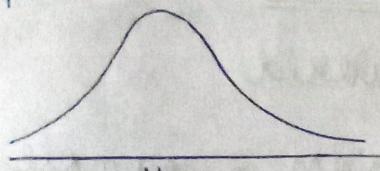


STATISTIC FOR DATA SCIENCE [DAY 4]

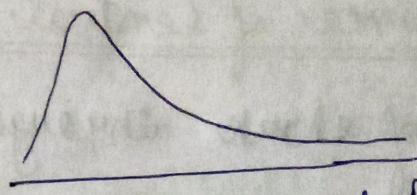
- (1) Central Limit Theorem
- (2) Probability
- (3) Permutation and Combination
- (4) Covariance, Pearson Correlation, Spearman Rank Correlation
- (5) Bernoulli's Distribution
- (6) Binomial Distribution
- (7) Power Law (Pareto Distribution)

(1) Central Limit Theorem

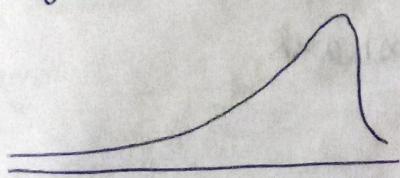
Population Data (N)



[Normal Distribution]



[Log Normal Distribution]



[Right skewed]

* In the above distribution data, if we consider sample data (n) 

Let ^{1st} sample data $\{x_1, x_2, x_3, \dots, x_n\}$ of mean \bar{x}_1 ,

^{2nd} sample data $\{x_1, x_2, x_3, \dots, x_n\}$ of mean \bar{x}_2

^{3rd} sample data $\{x_1, x_2, x_3, \dots, x_n\}$ of mean \bar{x}_3

m^{th} sample data $\{\dots\}$ of mean \bar{x}_m

n = size of sample, m = no. of sample

Let, $n \geq 30$ if we plot all the sample data then the central limit theorem then we will get a gaussian / normal distribution.

Let us consider a population data which may be normally distributed or not be then the central limit theorem say that if we take sample of size $n \geq 30$ and selected m no. of samples of respective mean and if we plot the data then we will get normal distribution ~~of total~~ of the sample mean.

Importance of Central Limit Theorem

Size of shark through the world

→ We will take 10 different region, sample population of size $n \geq 30$. and make an assumption on the size of shark.

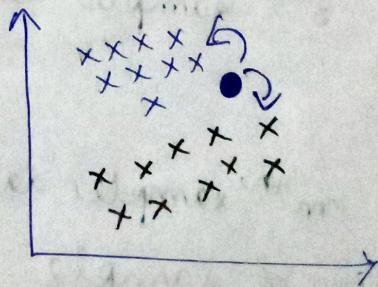
PRACTICAL APPLICATIONS

PROBABILITY :- Probability is a measure of the likelihood of an even

Eg; Tossing a fair coin $p(H) = 0.5$ $p(T) = 0.5$

Rolling a dice., $p(1) = \frac{1}{6}$, $p(2) = \frac{1}{6}$, $p(3) = \frac{1}{6}$

To check whether '●' data belongs to 'x x x' or 'x x x' ← we will use probability

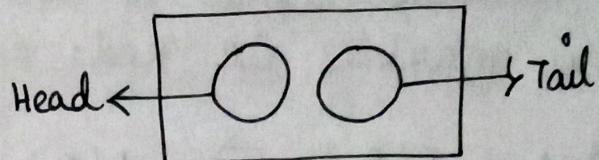


(i) Mutual Exclusive Event

Two events are mutually exclusive if they cannot occur at the same time

Example (i) Tossing a Coin (ii) Rolling a dice

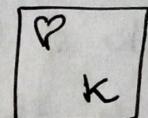
(ii) Non-Mutual Exclusive Event



Two events can occur at the same time

Example (i) Winning and Loosing game

(ii) Picking Randomly a card from a deck of cards, two events "heart" and "king" can be selected



PROBLEM STATEMENT OF MUTUAL EXCLUSIVE EVENT

Q1 What is the probability of coin landing on heads or tails

$$\text{Ans:- } P(A \text{ or } B) = P(A) + P(B) = \frac{1}{2} + \frac{1}{2} = 1$$

Q2 What is the probability of getting 1 or 6 or 3

$$\text{Ans:- } P(1) = \frac{1}{6}, \quad P(6) = \frac{1}{6}, \quad P(3) = \frac{1}{6}$$

$$P(1 \text{ or } 6 \text{ or } 3) = P(1) + P(6) + P(3) = \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{3}{6} = \frac{1}{2}$$

