

## Statistics | Basics

Mathematical Science including methods of ~~the handling~~ collecting, organising & analysing data in such a way that meaningful conclusions can be drawn from it.

Data  $\Rightarrow$  facts or pieces of info. that can be stored, measured & re-accessed.

ultimate (in case of business)

goal To make revenue to grow, from data insights.

Collecting data  $\Rightarrow$  manual, survey, internet clicks, ratings.

Organising data  $\Rightarrow$  SQL & NoSQL DB.

Analysing data  $\Rightarrow$  Python libraries (Pandas, NumPy, matplotlib, pyplot, bokeh).

Ex., for data

- 1) Heights of students in class 12 'A' section.
- 2) No. of units a product has sold ~~for~~ in each month over the year.
- 3) No. of products units produced by the company over several years.

Types of Statistics :-

1) Descriptive statistics

$\hookrightarrow$  It consists of organising & summarising the complete data (Population).

Ex., Avg. weight of students in classroom, strike rate of N.S. Dhoni.

## 2) Inferential statistics

You can't count all the trees present in all of Indian forests or calc. avg. weight of people by taking avg. of all 150cr+ weights. It uses sample data measured to form conclusion of population.

Sample is a subset of population.

Ex., calculating <sup>avg.</sup> weight of entire pop. of a country.

### Techniques of Descriptive statistics

- 1) Central Tendency (Mean, Median, Mode).
- 2) Symmetry (Skewness, Kurtosis).
- 3) Dispersion (Std. Deviation, Variance)

### Types of Sampling (Inferential statistics)

- 1) Simple Random Sampling
  - 2) Stratified Sampling
  - 3) Cluster Sampling
  - 4) Systematic Sample
- Most commonly used
- used mostly in research studies

## Sampling Types

### INFERENTIAL

### STATISTICS

#### 1) Simple Random Sampling :-

→ Every member of population ( $N$ ) has equal chances of getting selected.

Ex:- If  $N = 1000$ , the prob. of each every member getting selected is  $1/1000$ .

Disadvantage: If a sample is collected from all the states ~~of India~~ of India, the prob. of people getting selected from less populated states is less to none.

So, Not Equal Representation.

#### 2) Stratified Sampling :-



Strata → layers/groups

Diff. distinct categories are separated.

A Simple Random sample would be chose from each category (strata / layer / group).

Adv.: Equal Representation.

✓	○	✓	○	○	○	○	○
○	✓	○	○	○	○	○	○
○	○	○	○	○	○	○	○

○	○	○	○
○	○	○	○

#### 3) Cluster Sampling :-

clustering the population on grouping the population into clusters & then the sample is obtained / formed from elements of each cluster randomly.

Ex:- States with similar population are made into one cluster & sample is formed by random selection <sup>of states</sup> from the cluster.

#### 4) Systematic Sampling

→ Every  $n^{\text{th}}$  element will be selected.

Ex., selection based on

- odd roll no. (3, 5, 7, ...)

- odd year of birth (2001, 2003, 2005, ...)

### Types of Data

Quantitative

(Numerical) (Measurable)

↳ Discrete

↳ Continuous

Discrete (whole no.)

No. of children playing  
in park. Ex.,

Continuous (range)

(Real no.)

distance b/w. two  
places. Ex.,

Qualitative

(Categorical)

Nominal

Ordinal

Nominal

(No order is  
there, naturally)

Ex., gender,  
blood group.

Ordinal

(Order is there)

Ex., grades,  
ranks.

Sometimes, conversion from categorical  
data to Numerical data.

# Scales of Measurement

- i) Nominal scale } Qualitative
- ii) Ordinal scale }
- iii) Interval scale } Quantitative
- iv) Ratio scale }

## i) Nominal Scale Data

→ Basic Counting using numbering or percentage.

→ Ex., No. of females in Room A - 350. Pie,

No. of males in Room A - 209. Bar Plot

## ii) Ordinal Scale Data

→ Order & rank has a meaning.

→ Difference can't be measured. ( $1^{\text{st}}$  rank -  $2^{\text{nd}}$  rank  $\Rightarrow$  not meaningful)

→ Pie chart, Bar plot.

## iii) Interval Scale Data

→ Order & rank has meaning

→ Diff. can be measured (Excluding Ratios)

→ "0" need not be the starting value

Ex., Length, weight, height.

Histogram, Scatter plot, Line graph.

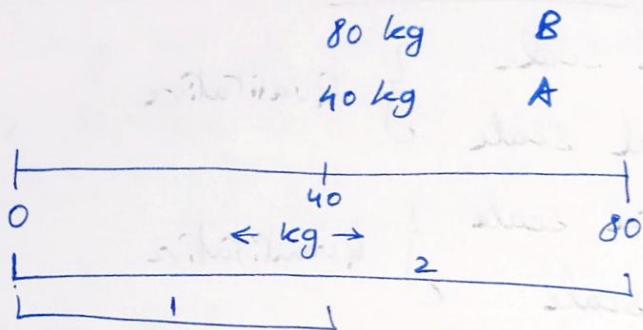
## iv) Ratio Scale Data

→ Order & rank has meaning.

