

Probability and Statistics

Introduction to Statistics

(Decision making process)



def: Statistics is the science of collecting, organizing and analyzing data.

Data: "facts to pieces of information"

e.g. Height of students in classroom.

IQ of students in classroom.

Types of Statistics

Descriptive Statistics

It consists of organizing and summarizing data.

(1) Measure of central tendency.

(i) Mean

(ii) Median

(iii) Mode

(2) Measure of dispersion

(i) Variance

(ii) Standard deviation

(3) Diff type of distribution of data.

e.g. Histogram

Probability distribution function.

Probability Mass function.

Inferential Statistics

It consists of using data you have measured to form conclusion.

from
Sample data
concluding

Population data

(i) Z-test
(ii) t-test

(iii) Chi Square test

Hypothesis
Testing
 H_0 , H_1
P-value
Significance value

g: There are 20 Statistics Classes at your University and you have collected the heights of student in the class. Heights are recorded as

{ 175, 180, 140, 140, 125, 160, 135, 190 }

Descriptive Question

"What is the ~~common~~ ^{avg} height of the entire classroom"

Inferential Question

"Are the height of the students in the classroom similar to what you expect in the entire University?"

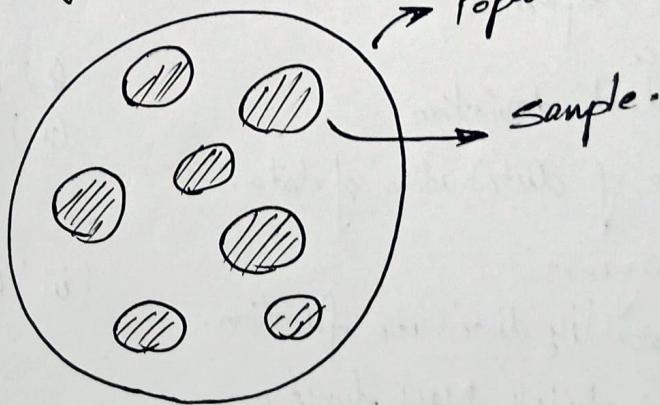
↓
↳ population data

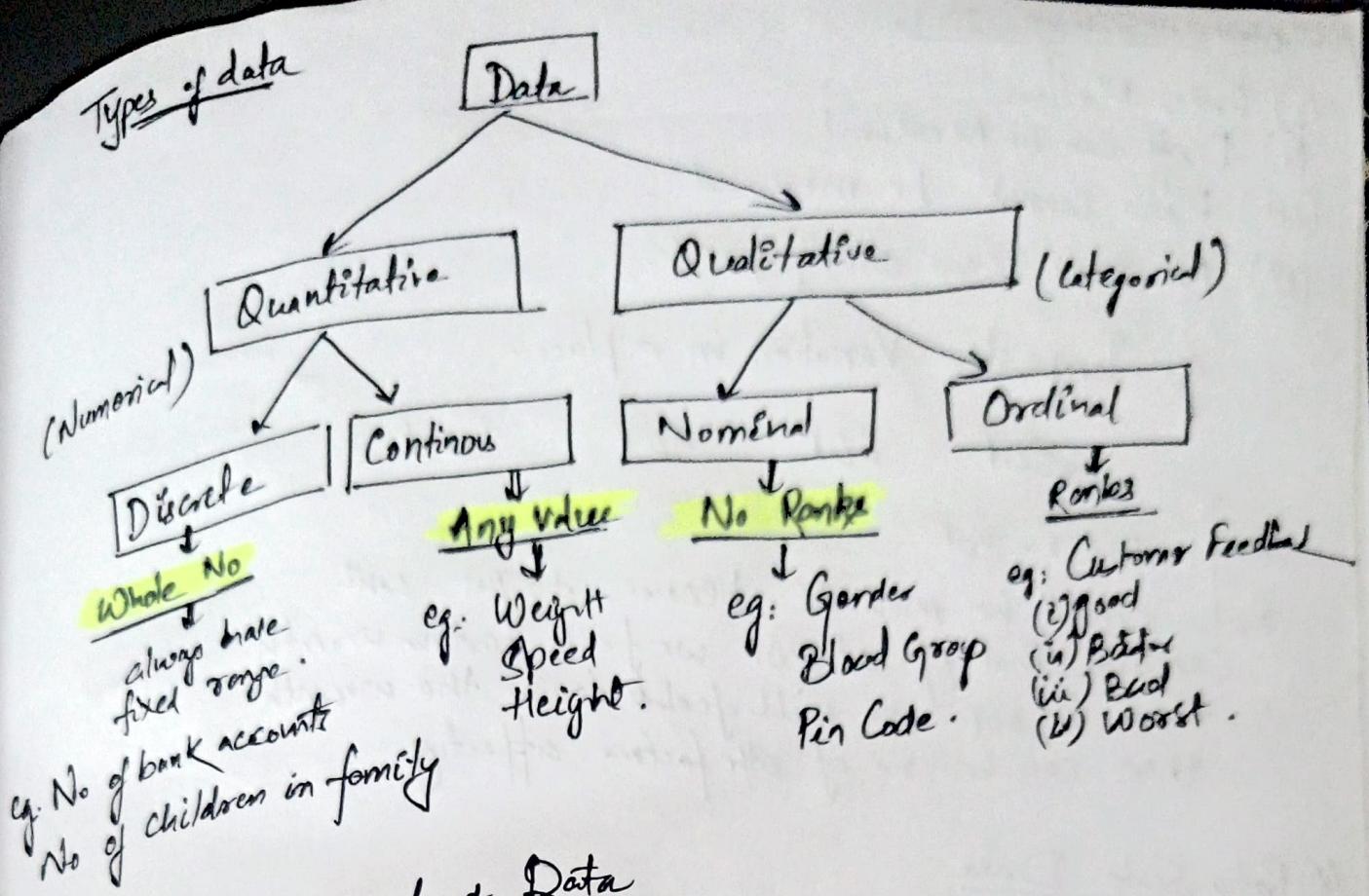
Population data & Sample data

Population : The group you are interested in studying.

Sample : A subset of Population.

e.g: Exit Poll.





Scale of Measurement of Data

- (1) Nominal Scale Data
- (2) Ordinal Scale Data
- (3) Interval Scale Data
- (4) Ratio Scale Data

(1) Nominal Scale Data

- (i) Qualitative
- (ii) Categorical Data
- (iii) Order/Rank does not matter

