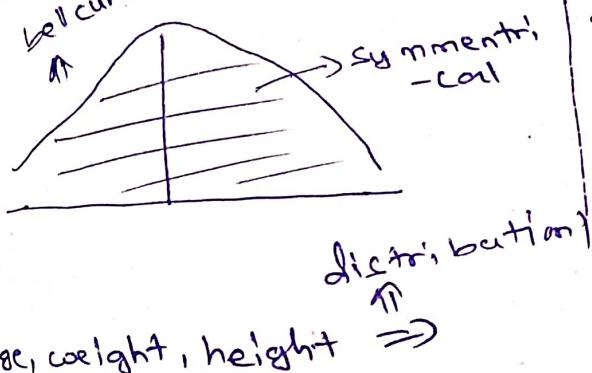


Day - 3

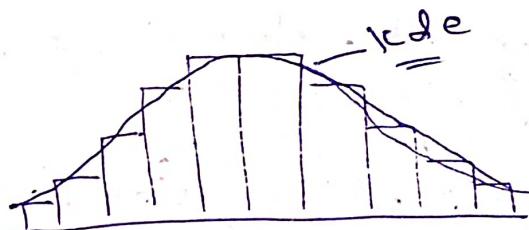
Statistics

- ① Normal Distribution
- ② Standard Normal distribution
- ③ Z-score
- ④ Standardization and Normalization

① Gaussian / Normal distribution is a probability distribution that is symmetric about mean, showing the data near the mean are more frequent in occurrence than data far from the mean.



In graphical form, the normal distribution appears as a "bell curve".



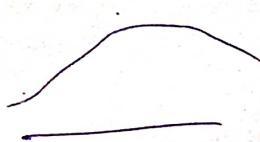
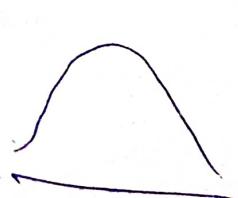
Example

IRIS DATASET

Petal length, Sepal width, Petal width, Sepal length

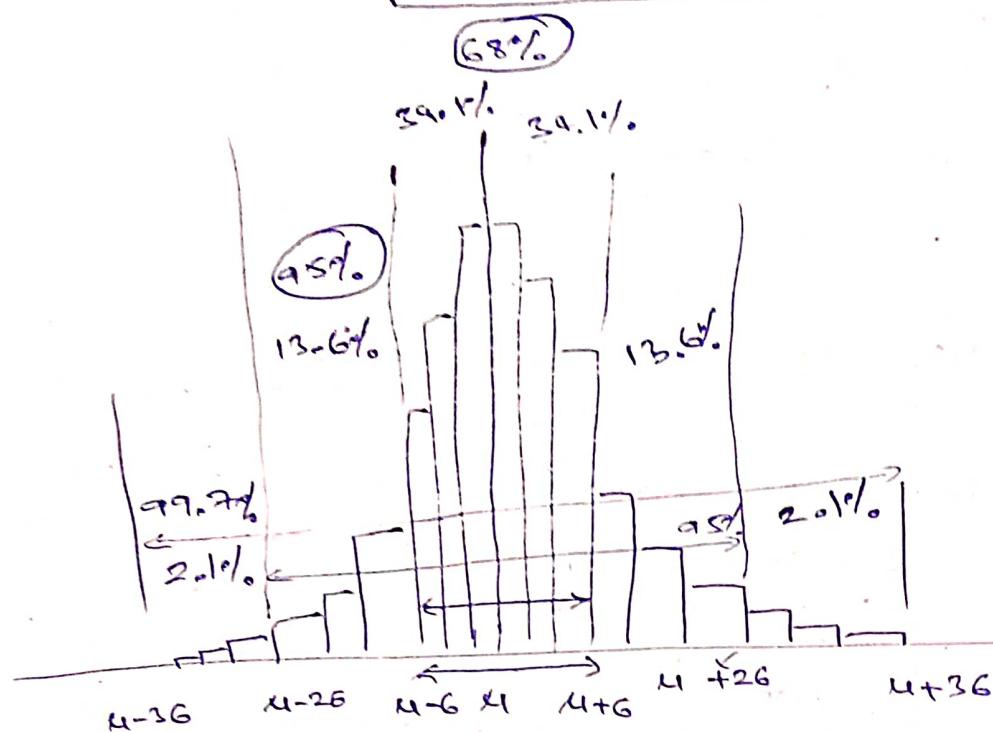


↑
Gaussian distribution



* Empirical Rule of Normal distribution

Empirical rule of $68-95-99.7\%$



Standard Normal Distribution

$X \approx$ Gaussian Distribution (M, σ)

