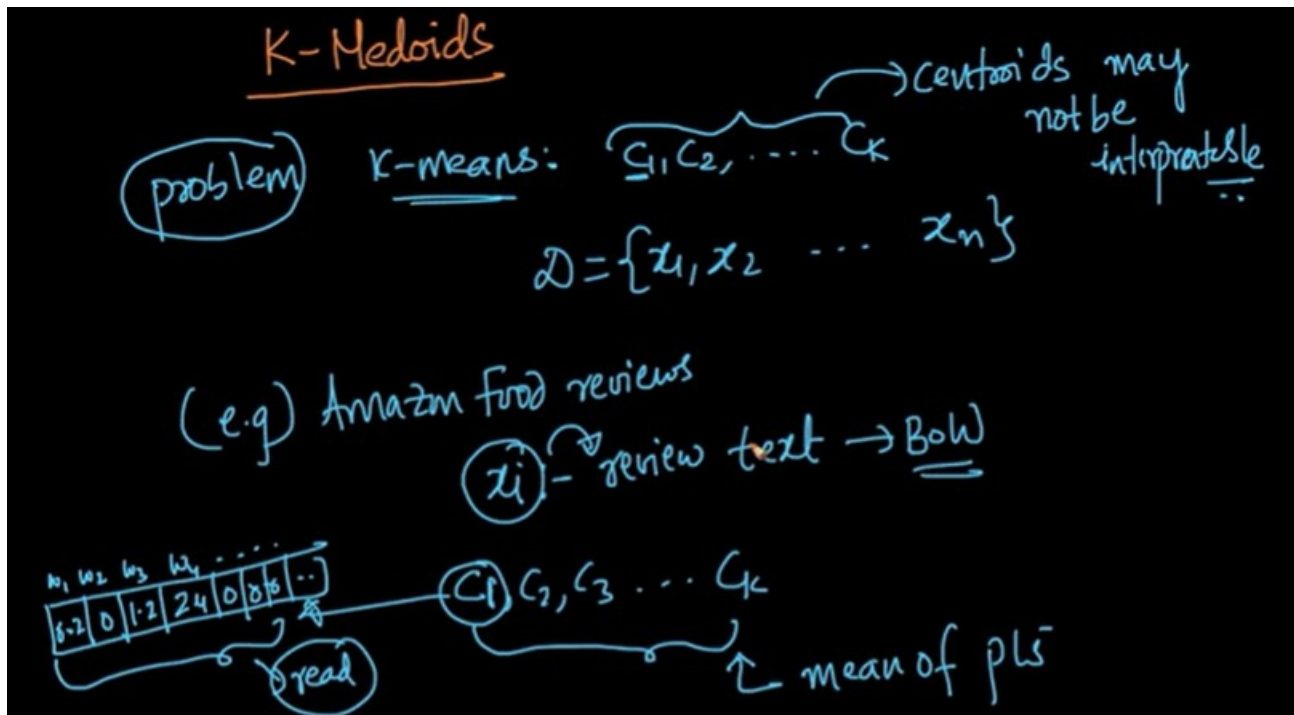In case of non – gobular structures, this is the best solution.
Evaluating clustering is not easy, there is no ground truth.



In the case of classification and regression there is no uncertainty in the results obtained.

K – Mediods: The centroids may not be interpret-able.



That means the centroids may be not interpretable in case of K-medoids.

When we represent the centroid as the data point in the dataset. Then it is called K-Medoids.

Especially in case of interpretation.

K – Medoids:

Partioning around mediods(PAM) algorithm:



We swap each medoid with a non-mediod point in the dataset.

