

Anomaly detection (USL)

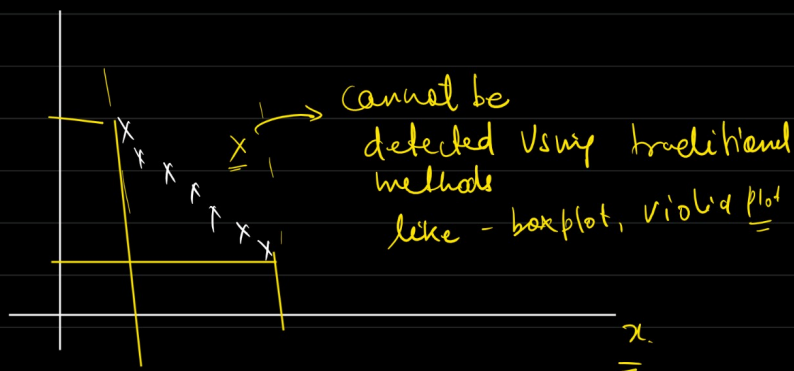


Abnormality

Ex1

transaction Amt	Acc no
1000	-
2000	-
3000	-
2Lakhs	-

→ 2Lakhs → this customer is a outlier.



Ex2

fraud transaction 3 trans → 15000

t₁ - 11 AM

t₂ - 1 PM

t₃ - 2 PM

t₄ } → night
t₅
t₆

DD
USA

Ex3 gmail or any portal?

✓ Ex4 0, 1, 0, 9, 10, 35, 76

Ex5 Fraud IP:

- ① Isolation Forest Anomaly detection
- ② DBSCAN
- ③ Local outlier factor Anomaly detection

* use case → specifically solving an outlier Anomaly problems.

* Use Case

ABC

→ 1 lakh transactions

SL → fraud / not fraud

→ 2004

→ 20 years

→ 3-5%

UK bank

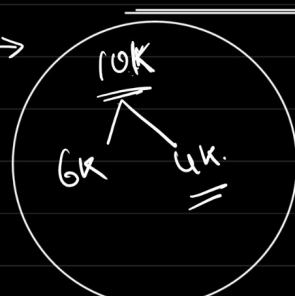
DBSCAN

data:

1 lakh

only the transaction which are abnormal / outliers

2010
2015
2019

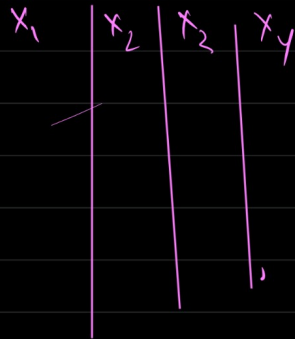


Classification Problem

① Isolation Forest (USL)

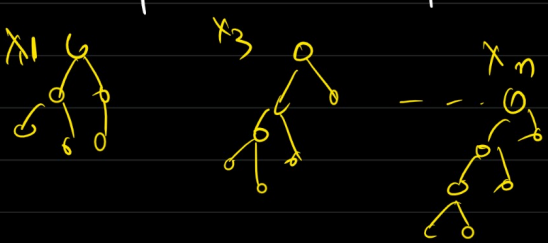


- multiple Isolation trees
- Isolation is like a decision tree.



Construction of Isolation tree.

- Select a feature randomly x
- Randomly choose a split value with range of x
- Repeat the process recursively to build the tree.



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

* How you will get outliers

Anomaly/outlier due to their distinctiveness, tend to end up in leaf node at shortest path.

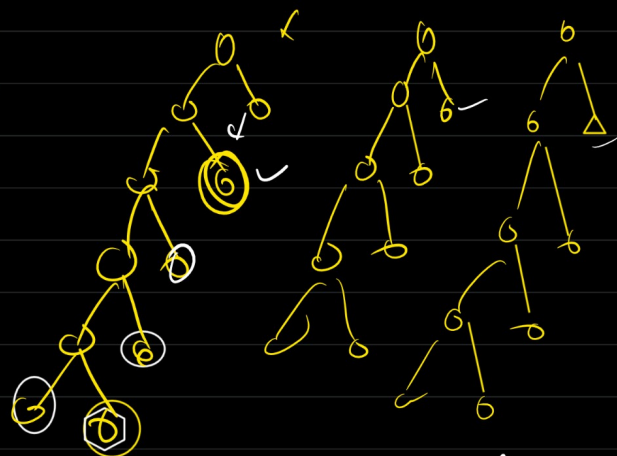
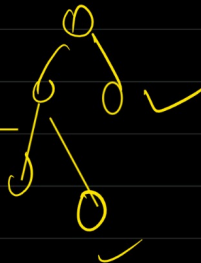


1, 2, 3, 4, 5, 100

≤ 100 , > 100

Isolation

x



----- n

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

How much confidence?

* Anomaly score

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

m is total no of dp's

(one dp)

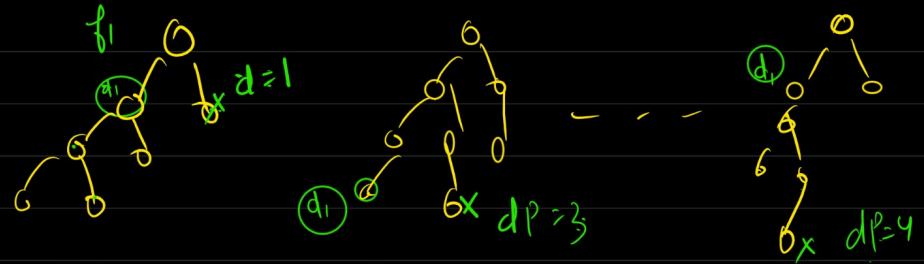
$E(h(x))$ → Average search depth of x in all isolated tree.

x → datapoint for which you want to check anomaly

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

$$E(h(x)) \ll C(m)$$

Since $C(m)$ is avg
depth of all dp's so
 $C(m)$ be far greater
than $E(h(x))$



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

very small
very large

very very small.

$$2^{-\text{very small no}} = \frac{1}{\text{very small sm}} \approx 1$$

generally if $S(x, m) \geq 0.5$ \Rightarrow outlier

threshold

$< 0.5 \Rightarrow$ Normal dp

$$E(h(x)) \gg C(m)$$

$$2^{-\frac{E(x)}{C(m)}} = \frac{1}{2^{\log_2}} = 0.5$$

DBSCAN

\rightarrow Studied VSL

\rightarrow Also detects the outlier \rightarrow If a dp is not
core / boundary dp \Rightarrow That's an
outlier.

* Local Outlier factor



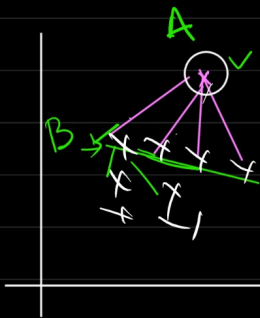
How to detect local outlier.

↓
Local outlier factor.

* Calculate distance of a dp with its k-nearest neighbour

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

* Higher distance \Rightarrow density less \Rightarrow outlier.



$$\underline{LOF}_{score} = \frac{\text{Interested dp. distance}}{\text{all dp's distance of neighbourhood dp}}$$

$> \underline{\text{Threshold}} \Rightarrow \underline{\text{outlier}}$

LOF - (0-∞)

$\sim 1 \rightarrow \text{Normal}$

$> 1 \rightarrow \text{outlier}$

$< 1 \rightarrow \text{Normal}$