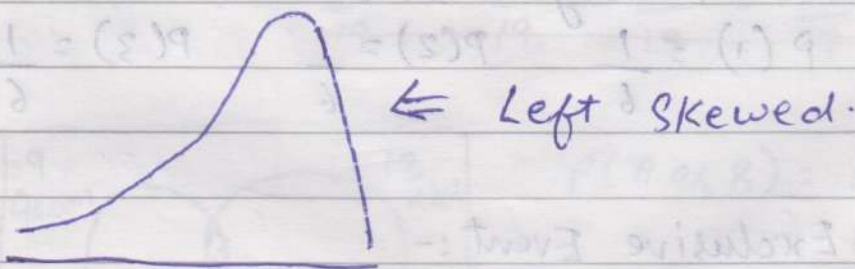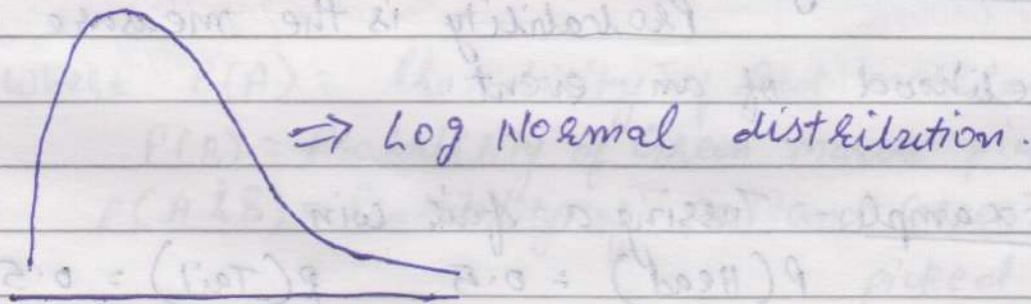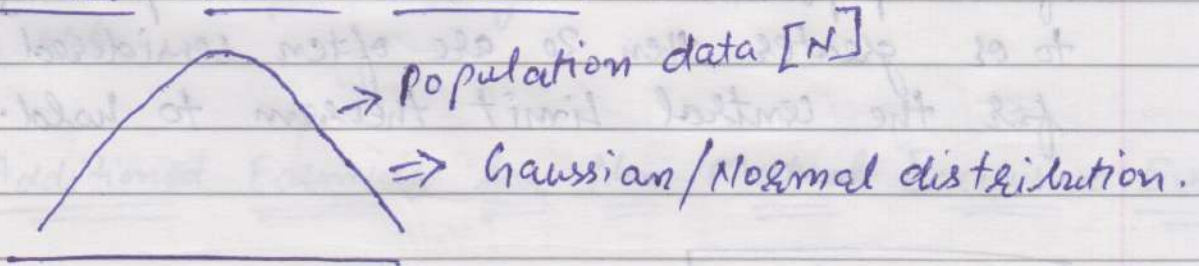① Central Limit Theorem
② Probability
③ Permutation And Combination
④ Covariance, Pearson Correlation, Spearman Rank Correlation.
⑤ Bernoulli's Distribution.
⑥ Binomial Distribution.
⑦ Power Law [Pareto Distribution]

★ Central Limit Theorem :-



→ Population data [N]

⟹ Gaussian/Normal distribution.



⟹ Log Normal distribution.



⟸ Left Skewed.

Sample data $\boxed{n}$

$$\{x_1, x_2, x_3 \ ----- \ x_n\} \rightarrow \bar{x}_1$$
$$\{x_1, x_2, x_3 \ ----- \ x_n\} \rightarrow \bar{x}_2$$
$$\{x_1, x_2, x_3 \ ----- \ x_n\} \rightarrow \bar{x}_3$$
$$\vdots \qquad \qquad \vdots \qquad \vdots$$
$$\{x_1, x_2, x_3 \ ----- \ x_n\} \rightarrow \bar{x}_n$$

The Central Limit Theorem (CLT) states that the distribution of sample means approximates a normal distribution as the sample size gets larger, regardless of the population's distribution. Sample size equal to or greater then 30 are often considered sufficient for the central limit theorem to hold.