PROBLEM STATEMENT : Non-MUTUAL Exclusive Event

Q1 Bag of Marbles : 10 Red, 6 Green, 3 (R & G)

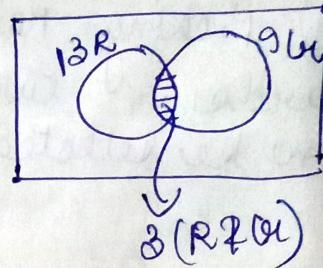
When picking randomly from a bag of marbles what is the probability of choosing a marble i.e. Red or green

$$\text{Ans: } P(R \text{ or } G) = P(R) + P(G) - P(R \text{ and } G) = \frac{10}{19} + \frac{6}{19} - \frac{3}{19} = \frac{13}{19}$$

$$= \frac{13}{19} + \frac{9}{19} - \frac{3}{19} = \frac{19}{19} = 1.$$

Q2 Deck of Cards

What is the probability of choosing heart ~~and~~ Queen



$$\text{Ans: } P(H \text{ or } Q) = P(H) + P(Q) - P(H \text{ and } Q) = \frac{13}{52} + \frac{4}{52} - \frac{1}{52} = \frac{13}{52} + \frac{3}{52} = \frac{16}{52}$$

MULTIPLICATION RULE

(i) Dependent Events : Two events are dependent if they effect one another.

e.g. Bag of Marbles { White - 4 }
Yellow - 3 }

$$P(W) = \frac{4}{7} \quad \rightarrow \quad P(Y) = \frac{3}{6}$$

↑ 1 white marble taken out

here white marble affecting yellow marble's probability of occurring

Q1. what is the probability of rolling a "5" and then a "3" with a normal 6 size sided dice?

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

[Independent Events]

* Independent Events Eg:- Tossing a coin

$$P(A \text{ and } B) = P(A) * P(B)$$

Q2. There are bag of marbles (4 orange, 3 yellow)
what is the probability of drawing a "orange" and then drawing a "yellow" marble from the bag [Dependent Event]

Ans:- Total no. of marble = $4 + 3 = 7$

Probability of orange, $P(O) = \frac{4}{7}$

After taking out 1 marble orange.

then probability of taking out yellow marble $P(Y) = \frac{3}{6} = \frac{1}{2}$

[Conditional Probability]

$$\therefore P(O \text{ and } Y) = P(O) * P(Y) = \frac{4}{7} * \frac{1}{2} = \frac{4}{14} = \frac{2}{7}$$

PERMUTATION

* All the possible arrangements (Dairy Milk, Kit Kat, Milky Bar, Smeakers, 5 stars)

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

= 60 ways chocolates can be chosen

* with permutation order matters

$$P_n^r = \frac{m}{(m-r)}$$

n = total no. of objects = 5
 r = no. of selection = 3.

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

COMBINATION

* Repetition will not occur

* unique combination possible

Formula

$${}^n C_r = \frac{m}{r(r-1)(r-2) \dots (r-r)} = \frac{5}{1 \times 2 \times 3 \times 2} = \frac{5 \times 4 \times 3}{1 \times 2 \times 3 \times 2}$$

\therefore the repeated one is removed = 10.

\therefore Permutation $>$ Combination.

[value]

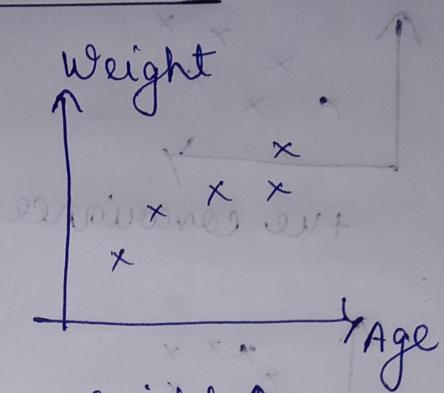
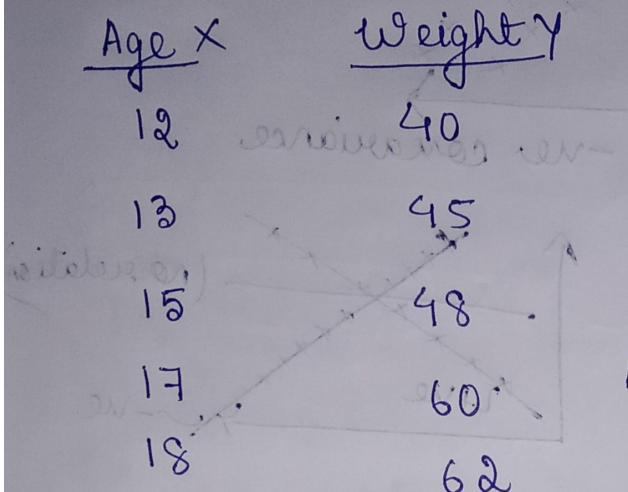
[value]

