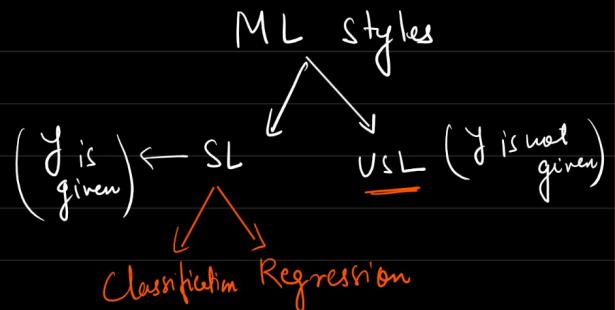
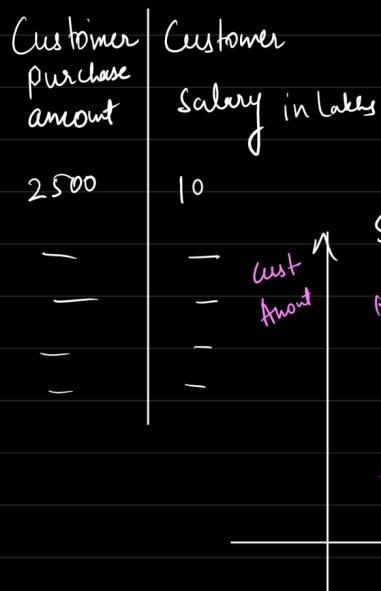


Introduction to Unsupervised learning

- * Target variable is Not given (y)
 - + VSL groups / segments your data
(finds pattern in the data)



Examples Zara store



To increase revenue

→ Group 2 customers should be focused more.

→ Group 1 → loyalty card & Extra discount.

* depends on you how you want to interpret the groups and bring business.

Motivation

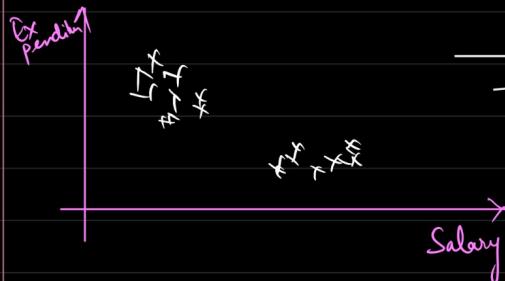
- Need to launch a campaign based on income | sales.
 - Different ways of fraud | money laundering.
 - Grouping of images | patterns of specific image
 - Document / Article Analysis.
 - Content Analysis.

* In higher dimensions, you cannot identify groups/patterns manually. Therefore you need VSL.

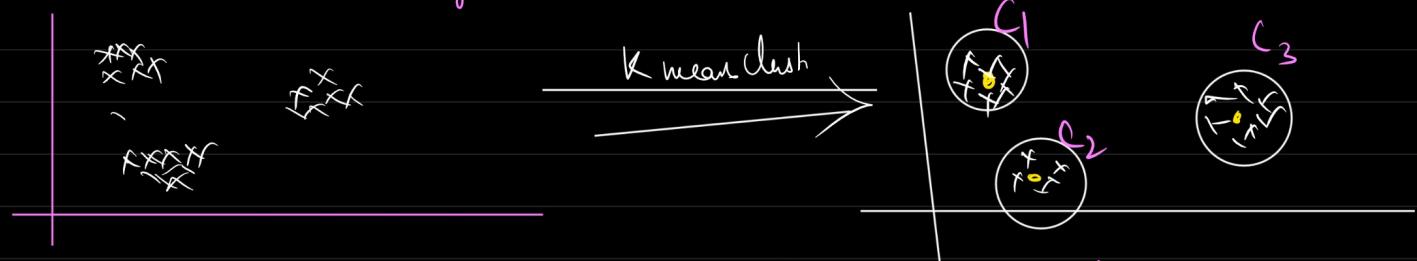
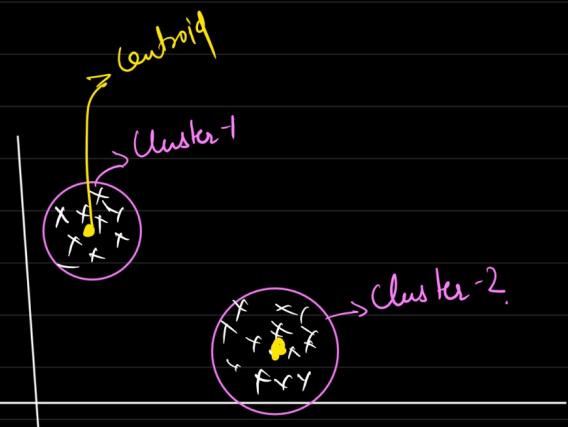
USL