→ Differences & ratio are measurable.

→ Starting pt. is zero.

Ex., Pehap. weight



$$\frac{B}{A} = \frac{80}{40} = \frac{2}{1}$$

$$B:A = 2:1$$

"0" is present

B is twice as heavy as A.

DATA	Nominal	ordinal	Interval	Ratio
Labelled	✓	✓	✓	✓
Meaningful order	X	✓	✓	✓
Measurable difference	X	X	✓	✓
Zero as starting point	X	X	X	✓
Example	Gender, Religion, Location	Rating, Grade, Rank	IQ, Temp, height, score, weight	height, weight, length.

## I) Measures of central Tendency

DESCRIPTIVE  
STATISTICS

"central"  $\Rightarrow$  what is the one value around which all data is revolving?

Mean,  
Median, mode }  $\rightarrow$  EDA, Data preparation,  
  feature engineering.

### ① Mean (Avg. or Arithmetic Mid Value)

$$\mu = \frac{\sum_{i=1}^n x_i}{N}$$

$\sum \Rightarrow$  summation

Summing up all observations.

No. of observations.

### ② Median (Physical Mid. point of data)

$\rightarrow$  sort the data points

$\rightarrow$  Divide the No. of data points by 2

$\rightarrow$  The output no. is the position of mid-value.

The above is ok for odd no. of values.

Even no. of observations : Case

$$P = \{1, 2, 3, 4\}$$

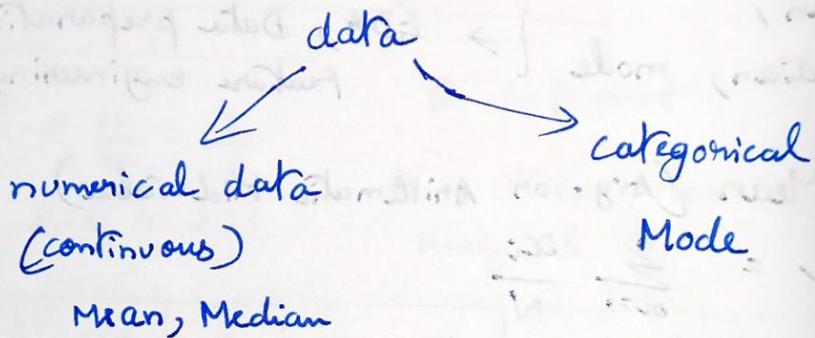
$$\frac{2+3}{2} = 2.5 \text{ (median)}.$$

Median is not affected by outliers whereas Mean is affected by outliers.

③ Mode (Highest occurring data point)

1, 2, 9, 10, 10, 8, 7, 7, 8, 8

Here  $\Rightarrow$  8 is the mode.



USE CASE OF CENTRAL TENDENCY

(Missing Value Imputation)

<u>AGE</u>	<u>GENDER</u>	<u>WEIGHT(kg)</u>
25	M	80
26	M	70
<u>24.75</u>	M	30
23	M	<u>75</u>
25	F	1000

Age is continuous variable

$\hookrightarrow$  Impute / Replace the missing / null value  
with mean if there is no outlier

$$\frac{25 + 26 + 23 + 25}{4} = 24.75$$

GENDER is categorical variable

↪ Impute with Mode

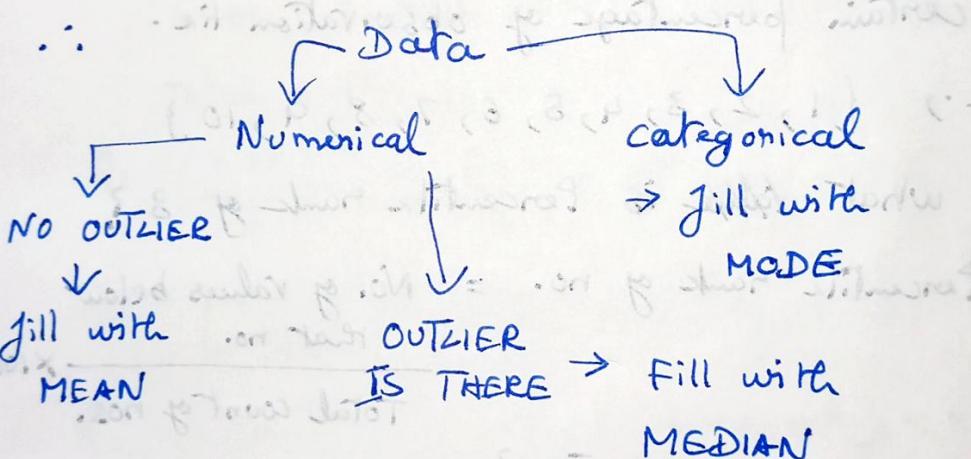
(M, M, M), F  
Mode = M

WEIGHT is continuous variable

↪ Impute with Median as outlier (1000) is present

30, 70, 80, 1000

$$\frac{70+80}{2} = 75$$



## II) Measures of Dispersion

How the data is spread out?

