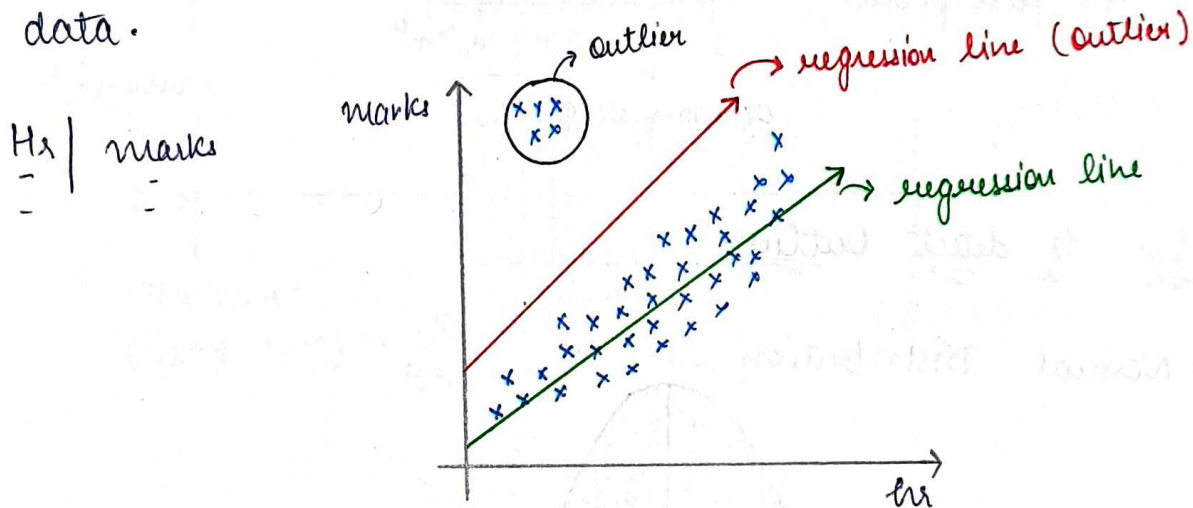# Outliers

## What are Outliers?

An outlier is a data point that is significantly different from the rest of the data. Outliers can be much higher or lower than the other data points. They can be caused by measurement or execution errors, bad data collection, or simply show variables not considered when collecting the data.

$\frac{H_s}{-} \Big| \frac{marks}{-}$



## When is outlier dangerous?

Age

300 →    Age 300 not possible

## Effect of Outliers on ML algorithms

Linear Regression          Ada boost
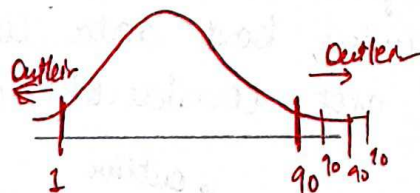
Logistic Regression        Deep learning

Weights algo.

# How to treat Outliers?

## Trimming
→ Data thin
→ fast process

## Capping



Outlier ←     Outlier →
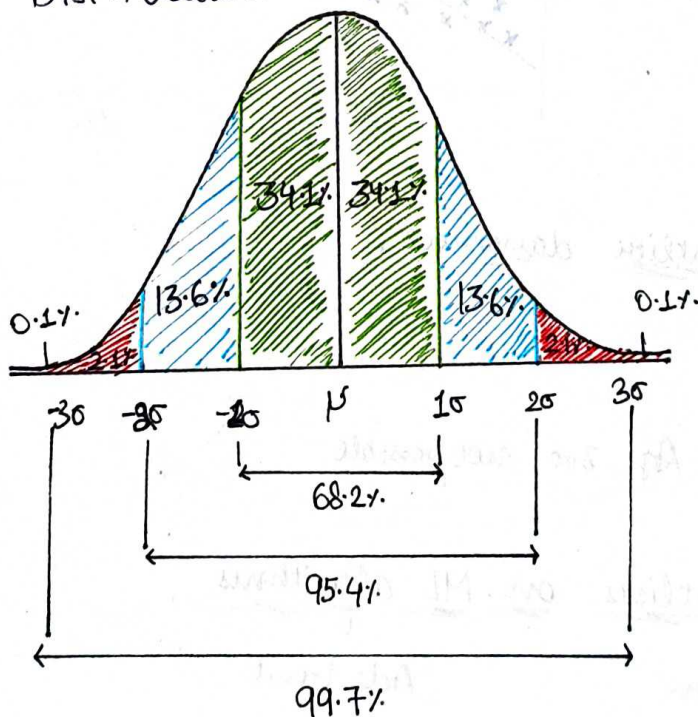
1        90 70 90 70

often 90 → all 90

→ missing value
→ Discritization
90-10
Numerical.