- ① K-means clustering
- ② Hierarchical clustering
- ③ DBSCAN

* K-means clustering



Kmeans clustering



Steps of K-means clustering

Step-1 initialize Centroid (K)

- ↳ Can be random dp
- ↳ Centroid of all the dp
- ↳ A random new point

Step-2 → Points nearer to the centroid will be labelled as that group

Step-3 Re-Calculate the centroid Again repeat ↳ Avg. Step 2 until Centroid doesn't change

Exp	Salary
2	5
3	6
5	7

$$\left(\frac{2+3+5}{3}, \frac{5+6+7}{3} \right)$$

$$\text{Distance}$$

$$\textcircled{1} \text{ Euclidean distance } c = \sqrt{a^2 + b^2}$$

$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

② Manhattan distance



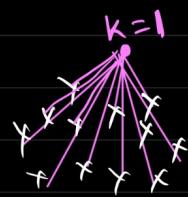
$$\text{Manhattan distance} = |x_2 - x_1| + |y_2 - y_1|$$

$$= a + b$$



How to Select the k value

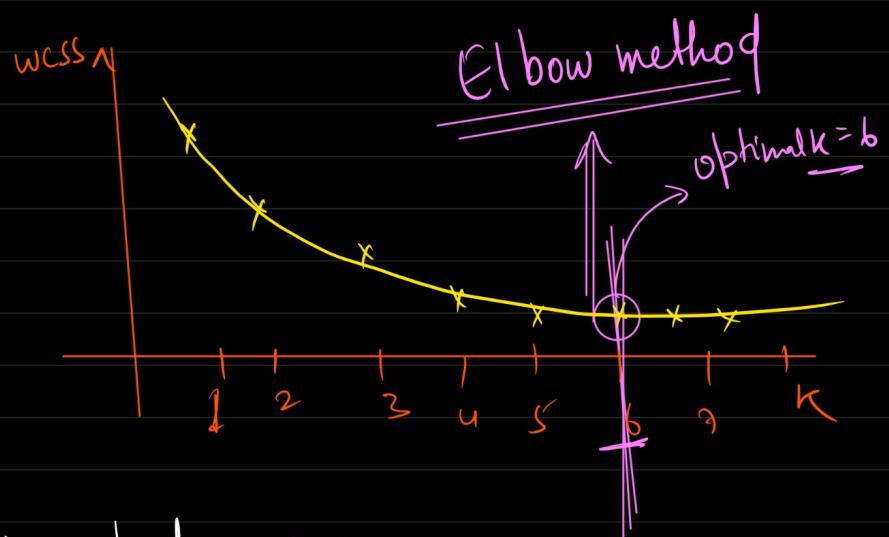
WCSS → Within Cluster Sum of squared distance



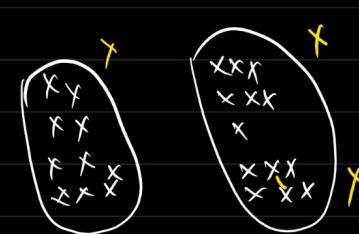
$$WCSS = \sum_{l=1}^n \left(\frac{\text{distance b/w points to nearest centroid}}{\text{Centroid}} \right)^2$$



$WCSS_1 + WCSS_2 \rightarrow$ for $k=2 \rightarrow$ As you increase k ,
WCSS decreases
→ And at optimal k ,
WCSS stops changes.