I) Range → Diff. between max. & min. value.

Ex: {1, 2, 3, 4}

$$\text{Range} = 4 - 1 = 3$$

Ex: {1, 2, 3, 1000}

$$\text{Range} = 1000 - 1 = 999 \quad \Delta$$

★ OUTLIER AFFECTS RANGE ★

## 2) Percentage

Ex., 1, 2, 3, 4, 5

what is % of no. of odd numbers?

$$\frac{3}{5} \times 100 = 60\%$$

60% of the numbers in the given data is odd.

## 3) Percentile

A percentile is a value below which a certain percentage of observation lie.

Ex., {1, 2, 3, 4, 5, 6, 7, 8, 9, 10}

i) what is Percentile rank of 3?

Percentile rank of no. =  $\frac{\text{No. of values below that no.}}{\text{Total count of nos.}} \times 100$

$$= \frac{2}{10} \times 100$$

= 20th Percentile.

ii) what value exists at 75th Percentile?

$$\text{Value} = \frac{\text{Percentile}}{100} \times (n+1)$$

$$= \frac{75}{100} \times (10+1)$$

$$= \frac{3}{4} \times 11 = 8.25$$

#### 4) Quantiles

Quantiles are values that divides a list of nos. into quarters (4 divisions).

Ex.1.  $\{5, 9, 13, 14, 21, 22, 29\}$

$$Q_1 = \frac{9}{4} = 2.25 \\ Q_2 = \frac{14}{4} = 3.5 \\ Q_3 = \frac{22}{4} = 5.5$$

Ex.2.  $\{1, 1, 1, 1, 2, 2, 2, 3, 3, 3, 4\}$

Case 1 : odd count of nos.

$$\text{Total} = 11$$

$$Q_1 = \left(\frac{n+1}{4}\right)^{\text{th}} = \frac{11+1}{4} = 3^{\text{rd}} \text{ no.} = 1$$

$$Q_2 = \left(\frac{n+1}{2}\right)^{\text{th}} = \frac{11+1}{2} = 6^{\text{th}} \text{ no.} = 2$$

$$Q_3 = \left(\frac{3(n+1)}{4}\right)^{\text{th}} = \frac{3 \times (11+1)}{4} = 9^{\text{th}} \text{ no.} = 3$$

Ex.3.  $\{25, 103, 921, 2517, 319, 31\}$

case 2 : even count of nos.

$$\text{Total} = 6$$

$$Q_1 = \left(\frac{n}{4}\right)^{\text{th}} = \frac{6}{4} = 1.5^{\text{th}} = \frac{25+103}{2} = 64$$

$$Q_2 = \frac{\left(\frac{n}{2}\right)^{\text{th}} + \left(\frac{n}{2}+1\right)^{\text{th}}}{2} = \frac{\frac{6}{2} + \left(\frac{6}{2}+1\right)}{2} = \frac{3^{\text{rd}} + 4^{\text{th}}}{2} \\ = \frac{921 + 2517}{2} = 1719: Q_2$$

$$Q_3 = \left(\frac{3n}{4}\right)^{\text{th}} = \frac{3 \times 6}{4} = \frac{18}{4} = 4.5^{\text{th}}$$

avg. of 4<sup>th</sup> & 5<sup>th</sup>  
numbers  
~~=  $\frac{3+4}{2}$~~  ~~=  $\frac{7}{2}$~~

$$= \frac{2517 + 319}{2}$$

$$= 1418$$

### Five pt. Summary

$Q_0 \Rightarrow \text{min. no.} \Rightarrow 1 \Rightarrow 0^{\text{th}} \text{ percentile}$

$Q_1 \Rightarrow 25^{\text{th}} \text{ p. ile} \Rightarrow 2.5$

$Q_2 \Rightarrow \text{median} \Rightarrow 50^{\text{th}} \text{ p. ile} \Rightarrow 5.5$

$Q_3 \Rightarrow 75^{\text{th}} \text{ p. ile} \Rightarrow 7.5$

$Q_4 \Rightarrow \text{max. no.} \Rightarrow 10 \Rightarrow 100^{\text{th}} \text{ p. ile}$

Ex:-  
data, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10  
 ↓      ↓      ↓      ↓  
 $Q_0 \quad Q_1 \quad Q_2 \quad Q_3$

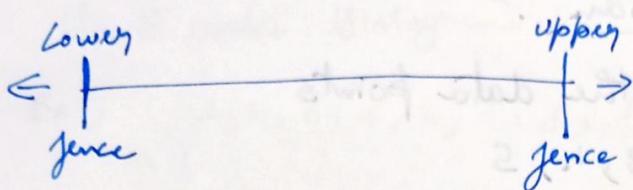
### Inter Quartile Range

$$\text{IQR} = Q_3 - Q_1$$

Not affected by range & outliers

Main usage of IQR is to find out outliers.

