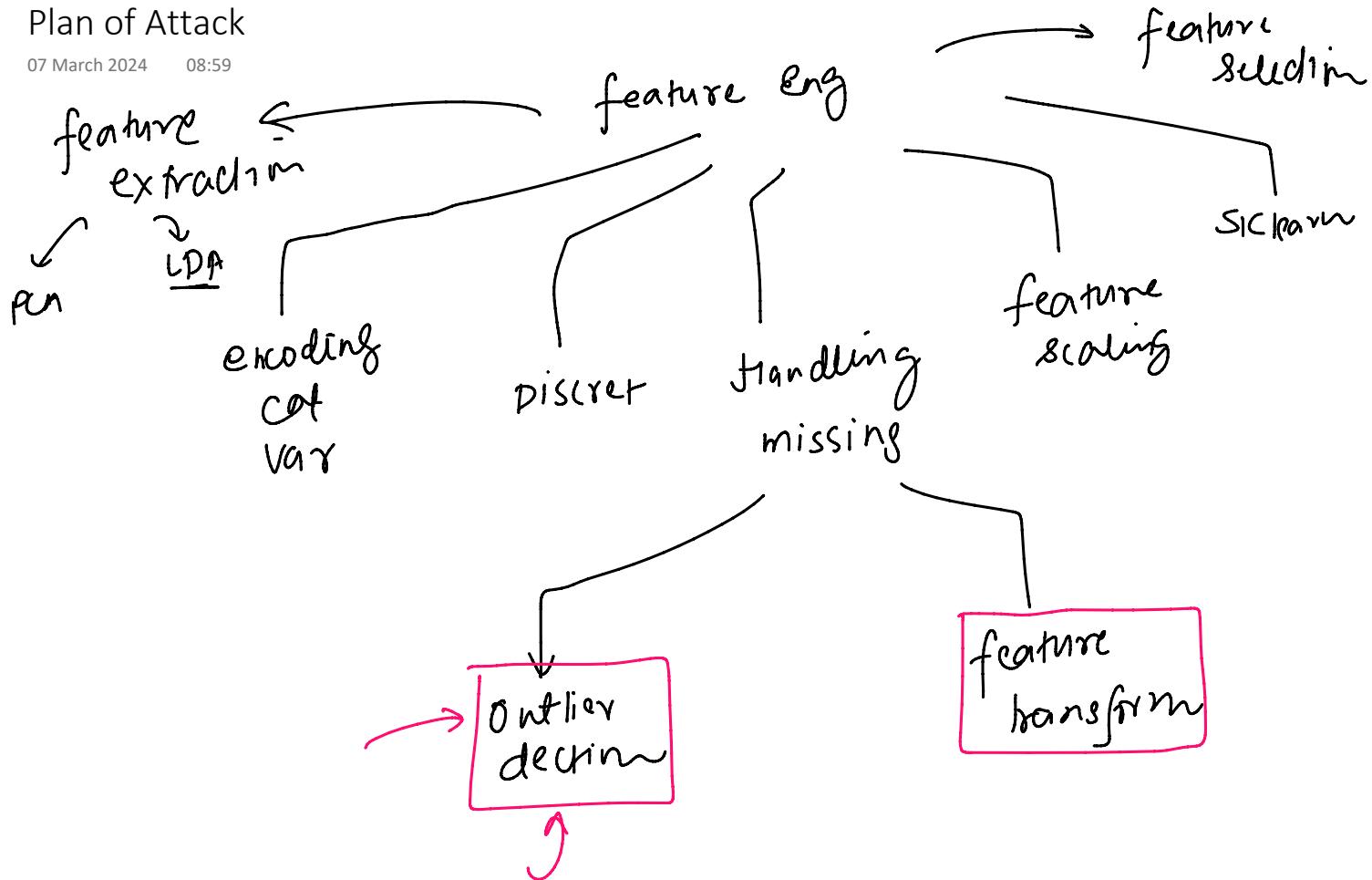


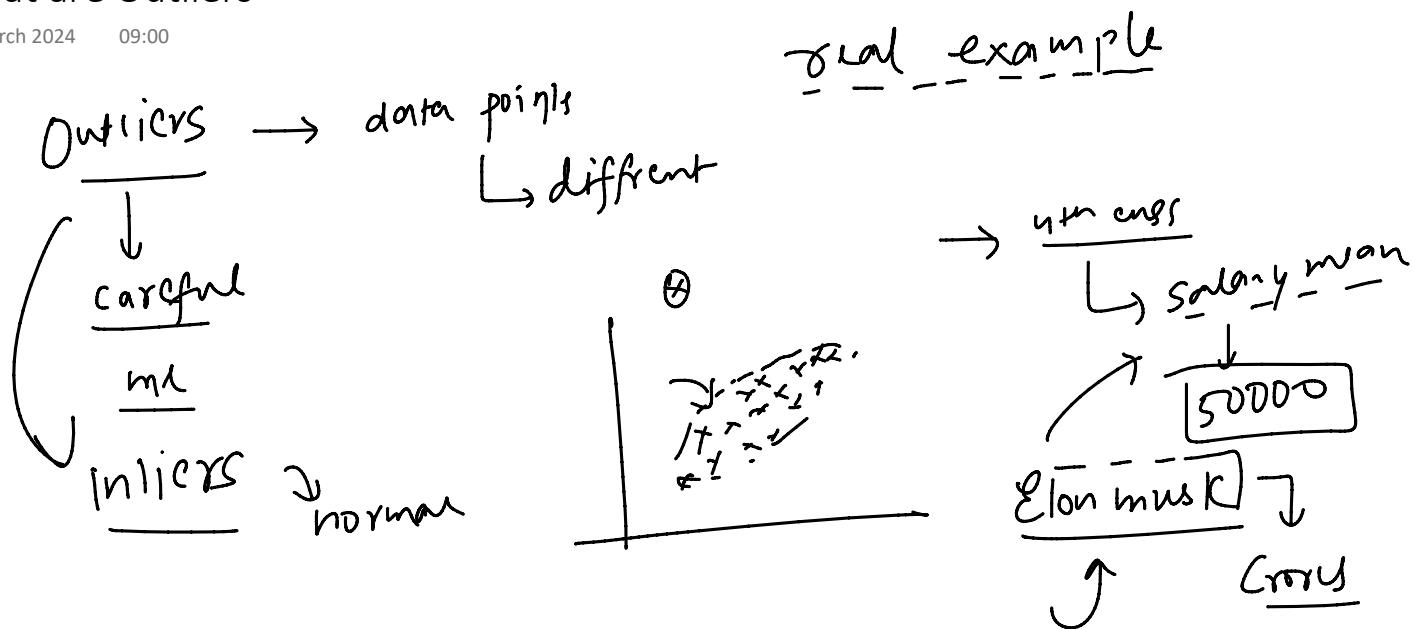
Plan of Attack

07 March 2024 08:59



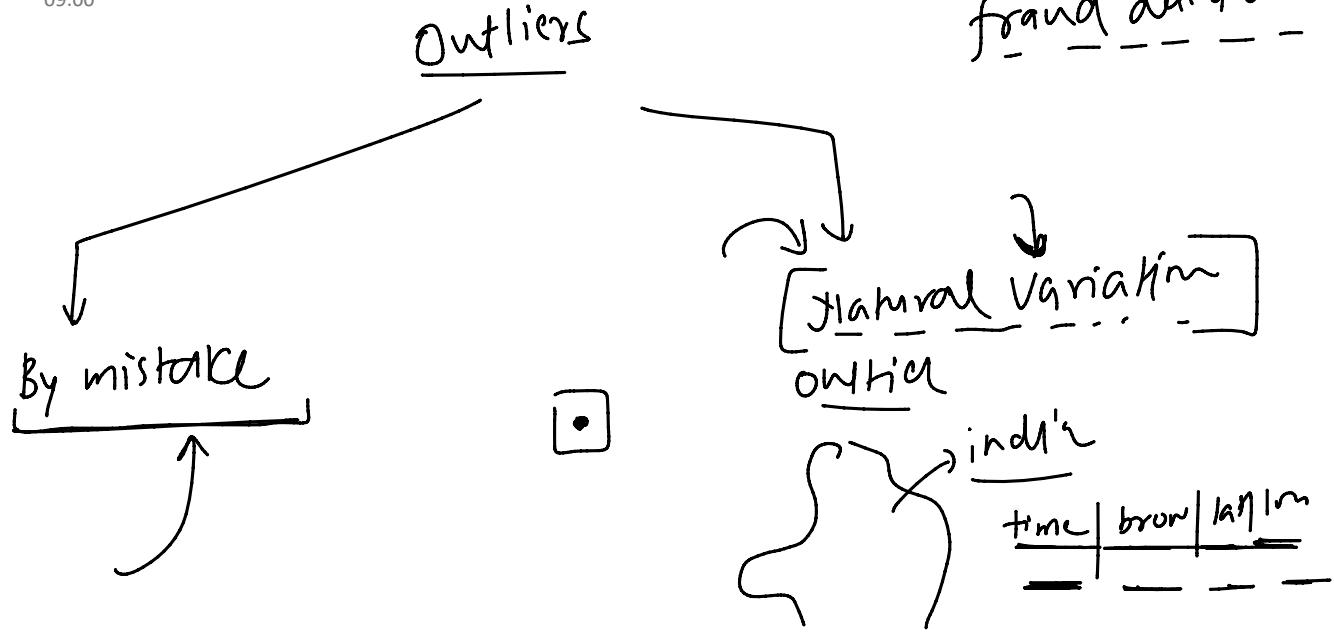
What are Outliers

07 March 2024 09:00



Types of Outliers

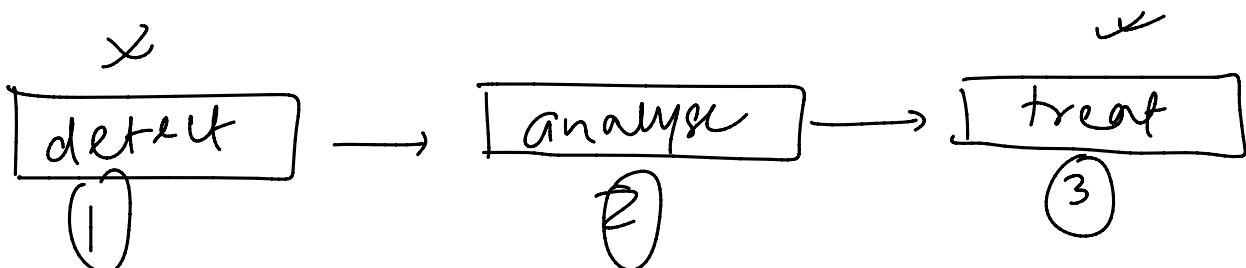
07 March 2024 09:00



Outlier

- 1) outlier detection →
- 2) By mistake) natural variation
- 3) outlier treatment

flow outlier

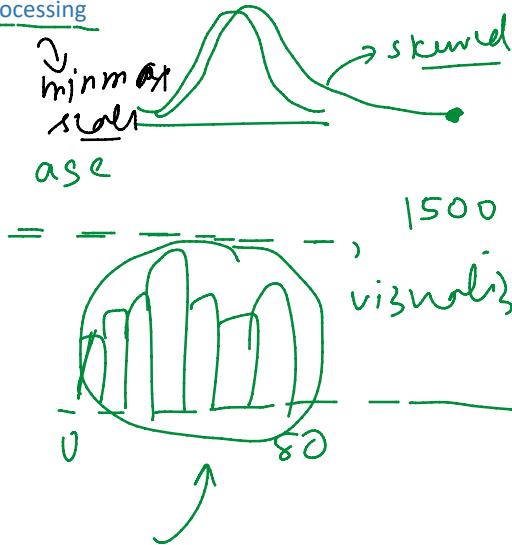