* Random Initialization trap $K=3$



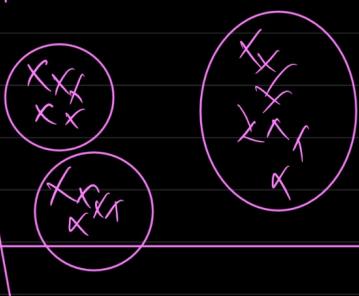
$\xrightarrow{K \text{ means}}$ output



Another possib

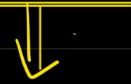


$\xrightarrow{K \text{ means}}$ output

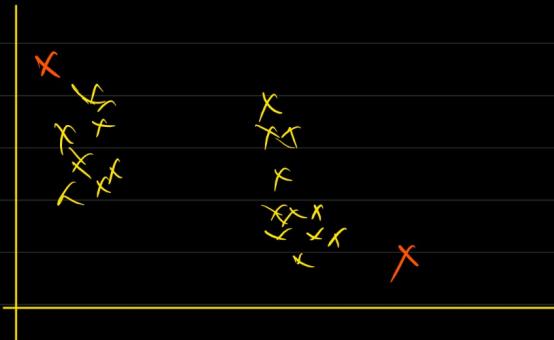


* The final cluster depends on the initialization of k .

Kmeans++



initialize the k 's as much far as possible



K-means clustering is an unsupervised machine learning algorithm used to partition data into K distinct groups based on similarity. It works by randomly initializing K centroids, then assigning each data point to the nearest centroid to form clusters. The centroids are recalculated as the mean of the points in each cluster, and the process repeats until convergence, meaning the centroids no longer move significantly. K-means is commonly used for tasks like customer segmentation, image compression, and pattern recognition. It requires the number of clusters (K) to be predefined, which can be a limitation in some cases.

Selecting the optimal value of K in K-means clustering can be done using several methods:

Elbow Method: Plot the sum of squared distances (inertia) between points and their assigned cluster centroids for various values of K . The optimal K is at the "elbow" point, where the reduction in inertia slows down significantly, indicating diminishing returns from adding more clusters.

Silhouette Score: This metric measures how similar points in a cluster are to each other versus points in other clusters. The value of K that maximizes the silhouette score is considered optimal.

The random initialization trap in K-means clustering refers to the problem where the algorithm's outcome depends heavily on the initial random placement of the centroids. Since K-means is sensitive to initial centroid positions, poor initialization can lead to:

Suboptimal Clusters: The algorithm may converge to a local minimum, resulting in inaccurate clustering.

Slow Convergence: Poorly chosen initial centroids may take longer to converge.

Inconsistent Results: Different runs of the algorithm can yield different clusters for the same data set due to random initialization.

Solution:

K-means++ Initialization: This technique strategically selects initial centroids by spreading them out, reducing the chances of poor clustering and improving both convergence speed and quality.

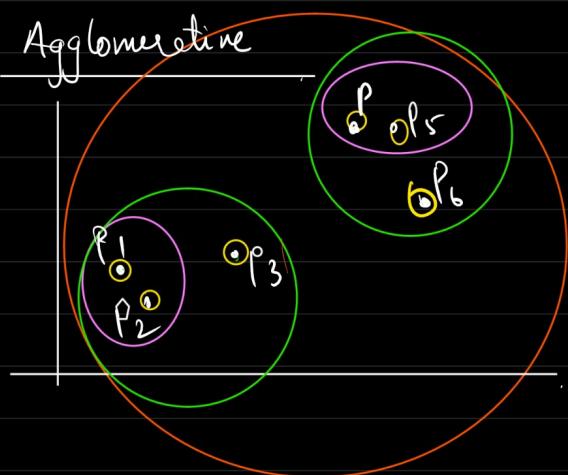
* Hierarchical clustering



* Hierarchical clustering.