Outlier is extreme values



$$\text{Lower fence} = Q_1 - 1.5 \times \text{IQR}$$

$$\text{Upper fence} = Q_3 + 1.5 \times \text{IQR}$$

If any value is present & is greater than the upper fence or lesser than the lower fence, it is an OUTLIER.

For the previous example,

$$\text{Lower fence} = 2.5 - 1.5 \times \text{IQR}$$

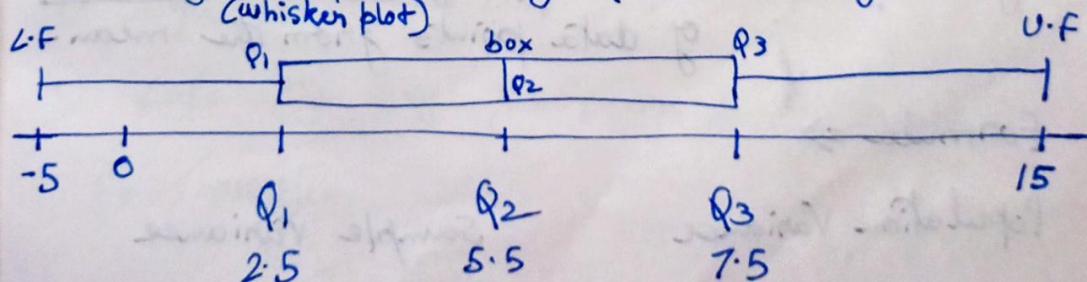
$$(\text{IQR} = Q_3 - Q_1 = 7.5 - 2.5 = 5)$$
$$= 2.5 - 1.5 \times 5$$

$$\text{Lower fence} = -5.$$

$$\text{Upper fence} = 7.5 + 1.5 \times 5$$
$$= 7.5 + 7.5$$

$$\text{Upper fence} = 15.$$

Drawing box plot using Quartiles & fences . . .



\* NOT TO SCALE

L.F  $\Rightarrow$   
Lower fence/whisker  
U.F  $\Rightarrow$   
Upper fence/whisker

## Main Concepts of Measures of Dispersion

### 6) Mean Deviation :-

Ex., Let's take the data points

$$1, 2, 3, 4, 5$$

Here, the mean is 3.

1 is 2 points away from 3

2 is 1 point away from 3

4 is 1 point away from 3

5 is 2 point away from 3

$$\therefore \frac{2+1+0+1+2}{5}$$

$$= 6/5 = 1.2 \Rightarrow \text{Mean Deviation}$$

$$\text{Formula} \Rightarrow \frac{\sum_{i=1}^n |x_i - \mu|}{N}$$

Definition  $\Rightarrow$  The average of difference of each data point from the mean.

### 7) Variance :-

Definition  $\Rightarrow$  The average of squared differences of data points from the mean.

Formula  $\Rightarrow$

Population Variance

$$\sigma^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{N}$$

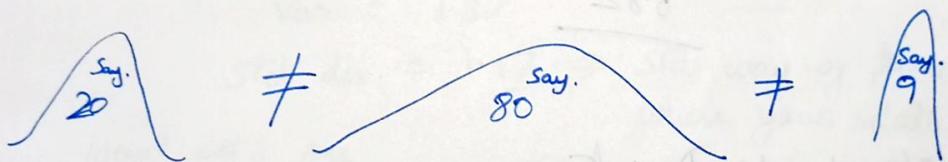
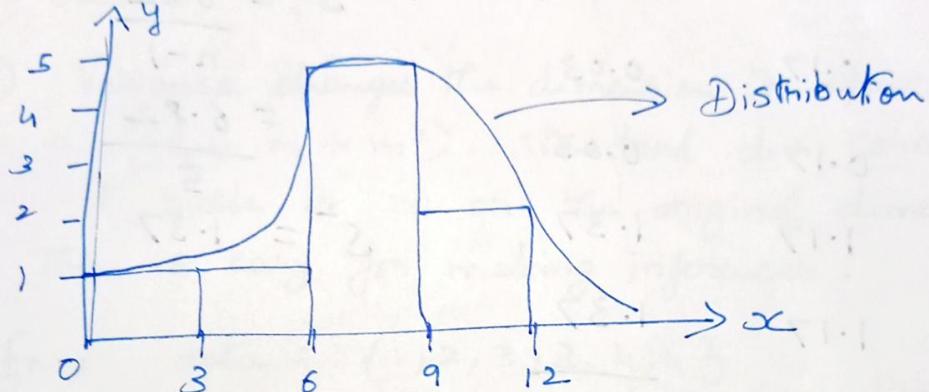
Sample Variance

$$S^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{N-1}$$

where,  $\mu \Rightarrow$  Population mean  
 $\bar{x} \Rightarrow$  Sample mean

How to make Histogram?