Impact of Outliers

07 March 2024 09:00

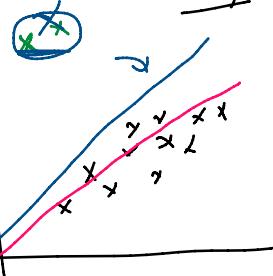
- 1. Impact on statistical measures
- 2. Impact on model performance
- 3. Assumption Violation
- 4. Data Visualization
- 5. Data Preprocessing

Robust
standard
error

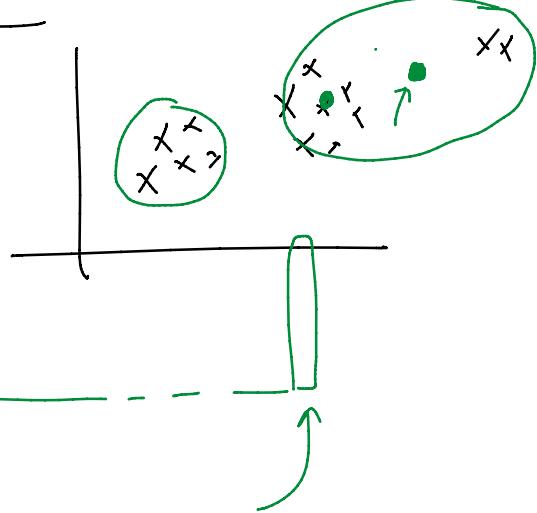


mean / std
data analysis

Why

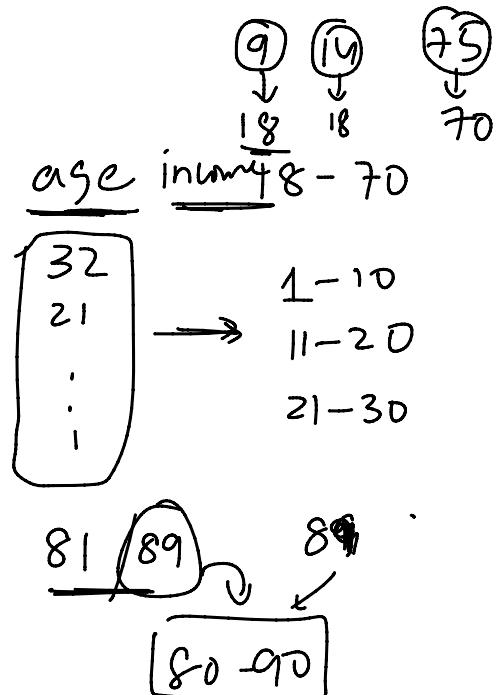
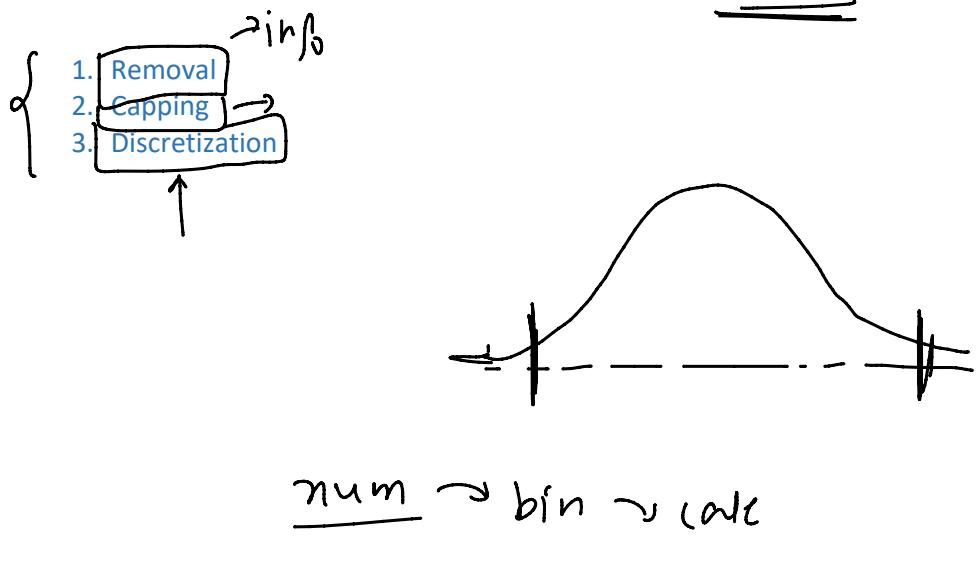


