

Unsupervised Machine Learning

⇒ Class of machine learning technique used to find patterns in the data.

⇒ The data given unsupervised learning algorithm is not labelled.

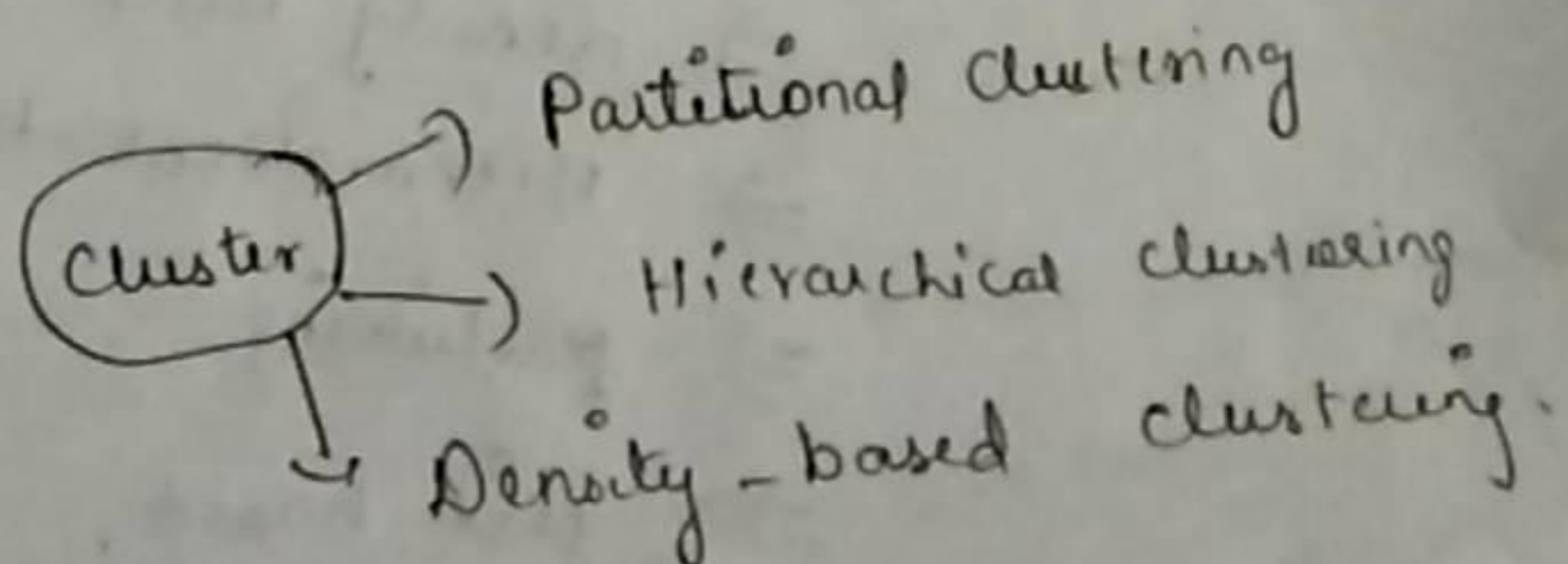
⇒ which means the only input variables (x) are given with no corresponding output variables.

⇒ The most prominent methods of unsupervised learning are cluster analysis & principal component analysis.

Clustering:

The data is divided into several groups with similar traits.

Three popular categories of clustering algorithm.



why clustering?

find hidden relationship between
the datapoints in the dataset

Ex: Marketing: Characterise and discover
customer segment

Biology: classify among dif types
of plants & animals

Library: It is used in clustering
dif books on the basis of topic &
information

Clustering method:

four types of clustering methods

⇒ density based

⇒ Hierarchical based

⇒ partitioning

⇒ grid based.