Ex: 2, 4, 6, 6, 6, 8, 8, 10, 11



Variance represents the spread.

Spread ↑, Variance ↑

How To Calc. Variance?

- 1) Calculate Mean
- 2) For each no. in data, Subtract the mean & the no.
- 3) Square the (difference)
- 4) Calc. the avg. of square of difference.

Ex: problem.

$$\text{data} = \{1, 2, 3, 3, 4, 4\} \quad \bar{x} = 17/6 = 2.83$$

$x$	$x - \bar{x}$	$(x - \bar{x})^2$	Sample Formula
1	-1.83	3.34	
2	-0.83	0.68	$S^2 = \frac{6.82}{n-1}$
3	0.17	0.03	
3	0.17	0.03	$= \frac{6.82}{5}$
4	1.17	1.37	$S^2 = 1.37$
4	1.17	1.37	
<u>17</u>		<u>6.82</u>	

### 8) Standard Deviation :-

S.D is a measure of how spread out numbers are.

↳ square root of Variance.

$$S = \sqrt{\text{Variance}}$$

From previous example,

$$S = \sqrt{1.37} = 1.17.$$

formula :

Population

$$\sigma = \sqrt{\text{Var}_p}$$

sample

$$S = \sqrt{\text{Var}_s}$$

why Std. Dev.?

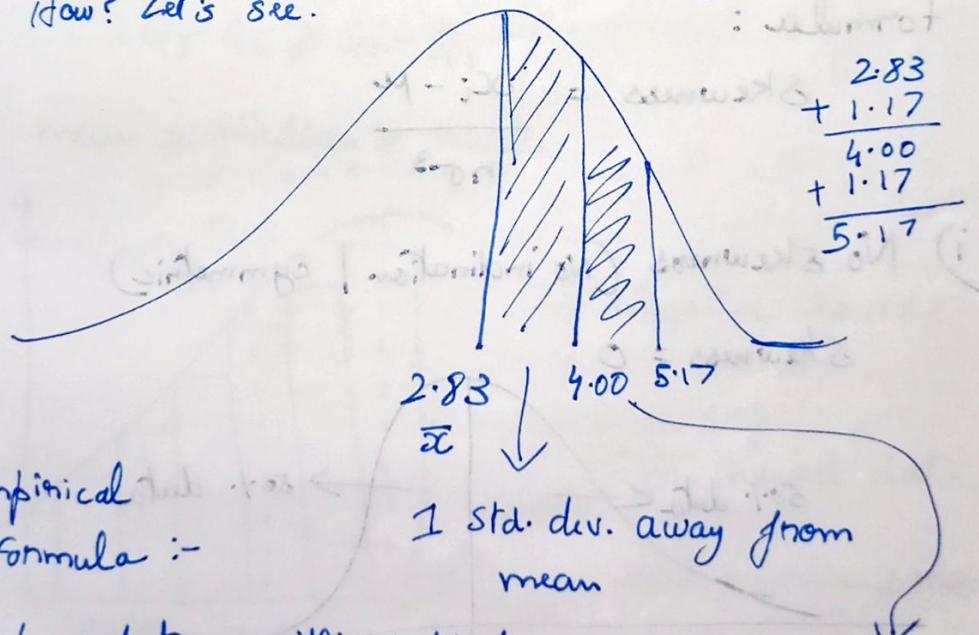
- i) Variance can be a huge no. because it talks about spread at an overall level. Comparison of each data point w.r.t to variance becomes difficult.
- ii) Variance changes the dimension to square (Ex.,  $m \rightarrow m^2$ ). Standard dev. converts it back to  $m$  on the original dimension. This is easy for making inferences.

Ex., data = {1, 2, 3, 3, 4, 4, 3}

$$\text{Var.} = 1.37$$

std. dev. =  $1.17 \Rightarrow$  std. way of knowing where your data lies.

How? Let's see.



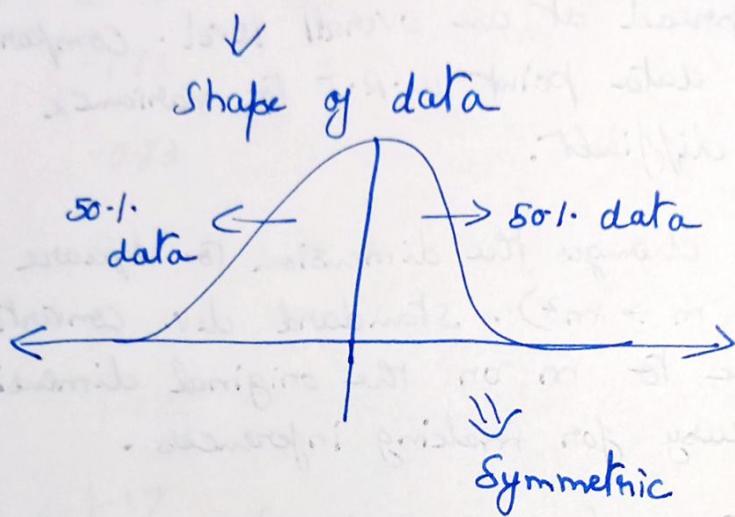
$$\begin{array}{r} 2.83 \\ + 1.17 \\ \hline 4.00 \\ + 1.17 \\ \hline 5.17 \end{array}$$

