

## \* Iterative Imputation (MICE)

Used when data is missing at random (MAR)

Advantage

↳ quite accurate.

Disadvantage

↳ slow

↳ memory problem

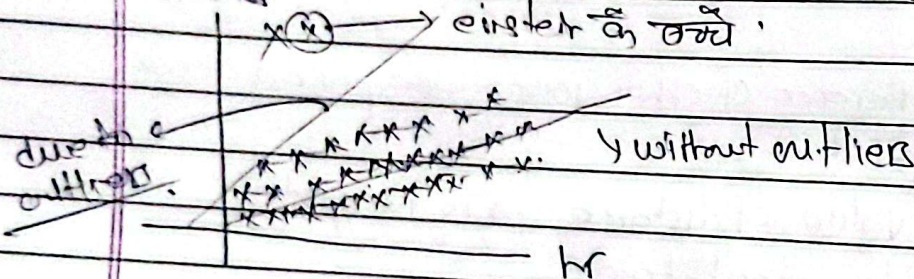
Main Motive

↳ predicting missing value using MICE

## \* Outliers

What are outliers?

यहाँ जिका बैठा - 99.1% exam में but still crying.  
marks



When is outlier dangerous?

eg.

Age  $\rightarrow 300$

But,

outliers are helpful in anomaly detection.  
like credit card anomaly detection.  
↳ should be used in anomaly detection

outlier detection  $\rightarrow$  simple

But

what to do with outliers is very tedious job



## Effect of outliers in ML algorithms?

- ↳ Linear Regression
  - ↳ Logistic Regression
  - ↳ Adaboost
  - ↳ Deep learning.
- } → weight based algorithm

But don't effect tree based algorithm

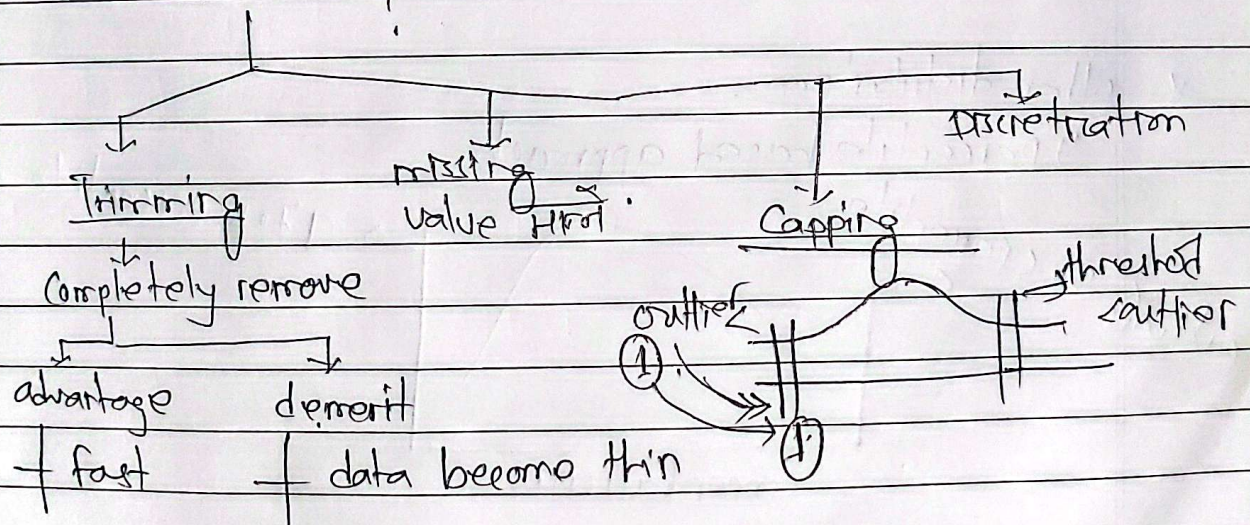
like

decision tree

Random forest

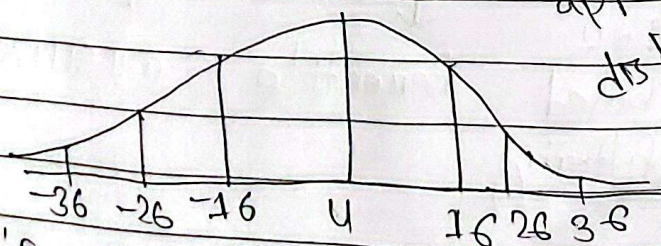
Xgboost algorithm.