- ① Agglomerative Clustering (combining)
- ② Divisive Clustering (dividing)

* Agglomerative



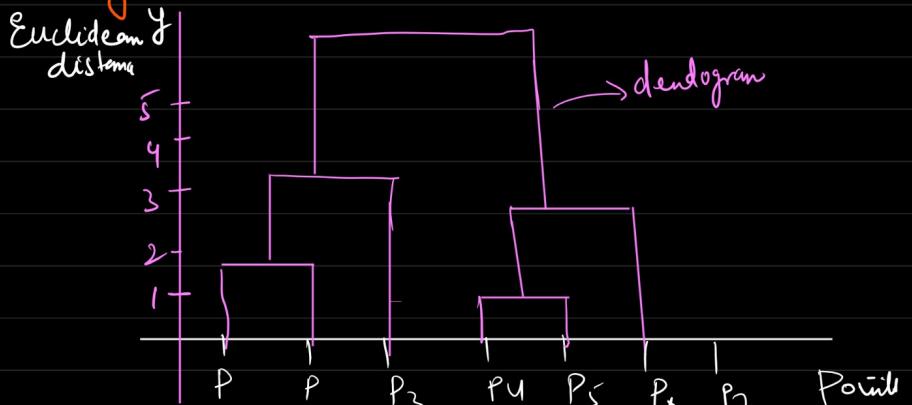
→ Euclidean dist

→ Manhattan distance

→ Cosine Similarity → To calculate distance between two vector (Categorical data / strip variable)

Still how to select K?

Using dendrogram



To determine K

→ threshold on Euclidean distance. (Can be tricky Job \Rightarrow business team)

→ trick \rightarrow we can take longest vertical line of dendrogram where none of the horizontal line of dendrogram is passing.

① Using threshold on Euclidean distance

→ Say threshold = 4
 \Downarrow
 $k=2$

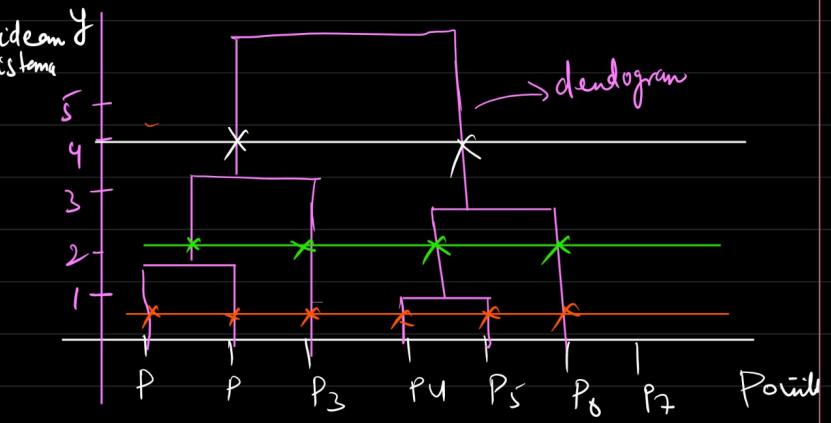
→ Say threshold is 2
 $k=2$

→ say threshold is 1
 $k=6$ (all the dP's is cluster in its own)

②

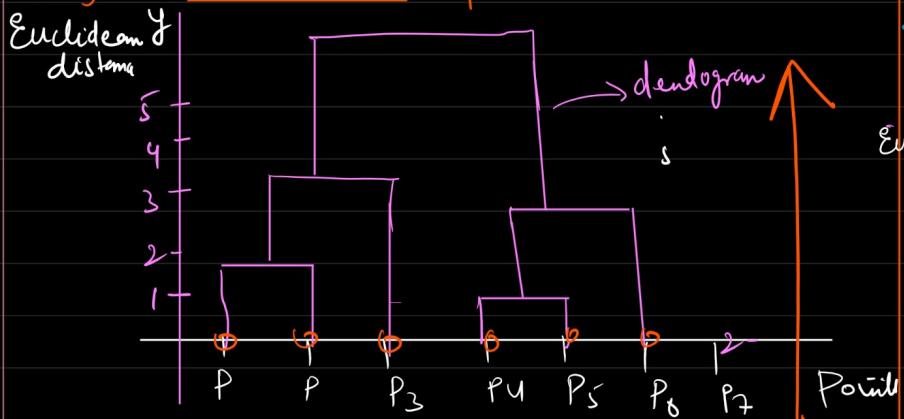
Apart from this horizontal line no any other horizontal line is cutting this vertical line of dendrogram

