

# STATISTICS

①



ATM



NEW  
ATM



ATM

→ How do you decide this ATM should open here or not.

② Find the avg size of shark in throughout world.

③ Amazon Big Billion Day Sale (which month?)  
{Intuit?}

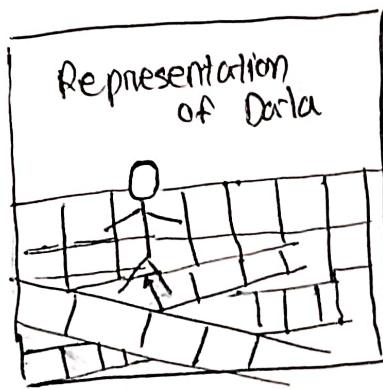


Life cycle of Data science project



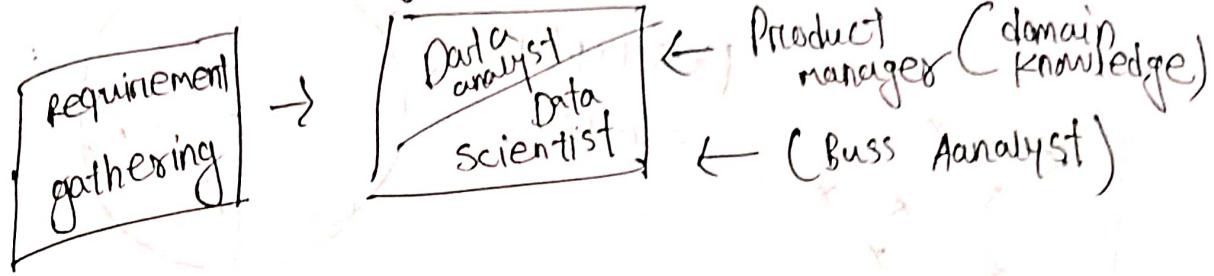
↓  
product manager

"Data is more useful when we arrange data in a perfect manner."

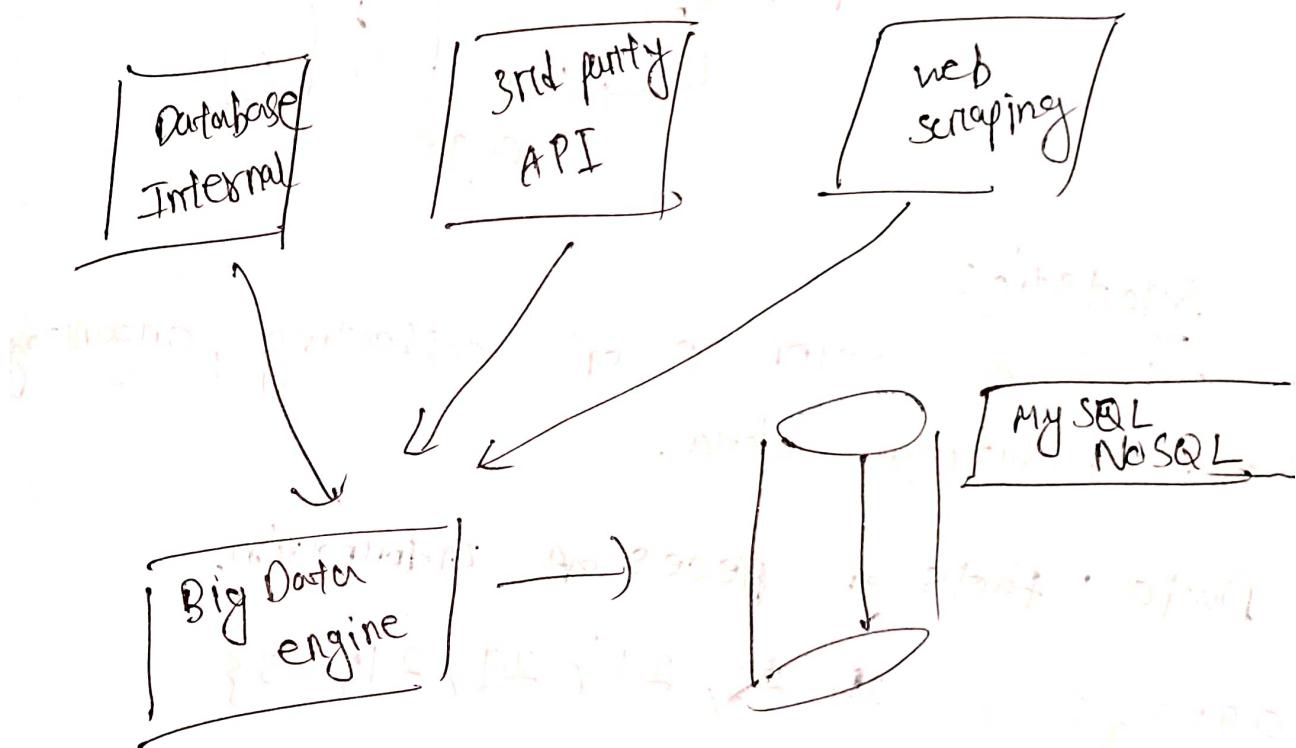


(use of ladder)