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

$$\text{z-score} = \frac{x_i - M}{\sigma}$$

$$M = 3$$

$$\sigma = 1.41$$

$y \approx$ SND ($M=0, \sigma=1$)

$$\text{z-score} = \frac{x_i - M}{\sqrt{\frac{\sigma^2}{n}}} = \frac{x_i - M}{\sqrt{\frac{1}{n}}} = \frac{x_i - M}{\sqrt{n}}$$

\Rightarrow Standard Error

\Rightarrow Inferential stats

$$\text{z-score} = \frac{x_i - M}{\sigma} \in \text{Simple}$$

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

$$\Rightarrow \frac{1-3}{1.414} = -1.414$$

$$N=3, \sigma = 1.414$$

$$\Rightarrow \frac{2-3}{1.414} = -0.707$$

$$y = \{-1.414, -0.707, 0, 0.707, 1.414\}$$

$$\frac{3-3}{1.414} = 0$$

$$\frac{4-3}{1.414} = 0.707$$

$$\frac{5-3}{1.414} = 1.414$$

$$\overbrace{-\frac{1}{3} \frac{1}{2} \frac{1}{1} \frac{1}{0} \frac{1}{1} \frac{1}{2} \frac{1}{3}}$$

why? [standardization] $\Rightarrow N=0, \sigma=1$

Age (years)	weight (kg)	height (cm)
24	72	150
26	78	160
27	84	165
28	92	170
29	80	150
34	82	180
28	83	175
29	80	170



feature scaling

$$\Downarrow \frac{x-\bar{x}}{\sigma}$$

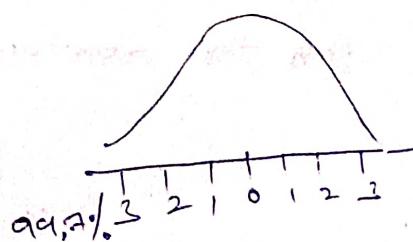
Normalization

standardization $\{z\text{-score}\}$

[0-1]

[0-5]

E3, ③



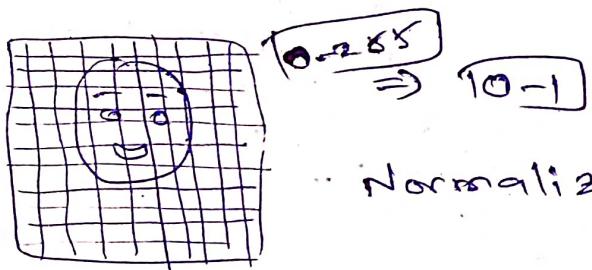
Normalization: [lower scale \leftrightarrow high scale]

① Min Max Scaler [0-1]

$$x_{scaled} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

x	y	$y_1 - \text{standard}$
1	0	-1.414
2	0.25	-0.707
3	0.5	0
4	0.75	0.707
5	1	1.414

* we mainly used in
deep learning, image process
(CNN, ANN, RNN)



Normalization,

Standardization

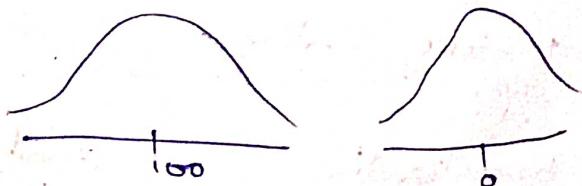
$$\text{z-score} = \frac{x_i - \bar{x}}{\sigma}$$

$x \rightarrow$ Normal distribution (μ, σ)

\Downarrow z-score

$y \rightarrow$ SND ($\mu=0, \sigma=1$)

why do we do this?