1 std. dev. away from mean

2 std. dev.

away from mean

### III) Measures of Symmetry



#### D Skewness:-

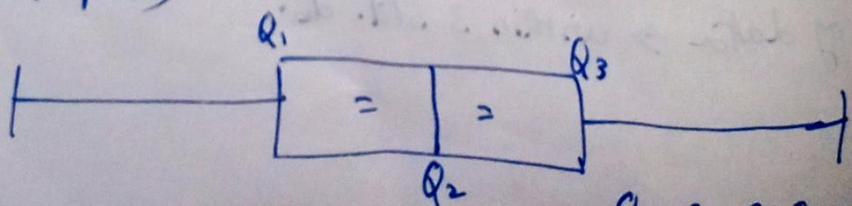
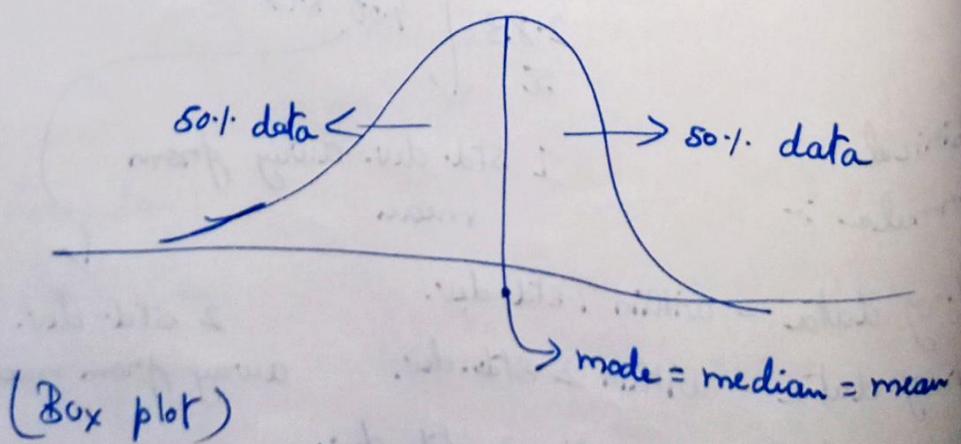
Inclined to some side

Formula :

$$\text{Skewness} = \frac{\bar{x}_i - \mu}{n\sigma^3}$$

i) No skewness (No inclination / symmetric)

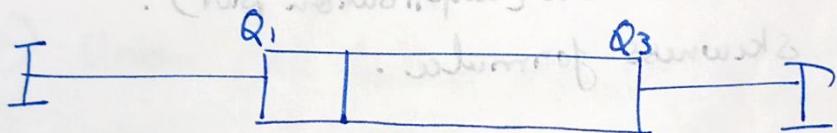
$$\text{Skewness} = 0$$



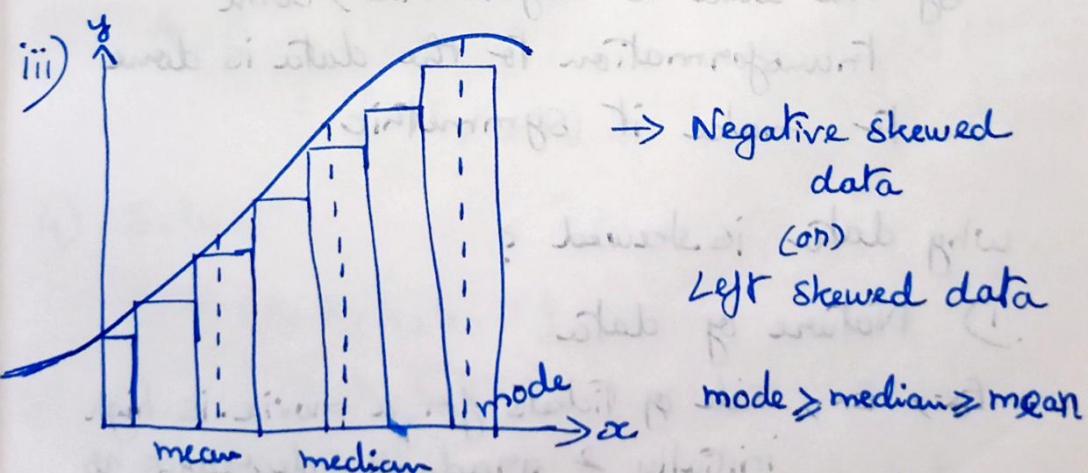


→ Tail is on right side of distribution  
→ Most of data is in the right side of the distribution.