What is loss in K-Means?

loss - in K-means

$$\min \quad \sum_{i=1}^{K} \sum_{x \in S_i} \| x - m_j \|^2$$

↑ medoid j

After swapping again compute the loss values.

$$M_1 \quad\quad\quad\quad M_2$$
$$\boxed{x_1} \; x_2, x_3, x_4, x_5, \boxed{x_6} \; x_7, x_8, x_9, x_{10}$$

(a) loss - value $\quad x_1 = M_1 \; ; \; x_6 = M_2 \; \rightarrow \boxed{l_1}$

$$M_1 \boxed{x_1} \leftrightarrow x_2$$

(b) swap $\quad M_1 = x_2 \; ; \quad x_1$ as a non-medoid pt
$$M_2 = x_6$$

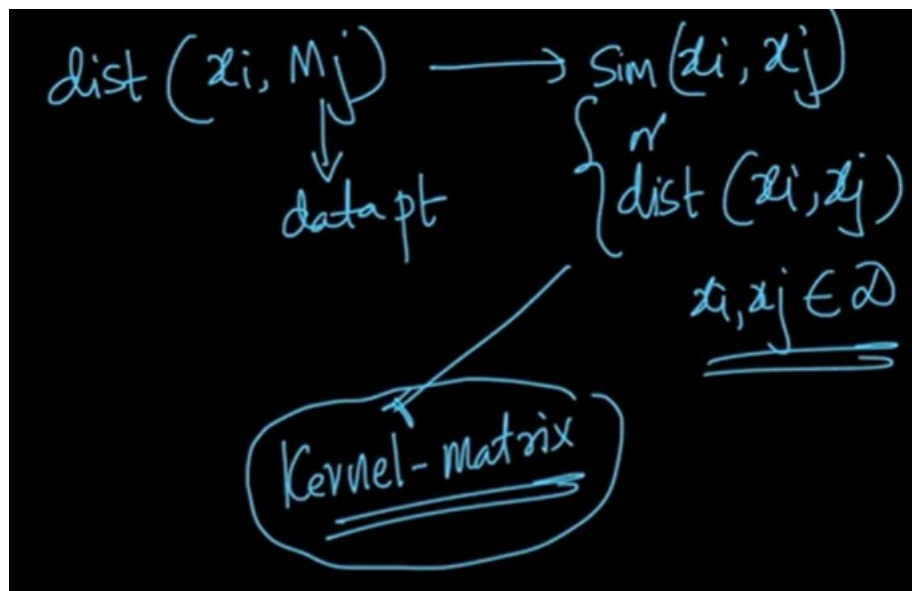↳ loss - value $\left. \begin{array}{l} M_1 = x_2 \\ M_2 = x_6 \end{array} \right\} \rightarrow \boxed{l_2}$

Is loss decreases by swapping then we keep the point that decreases the loss. There are lot of swaps possible.

$$\text{if } l_2 < l_1$$
$$M_1 = x_2$$
$$M_2 = x_6$$

$$\text{else}$$
$$M_1 = x_1$$
$$M_2 = x_6$$

$$\{ \text{lots of swaps that are possible} \checkmark$$

The swap is successful, when the loss decreases. Medoid is also a data point.

We can use the kernel matrix or distance matrix and apply $K - Mediods$. The massive advantage is

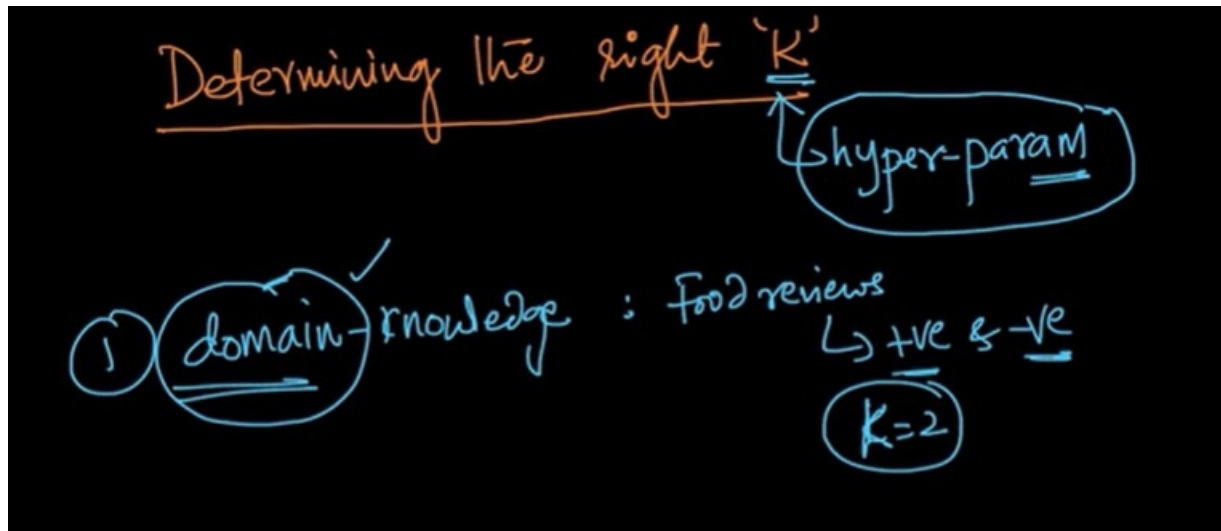1. More interpretation.
2. Kearnelization.
It is trivially kernalizable.