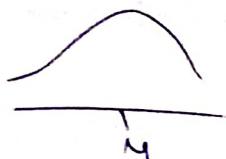


\rightarrow bring the feature in the same scale

Normalization (0-1)

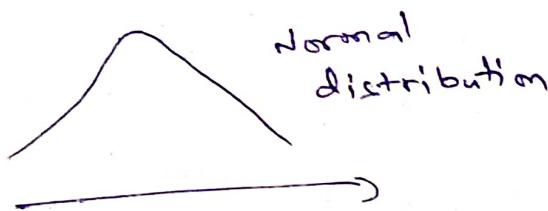
① min max scalar \Leftrightarrow standardization

$$\frac{y - y_{\min}}{y_{\max} - y_{\min}}$$

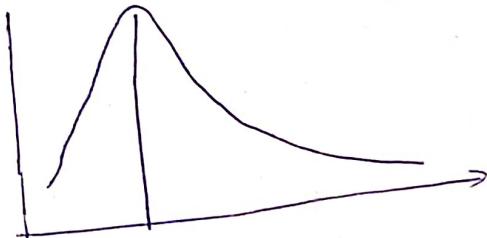


min max scale

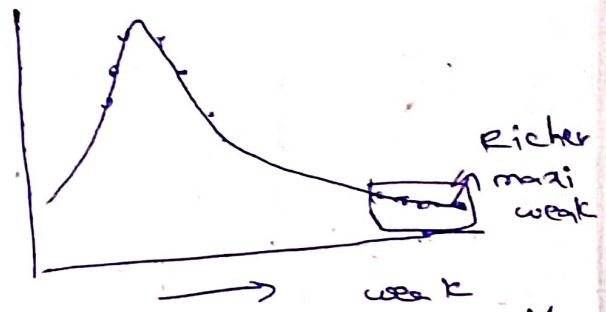
① Log Normal Distribution



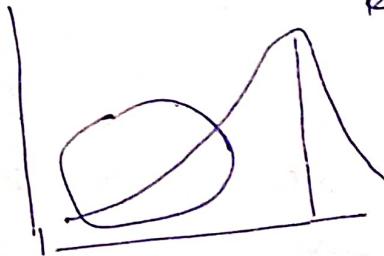
mean will be higher



log normal distribution

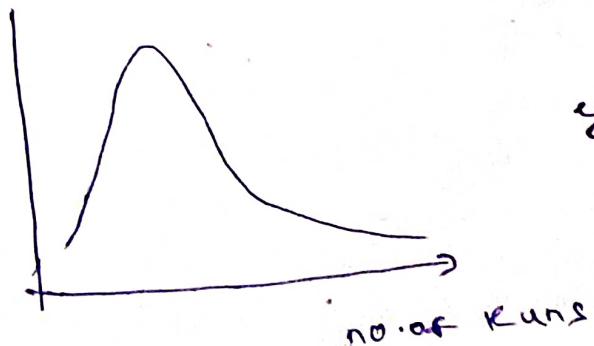


relation of mean median mode

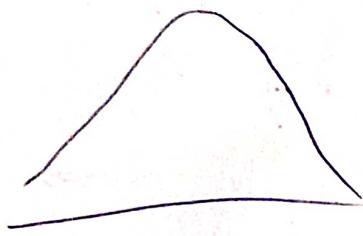


from ascending order given the relationship of
mean median & mode.

* log normal distribution to Normal distribution



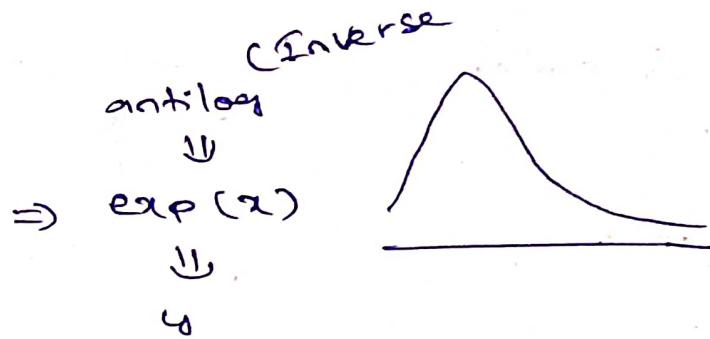
$$y = \ln(x)$$



$x \approx \log$ Normal distribution (μ, σ)