A  |Probability.| =
                    Probability is the measure of the likelihood of an event.

Example:- Tossing a fair coin
$$P(Head) = 0.5 \qquad P(Tail) = 0.5$$

Example 2 :- Rolling a Dice
$$P(1) = \frac{1}{6} \qquad P(2) = \frac{1}{6} \qquad P(3) = \frac{1}{6}$$

① ⟹ Mutual Exclusive Event :-
                    Two Events are mutually exclusive if they cannot occur at the same time.
$$P(A \ or \ B) = P(A) + P(B)$$
Eg:- Tossing a Coin        Eg :- Rolling a dice.

② ⟹ **Non - Mutual Exclusive Events:-**

    Two Events can occur at the same time.

Eg:- Picking randomly a card from a deck of cards, two events "Heart" and "King" can be selected.

Eg:- Bag of Marbles:- 10 Red, 6 Green; 3 (R & G)

      a              b            or

   Red Marble      Green Marble     Red and Green
                                           Marbel

   There is a probability that we choose red and Green marble.

**\* Additional Formula for Non Mutual Exclusive Event**

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

  where  $P(A)$ = Probability of Red marble picked up
         $P(B)$ = Probability of Green marble picked up.
       $P(A \& B)$ = Probability of Red and Green marble
                                picked up.

$$\therefore P(A \text{ or } B) = \frac{10}{19} + \frac{6}{19} - \frac{3}{19} = \frac{13}{19}$$

Red & Green Marble

$$P(A \text{ or } B) = P(A) + P(B) - P(A \& B)$$
$$= \frac{13}{19} + \frac{9}{19} - \frac{3}{19}$$
$$= \frac{19}{19} = 1$$

Another Example → What is the probability of choosing Heart ♡ or Queen.

$$P(♡ \text{ or } Queen) = P(♡) + P(Queen) - P(♡ \text{ or } Queen)$$
$$= \frac{13}{52} + \frac{4}{52} - \frac{1}{52} = \boxed{\frac{16}{52}}$$

* **Multiplication Rule for Non Mutual Exclusive Event.**

**\* Dependents Evens:-** Two Events are dependent if they affect one another.

White Marble

Bag of marble $\left\{ \begin{array}{c} 0\ 0\ 0\ ✕ \\ 0\ 0\ 0 \end{array} \right\}$

Yellow marble

$$P\left(\begin{array}{c} \text{white} \\ \text{marble} \end{array}\right) = \frac{4}{7} \longrightarrow P\left(\begin{array}{c} \text{yellow} \\ \text{marble} \end{array}\right) = \frac{3}{6}$$

White 7 marble

**\* Independent Events:-**

**Question:-** What is the probability of rolling a "5" and then a "3" with a normal 6 sided dice ?
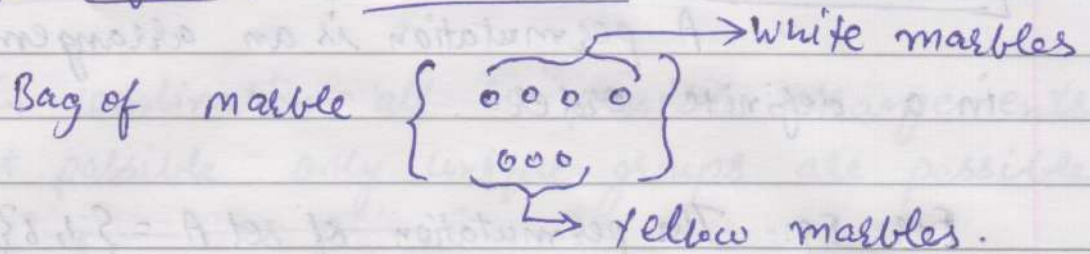
**Ans:-** $P(1) = \frac{1}{6}$   $P(2) = \frac{1}{6}$   $P(3) = \frac{1}{6}$   $P(4) = \frac{1}{6}$

Multiplication Rule For Independent Events.

$$P(A \text{ and } B) = P(A) * P(B)$$
$$= \frac{1}{6} * \frac{1}{6} = \boxed{\frac{1}{36}}$$

**\* Example of Dependent Events:-**

$$Bag\ of\ marble \left\{ \underset{\underbrace{\large\circ\circ\circ}}{\overbrace{\circ\circ\circ\circ}} \right\}$$ → white marbles

→ yellow marbles.