## How to treat outliers?



## \* How to detect outliers?

### 1. Normal distribution



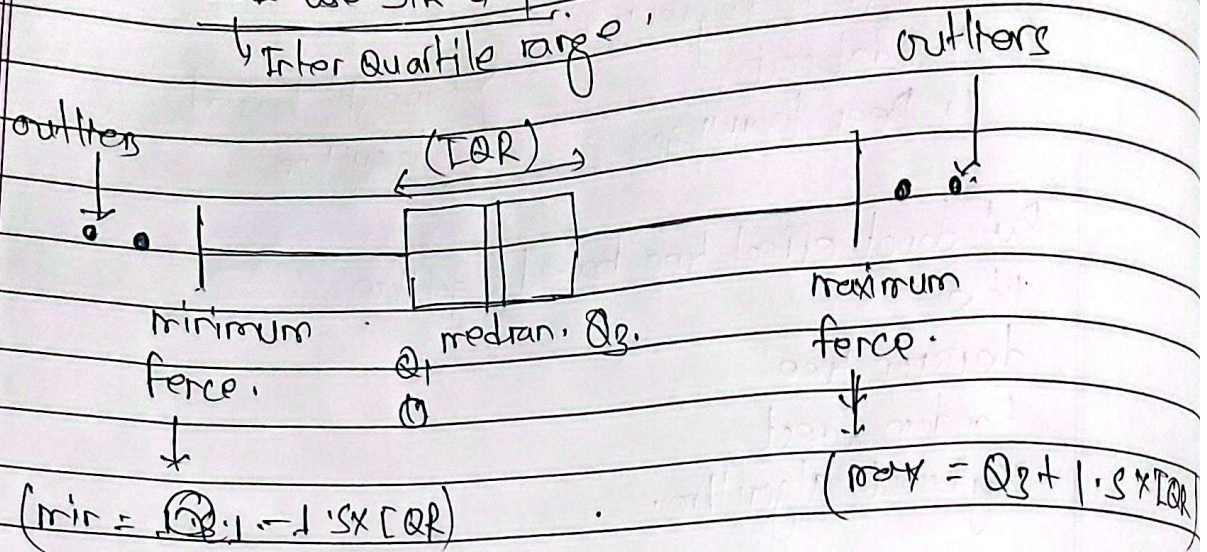
→ but only applied on normally distributed data

if  $-3(\mu - 3\sigma) < \text{Observation} < (\mu + 3\sigma) \Rightarrow \text{no outlier}$   
else outlier.



## 2) Skewed Distribution

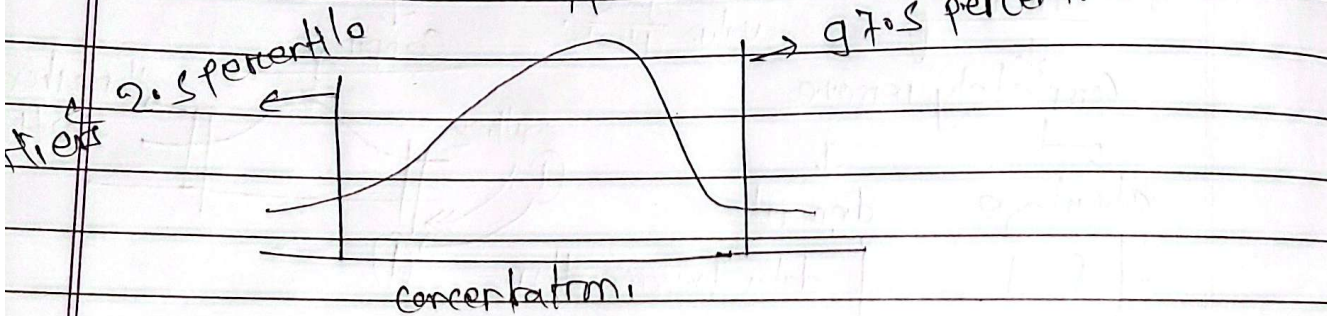
↳ IQR use sikre  
↳ Inter Quartile range