# How to detect Outliers?

## 1. Normal Distribution



0.1%    13.6%    34.1%    34.1%    13.6%    0.1%

$3\sigma$   $-2\sigma$   $-1\sigma$   $\mu$   $1\sigma$   $2\sigma$   $3\sigma$

68.2%

95.4%

99.7%

$(\mu + 3\sigma) >$

$(\mu - 3\sigma) <$

## 2. Skewed distribution



→ right skewed

→ left skewed

Right Skewed
of
Left Skewed Data

Interquartile Range
(CIQR)

Outlier

Outlier

"Minimum"
$(Q_1 - 1.5 * IQR)$

Median

$Q_1$
$(25\%)$

$Q_3$
$(75\%)$

$(Q_3 + 1.5 IQR)$

## 3. Other Distribution

$(2.5\%)$

$(97.5\%)$

Percentile    1    5  10 15 20 25 30 35 . . . . .    95    99
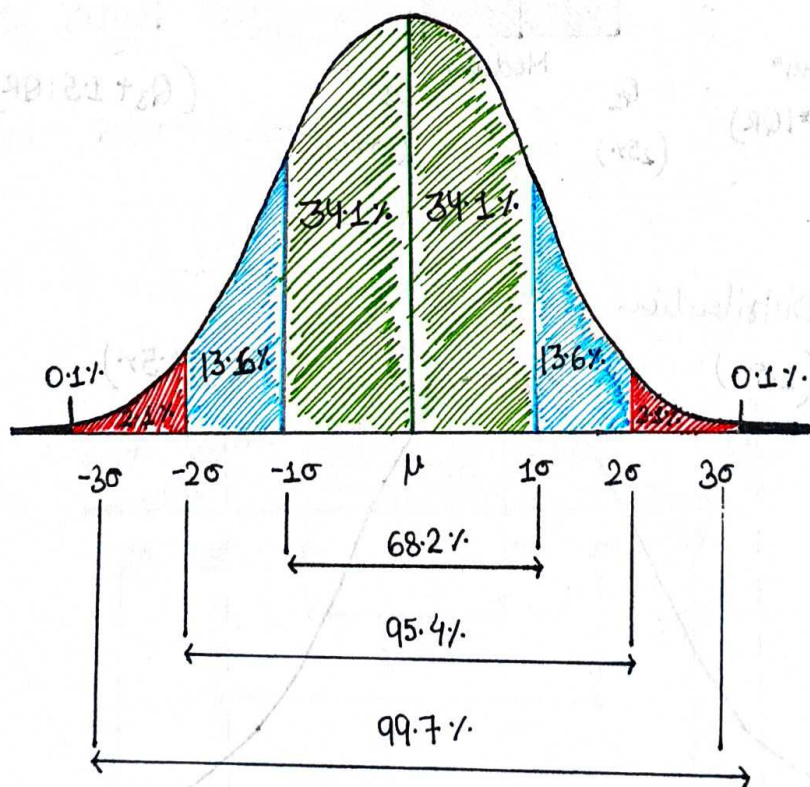
# Techniques for outlier Detection and Removal

1. Z-score

2. IQR

3. Percentile

4. Winsorization

## Outlier removal using Z score



$$\left\{ \begin{array}{c} \mu + \sigma \\ \mu - \sigma \end{array} \right\} 68\%$$

$$\left\{ \begin{array}{c} \mu + 2\sigma \\ \mu - 2\sigma \end{array} \right\} 95\%$$

Z score $\longrightarrow$ Age $\uparrow$    $\boxed{-3 \text{ to } 3}$

$$x' = \frac{x_i - \mu}{\sigma}$$

27
32

# Outlier Treatment

1. Trimming      2. Capping →      3 outlier

↓

1000 dataset

5 row → outlier
↳ Trim from original data

995 → data

$\mu + 3\sigma = 80$
↓
upper limit

$\mu - 3\sigma = 5$
↳
lower limit

Outlier → 85, 3, 90
Caping → 80, 5, 80

Steps:- Sure Data should be Normally distributed
Step2:- find upper limit and lower lower
Step3:- Decide → Trimming
↳ capping

## Percentile

max - 95 → 100%
No age more than 95

min - 10 → 0%
No age less than 10

50% → median

Age

1% → 99%

### Percentile

remove      Capping
                ↳ Winsorization