→ classification
clustering
2 clusters



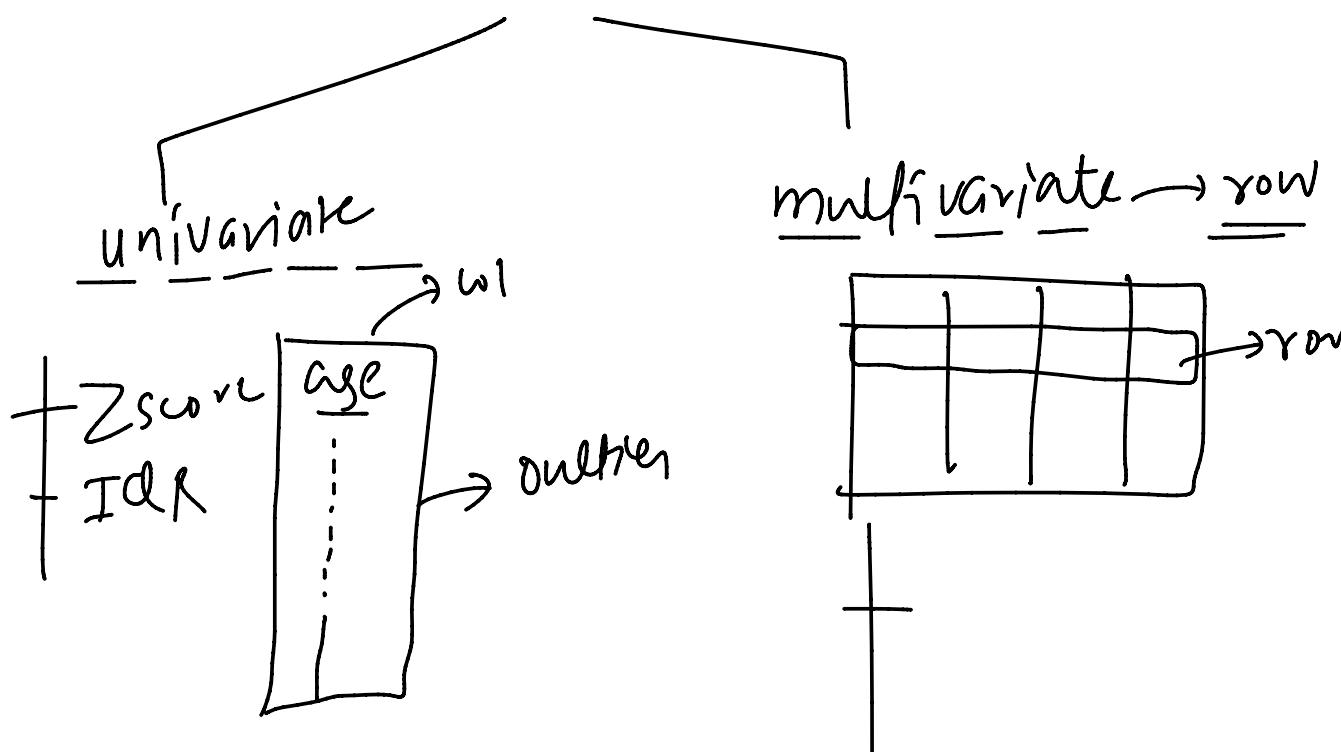
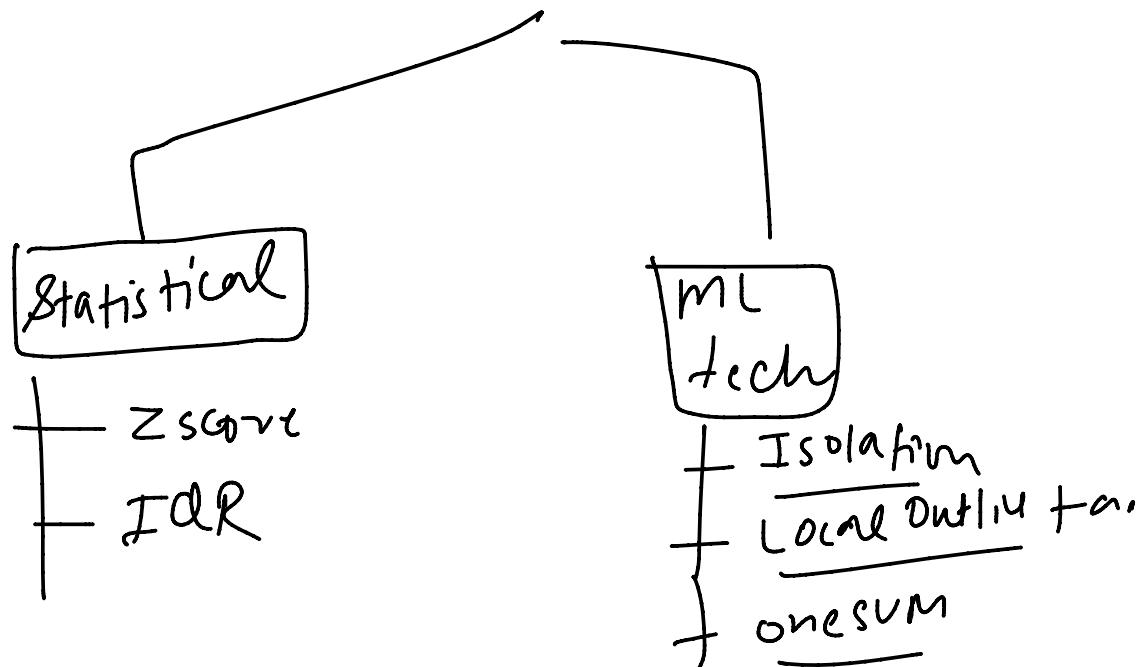
How to deal with outliers

07 March 2024 10:25



Outlier Detection Techniques

07 March 2024 09:41



Z-Score → univariate outlier detection

07 March 2024

09:41

	age	income	purchase
1	32	32000	0
2	64	150000	1

Big assumption

$$\begin{bmatrix} \mu \rightarrow \\ \sigma \rightarrow \end{bmatrix}$$

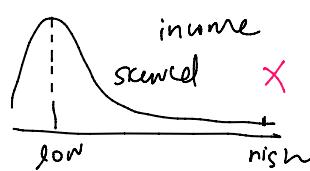
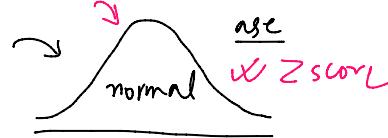
→ normal / normal like

$$[\mu + 3\sigma / \mu - 3\sigma] \quad [82, 1, 3]$$

outliers

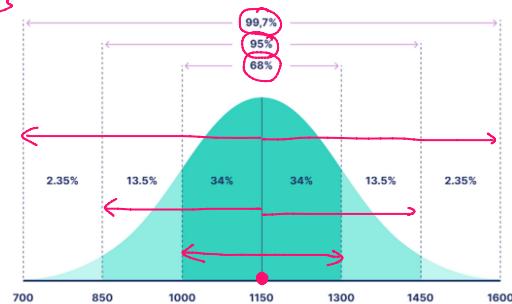


Z score



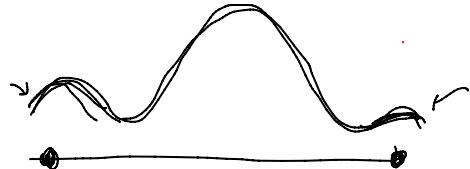
$$\frac{x_i - \mu}{\sigma}$$

$[-3 \text{ to } 3]$



$$\begin{bmatrix} \mu - 3\sigma \\ \mu + 3\sigma \end{bmatrix}$$

0.03 %
outlier threshold



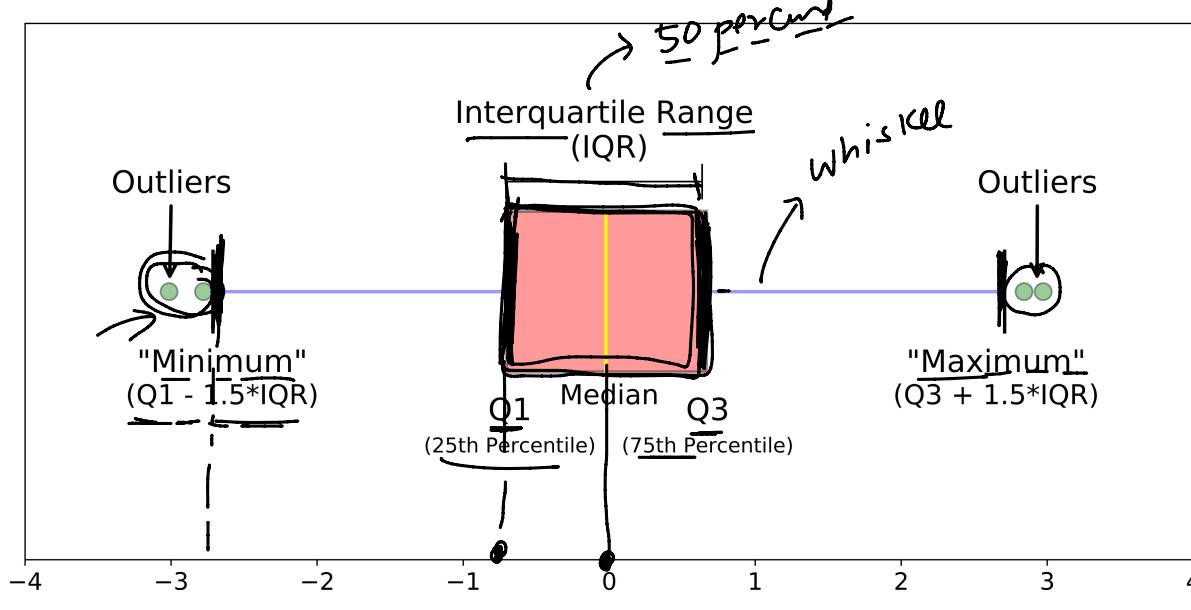
[IQR and BoxPlot]