**[Q:-]** What is the probability of drawing a "white" and then drawing a "yellow" marble from the bag?

**Ans**

$$\boxed{\begin{array}{c} \circ\circ\circ\circ \\ \circ\circ\circ \end{array}} \longrightarrow P\left( \circlearrowleft \substack{white \\ marble} \right) = \frac{4}{7}$$

Now If we remove 1 white marble then.

$$\boxed{\begin{array}{c} \circ\circ\circ \\ \circ\circ\circ \end{array}} \longrightarrow P\left( \frac{Yelow\ marble}{white\ marble} \right) = \frac{3}{6} \rightarrow \boxed{\begin{array}{c} Conditional \\ Probability \end{array}}$$

$$\boxed{P(White\ and\ Yellow) = P\left( \substack{White \\ marble} \right) * P\left( \frac{Yellow}{white\ marble} \right)}$$

$$= \frac{4}{7} * \frac{3}{6} = \frac{2}{7}$$

**Note:-** Naive Baye's ML Algorithm is derived from conditional Probability.

**\* Permutation :-**

A permutation is an arrangement of objects in a definite order.

For Eg:- The permutation of set A = {1, 6} is 2

Such as {1,6} , {6,1} . As you can see, these are no other ways to arrange the elements of set A.

ANOTHER EXAMPLE.

$$\underline{5} * \underline{4} * \underline{3}$$

$$= 60 \text{ ways}$$

⇓

permutation.

Dairy → milk

School of childrens.

← 5 Star

Kit Kat      milky Bar      Sneakers

means there will be 60 possible arrangements like.

{ DM   KK   MB } { ⋯ } { ⋯ }
{ KK   DM   MB } { ⋯ } { ⋯ }
{ KK   MB   DM } { ⋯ } { ⋯ }
{ ⋯ } { ⋯ } { ⋯ } and so on.

$$\boxed{{}^{n}P_{k} = \frac{n!}{(n-k)}} = \frac{5!}{(5-3)!}$$

$$= \frac{5 \times 4 \times 3 \times \cancel{2} \times \cancel{1}}{\cancel{2} \times \cancel{1}} = 60$$

where n = Total no. of objects
      k = No. of Selections.

* **Combination:-** {Repeatation will not occur}

In combination all 60 possible arrangements are not possible only unique groups are possible.

$$\boxed{^nC_k = \frac{n!}{k!(n-k)!}} = \frac{5!}{3!(5-3)!}$$

$$= \frac{5 \times 4 \times 3 \times 2 \times 1}{3!(2!)} = \frac{5 \times \overset{2}{\cancel{4}} \times 3 \times 2 \times 1}{3 \times 2 \times 2}$$

$$= \boxed{10}$$

Where n = Total no. of objects.
k = No. of selections.

* **Covariance:-** { Feature Selection }

| x | y |
|---|---|
| Age | Weight |
| 12 | 40 |
| 13 | 45 |
| 15 | 48 |
| 17 | 60 |
| 18 | 62 |

Age ↑    weight ↑

Age ↓    weight ↓

⇈

Quantify the relationship of x & y using mathematical question.

$$\boxed{Cov(x,y) = \frac{\sum[(x_i - \bar{x}) * (y_i - \bar{y})]}{n-1}}$$

This formula is derived from variance.

$$S^2 = \frac{\sum(x_i - \bar{x})^2}{n-1} \quad \therefore \quad S^2 = \frac{\sum(x_i - \bar{x}) * (x_i - \bar{x})^2}{n-1}$$

∴ We can say that

$$Cov(x, x) = Var(x)$$

Now According to our Age and weight data.

