

Anomaly detection (To detect outliers)

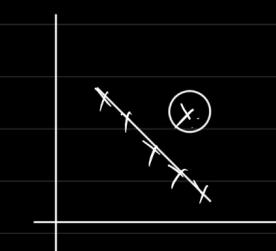
↓
Abnormality

→ Anomaly detection algos are Unsupervised.

Why?

→ In few case Anomaly plays an important role.

Ex1	transaction amount	Acc.no
1000	—	—
2000	—	—
1500	—	—
3000	—	—
2 Lakh	—	—



→ sometimes to detect outliers traditional methods like Box plot, violin plot fails.

This customer being outlier, holds an important role. ⇒ Bank will target this customer for the loan.

∴ Outliers sometimes makes sense.

Ex2 → Fraud transaction.

tr₁ - Delhi
tr₂ - Delhi
tr₃ - Delhi } Normal

tr₄ - USA → You receive a msg ⇒ anomaly.

Ex3 → When you login a gmail or any Portal.

Ex4 VK. → 5, 10, 5, 4, 3, 5, 110

Ex5 Cancer prediction ⇒ Normally people don't have Cancer. This outlier but important.

Ex6 → Fraud IP.

* Detection of Outliers

- ① Isolation Forest Anomaly detection
- ② DBSCAN
- ③ Local Outlier Factor Anomaly detection.

① Isolation Forest (VSL algorithms)



→ multiple Isolation trees

→ Isolation tree is like a decision tree

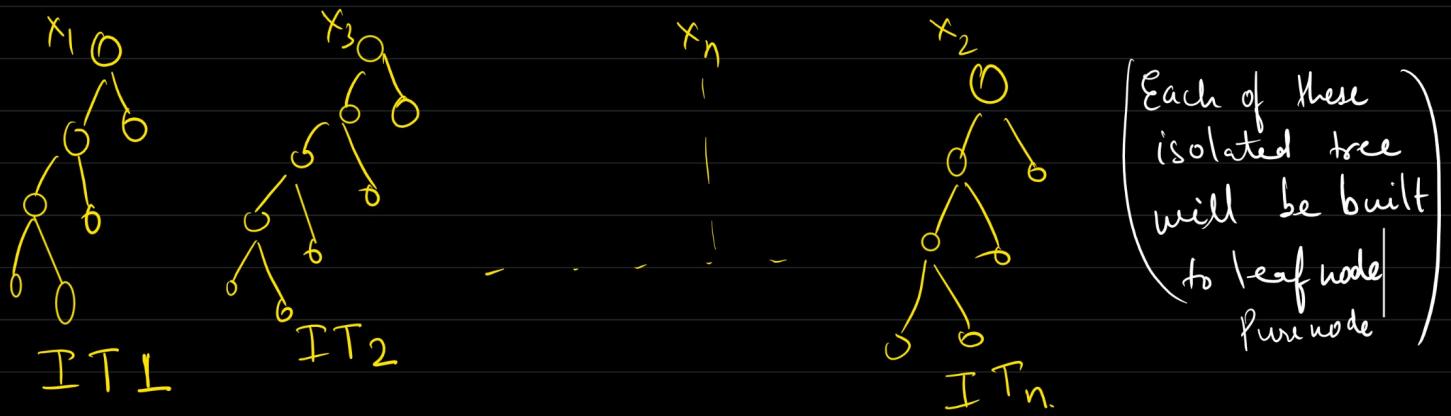
Construction of Isolation tree

→ Select a feature randomly X_1

→ Randomly choose a split value within the range of X_1

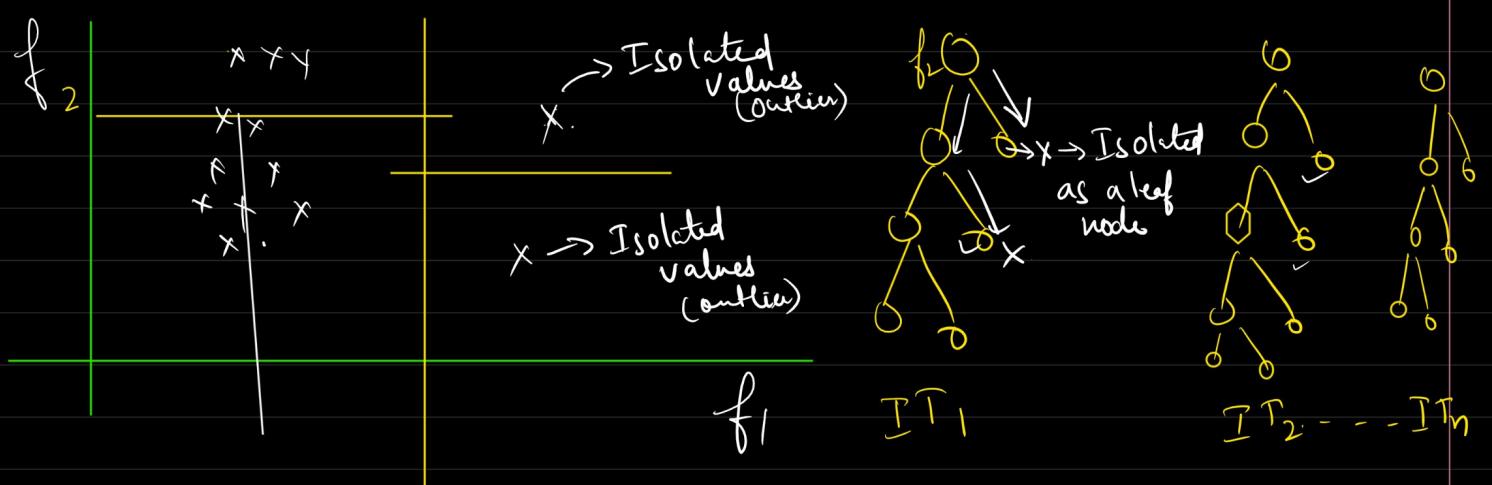
→ Repeat the process recursively to build the tree.

X_1	X_2	X_3	X_4
-	5		
5	6		
6	7		
7	8		
8	9		
9	10		
10	-		



* How you will get outlier?

→ Anomaly / outlier due to their distinctiveness, tend to end up in leaf nodes at shorter path.



* Since outliers are different from normal dp, during construction of isolated trees it will be isolated as a separate leaf node.

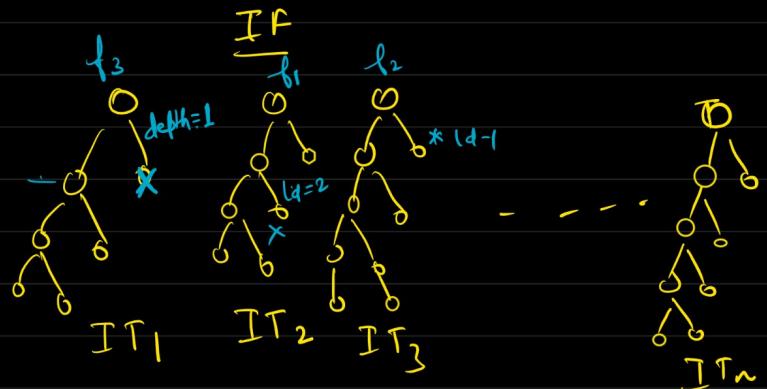
Anomaly Score

$$S(x, m) = 2 \frac{-E(h(x))}{C(m)}$$

where m is total no of datapoints
 $x \rightarrow$ datapoint for which you want to check anomaly score.

$E(h(x)) =$ Average search depth
 ↓
 for one dp of x in all of the isolated tree.

$C(m) =$ Average depth of all the dp's in all Isolation tree



$$\underline{E(h(x)) \ll C(m)}$$

Since $C(m)$ is avg depth for all the dp's so $C(m)$ will be far greater than $h(x)$.

$$S(x, m) = \frac{-E(h(x))}{C(m)}$$

will be very small value

$S(x, m) \approx 1 \Rightarrow$ Anomaly Score \Rightarrow outliers

Generally If $S(x, m) > \frac{0.5}{\text{threshold}} \Rightarrow \text{Outlier.}$

$E(h(x)) \gg C(m) \Rightarrow S(x, m) \approx 0.5 \Rightarrow \text{Normal dp's}$

≤ 0.5

Isolation Forest is an anomaly detection algorithm that isolates outliers by randomly partitioning data. It constructs multiple isolation trees by recursively splitting data using random features and split values. Anomalies, being distinct and few, tend to be isolated with fewer splits, leading to shorter paths in these trees. The algorithm calculates an anomaly score based on the average path length, where a higher score indicates a higher likelihood of the point being an anomaly. Isolation Forest is efficient, does not assume any data distribution, and performs well with high-dimensional data, making it ideal for large-scale anomaly detection tasks.

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) can be used for anomaly detection by identifying points that do not belong to any dense cluster. It groups data points based on their density, using parameters eps (neighborhood radius) and minPts (minimum points in a neighborhood). Points that are isolated and do not meet these criteria are classified as noise, which are considered anomalies. DBSCAN is effective because it doesn't require specifying the number of anomalies, can detect outliers in clusters of arbitrary shapes, and is robust to noise, making it suitable for complex, real-world datasets.

Local Outlier Factor (LOF) is an anomaly detection algorithm that identifies outliers by comparing the local density of a data point with that of its neighbors. It calculates a LOF score, where a score significantly greater than 1 indicates the point is an outlier. The algorithm considers the distance to a point's k -nearest neighbors to determine its local reachability density. LOF is effective in datasets with varying density, as it detects outliers based on the local context rather than global distribution. It's commonly used in fraud detection, network security, and other applications requiring the identification of subtle anomalies.

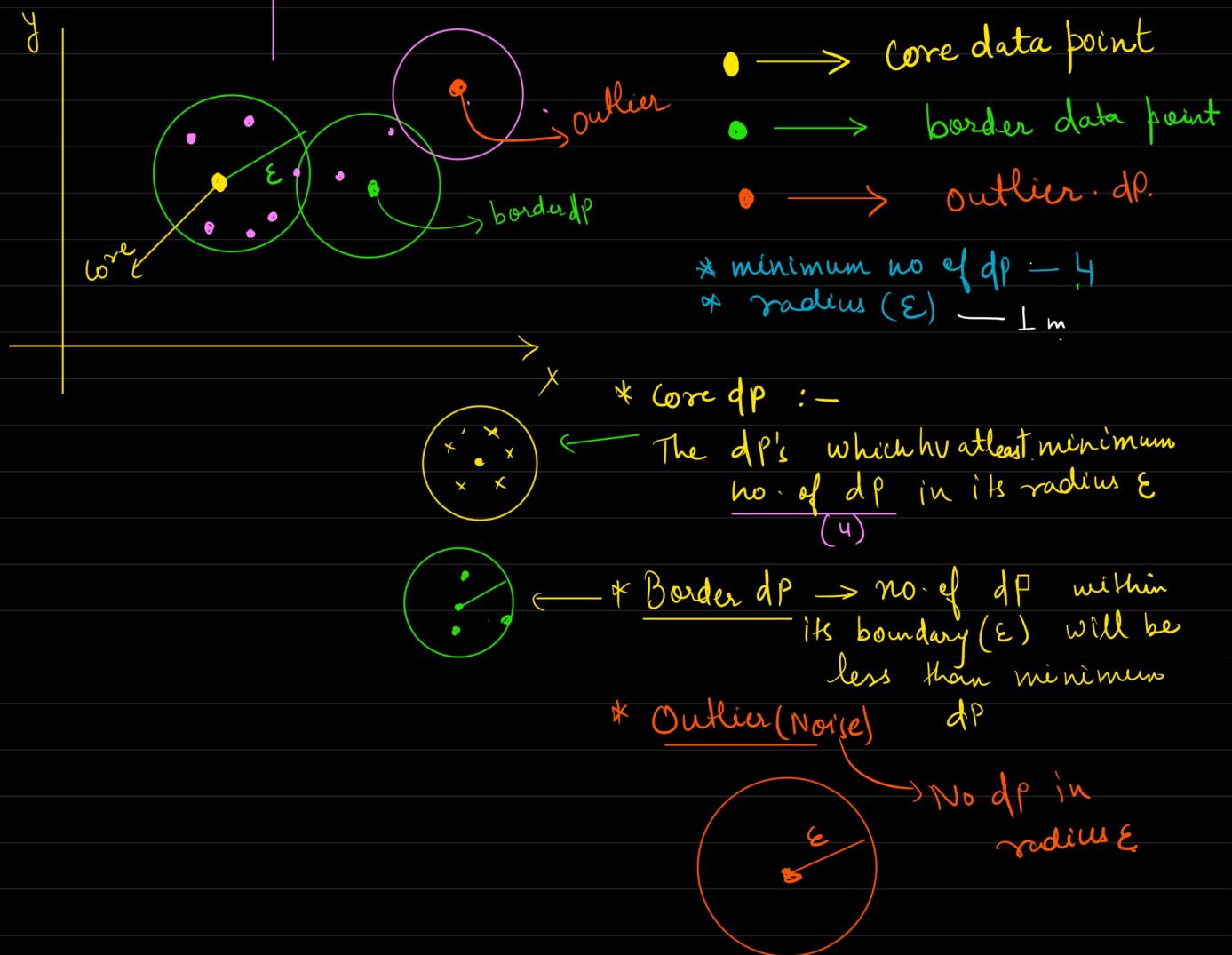
DBSCAN (density based spatial clustering of Application)

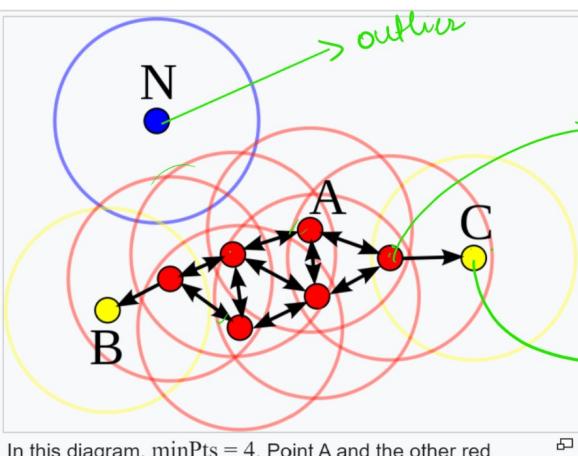
* Characteristics of DBSCAN

- 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)

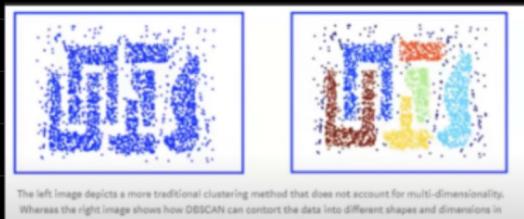
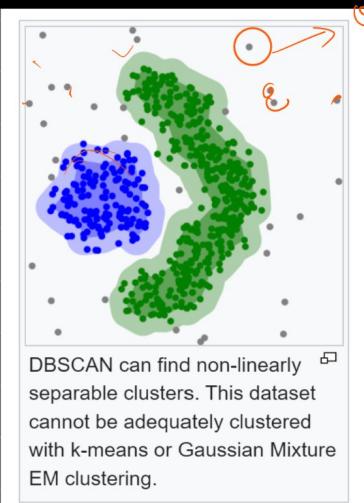




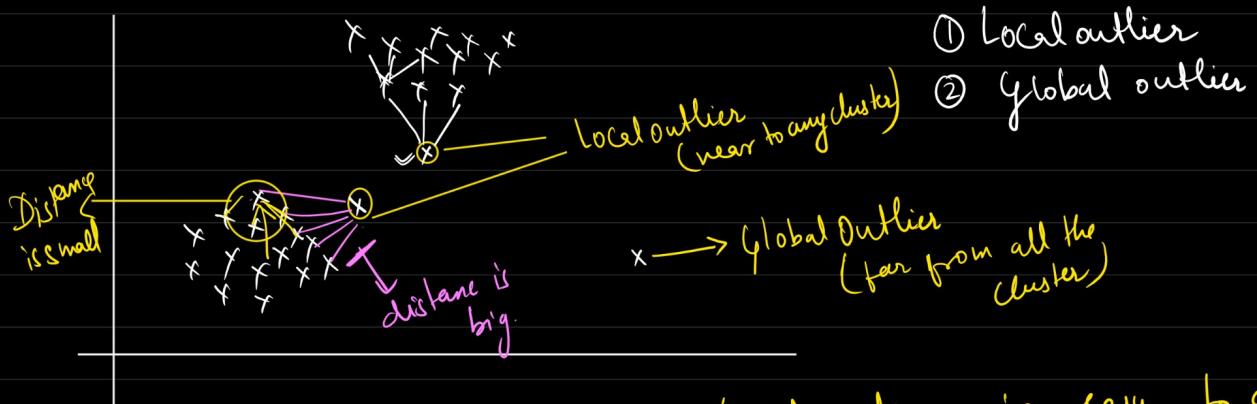
In this diagram, $\text{minPts} = 4$. Point A and the other red points are core points, because the area surrounding these points in an ϵ radius contain at least 4 points (including the point itself). Because they are all reachable from one another, they form a single cluster. Points B and C are not 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.

outlier
core p
border p

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



Local Outlier factor Anomaly detection



→ Global outlier is easy to detect using Isolation forest / DBSCAN.

→ How to detect local outlier?

→ Using Local outlier factor.

* Calculate distance of a dp with its k-nearest neighbour. ↑
local density - Using kNN.

* dp's which are far will have more distance from its neighbours. and vice-versa

* Higher distance → density less \Rightarrow Outlier

$$\text{LOF score} \Rightarrow \frac{\text{# of Interest distance}}{\text{All dp's distance}}$$

\geq threshold \Rightarrow outlier.