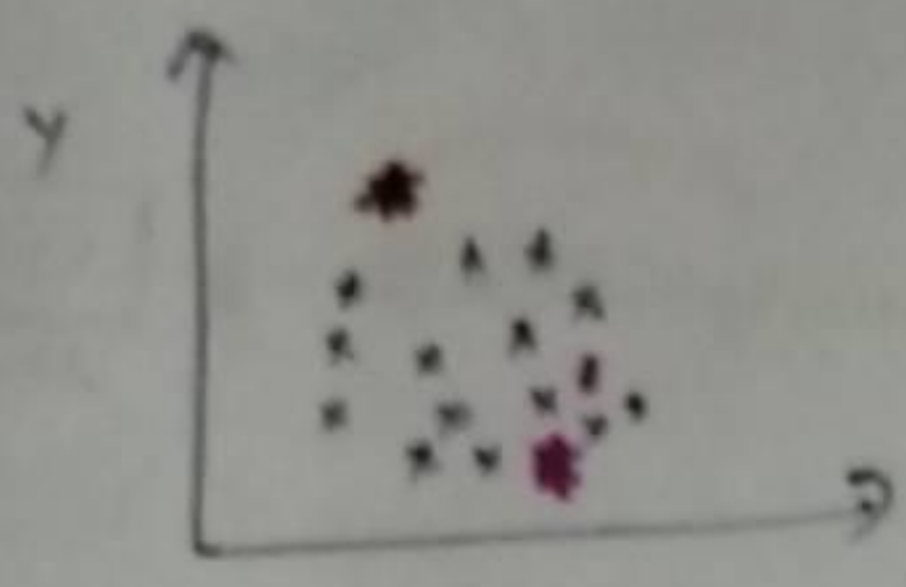
Partitioning methods:

- ⇒ find mutually exclusive clusters of spherical shape.
- ⇒ distance based
- ⇒ May use mean or median (etc.) to represent cluster center
- ⇒ effective for small to medium size data.

K-means cluster:

The K-means algorithm defines the centroid of cluster as the mean value of the points within the cluster.

- ⇒ K-value = 2 (centroids)
- ⇒ initialize K-value (the centroid randomly)
- ⇒ distance find $\begin{cases} \text{Euclidean} \\ \text{Manhattan} \end{cases}$
- ⇒ Select the group & find average
- ⇒ After that centroid moved again unless until when the point can't change.



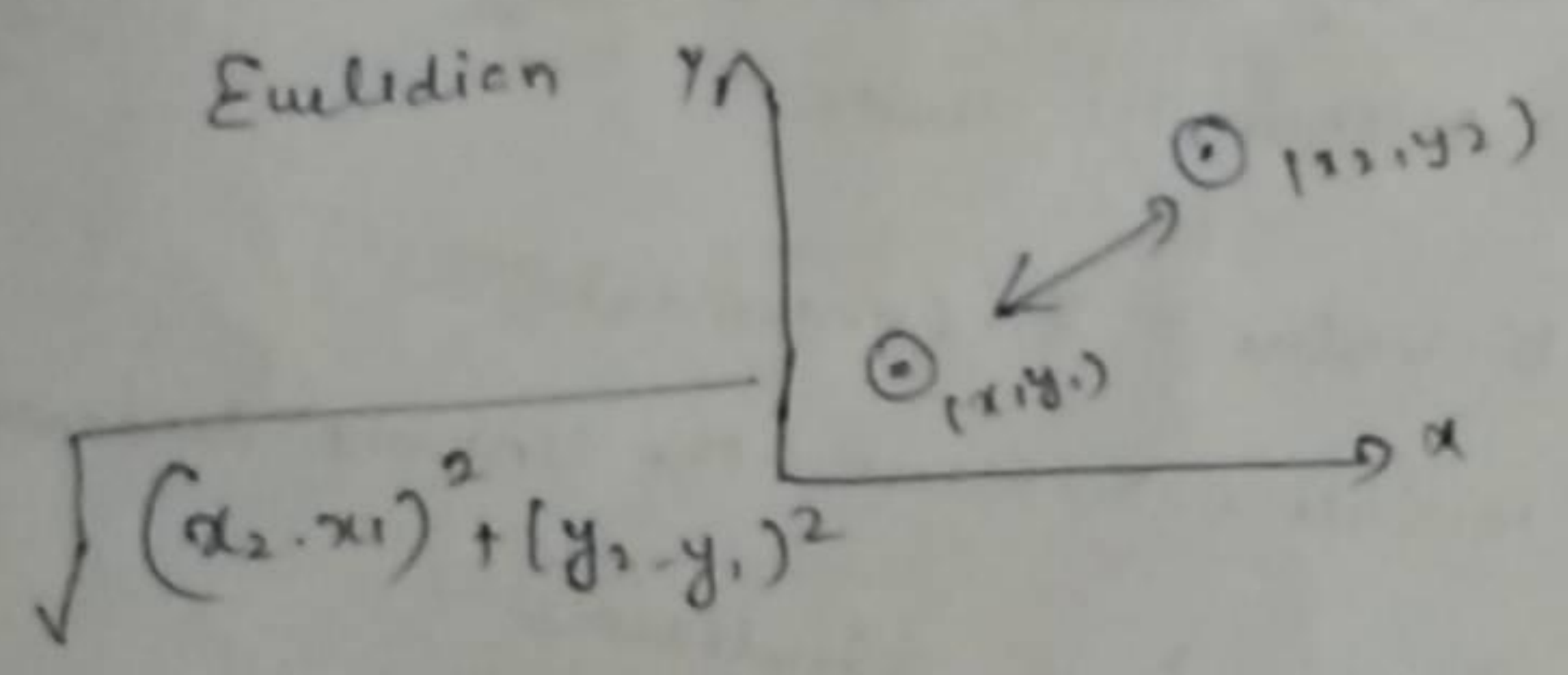
→ This is my data points

→ Randomly select centroids

→ my randomly cluster is blue, red...