* Permutation used in

Dream 11

Combination [Eg; choosing 3 deserts from 10 of menu dishes]

COVARIANCE [Feature Selection]



Age ↑ weight ↑

Age ↓ weight ↓

* Quantify the relationship x & y using
Mathematical question

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

$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n-1} = \frac{\sum (x_i - \bar{x}) \times \sum (x_i - \bar{x})}{n-1} = \text{Cov}(x, x)$$

$$\text{Cov}(x, x) = \text{Var}(x)$$

[Interview Question]

Now,

$$\bar{x} = 15, \bar{y} = 51$$

$$\begin{aligned} \sum (x_i - \bar{x})^2 &= (12-15)^2 + (13-15)^2 + (15-15)^2 + (17-15)^2 + (18-15)^2 \\ &= (-3)^2 + (-2)^2 + 0 + (2)^2 + (3)^2 \end{aligned}$$

+ve Covariance

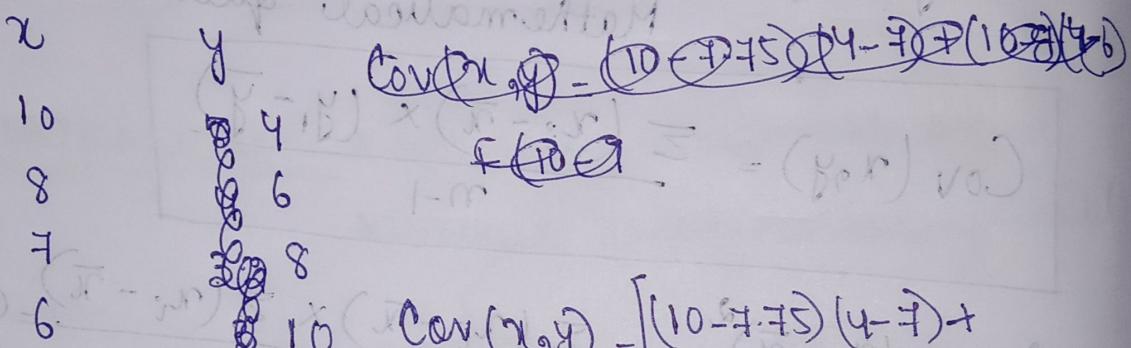
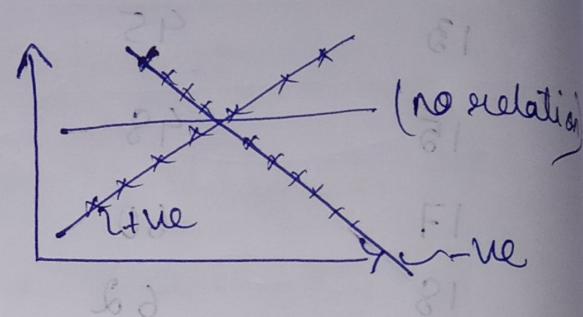
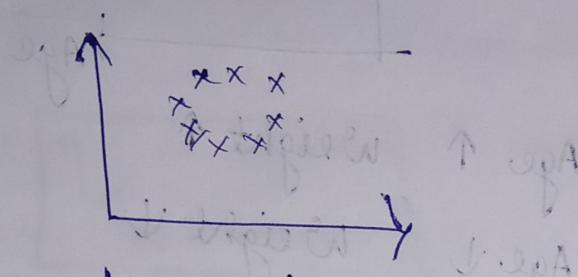
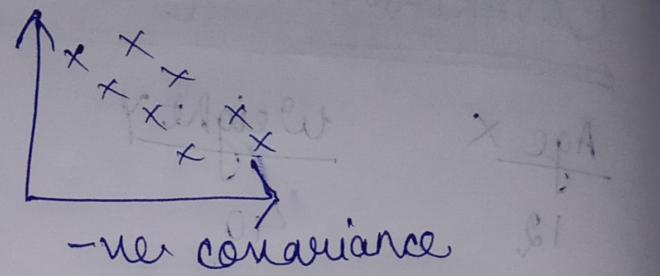
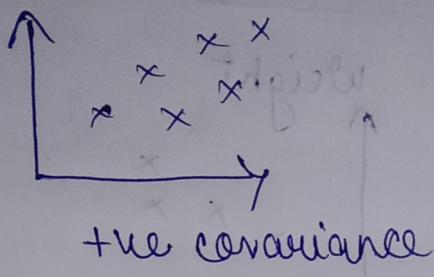
x ↑ y ↑
x ↓ y ↓