$$Cov(x, y) = \frac{\sum (x_i - \bar{x}) * (y_i - \bar{y})}{n-1}$$

Here

| $x_i$ | $x_i - \bar{x}$ | $y_i - \bar{y}$ |
|---|---|---|
| 27 | 12-15 = -3 | 40-51 = -11 |
| 26 | 13-15 = -2 | 45-51 = -6 |
| 30 | 15-15 = 0 | 48-51 = -3 |
| 32 | 17-15 = 2 | 60-51 = 9 |
| 33 | 18-15 = 3 | 62-51 = 11 |

$\bar{x} = 15$
$\bar{y} = 51$
$n = 5$

Now $\sum (x_i - \bar{x}) * (y_i - \bar{y}) = (-3 \times -11) + (-2 \times -6) +$
$(0 \times -3) + (2 \times 9) +$
$(3 \times 11)$
$= 96$

∴ $Cov(x, y) = \frac{96}{n-1} = \frac{96}{4} = \boxed{24}$

24 is +ve covariance.

If +ve covariance
| $x \uparrow$ | $y \uparrow$ | $x$ increase |
|---|---|---|
| $x \downarrow$ | $y \downarrow$ | $y$ increase. |

If -ve covariance
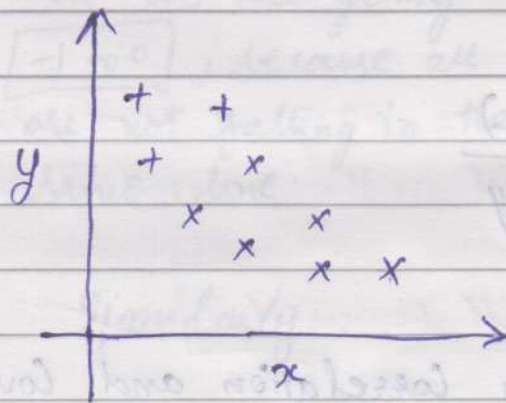| $x \uparrow$ | $y \downarrow$ |
|---|---|
| $x \downarrow$ | $y \uparrow$ |

↑ $x$ increase then $y$ decrease.

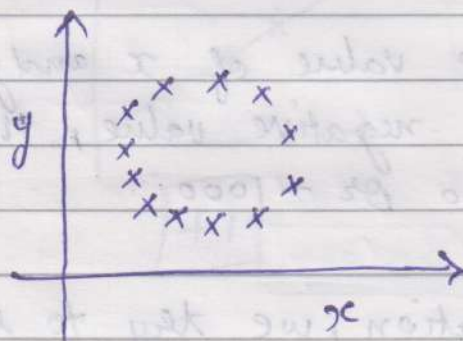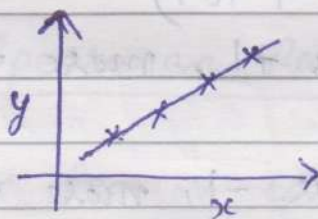If covariance is 0 | No relation with $x$ & $y$
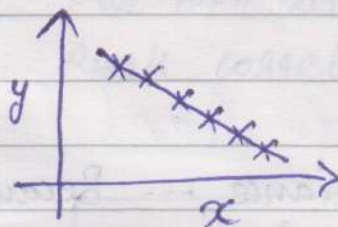
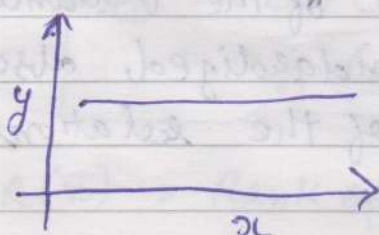$\Rightarrow$ +ve covariance.



$\Rightarrow$ -ve covariance.



$\Rightarrow$ 0 covariance, No relation with x & y



$\Rightarrow$ +ve covariance.



$\Rightarrow$ -ve covariance.



$\Rightarrow$ 0 covariance, No relation with x & y.

