# five number summary

1) minimum

2) first Quartile $(Q_1)$ 25%

3) median $(Q_2)$ 50%

4) Third Quartile $(Q_3)$ 75%

5) maximum.

Note: Choose these 5 numbers after removing the outlier from the data by finding boundary values.

[ Lower fence      upper fence ].

$LF = Q_1 - 1.5 \ (IQR)$

$UF = Q_3 + 1.5 \ (IQR)$

$IQR =$ (inter Quartile range)

$IQR = Q_3 - Q_1$

$\{ 1, 1, 2, 3, 4, 4, 4, 5, 5, 6, 7, 7, 8, 8, 9 \} \{ 28, 36 \}$ = 17 No.

outlier

$LF \quad \underline{\quad\quad} \quad Q_1$

$Q_1 = \dfrac{25}{100} \times 18$         $IQR = Q_3 - Q_1$

$= 4.5$ index               $= 8 - 3 = 5$

$Q_1 = 3$

$Q_3 = \dfrac{75}{100} \times 18$

$= 13.5$ index

$Q_3 = 8$

$LF = Q_1 - 1.5 (IQR)$         $UF = 8 + 1.5 (5)$

$= 3 - 1.5 (5)$                 $= 15.5$
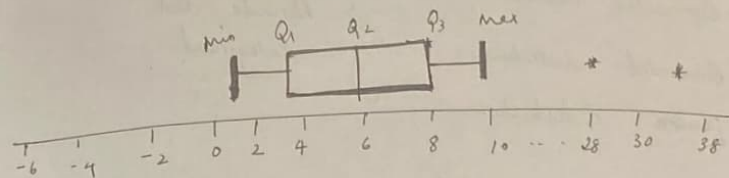
$= 4.5$

y < 4.5 and > 15.5 is outlier

---

minimum = 1

$Q_1 = 3$

$Q_3 = 8$

median = 5

max = 9

## Boxplot



The graph is used to find the outlier

## Different types Distribution

- To understand data patterns

- To summarize the data easily.

- To calculate probability

- To make prediction and decision
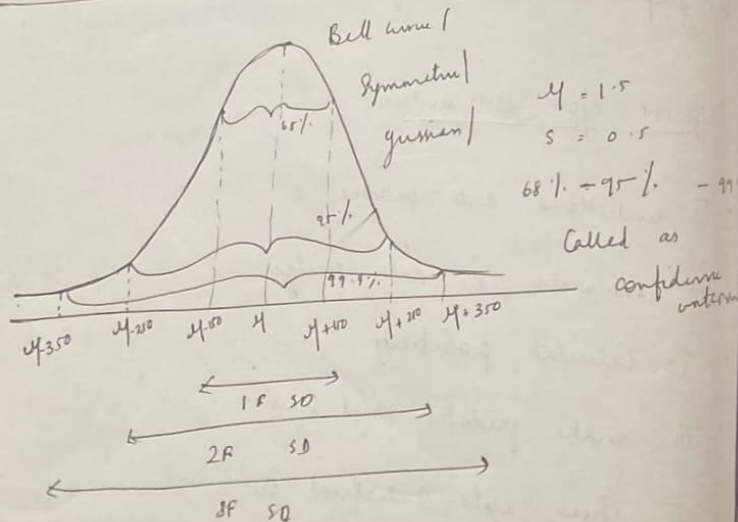
- To choose right statistical out

There are 2 category of

1) Continuous distribution ( Number)
   ↳ any value including decimals (meaning)
     eg. Weight = 52.8 kg, time = 4.36 sec
2) Discrete distribution ( Categorical distribution)
   ↳ specific values - you count them
     eg. No. of students in class = 20, 21, 22 (not 20.5)

1) Normal distribution          } continuous
                                    dist
2) Standard normal distribution }   numerical

3) Bernoulli distribution       } Discrete dist
4) Binomial distribution        }   Categorical.
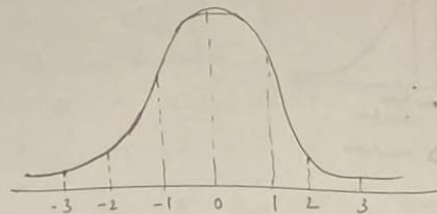5) poison distribution

## 1) Normal distribution



Bell curve
Symmetrical
gaussian

$\mu = 1.5$
$S = 0.5$
$68\% - 95\% - 99$

Called as confidence interval

1F SD
2F SD
3F SD

## Empirical rules

- 68% of data will be present in 1SD
- 95% of data will be present in 2SD
- 99.7% of data will be present in 3SD

## 2) Standard Normal distribution.

$\mu = 0$ } (always)
$SD = 1$ }



$$Z \text{ score} = \frac{xi - \mu}{\sigma} ; \frac{2 - 3.86}{2}$$

| Normal Dist data | Standard Normal dist data |
|---|---|
| 2 | -0.93 |
| 7 | 1.57 |
| 5 | 0.57 |
| 4 | 0.07 |
| 1 | -1.43 |
| 3 | -0.43 |
| 5 | 0.57 |