$$dist(x_i, M_j) \longrightarrow Sim(x_i, x_j)$$
$$\downarrow \qquad\qquad \{ \text{or}$$
$$data\ pt \qquad \{ dist(x_i, x_j)$$
$$\qquad\qquad x_i, x_j \in d$$

Kernel - matrix
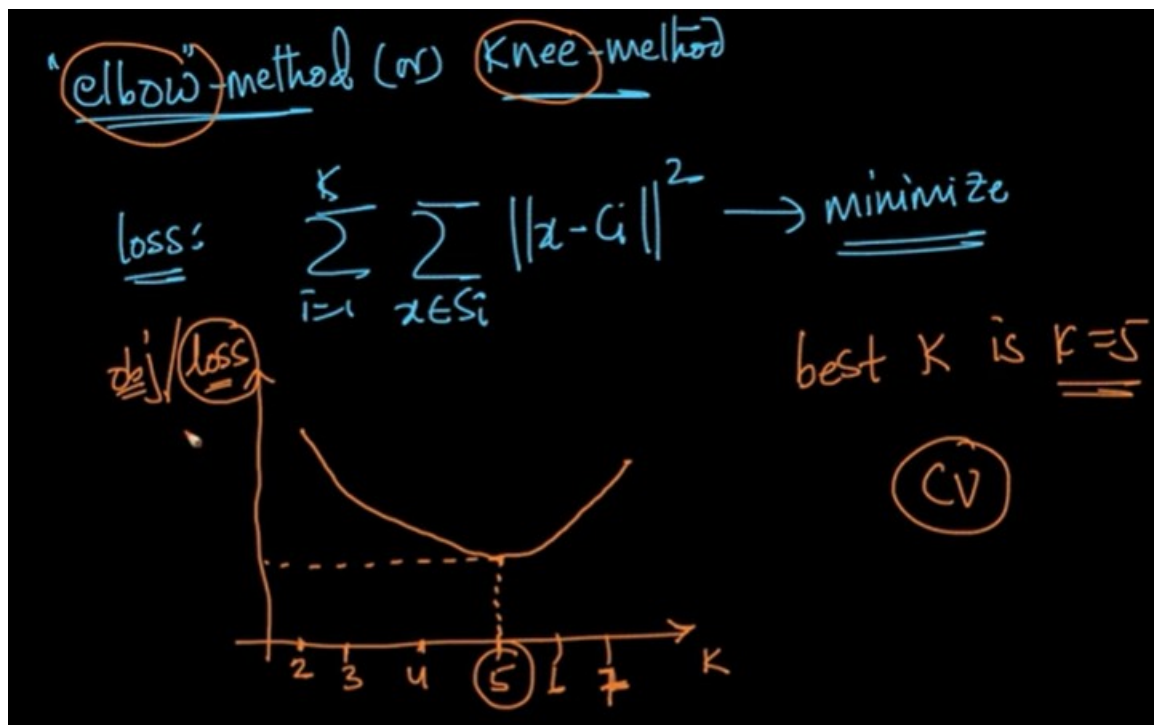
Determining the right K:

In K – means, K is the hyper parameter,

1. Domain knowledge – we can know the clusters.



Then we use the elbow-method (or) Knee-method -

We want to minimize the loss function is k means.