K=2

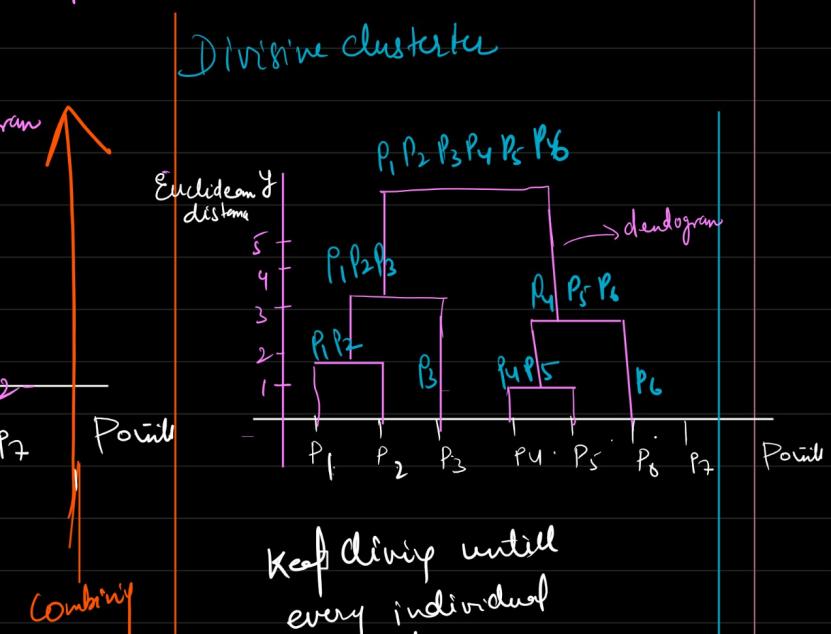


Conclusion → Select the longest vertical line such that no horizontal line passed through it.

Agglomerative clustering



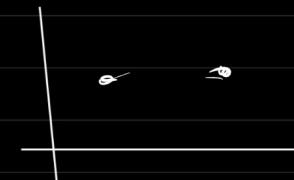
→ Here every dP's are getting combined



Keep dividing until every individual dP becomes cluster in its own.

dividing

K Means	VS	Hierarchical (Scalability & flexibility)
① Size of data	→ Huge → Small	→ K Means → Hierarchical
② K means	— Numerical data (Euclidean Manhattan distance) <small>(only)</small>	
Hierarchical clustering	— Variety of data (Euclidean, Manhattan, cosine similarity)	
③ K means	→ Centroid(k)	→ Elbow method
Hierarchical	— No centroid.	



Hierarchical clustering is an unsupervised machine learning algorithm used to build a hierarchy of clusters. It can be performed using two approaches: agglomerative (bottom-up) or divisive (top-down). In agglomerative clustering, each data point starts as its own cluster, and pairs of clusters are merged iteratively based on their similarity until all points are in a single cluster. Divisive clustering works in reverse, starting with one large cluster and splitting it. The results are often visualized using a dendrogram, which shows the cluster hierarchy. It doesn't require a predefined number of clusters, unlike K-means.