e.g.: Survey on favorite color.
 (These order does not matter on what color people like.)

(2) Ordinal Scale Data:

- (i) Ranking is Important
- (ii) Order Matters
- (iii) Difference cannot be measured

e.g.: A feedback form
 (i) Best
 (ii) Good
 (iii) Bad
 (iv) Worst

→ Just on the basis of ranks we cannot calculate difference

like in this case we don't know the reason for this ranking

(3) Interval Scale of Data

- (i) Order Matters
- (ii) Diff can be Measured
- (iii) Ratio cannot be measured
- (iv) No true '0' startig point.

e.g. Temperature Variation in a place.

30°F 60°F 90°F 120°F

$$\text{Diff: } 60 - 30 = 30^{\circ}\text{F}$$

Ratio cannot be measured because it's just cont
 conclude that if at 30°F we feel certain warmth
 then at 60°F we will feel double the warmth
 there can be no of other factors effecting.

(4) Ratio Scale Data

- ① The Order matters
- ② Diff are measurable (including ratio)
- ③ Common '0' starting point.

Grades of students

0, 90, 60, 30, 75, 40, 50

- Measure of Central Tendency
- (i) Mean or Average
 - (ii) Median
 - (iii) Mode

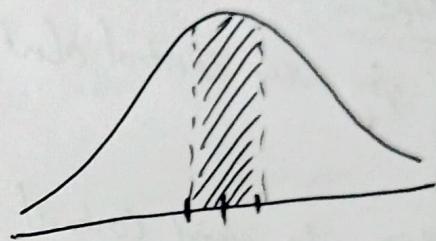
Mean Population (N)

$$L = \{1, 1, 2, 2, 3, 3, 4, 5, 6\}$$

$$\text{Population Mean } (\mu) = \frac{\sum_{i=1}^N x_i}{N} = \frac{[1+1+2+2+3+3+4+5+6]}{10} = 3.2$$

Sample Mean (n)

$$\text{Sample Mean } (\bar{x}) = \frac{\sum_{i=1}^n x_i}{n}$$



Median $x = \{4, 5, 3, 2, 1\}$

$$\text{Step 1: Sort the random variable } x. : \{1, 2, 3, 4, 5\}$$

Step 2: No. of elements count: 5

Step 3: if Count % 2 == 0
find centrd element

$$\text{if count} = 6 \quad \{1, 2, \boxed{3}, 4, 5\}$$

$$\frac{2+3}{2} = 2.5 \quad \text{Median}$$

or if (Count % 2 != 0)
find centrd element
 $\frac{2+5}{2} \rightarrow 3$
Median

Why Median?