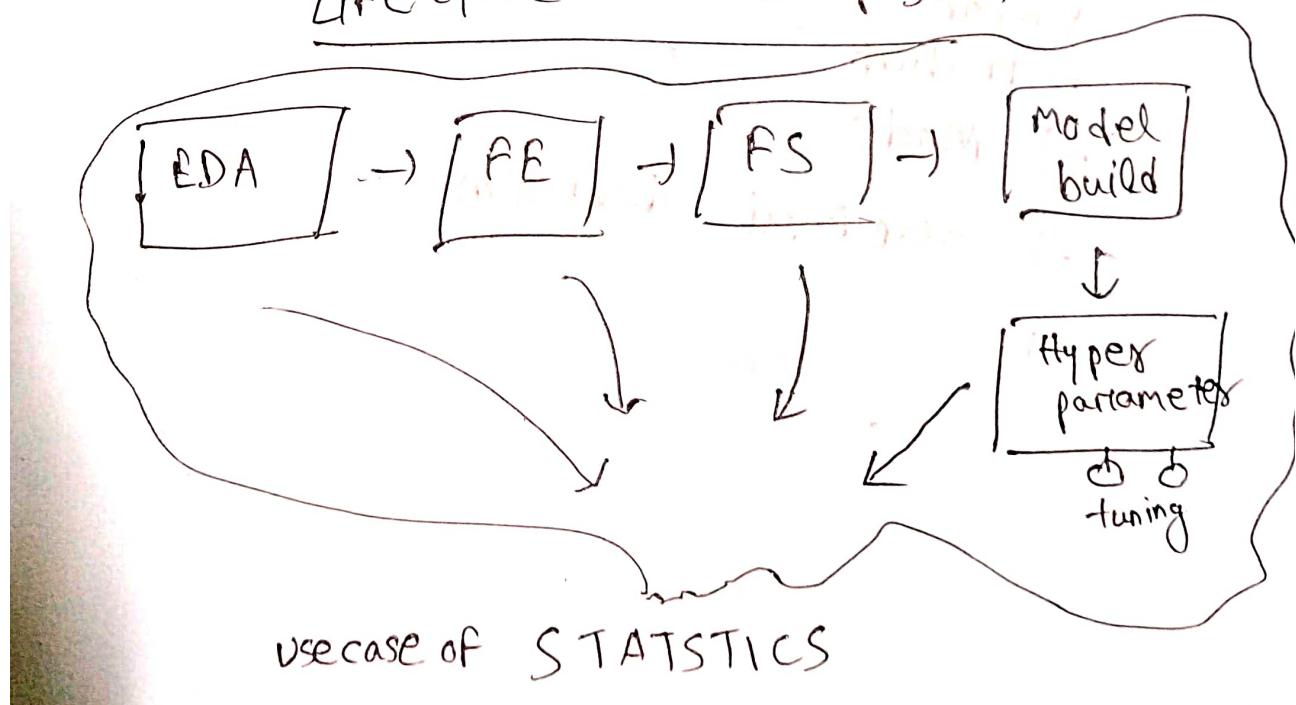
## Q. what companies do ?



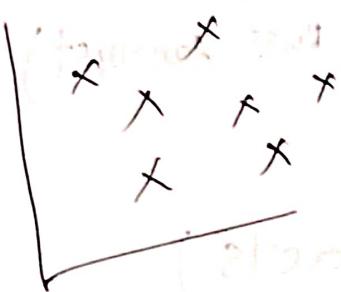
## B Data 2. (collection of facts)