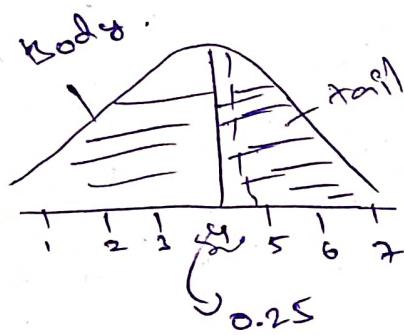
$$x \approx \dots$$



* $x = \{1, 2, 3, 4, 5, 6, 7\}$

$$\mu = 4$$

$$\sigma = 1$$



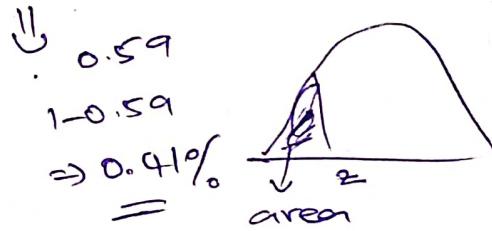
$$z\text{-score} = \frac{x_i - \mu}{\sigma} \Rightarrow \frac{4.25 - 4}{1} = 0.25$$

z-table cover the area

↳ +ve table

↳ -ve table

Question: what is the PCTL score that falls above $\underline{+0.25\%}$ fall below $\underline{3.75\%}$?



$$z = \frac{x_i - \mu}{\sigma} = \frac{3.75 - 4}{1} = -0.25$$

$\underline{40\%}$



$$\Rightarrow \frac{+0.25 - 4}{1} \Rightarrow \underline{0.75} \Rightarrow$$

$$\Rightarrow \frac{5.25 - 4}{1} \Rightarrow$$

* In India the average IQ is 100 with a standard deviation of 15. what is the percentage of population would you expect to have an IQ

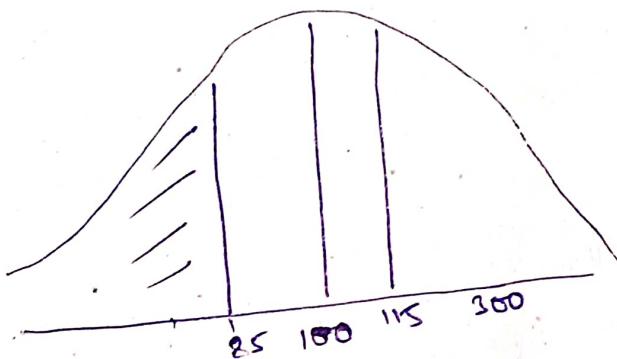
① lower than 85 = 0.1587

② Higher than 85 = 0.8413

③ Between 85 and 100 = 0.8413

$$z\text{ score} = \frac{85-100}{15} = -1$$

$$0.5 - 0.1587$$



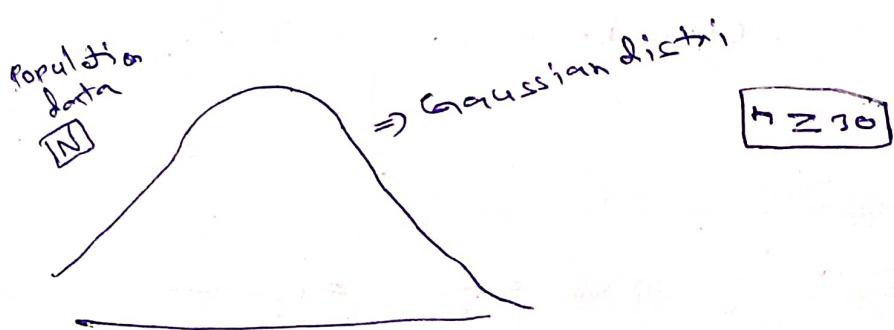
Day - 4

Stats

- 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

Central Limit Theorem

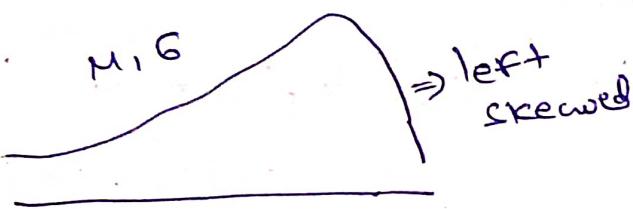
The CLT states that the distribution of sample means approximates a normal distribution as the sample size gets larger, regardless of population's distribution. Sample size equal to or greater than 30 are often considered sufficient for the CLT to hold.