$$X = \{1, 2, 3, 4, 5\}$$

$$\bar{x} = \frac{1+2+3+4+5}{5} = 3$$

if now introduce another sheet.

$$X = \{1, 2, 3, 4, 5, 100\}$$

(outlier)

because it
does not
belong to the
distribution.

$$\bar{x} = \frac{1+2+3+4+5+100}{6} = \frac{115}{6} = 19.1$$

Now Median for the central sheet in ① Col $\rightarrow 3$

$$\textcircled{2} \text{ Col} = \frac{3+4}{2} = 3.5$$

* Median to find central tendency when outliers present.

Mode

frequency \rightarrow Maxm frequency.

$$\{2, 1, 1, 1, 4, 5, 7, 8, 9, 9, 10\}$$

Maxm frequency of an sheet: 1 80, Mode = 1

Application of Mean, Median, Mode in Foatra Engineering

Age	Weight	Salary	Gender	Degree
29	70	40K	M	B.E.
25	80	70K	F	-
27	95	45K	M	-
24	-	50K	-	Ph.D
32	-	60K	-	Masters
-	60	-	-	B.Sc
-	65	55K	M	-
40	72	-	F	-

Numerical Value → We can find mean and fill null values.
 ↳ of outliers → median value.

Categorical Values: We can fill it with mode values.

Measure of Dispersion (Spread of the data)

- (1) Variance
- (2) Standard deviation.

Variance

Population Variance (σ^2)

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

x_i = Data points

μ = population mean

N = population size

$$\text{Sample Variance}$$

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

↓
Diff

x_i = Data points

\bar{x} = sample mean

n = sample size

Q Why we divide sample variance by $(n-1)$?

In order to create unbiased estimator of σ^2 .

Q What is unbiased and biased estimator?

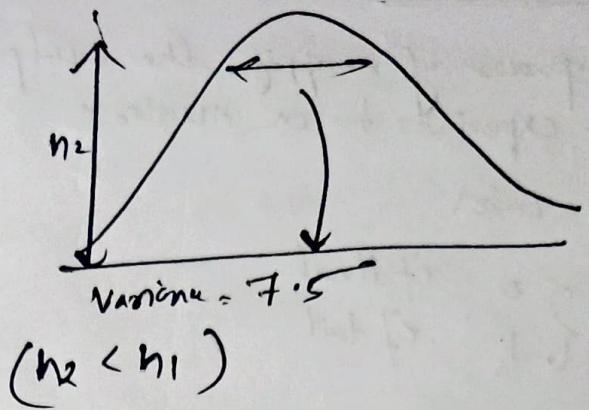
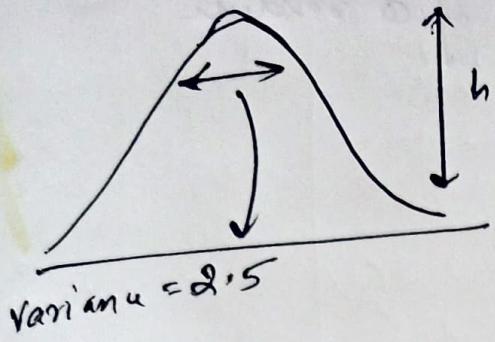
Unbiased: like expected value is equal to the true value of the parameter.

Biased: expected value is not equal to the true value due to bias.

eg. If we pick 20 random students from a school and find avg
there is a good chance our estimator will give avg height true to the actual height

Biased { If we pick up student only from the basketball team where students are typically taller our avg value will not be equal to true value.