-ve Covariance

x ↑ y ↓

Covariance 0

[No relation with x & y]



$$\bar{x} = 7.75 \quad \bar{y} = 7.5$$

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

$$\text{Cov}(x, y) = \frac{1}{3} [(-10 - 7.75)(4 - 7.5) + (-8 - 7.75)(6 - 7.5) + (-7 - 7.75)(8 - 7.5) + (-6 - 7.75)(10 - 7.5)] = -4.33$$

Pearson Correlation Coefficient [-1 to 1]

$$r(x, y) = \frac{\text{Cov}(x, y)}{\sqrt{s_x^2 s_y^2}}$$

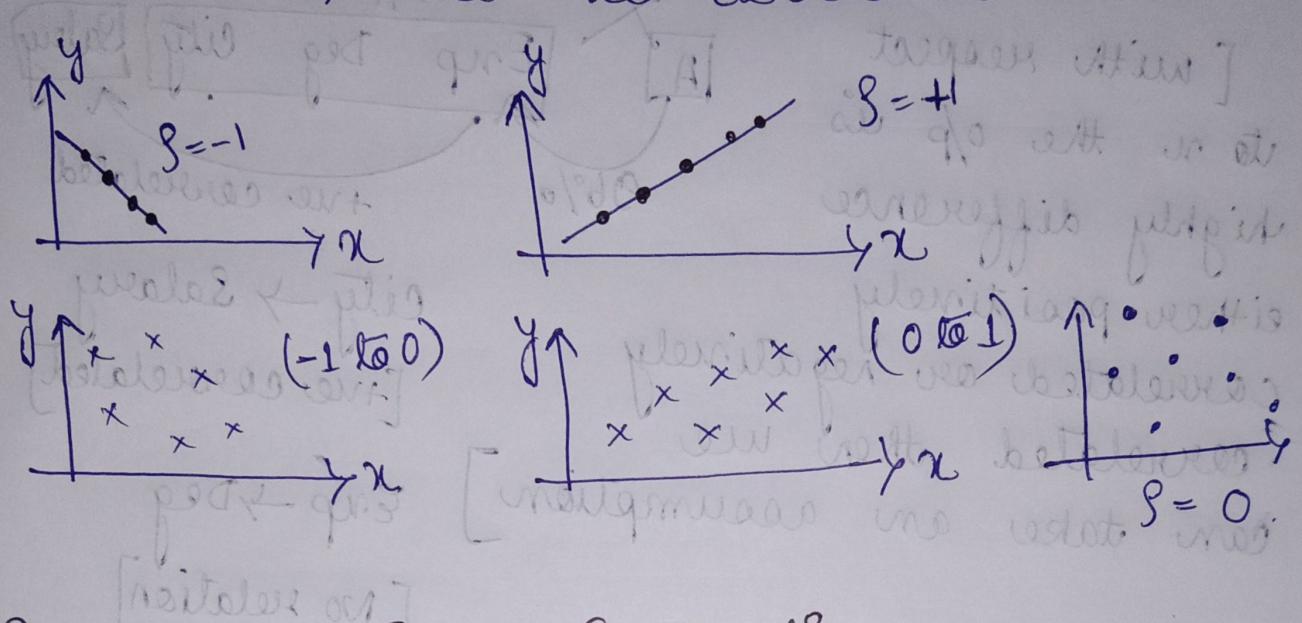
* with the help of this we are trying to restrict the value between -1 to 1.

Move the value towards +1

→ More +ve correlated it is

And, Move the * value towards -1

→ More -ve correlated it is

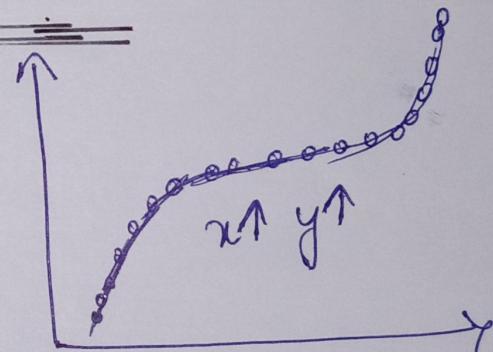


Spearman Rank Correlation

For non-linear data
we use Spearman Rank
Correlation

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

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



Spearman Correlation = 1

Pearson Correlation = 0.88

* Rank = [Assigning value by no. in ascending order]