* Population data is Norm dis or logND or left skew we take $n \geq 30$ and n size of sample and m is no. of samples mean we can draw the Gaussian distribution.



$n \geq 30$ → size of sample
 \bar{x} → mean
 \bar{S}_n → No. of sample



$$\rightarrow S_1 \rightarrow \{x_1, x_2, \dots, x_n\} \rightarrow \bar{x}_1 = \bar{x}_1$$

$$\rightarrow S_2 \rightarrow \{x_3, x_4, \dots, x_n\} \rightarrow \bar{x}_2 = \bar{x}_2$$

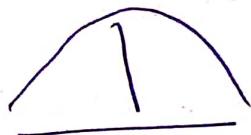
$$\rightarrow S_3 \rightarrow \{x_4, x_5, \dots, x_n\} \rightarrow \bar{x}_3 = \bar{x}_3$$

\vdots

$$\vdots$$

$$\frac{1}{\bar{x}_m}$$

sampling with replacement



≡ Gaussian distribution

size of shark through out the world \Rightarrow A solution

$$n \geq 30$$

* Probability: Probability is a measure of the likelihood of an event

Eg. Tossing a fair coin $P(H) = 0.5$ $P(T) = 0.5$

ii)

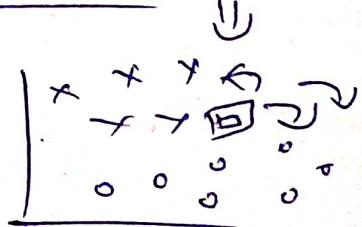
shallow \rightarrow coin

ii)

unfair coin

$$P(H) = 1$$

using probability



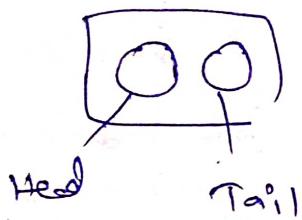
Rolling a Dice $P(1) = \frac{1}{6}$, $P(2) = \frac{1}{6}$, $P(3) = \frac{1}{6}$

① mutually exclusive events:

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

① Tossing coin

② Rolling dice



② Non mutual Exclusive Events

Two events can occur at the same time

- * Picking randomly a card from a deck of cards
- * Two events "heart" and "king" can be selected

Head



Bag of marbles

mutual exclusive events

① what is the probability of coin landing on head or tails

Addition Rule for mutual exclusive events

$$\begin{aligned} P(A \text{ or } B) &= P(A) + P(B) \\ &= \frac{1}{2} + \frac{1}{2} = 1 \end{aligned}$$

$$\frac{1 \times 2 + 2 \times 1}{2+2} = \frac{4}{4} = 1$$

② what is the probability of getting 1 or 6 or 3 while rolling a dice?

$$\begin{aligned} P(1 \text{ or } 6 \text{ or } 3) &= P(1) + P(6) + P(3) \\ &= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \underline{\underline{\frac{1}{2}}} \end{aligned}$$

Non Mutual exclusive events

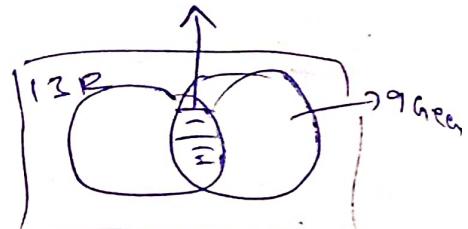
Bag of marbles \therefore 10 red, 6 green, 1 (R or G)

- * when picking randomly from a bag of marble what is the probability of choosing a marble that is red or green?



Non mutual exclusive.