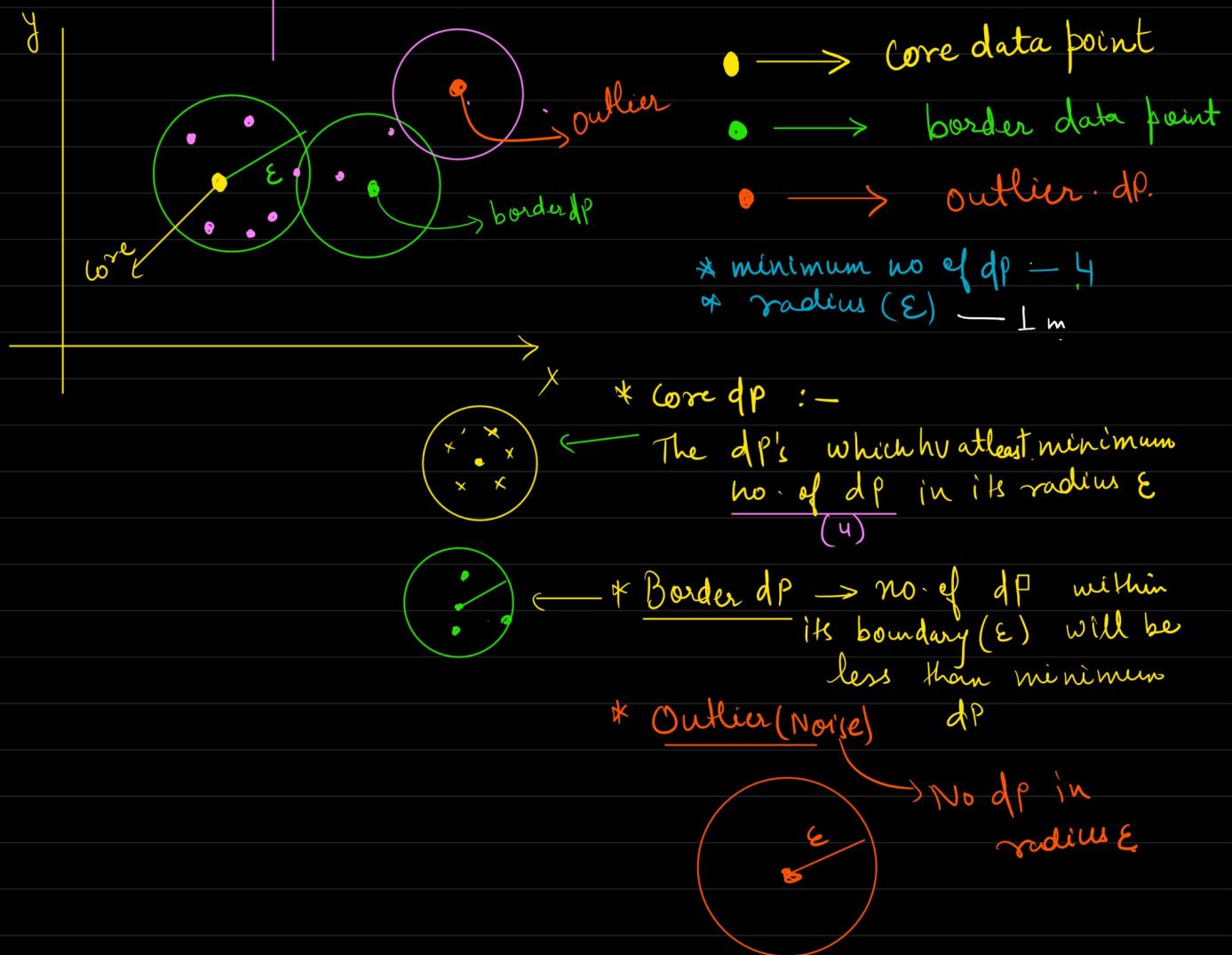
DBSCAN (density based spatial clustering of Application)

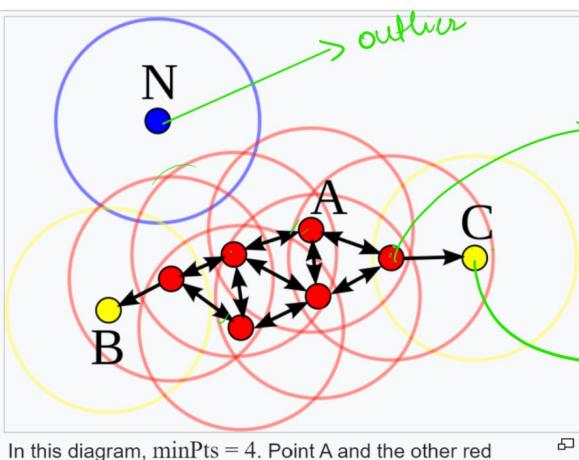
* Characteristics of DBSCAN

- finds groups/pattern
- finds outliers/noise

if already we had Kmean/hierarchical, then why DBSCAN?