07 March 2024

09:41

→ normal → col → skewed

numerical
col

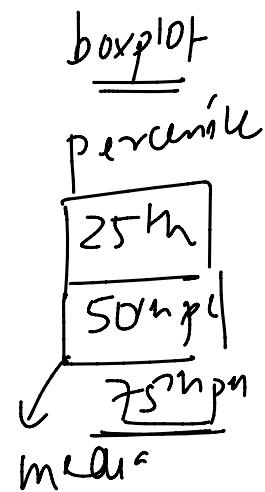


$$\underline{Q_3 + 1.5 IQR}$$

$$\underline{99\%} \rightarrow \underline{360}$$

$\underline{5\%} \rightarrow \underline{10}$

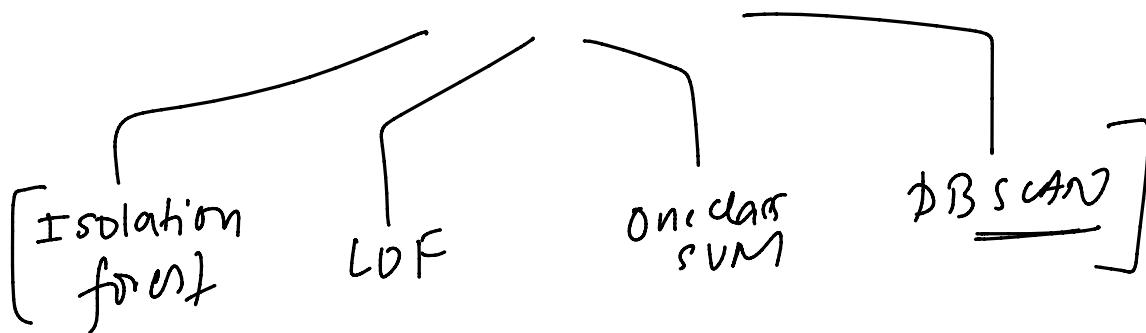
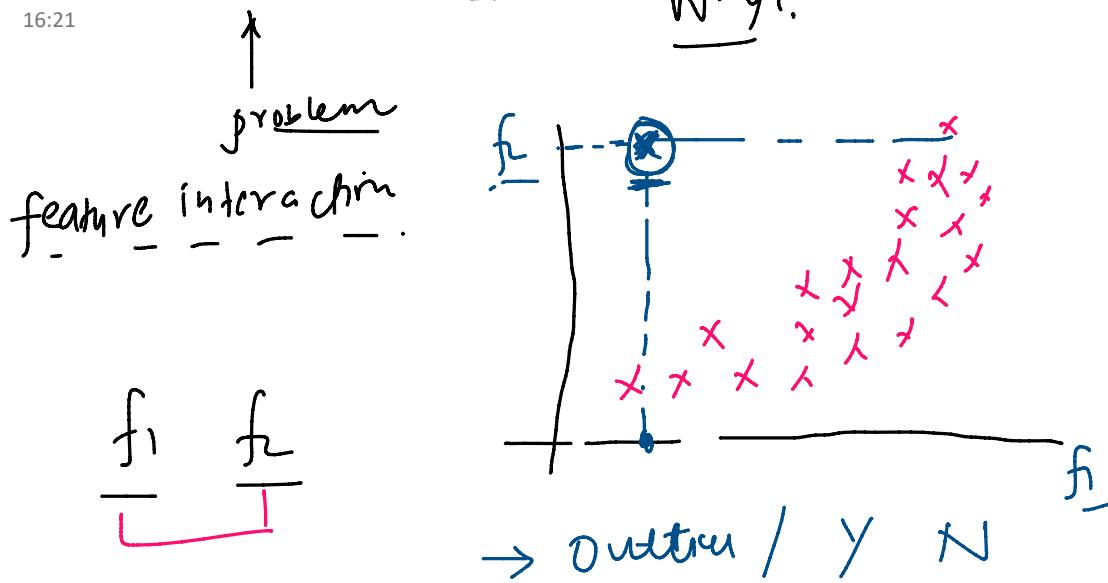
$$\underline{\overline{Q_1} - 1.5 \overline{IQR}}$$