Dispersion or Spread (Variance)



Standard Deviation

Population S.D

$$\sigma = \sqrt{\text{variance}} \quad (\sigma^2)$$

Sample S.D

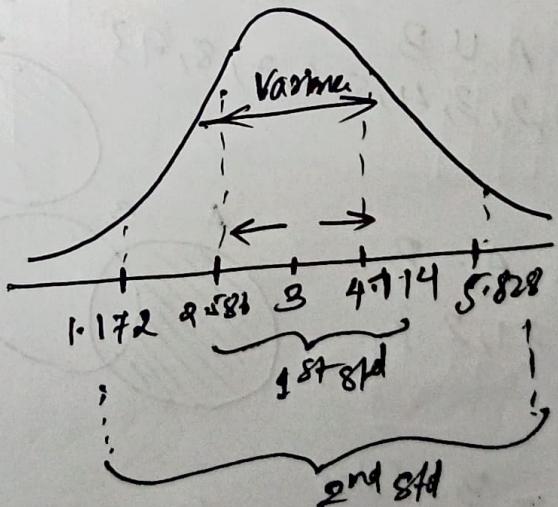
$$\text{std} = \sqrt{s^2}$$

Sample variance

e.g. $\{x = 1, 2, 3, 4, 5\}$

$$\bar{x} = 3$$

$$\sigma = 1.414 \cancel{2.23}$$



Random Variables

If it is the process of mapping the output of a random process or experiment to a number.

e.g. (1) Tossing a coin

$$X \begin{cases} 0 & \text{if Head} \\ 1 & \text{if tail} \end{cases}$$

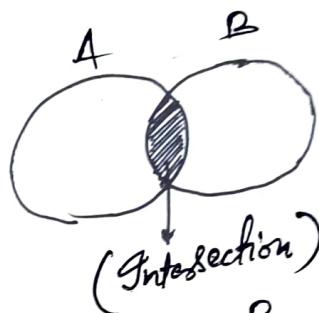
(2) Rolling a dice

Sets

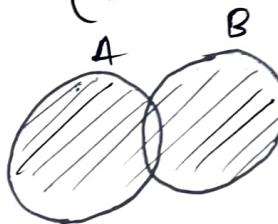
$$A = \{1, 2, 3, 4, 5, 6, 7, 8\}$$

$$B = \{3, 4, 5, 6, 7, 8, 9\}$$

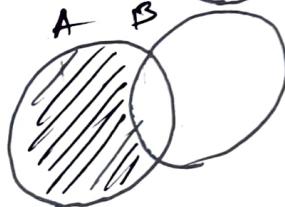
① Intersection: $A \cap B$
 $\{3, 4, 5, 6, 7, 8\}$



② Union: $A \cup B$
 $\{1, 2, 3, 4, 5, 6, 7, 8, 9\}$



③ Difference: $A - B$
 $\{1, 2\}$



④ Subset:

$$\begin{aligned} A \rightarrow B &\Rightarrow \text{False} \\ B \rightarrow A &\Rightarrow \text{False} \end{aligned}$$

Histogram & Skewness

Ages = {10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51}

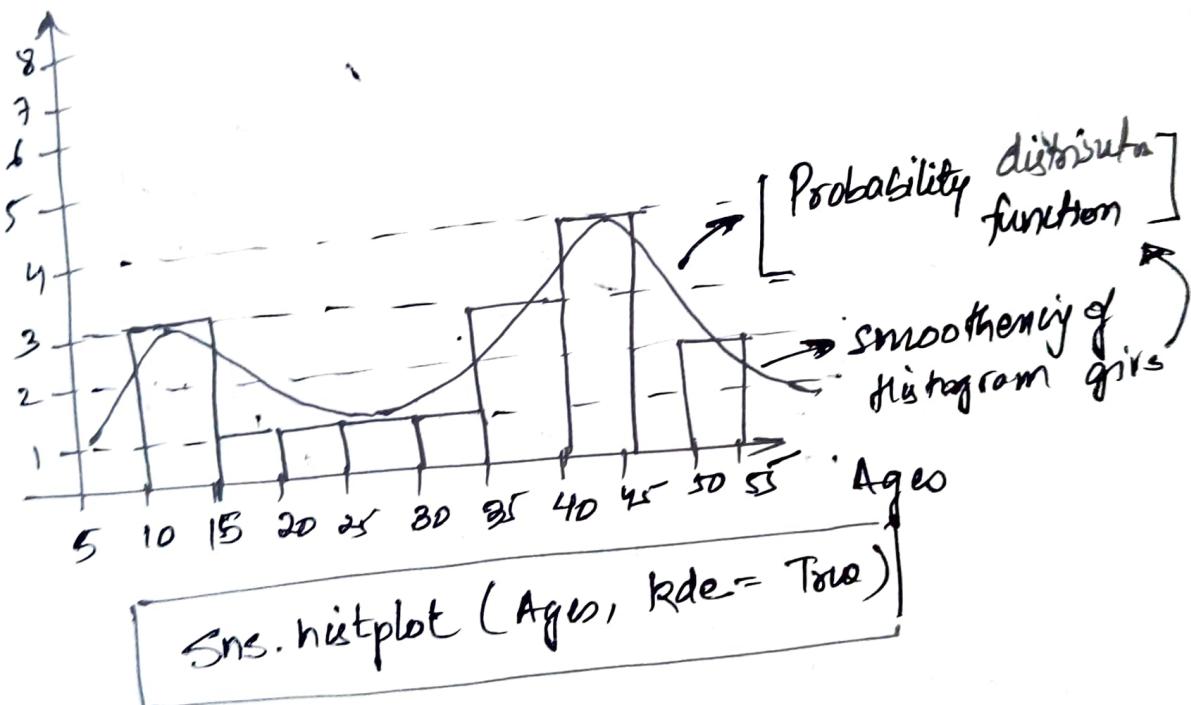
Consider, we have 60 values and we want bin off 10 bins.