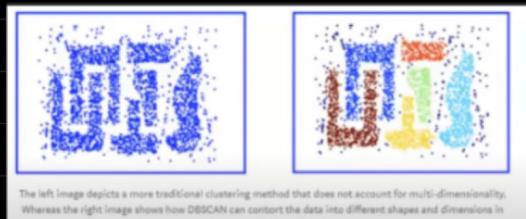
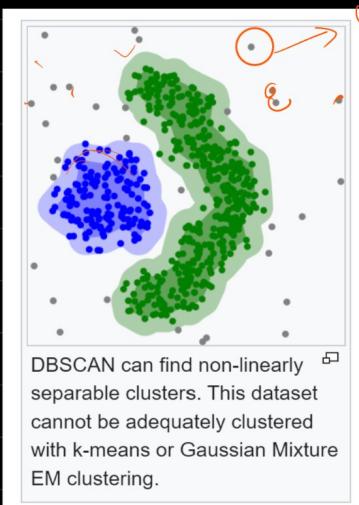
* Distance based Algorithm
such as Kmeans &
hierarchical
can not be used
here (non linear data/non
linear clustering)





In this diagram, $\text{minPts} = 4$. Point A and the other red points are core points, because the area surrounding these points in an ϵ radius contain at least 4 points (including the point itself). Because they are all reachable from one another, they form a single cluster. Points B and C are not core points, but are reachable from A (via other core points) and thus belong to the cluster as well. Point N is a noise point that is neither a core point nor directly-reachable.

* Some Examples of DBSCAN working very well with non linear data.



DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is a clustering algorithm that groups together points that are closely packed, while marking points in low-density regions as outliers. Unlike K-means, DBSCAN doesn't require specifying the number of clusters beforehand. It works by defining a neighborhood around each point using two parameters:

Epsilon : The maximum distance between points to be considered neighbors.

MinPts: The minimum number of points required to form a dense region (cluster).

DBSCAN is effective for identifying clusters of arbitrary shapes and is robust to noise, but its performance depends heavily on appropriate parameter selection.

Silhouette Score / Coefficient

→ Metrics for clustering algorithms.

Classification → accuracy
Regression → r² score

* When we will say, a given cluster is a good cluster?

→ Within cluster datapoints are compact

↑ tightly packed.

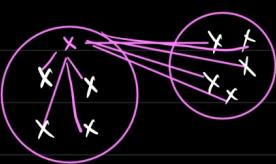


WCSS (within cluster sum of square distance should be minimum)



(intra cluster distance)

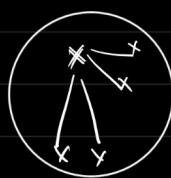
→ Outside cluster sum of square should be as maximum as possible (OCSS)



(inter cluster distance)

Separable clusters | good cluster

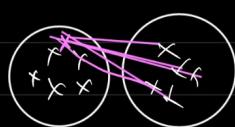
① WCSS



for $d_p(i)$ belonging a Cluster I is given by

$\Rightarrow a(i) = \text{distance within cluster wrt to each } d_p$

② OCSS → distance from nearest cluster of p_s'



$b(i)$

③ Silhouette Score

$$s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}$$

i is any dP

-1 to 1

{ More near to 1 better the
clustering model
is created }

The silhouette score is a metric used to evaluate the quality of clusters created by a clustering algorithm. It measures how similar a data point is to its own cluster (cohesion) compared to other clusters (separation). The score ranges from -1 to +1:

+1: The data point is well-clustered and far from neighboring clusters.

0: The data point is on or very close to the boundary between two clusters.

-1: The data point is likely misclassified, being closer to another cluster than its own.

A higher silhouette score indicates better-defined clusters. It is commonly used to assess the optimal number of clusters.