We only have data D={xi}.

Code examples:

Sklearn – sklearn.cluster.Kmeans

**Examples**

```
>>> from sklearn.cluster import KMeans
>>> import numpy as np
>>> X = np.array([[1, 2], [1, 4], [1, 0],
...               [4, 2], [4, 4], [4, 0]])
>>> kmeans = KMeans(n_clusters=2, random_state=0).fit(X)
>>> kmeans.labels_
array([0, 0, 0, 1, 1, 1], dtype=int32)
>>> kmeans.predict([[0, 0], [4, 4]])
array([0, 1], dtype=int32)
>>> kmeans.cluster_centers_
array([[ 1.,  2.],
       [ 4.,  2.]])
```
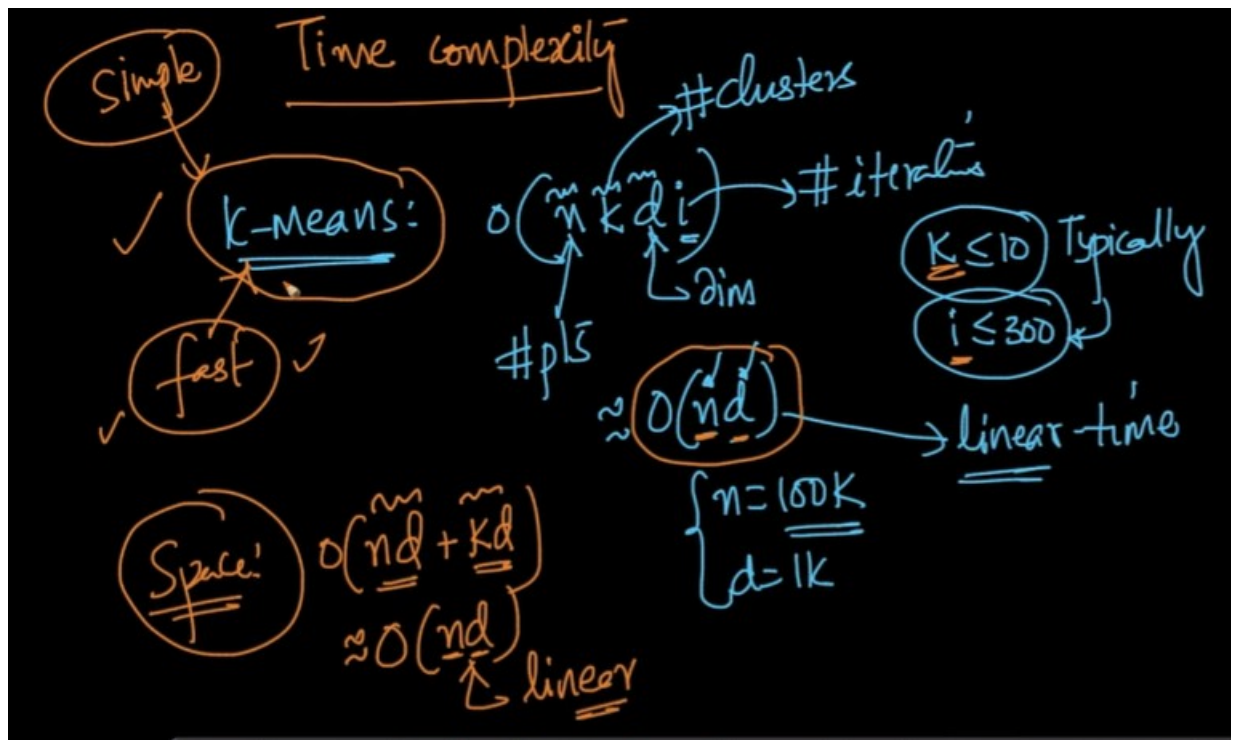
**Methods**

Time and space complexity:

K-Means: O(nkdi) This is still a linear in time.
As far as space concerned we need to store nd + kd. Which is also linear.

It is fast and simple to understand.

Cluster Amazon reviews: