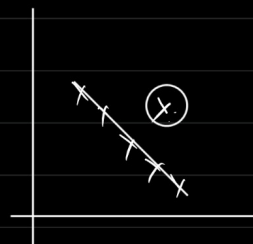


Anomaly detection (To detect outliers)

⇓
Abnormality

→ Anomaly detection algos are Unsupervised.



Why?

→ In few case Anomaly plays an important role.

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

Ex1

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

→ 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 (USL algorithms)

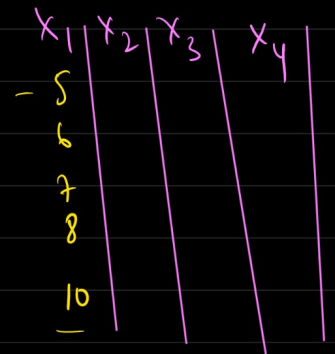
⇓

⇒ 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.



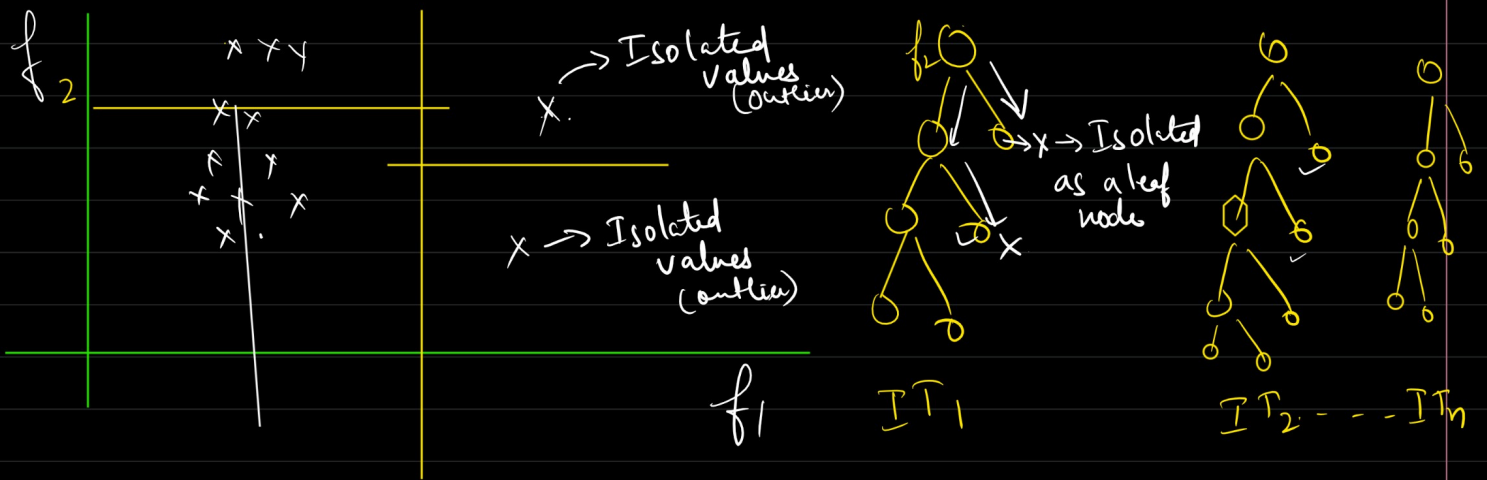
x_n



(Each of these isolated tree will be built to leaf node / pure node)

* 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.

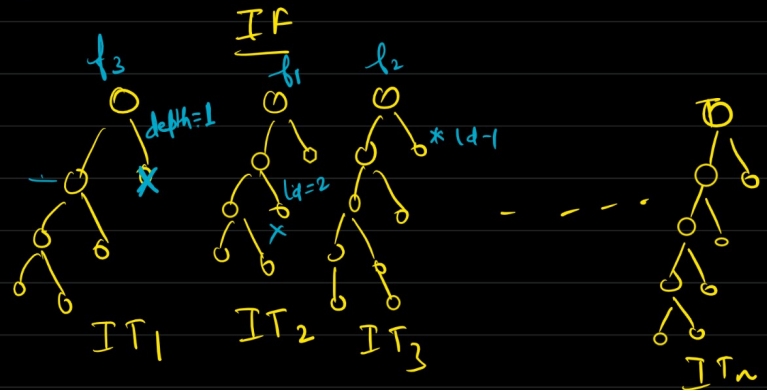
* Anomaly score

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

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

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

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



$E(h(x)) \ll C(m)$ Since $C(m)$ is avg depth for all the DP's so $C(m)$ will be far greater than $h(x)$

$$S(x, m) = \frac{-E(h(x))}{2^{C(m)}} \rightarrow \text{will be very small value}$$

$S(x, m) \approx 1 \Rightarrow$ Anomaly score \Rightarrow outlier

Generally $gf \leq (x, m) \geq \underbrace{0.5}_{\text{threshold}} \Rightarrow \text{Outlier.}$

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