Problem with Univariate Techniques

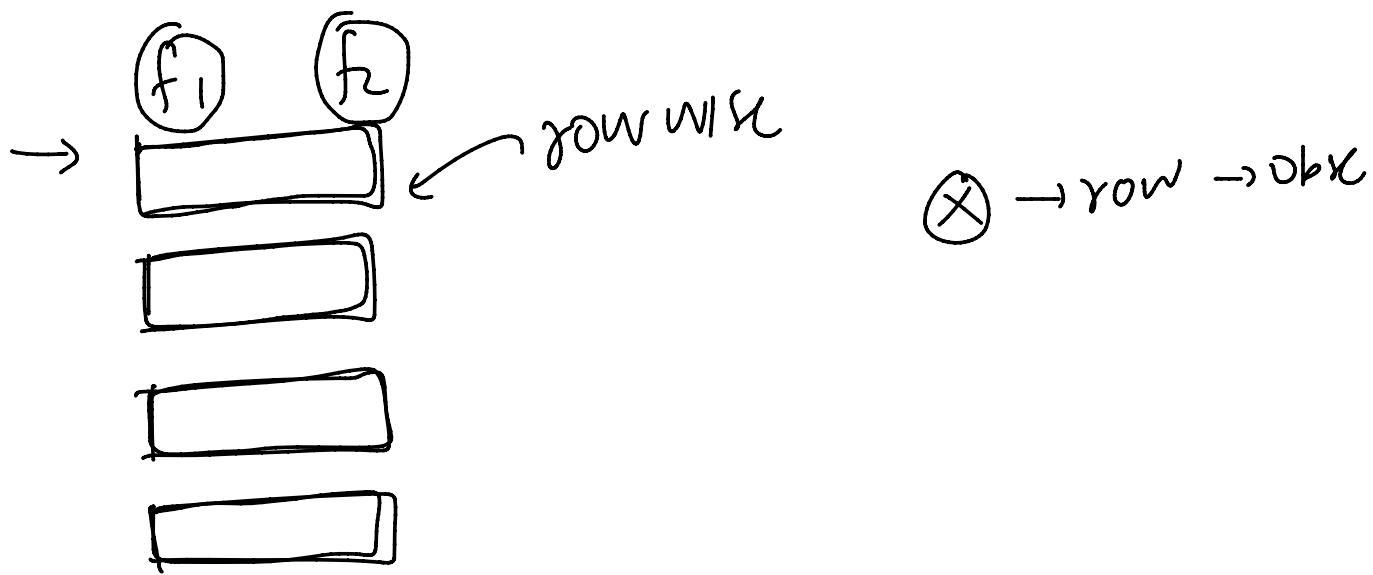
07 March 2024 16:21

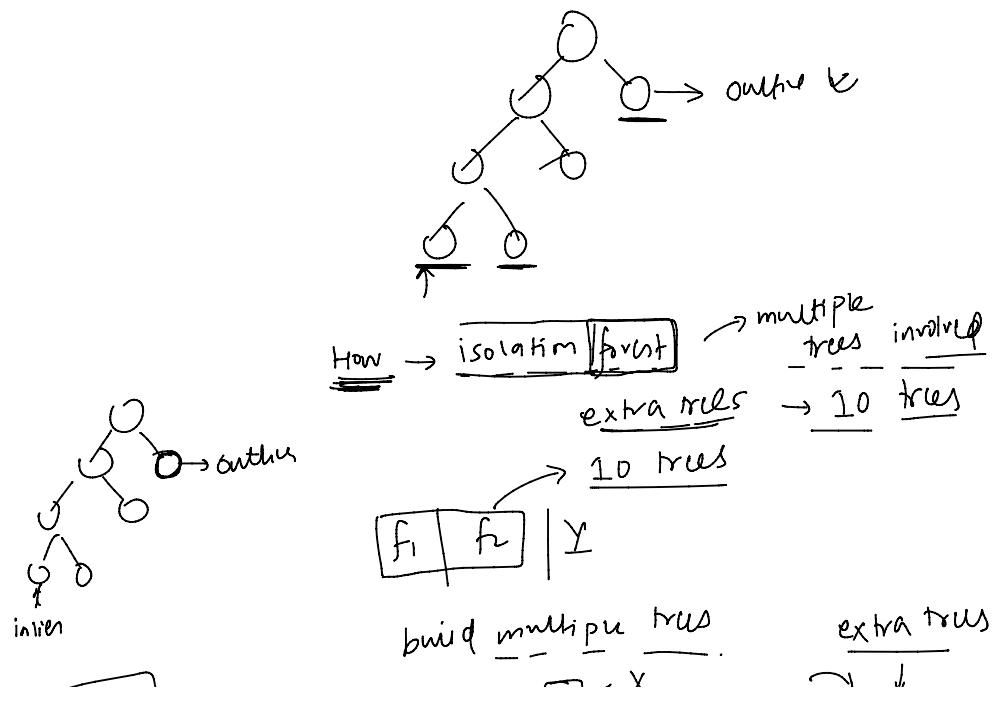
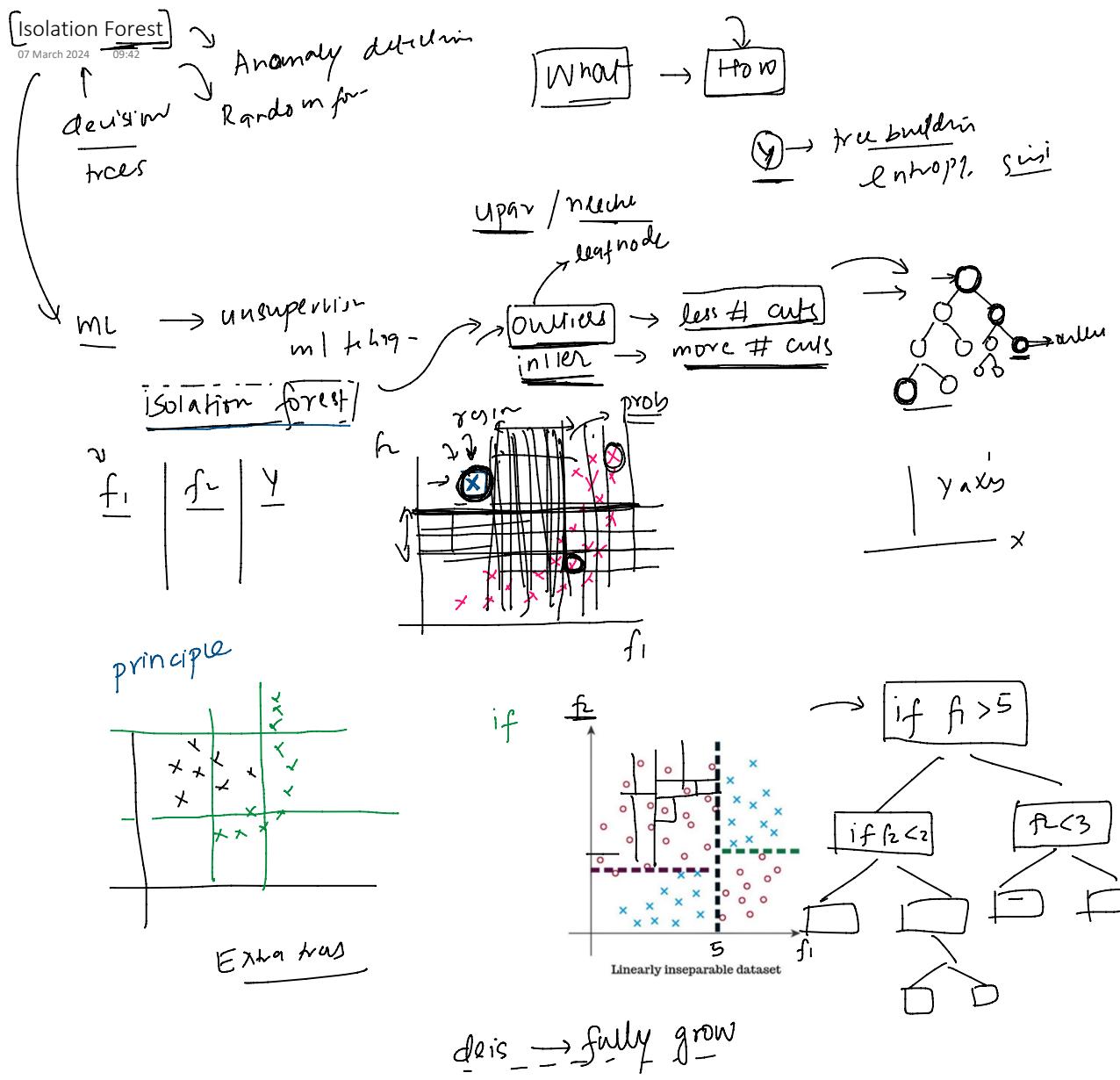
Why?