$P(R \text{ or } G)$



Addition Rule for non mutual exclusive events

events \therefore

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

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

Dice of cards \rightarrow what is the probability of choosing 7 or Queen

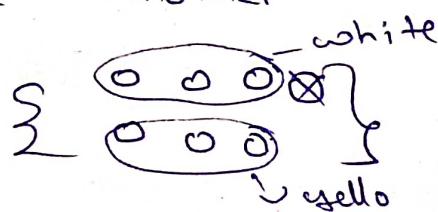
$$P(7 \text{ or Queen}) = P(7) + P(\text{Queen}) - P(7 \text{ and Queen})$$

$$= \frac{12}{52} + \frac{4}{52} - \frac{1}{52} = \frac{16}{52}$$

* multiplication rule

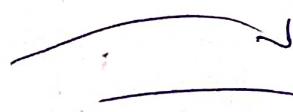
* Dependent Events: Two events are dependent if they affect one another

Bag of marble



$$\Rightarrow P(w) = \frac{4}{7}$$

↑ white marble



$$P(4) = \frac{3}{6}$$

Independent events:

what is the probability of rolling a "5" and then a 9 with a normal 6 six sided dice?

$$P(1) = \frac{1}{6} \quad P(2) = \frac{1}{6} \quad P(3) = \frac{1}{6} \quad P(4) = \frac{1}{6}$$

multiplication Rule for Independent Events

$$\begin{aligned} P(A \text{ and } B) &= P(A) * P(B) \\ &= \frac{1}{6} * \frac{1}{6} = \frac{1}{36} \end{aligned}$$

$$P(A \text{ or } B) = \begin{cases} \text{mutual Events} \\ \text{Non mutual Events} \end{cases}$$

$$\begin{aligned} P(A \text{ or } B) &= P(A) + P(B) - P(A \text{ and } B) \rightarrow \text{Non mutual} \\ P(A \text{ or } B) &= P(A) + P(B) \text{ [mutually exclusive]} \end{aligned}$$

Dependent and Independent Events

Events A and B

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

Tossing a coin

$$P(H) = 0.5 \quad P(T) = 0.5$$

②



Dependent Event

Probability of drawing a "orange" and then drawing a "yellow" marble from the bag?

$$\begin{aligned} \boxed{\text{oooo}} \quad P(O) = \frac{4}{7} &\rightarrow \boxed{\text{o ooo}} \rightarrow P(Y/O) \rightarrow \text{conditional probability} \\ \downarrow \quad \text{orange} \\ \text{marble} \end{aligned}$$

$$\therefore P\left(\frac{3}{6}\right) = \frac{1}{2}$$

$$P(O \text{ and } Y) = P(O) * P(Y/O)$$

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

* Permutation :-

The term permutation refers to a mathematical calculation of the number of ways a particular set can be arranged.

* ~~Unique~~ combination

$n P_r \Rightarrow$ Permutation

$n =$ Total no. of objects

$r =$ no. of objects taken at a time

$$P(n,r) = {}^n P_r = \frac{n!}{(n-r)!}$$

* Possible arrangements

Example :- How many ways you can arrange the letters in the word MATH?

$$n=4 \quad r=1$$

$$\text{Sol} \quad n P_r = \frac{n!}{r!(n-r)!} = \frac{4!}{1!(4-1)!} = \frac{24}{6} = 4$$

$$\frac{3!}{0!(3-1)!} = \frac{6}{2} = 3 \quad \frac{2!}{1!(2-1)!} = 2 \quad \frac{1!}{0!(1-1)!} = 1 = 4 \times 3 \times 2 = 1 \\ = 24 \text{ ways}$$

Combination :-

* Unique elements

* Repeat will not occur

$$n C_r = \frac{n!}{r!(n-r)!} = \frac{5!}{2!(3!)!} = \frac{5 \times 4 \times 3 \times 2 \times 1}{2 \times 1 \times 2 \times 1} \\ = 10.$$

IX
12
4
2
1
X
23
X

Covariance

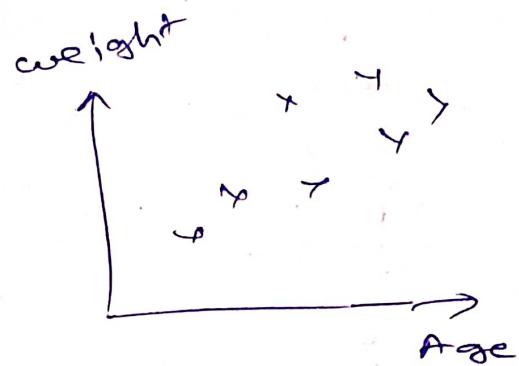
In statistics and mathematics, covariance is a measure of the relationship between two random variables.