$$\frac{60}{10} = \boxed{6} \rightarrow \text{bin size} \quad (\text{No of bins} = 10)$$

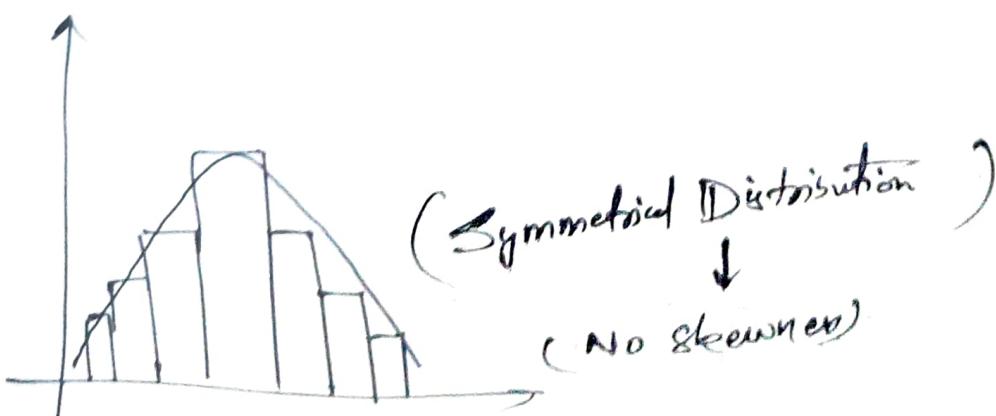
$$\frac{60}{20} = \boxed{3} \text{ bin size} \quad (\text{No of bins} = 20)$$

Histogram

Count.