$\mu = 0$
$\sigma = 0.94$

## Symmetric



mode , mean , median
Symmetric



mean

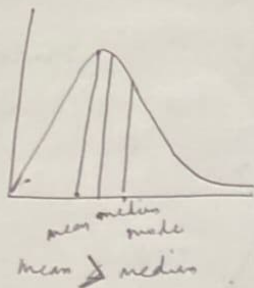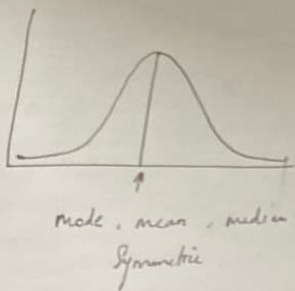median  mode
Skewed
left

mean < median



mean  median
mode

mean > median

## Positive Skew (Right Skew).

Tail on the right side is longer most data are the left

## Negative Skew (Left Skew)

Tail on the left side is longer most data are on the right.

## Zero Skew (Symmetric)

The data is evenly distributed around the mean (like a normal distribution)

$$Skewness = \frac{3(mean - median)}{Standard\ deviation}$$

If Value near to -1 then it is -Ve Skew

If Value near to +1 then it is +Ve Skew

If Value lies in -0.5 to 0.5 then it is Zero Ske

## Kurtosis

$$k = \frac{1}{n} \sum_{i=1}^{n} \left( \frac{x_i - \mu}{\sigma} \right)$$

It measures the tailedness or peakness of distribution

Types of kurtosis

1) Mesokurtic $(k=3)$
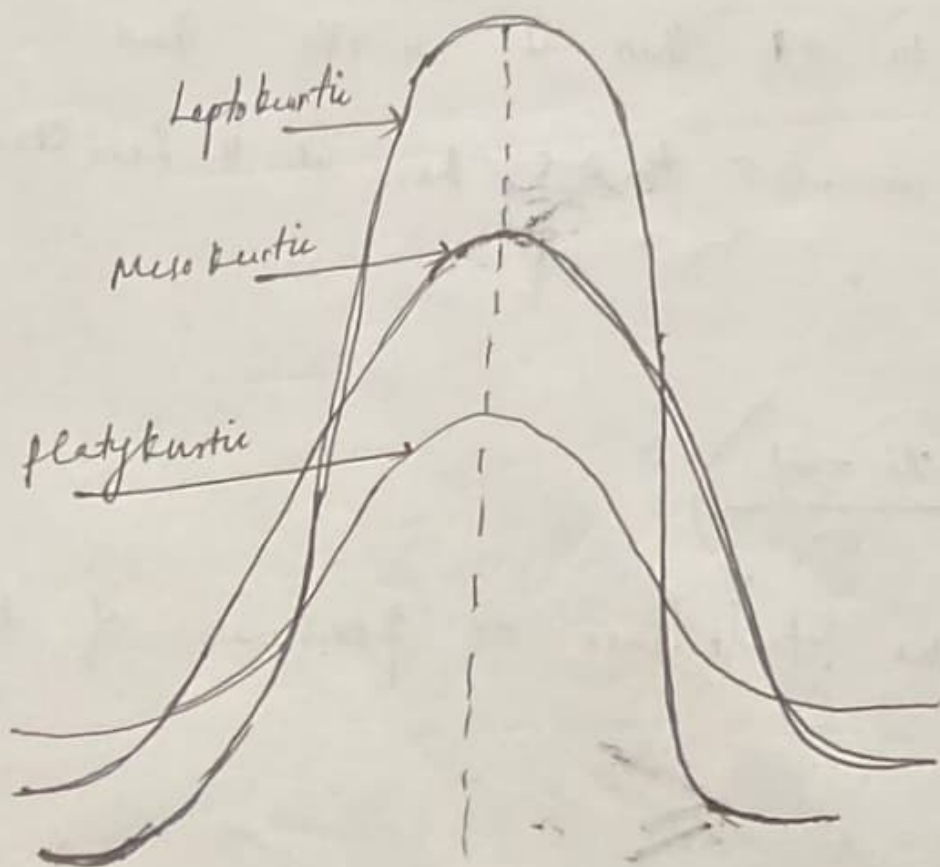 - Normal distribution
 - No outliers
 - Moderate tail and peak

2) Leptokurtic $(k>3)$
 - Heavy tails and sharp peak
 - More outliers

$(x=k) = c(n,k) \, p^k (1-p)^{n-k}$

3) platy kurtic $(k<3)$
 - light tail and flat peak
 - fewer outliers

# 5) poison distribution

It is used to model the number of event tha[t]
occur in a fixed time interval or space and
occur independently the parameter $\lambda$ represents th[e]
avg number of event in the interval

$$\boxed{P(x=k) = \dfrac{\lambda^k \, e^{-\lambda}}{k!}}$$

how many times something
[beca]... in a fixed time or area

$\lambda = u$ the avg no of ev[ent]

$k = $ no of argumen[t]