X	Y
age	weight
12	40
13	45
14	48
15	60
16	62
17	66
$\bar{x} = 15$	$\bar{y} = \underline{\quad}$



Quantity the relationship

x & y using mathematical questions



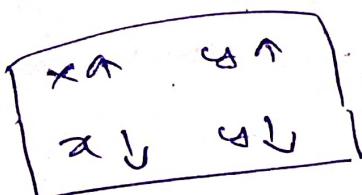
$$\text{Cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

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

$$\sigma^2 = \frac{\sum (x_i - \bar{x})(z_i - \bar{z})^2}{n-1}$$

$\text{Cov}(x, x)$ equal

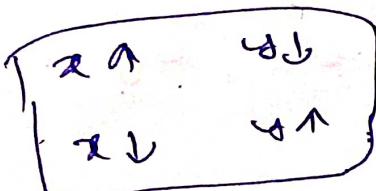
$$\text{Cov}(x, x) = \text{Var}(x)$$

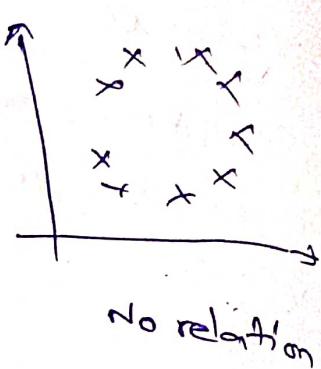
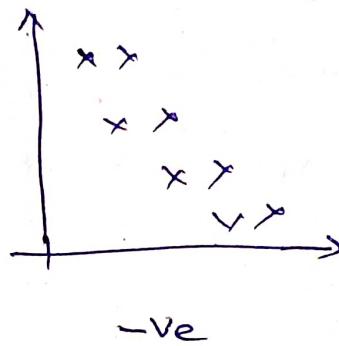


+ve Covariance

Covariance = 0 [No relationship
ship x & y]

-ve Covariance

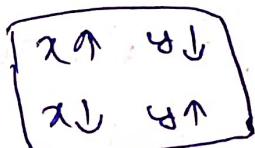




X	Y
10	4
8	6
7	8
6	10
<hr/>	
7.75	7

$$\text{cov}(x, y) = -\text{ve}$$

$$\begin{aligned}
 &= [(10-7.75)(4-7) + (8-7.75)(6-7) \\
 &\quad + (7-7.75)(8-7) + (6-7.75)(10-7)] \\
 &\Rightarrow -3.25
 \end{aligned}$$



* Person correlation coefficient (-1 to 1)

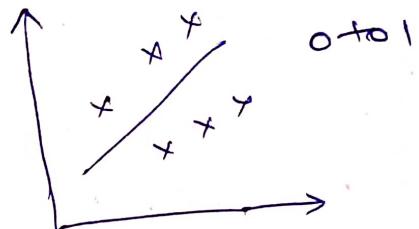
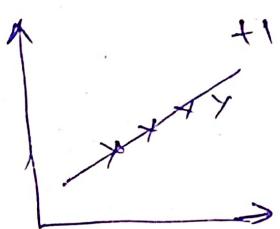
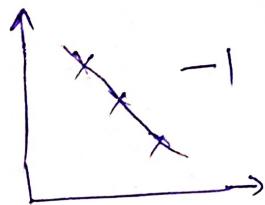
is the most common way of measuring a linear correlation. It is a number between -1 to 1 that measures the strength and direction of the relationship between two variables.

$$\rho(x, y) = \frac{\text{cov}(x, y)}{\sigma_x \cdot \sigma_y}$$

scale
+ covar
- covariate



-1 negative



Spearman's rank correlation coefficient
measures the strength and direction of association
between two ranked variables

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

<u>Grade</u>	<u>Ascending Order</u>		<u>Spearman rank correlation</u>
	$R(x)$	$R(y)$	
10	4	1	
8	3	2	
7	2	3	
6	1	4	