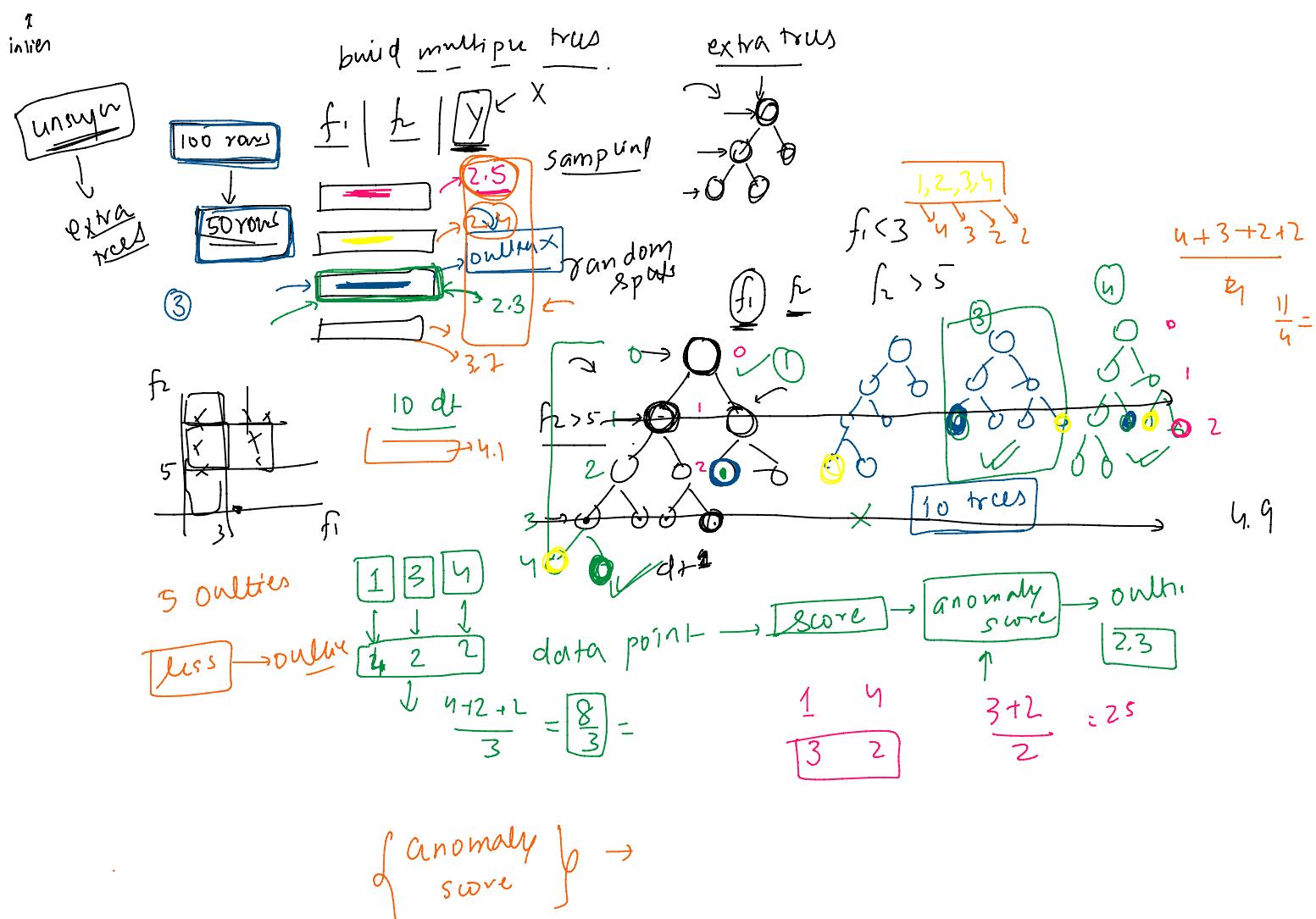


Multivariate Outlier Detection

07 March 2024 16:38

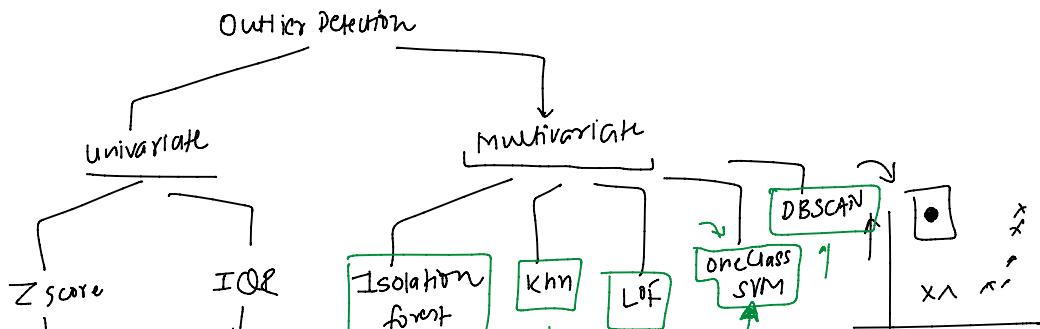


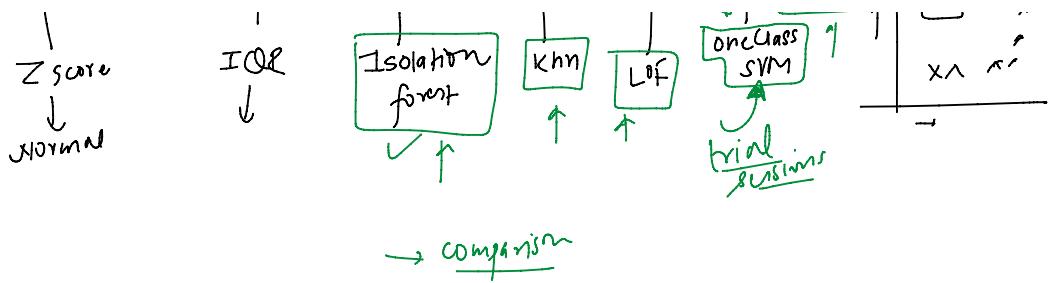




- deduce # extra trees
- max sample
- build all trees
- for every row
 callibrate the avg position of leaf
soullie in the tree in which they
 are present

1. anomaly / ȝow

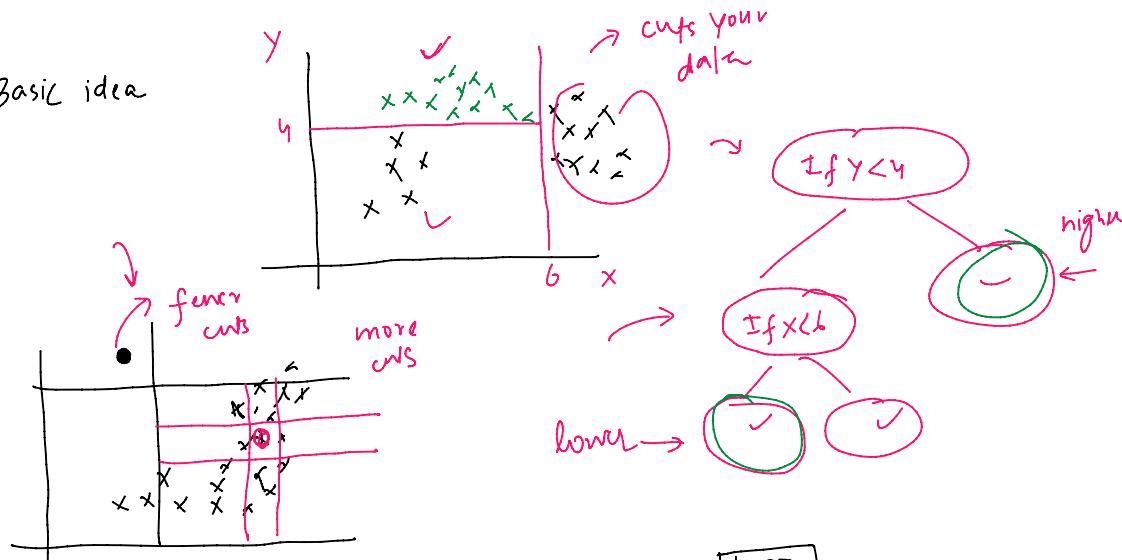




Quick Revision

Isolation forest

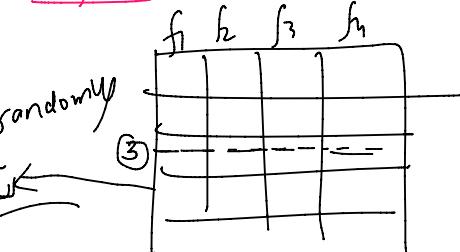
Basic idea



isolation \rightarrow unsupervised

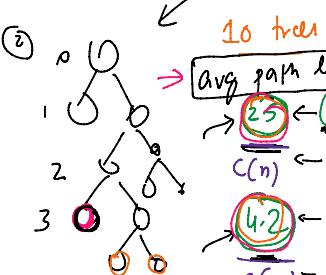
$\rightarrow 1000$

multiple trees



10 trees \rightarrow 700 data

randomly



path len 1 path len 2 ... path len 40

extra tree

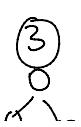
$$E(h(x_i))$$

(P3)

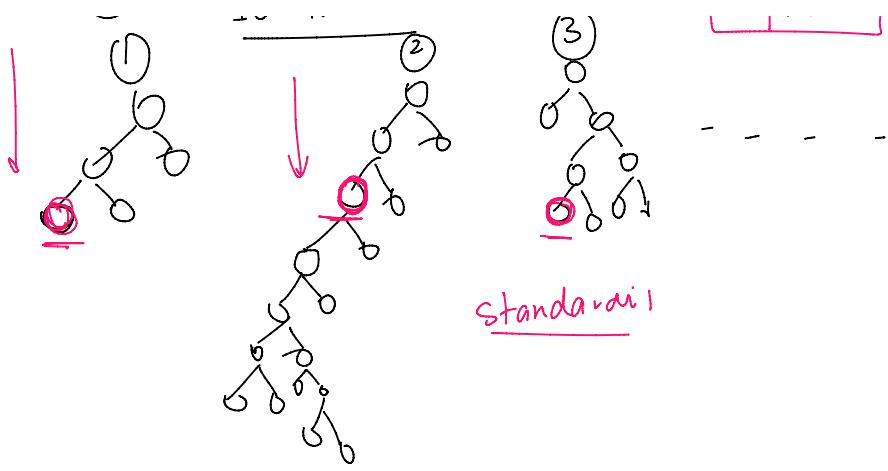
anomaly score \neq avg path length

(2)

10 trees



$\boxed{\text{depth}=2}$



Calculation of Anomaly Score

07 March 2024 17:50

$$\rightarrow \boxed{P3} \quad \frac{1000}{c(n)}$$

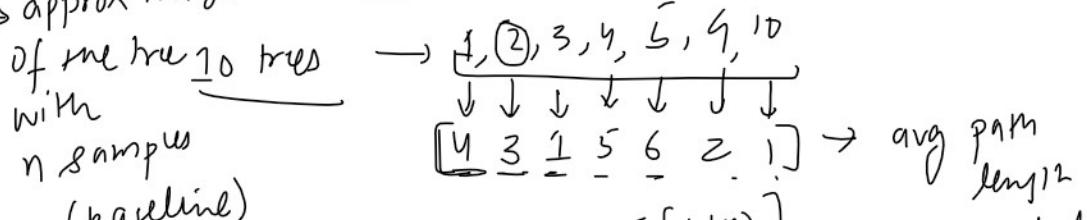
$$s(x, n) = 2^{-\frac{E(h(x))}{c(n)}} \quad \text{approx height}$$

$$\text{Anomaly: } Z^{-\frac{E(h(x))}{c(n)}}$$

$$E(h(x)) = \frac{1}{N} \sum_{i=1}^N h_i(x)$$

with
n samples
(baseline)

$h(x)$ → path length



$$c(n) = 2(\ln(n-1) + 0.5772156649) - \frac{2(n-1)}{n}$$

It's a baseline measure of how deep we would expect to go in a tree of n samples if the points were uniformly distributed (i.e., no anomalies).