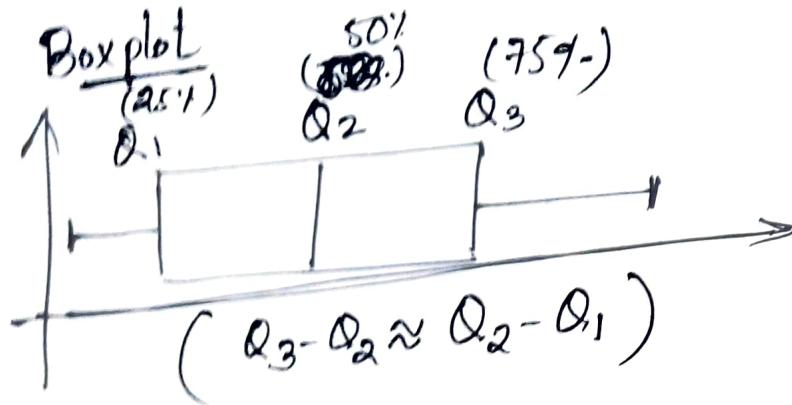


Symmetrical

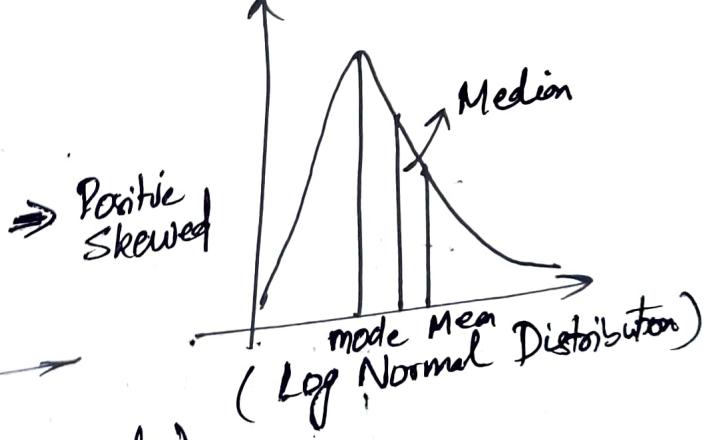
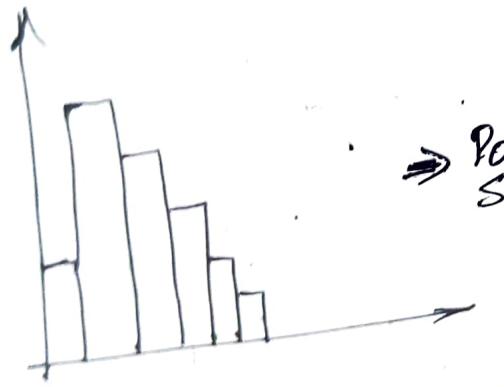


No skewness

- ① The mean, median, mode all are at perfect centre
 $\rightarrow (\text{Mean} = \text{Median} = \text{Mode})$

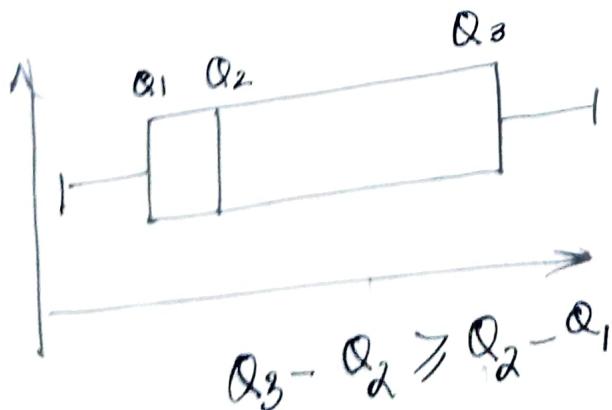


- ② Right Skewed data

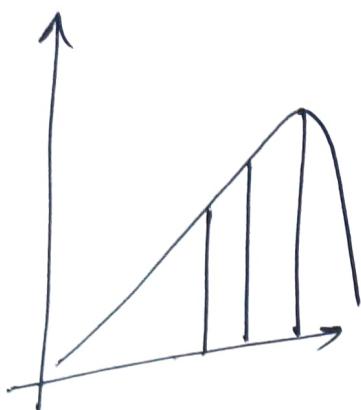


(Relationship b/w mean, median, mode)

$\boxed{\text{Mean} > \text{Median} > \text{Mode}}$

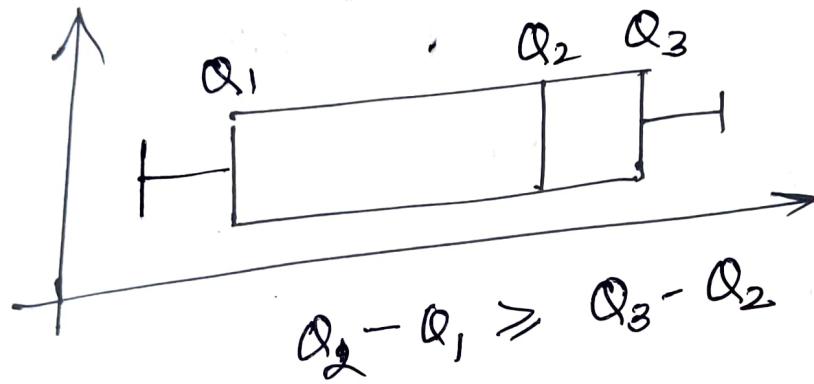


Left-skewed distribution



⇒ Negative skewed

Mean < Median < Mode



Covariance & Correlation

[Relationship b/w X and Y]