(Box plot)

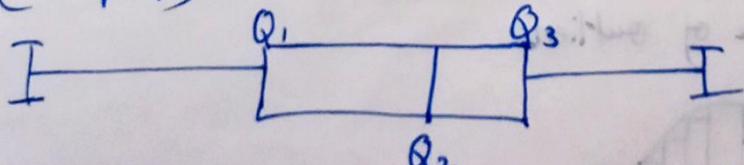


$\text{mean} \geq \text{median} \geq \text{mode}$



$\text{mode} \geq \text{median} \geq \text{mean}$

(box plot)

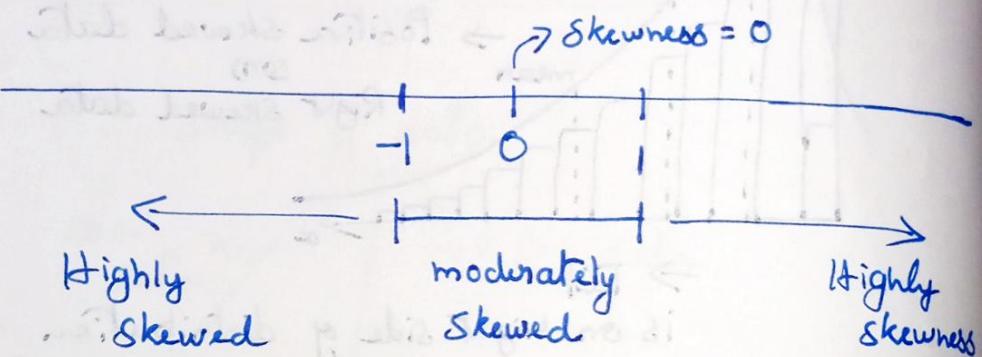


$Q_2 - Q_1 \geq Q_3 - Q_2$

Tail on the left side of distrib.

Most of data on left side of distribution.

Possible Values for skewness :



To know if data is skewed or not?

- Visualisation (distribution plot).
- Skewness formula.

Use Case : Some of ML algorithms requires Symmetric data to build the ML model.

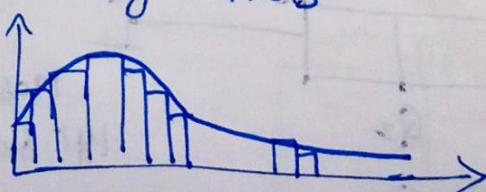
If the data is asymmetric, some transformation to the data is done to make it symmetric.

why data is skewed ?

i) Nature of data

Ex: Ex., sale of tickets for a movie is high initially & gradually decreases as days pass on.

ii) Presence of outliers



# SET

Collection of un-ordered unique elements.

Properties of set :

1) Intersection (common elements)

$$A = \{1, 5, 7, 9, 13\}$$

$$B = \{13, 15, 17\}.$$

$$A \cap B = \{13\}$$

2) Union (All distinct elements from both sets)

$$A \cup B = \{1, 5, 7, 9, 13, 15, 17\}$$

3) Difference

$$\begin{aligned} A - B &= \{1, 5, 7, 9, 13\} - \{13, 15, 17\} \\ &= \{1, 5, 7, 9\}. \end{aligned}$$

4) Subset

$$C = \{1, 2, 4, 5, 7\}$$

$$D = \{2, 4, 5\}$$

D is subset of C as all the elements of D is present in C .

## 5) Superset

From the C & D example, C is the Superset of D as

C contains all the elements of D.

## 6) Symmetric Difference (Opposite of Intersection)

$$E = \{1, 20, 200, 4000\}$$

$$F = \{20, 40, 80, 100, 200\}$$

$$E \Delta F = \{1, 200, 4000, 40, 80, 100\}$$

All the elements of both sets except the common ones.

## Covariance & Correlation

### Understanding the Relationships

#### i) Direct Relationship:

If one variable  $\uparrow$ , another variable  $\uparrow$ , then they are in a direct relationship.

Ex., Predict price of house based on area of house.

If the area of house increases, price of house also increases.

#### ii) Indirect relationship:

If one variable  $\uparrow$ , another variable  $\downarrow$ , or vice versa, then they are in indirect relationship.

Ex., As you drive longer, the distance covered will be increasing but the fuel level in the vehicle will be decreasing.

To Quantify / Measure The Relationships :-

#### ① Covariance

$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

where,  $x$  &  $y$   
are two features.