$$\frac{E[h(x)]}{c(n)} \rightarrow \text{not standn.} \quad Z^{-1}$$

avg path len = 0

$$E[h(x_i)] \approx \log \frac{x=0}{y=1} \quad x=0 \quad y=1 \quad \boxed{x=1} \quad \boxed{y=0}$$

comp → stand $[0-1]$

$[0-1]$ anomaly → $1 \rightarrow 0$ outlier

% → top 10 point $0.5 \rightarrow$ inlier/outlier

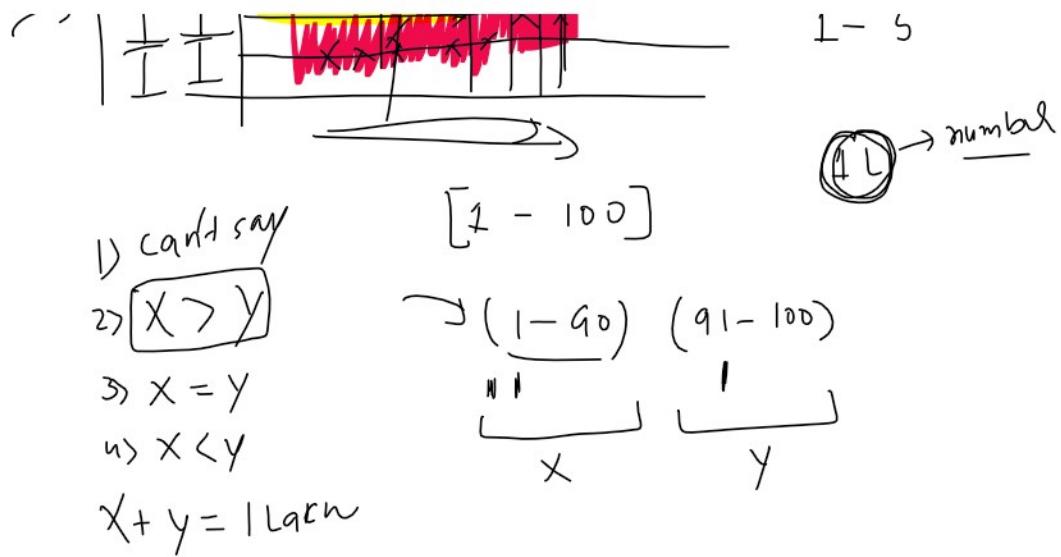
$0 \rightarrow$ inlier

contaminant → $\frac{10}{10} \rightarrow 5$

7 splits



$$3 - 72 \\ 2 - 92 \\ 1 - 5$$

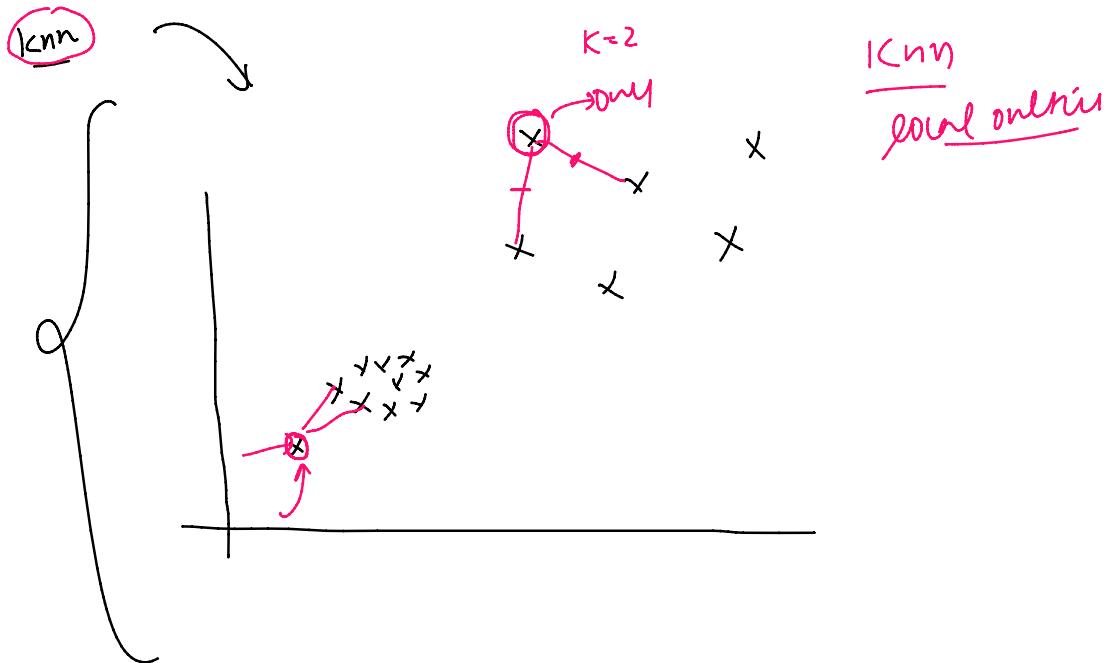
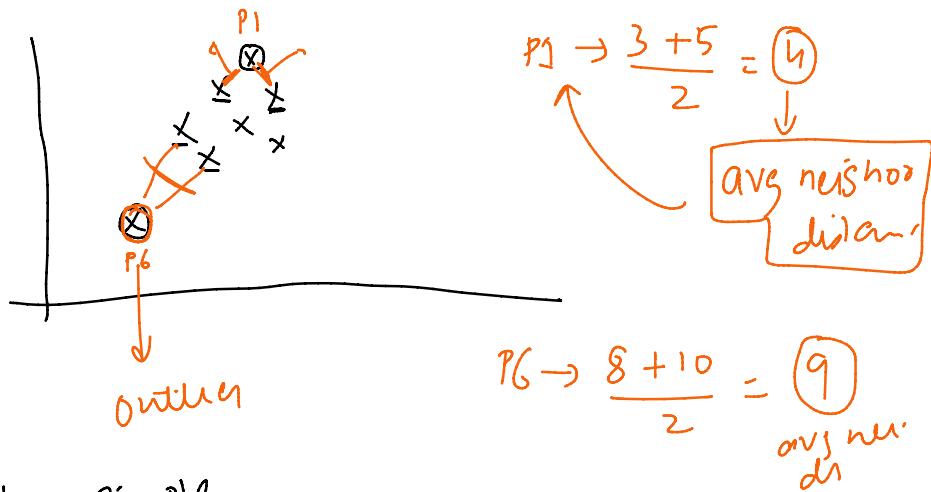