\* **Pearson Correlation Coefficient :-**

rho
$$p(x,y) = \frac{\sum [(xi - \bar{x}) * (yi - \bar{y})]}{\sigma x * \sigma y}$$

$$p(x,y) = \frac{Cov(x,y)}{\sigma x * \sigma y}$$

Difference between Pearson Correlation and Covariance.

In Covariance (x,y), the value of x and y can be any positive value or negative value, there is no limit, it can be + 1000 or - 1000.
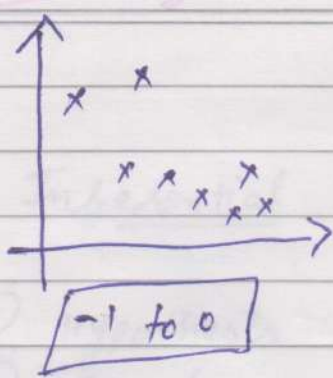
By using Pearson Correlation, we try to restrict value of x and y between (-1 to 1)

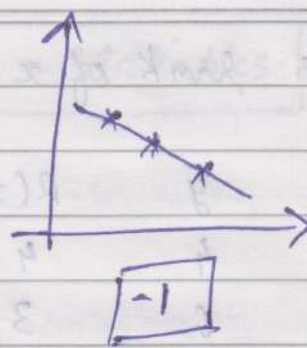[Note] :- - More the value towards +1, more +ve correlated it is.

- More the value towards -1, more -ve correlated it is.

Correlation is better then Covariance --- Because correlation removes the effect of the variance of the variables, it provides a standardized, absolute measure of the strength of the relationship, bounded by -1.0 and 1.0.
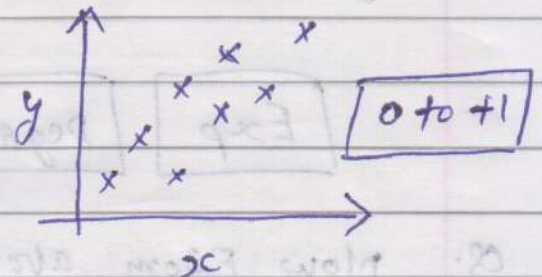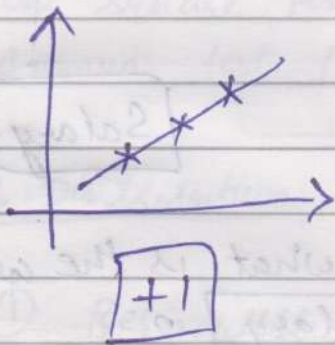
$$\boxed{-1 \ to \ 0}$$

$$\boxed{-1}$$

Since all the points are in
same line rho(p) is -1

Here we are getting
$\boxed{-1 \ to \ 0}$, because all
are not falling in the
same line.

Similarly



$$\boxed{+1}$$

$$\boxed{0 \ to \ +1}$$

* $\boxed{\text{Spearman Rank Correlation:-}}$

- Pearson Correlation is only good for Linear Data.
- For Non linear data we have to use Spearman
  Rank correlation.

$$\boxed{r_S = \frac{Cov(R(x), R(y))}{\sigma(R(x)) * \sigma(R(y))}}$$

$R(x) = $ Rank of $x$

What is rank of x ? R(x)

| x | y | R(x) | R(y) |
|----|----|------|------|
| 10 | 4 | 4 | 1 |
| 8 | 6 | 3 | 2 |
| 7 | 8 | 2 | 3 |
| 6 | 10 | 1 | 4 |

R(x) — 1,2,3,4 are the ranks assigned to x, and all the ranks are assigned in the ascending order.

[ Why this correlation will be used ? ]

[ Exp ]  [ Degree ]  [ City ]          [ Salary ]

Q.   Now From above example what is the correlation between [ Exp ] and [ Salary ] ?

Ans.   It is +ve correlated.

Q.   What is the correlation between [ City ] and [ Salary ] ?

Ans   +ve correlated as the city is good like banglore or girgaon the salary is much better.

Q.   What is the correlation between [ Exp ], [ Degree ] ?

Ans   No relation.

Q.   Correlation between [ Degree ] and [ Salary ] ?

Ans.   +ve correlated.