X	Y
2	3
4	5
6	7
8	9

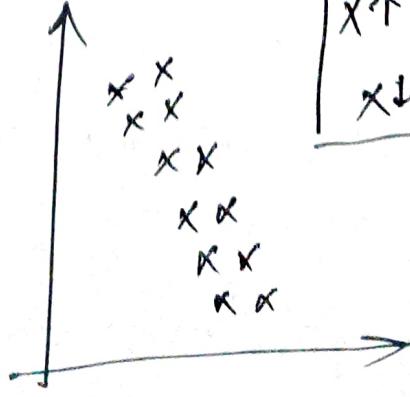
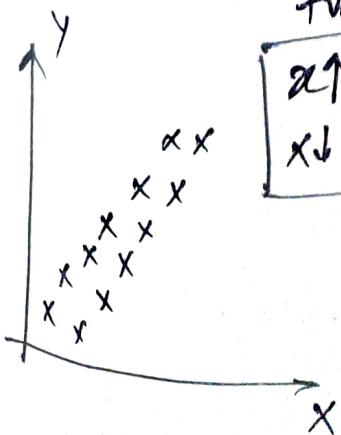
can be

X↑	Y↑
X↓	Y↑
X↓	Y↓
X↑	Y↓

-ve (Cov)

+ve (Cov)

X↑ Y↑
X↓ Y↓



Covariance

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

x_i = data point

\bar{x} = mean (sample)

y_i = data point

\bar{y} = mean (sample)

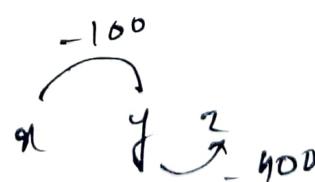
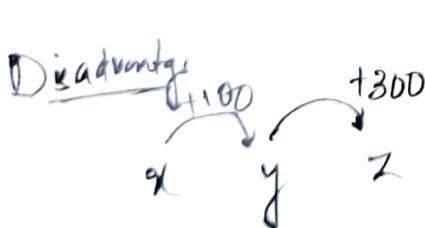
$$\begin{aligned}\text{Var}(x) &= \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \\ \downarrow \\ \text{Cov}(x, x)\end{aligned}$$

* ($\text{Cov}(x, x)$ is basically $\text{Var}(x)$ which is the spread of data on x .)

Advantages of Covariance

① find out the relationship b/w x and y .

②



- * No limit on the increment or decrement of value of random variables, so conclusion is difficult.
- * Covariance does not have limit value.

To fix this we use: Pearson Correlation Coeff.

Pearson. Correlation Coefficient [-1 to 1]

$$f(x,y) = \frac{gV(x,y)}{\epsilon_x \epsilon_y}$$

s_x : Std of x
 s_y : Std of y

(for linear relationships) or approx (Normal distribution) because of this Cov is limited

- (~~referred~~) (Now called) Cov.

The more the value towards +1, the more +vely correlated it is.
 The more the value towards -1, the more -vely correlated it is.
 Coefficient [-1 to 1]

Spearman Rank Correlation [-1 to 1]

$$\gamma_s = \frac{\text{Cov}(R(x), R(y))}{\sigma_{R(x)} \times \sigma_{R(y)}}$$

$$R(x) = \text{Rank } g^x$$

$$R(y) = \text{Rank } g^y.$$

x	y	$R(x)$	$R(y)$
1	2	5	5
3	4	4	4
5	6	3	3
7	8	2	1
0	7	6	2
8	1	1	6

(prefixed)

for data Not Normally distributed

- ① has Outliers
- ② Not strictly linear

Implementation

① Feature Selection

Size of
flame

No of Rooms ↑
(+ve)

(+ve)
Location ↑

(+ve)
Location ↑

(pp features) w/ relation

$$\approx 0 \text{ (No } r^c \text{)}$$

we can delete feature which do not effect the output

Price ↑

* Correlation and Covariance Relationship

Covariance : tells you how two random variables are connected to each other.
but does not give you a measure of strength.

Correlation : makes the relationship a standardized relationship which is a simple scale from -1 to 1.
It tells you the strength of a relationship.

Kendall Correlation Coeffic

It also tells the strength and direction, but it maps the value on basis of its agreement or disagreement.

formula

$$T = \frac{C - D}{\frac{1}{2} n(n-1)}$$

C-Concordant: if ranks of all pairs agree with each other.

D-Discordant: if ranks of all the pairs dont agree with each other.

If most pair are in same order = 1

If most pair are in opposite order = -1

If most pair are in random = 0

- * Use "Pearson" if you care about the actual values and assume a straight line relationship.
- * Use "Spearman" if you care about the ranks of the values
- * Use "Kendall" if you want to measure the agreement or disagreement of the order of data points

Syntax covariance
`df.cov()`

corr

- ① `df.corr(method='pearson', numeric_only=True)`
- ② `df.corr(method='spearman', numeric_only=True)`
- ③ `df.corr(method='kendall', numeric_only=True)`