only used when observation is skewed

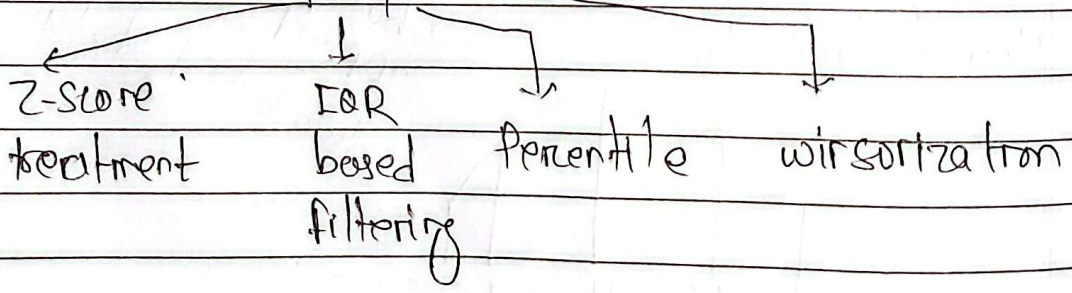
## \* other distribution

↳ percentile based approach



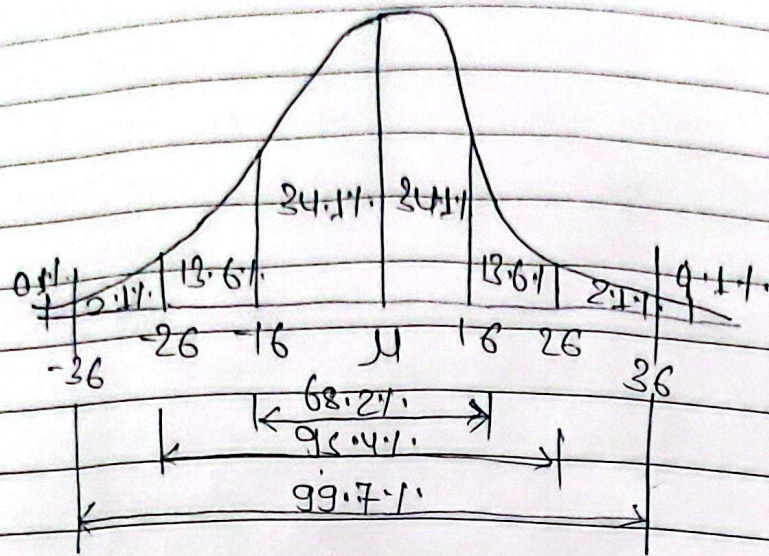
So,

Techniques





# outlier removal using z-score.



$$\left. \begin{matrix} \mu + 6 \\ \mu - 6 \end{matrix} \right\} \rightarrow 68.1\% \text{ data.}$$

$$\left. \begin{matrix} \mu + 26 \\ \mu - 26 \end{matrix} \right\} 95.1\% \text{ data.}$$

So, outliers are value out of  $\mu - 36$  &  $\mu + 36$ , or,  
 $[-3 < \text{Z-score} < 3] \rightarrow$  range for outlier detection.

Z-Score?

$$Z_i = \frac{x_i - \mu}{s} \quad \text{for detection}$$

outliers Treatment

Trimming  
 Ignoring

Capping

significant data loss,

\*capping

85, 290  
 ↓ ↓ ↓  
 80 5 80

→ should be normally distributed or close to normally normal distribution