Okay, After select random centroids we will find distance b/w actual points and randomly select centroids

How find distance? use Euclidean, Manhattan...

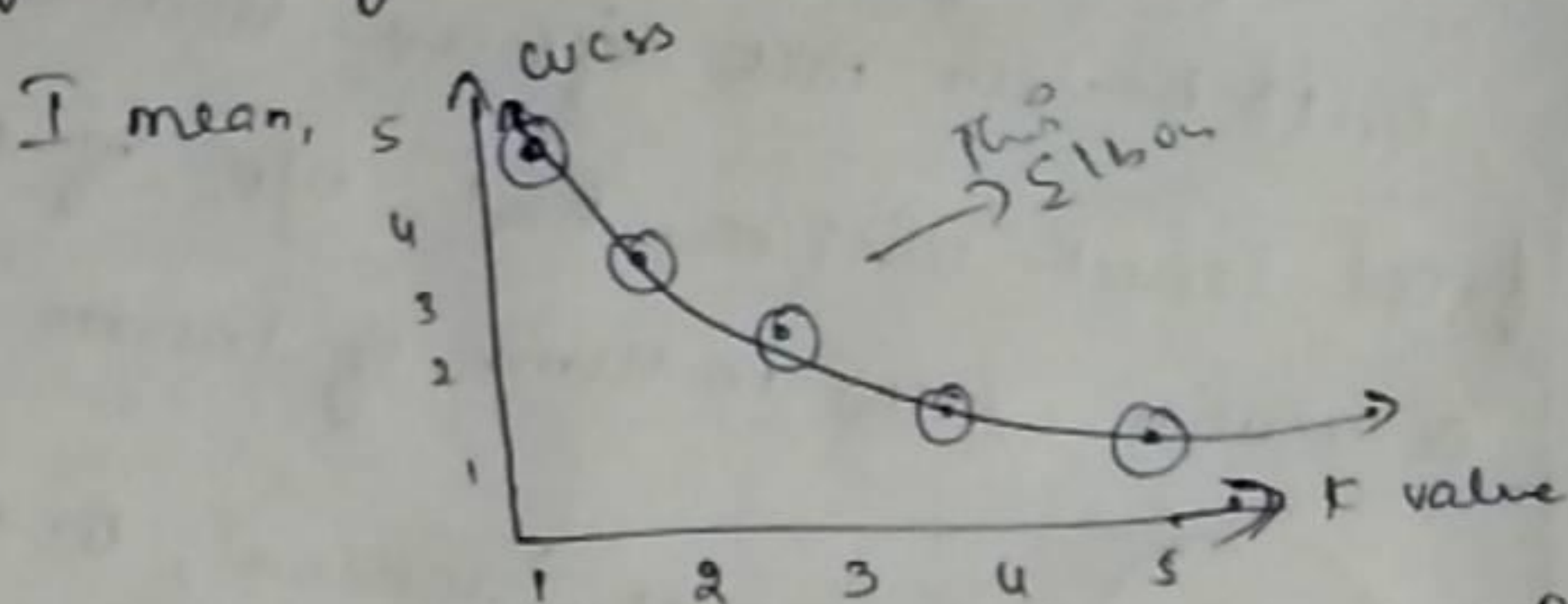


How select K value?

Elbow

Elbow method?

Now I am run $K=1$ to 20 through looping
when my $K=1$ my wcss is **high**



which is sudden decrease that is
correct and best value of K.

Okay, what is wcss?

within cluster sum of square

$$\sum_{i=1}^n (c_i + x_i)$$

Centroid \rightarrow datapoints.

Hyperparameter tuning :

⊗ $n_cluster$ → no of clusters that we want (Elbow)

init → initial cluster centroids

⊗ n_init → no of times k-means alg will be run with different centroid seeds.

final result will be best of n_init consecutive runs in terms of inertia

⊗ init (k-means++, 'random', or ndarray)

k-means++ (default) → select initial cluster

↓
small way to speed up convergence

random → choose k rows as initial centroids

ndarray ($n_cluster$, $n_features$)

max_iter (maximum no of iterations of k-means algorithm for a single run)