Outlier Detection using KNN

12 March 2024 08:24

K=2

Disadvantages

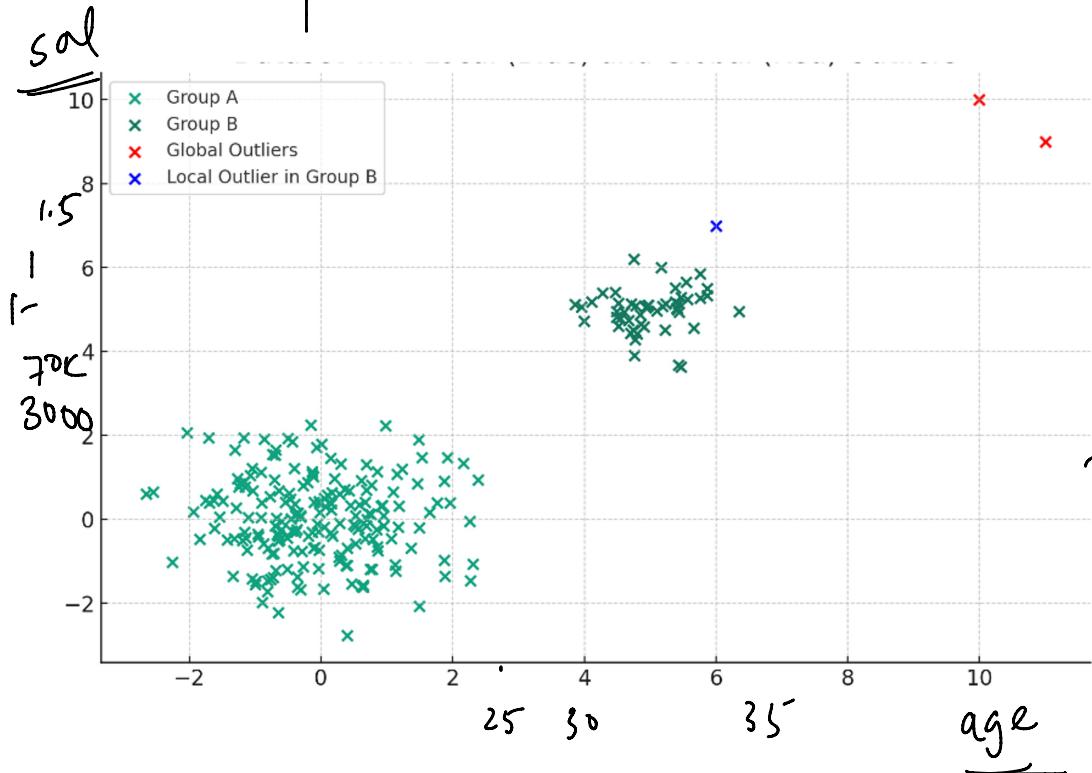


LOF

↳ Local outlier → Knn

Local vs Global Outliers

12 March 2024 08:35



global
local

local
outlier

knn → good algo

→ students

promise
↓

1.5 L

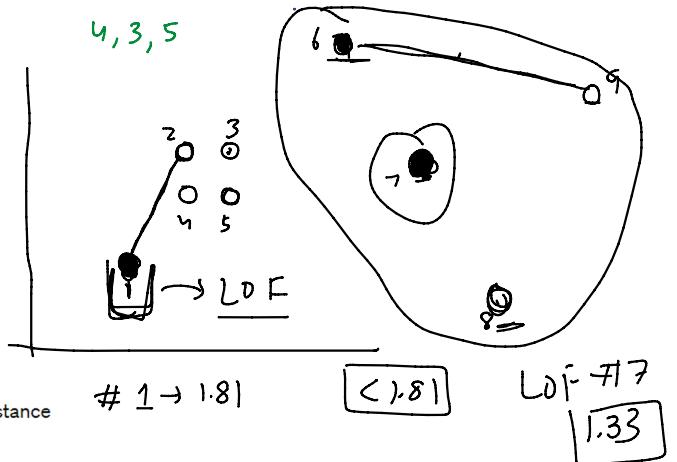
Local Outlier Factor

07 March 2024 09:42

$$\begin{aligned} 7 &\rightarrow \text{LRD} \rightarrow 0.06 \\ 6 &\rightarrow \text{LRD} \rightarrow 0.075 \\ 8 &\rightarrow \text{LRD} \rightarrow 0.075 \\ 9 &\rightarrow \text{LRD} \rightarrow 0.1 \end{aligned}$$

Step 1: Calculate k-distance for each point

- For every point p in the dataset:
 - Calculate the distance from p to every other point.
 - Determine the k -th smallest distance; this is the k -distance for point p .



$$\frac{0.075 + 0.075 + 0.1}{3}$$

Step 2: Compute reachability distance

- For every point p and each of its neighbors o :
 - The reachability distance from p to o is the maximum of the k -distance of o and the distance between p and o :

$$\text{reachability distance}_k(p, o) = \max(\text{k-distance}(o), d(p, o))$$

Step 3: Calculate local reachability density (LRD)

- For each point p :
 - Compute the LRD of p as the inverse of the average reachability distance from p to its k -nearest neighbors:

$$\text{LRD}_k(p) = \frac{1}{\sum_{o \in k\text{-neighbors}(p)} \text{reachability distance}_k(p, o) / |k\text{-neighbors}(p)|}$$

Step 4: Compute the Local Outlier Factor (LOF)

- For each point p :
 - Calculate the LOF by averaging the ratios of the LRD of p 's neighbors to the LRD of p :

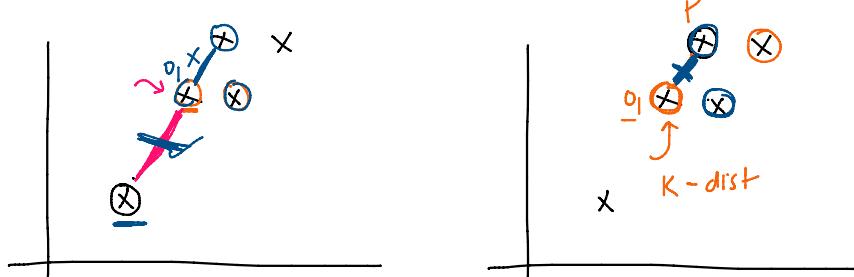
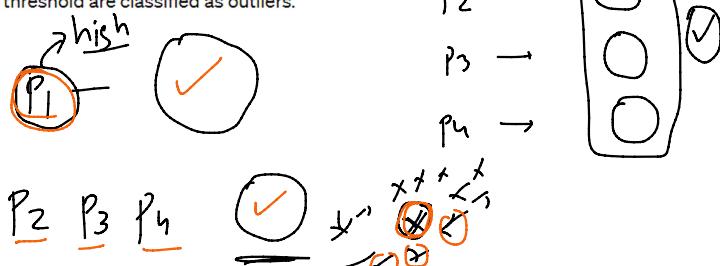
$$\text{LOF}_k(p) = \left(\frac{\sum_{o \in k\text{-neighbors}(p)} \text{LRD}_k(o)}{|k\text{-neighbors}(p)|} \right)$$

Step 5: Interpret the LOF scores

- Points with an LOF score around 1 are similar in density to their neighbors and likely not outliers.
- Points with higher LOF scores have lower local density compared to their neighbors and are considered outliers.

Step 6: Identify outliers

- Choose a threshold LOF score. Points exceeding this threshold are classified as outliers.



$$\text{reach} = \max$$

$$\text{reach} =$$

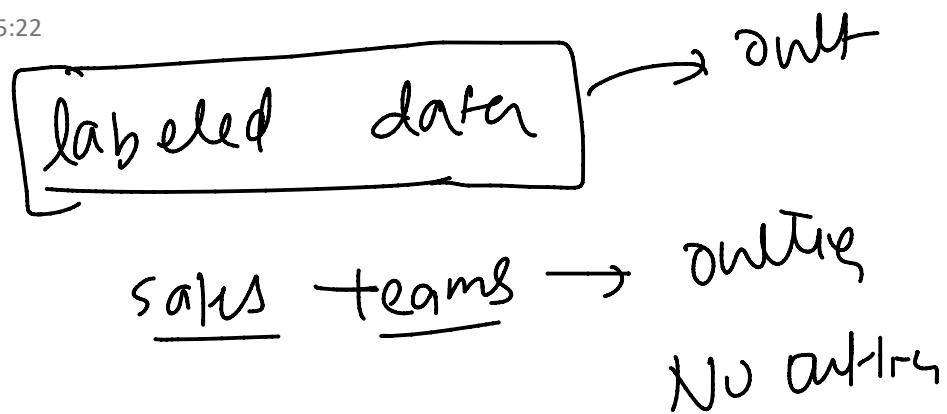
$$\text{reach} = \max$$
$$\text{reach}(p, \sigma) = \max(\text{cd}(\sigma), \underline{\text{act}}(p, \sigma))$$

DBSCAN

07 March 2024 09:42

How to access the accuracy?

12 March 2024 15:22



When to use which algo?

12 March 2024 16:41

Col → outliers

1. Z-Score Technique:

- Use when your data is approximately normally distributed.
- Ideal for univariate data analysis where you want to find outliers based on a single feature.
- Suitable for quick, preliminary outlier detection to identify extreme values.

2. IQR and Box Plot Method:

- Use for univariate data that may not be normally distributed or when the data is skewed or has outliers.
- Effective for exploratory data analysis to visually identify outliers.
- Suitable when you need a robust method that is less affected by extreme values.

3. Isolation Forest: → high dim

- Use for high-dimensional datasets where other algorithms might struggle due to the curse of dimensionality.
- Ideal when you have a mix of normal and anomalous instances, and the anomalies are 'few and different', which makes them more susceptible to isolation.
- Suitable for situations where you do not have a clear definition of what constitutes an outlier but want to isolate points that appear to be different.

4. KNN (k-Nearest Neighbors):

- X [
- Use for datasets where the notion of density and local neighborhoods is meaningful.
 - Ideal for moderate-sized datasets; computational complexity can be a concern for very large datasets.
 - Effective when outliers are expected to be points that deviate significantly from their local neighborhoods.

5. LOF (Local Outlier Factor):

- Use for datasets where you want to detect outliers that are not just globally anomalous but also locally anomalous in varying density regions.
- Ideal when it's important to consider the relative density of a point's neighborhood, not just the absolute distance to the nearest neighbors.
- Suitable for datasets where clusters have different densities, and you want to identify outliers relative to these local densities.

6. DBSCAN (Density-Based Spatial Clustering of Applications with Noise):

- Use for spatial data or when you want to identify outliers as part of a clustering process.
- Ideal when you have data that forms natural clusters of varying shapes and sizes, and you want to identify points that do not belong to these clusters.
- Suitable for datasets where you have some intuition about the density or the distance that defines clusters and noise (outliers).