Covariance  $\Rightarrow$  relationship b/w two variables

$x \uparrow y \downarrow$   
 $x \downarrow y \uparrow$

(on)

$x \uparrow y \uparrow$   
 $x \downarrow y \downarrow$

-ve covariance

+ve covariance

Ex:- Listed below following table

$x$	$y$	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})(y - \bar{y})$
2	3	-1	2.5	2.5
3	5	0	-0.5	0
6	6	3	0.5	1.5
1	8	-2	2.5	-5
$+$				<u>-1</u>
$\bar{x} = 3.5$				
$\bar{y} = 5.5$				
$\text{cov}(x, y) = \frac{-1}{4-1}$				
$= -1/3 = -0.33$				

Interpretation  $\Rightarrow x$  &  $y$  are negatively related.

### Disadvantage of Covariance

- $\rightarrow$  The range of Covariance is  $-\infty$  to  $\infty$ .
- $\rightarrow$  This means the comparison of strength of relationship is not possible.

→ There is no standardised scale to interpret the strength.

→ Covariance has dimension.

$$\text{Ex. } \text{Cov}(\text{price}, \text{height}) = \text{Rs. ft} = 450 \text{ Rs. ft}$$

$$\text{Cov}(\text{height}, \text{weight}) = \text{ft. kg.} = 600 \text{ ft. kg.}$$

We cannot

Say that  $\text{Cov}(\text{ht.}, \text{wt.})$  is larger than  $\text{Cov}(\text{price}, \text{ht.})$  as both are diff. units.

(The above example is for height, weight & Price of tiles (floor tiles))

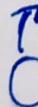
To overcome these disadvantages, we use,

## ② Pearson Correlation Coefficient (-1 to 1)

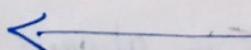
$$f_{(x,y)} = \frac{\text{Cov}(x,y)}{\sigma_x \cdot \sigma_y}$$

Dimensionless quantity as the dimensions are cancelled out by the standard deviations of the features.

No relation



-1



Increase in  
-vety of relations  
btw. features  
(Indirect)

Increase in positivity  
of relations btw.  
features (direct)

$$f_{(x,y)} = 0.4$$

$$f_{(a,b)} = 0.8$$

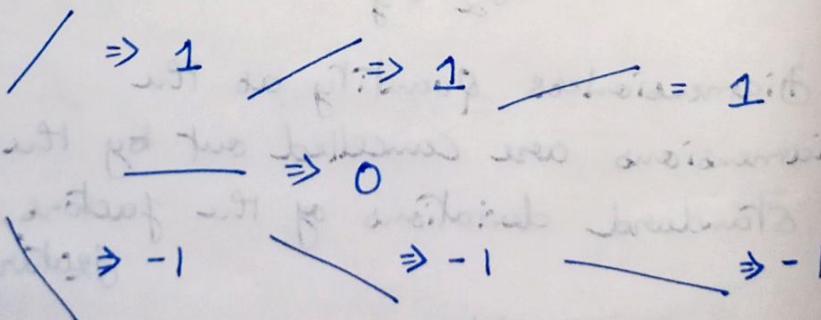
Features  $a, b$  is highly +vely correlated as compared to features  $x, y$ .

$$f_{(c,d)} = -0.2 \quad f_{(o,p)} = -0.5$$

Features  $o, p$  are more -vely correlated as compared to Features  $c, d$ .

IMP :: Pearson Correlation Coefficient always measures the linear relationships.

Correlation is not dependent on the angle of the slope.



For non-linear Relationships, the Pearson Correlation Coefficient is always 0.

To understand non-linear relationships,  
we use,

### ③ Spearman Rank Correlation

$$r_s = \frac{\text{Cov}(R(x), R(y))}{\sigma(R(x)) \cdot \sigma(R(y))}$$

where,  $R(x) \Rightarrow$  Rank of  $x$ .

$R(y) \Rightarrow$  Rank of  $y$ .

Ex.	$x \rightarrow$	5	7	8	1	2
	$y \rightarrow$	6	4	3	1	2
	$R(x) \rightarrow$	3	2	1	5	4
	$R(y) \rightarrow$	1	2	3	5	4

$x \rightarrow 5, 7, 8, 1, 2 \Rightarrow$  Sort the values in  
descending order

8, 7, 5, 2, 1 → 5<sup>th</sup>  
1<sup>st</sup> ← 2<sup>nd</sup> ↓ 3<sup>rd</sup> ↓ 4<sup>th</sup>

$y \rightarrow 6, 4, 3, 1, 2$

6, 4, 3, 2, 1 → 5<sup>th</sup>  
1<sup>st</sup> ← 2<sup>nd</sup> ↓ 3<sup>rd</sup> ↓ 4<sup>th</sup> ↓ 5<sup>th</sup>

Then, we use the formula

to find out the correlation.

Use case for correlation  $\Rightarrow$  used for feature selection in ML modelling.

Ex., Let's take factors

$x_1, x_2, x_3, x_4, x_5, \dots, x_{1000}$

&  
 $y$  (price)

Correlation is used to find which factors/features from  $x_1$  to  $x_{1000}$  correlate well with feature  $y$  (price) in-order to use it for machine learning modelling & predictive modelling.