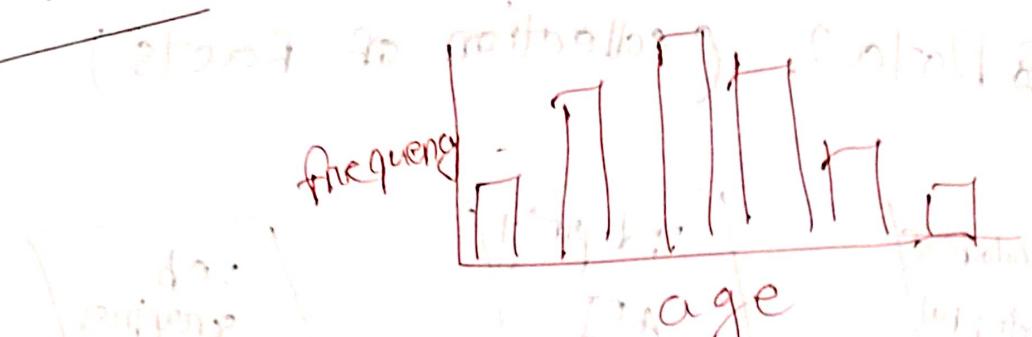
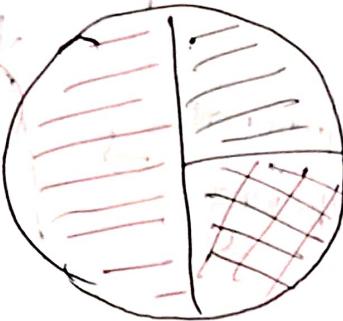
## Life cycle of DS project



## Analysis of Data



Descriptive  
starts



## Statistics

It is a science of collecting, organizing and analyzing data.

Data: facts or pieces of information.

age = {21, 22, 23, 21, 24, 21, 23}

→ mean

→ median

→ mode

→ standard deviation

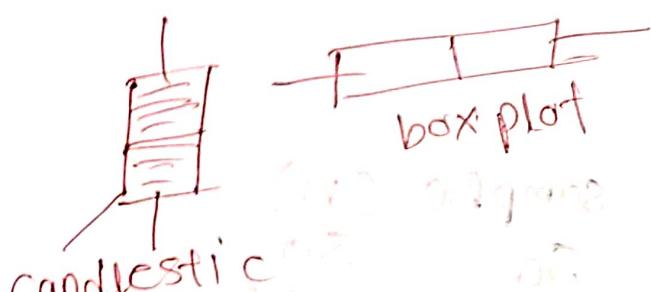
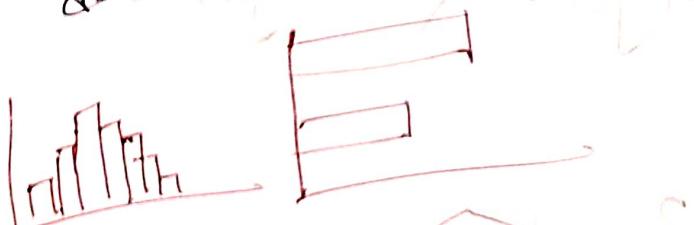
deviation

Avg

Starts

## Descriptive

- ① Organizing Summarizing  
data.



candlestick



Scatter

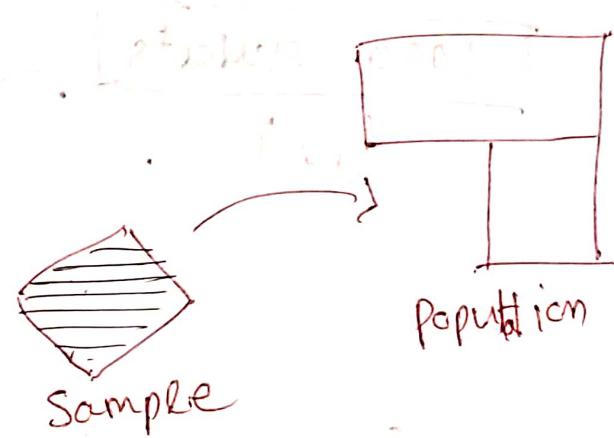
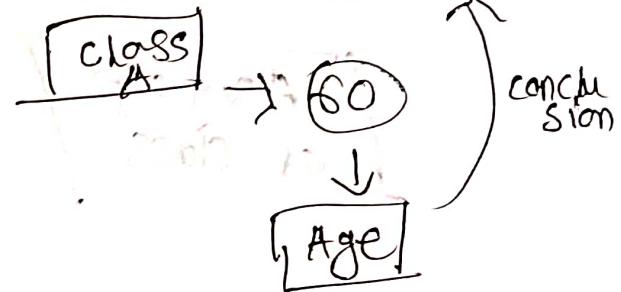
Sample Data

## Inferential

- ① It consist of collecting sample data and making conclusion about population.

## Hypothesis testing

University (500)



Population

Sample

Population Data

vs

Sample Data



Sample  
1000

Population  
100

Eg. Age = { 21, 20, 18, 24, 28, 19, 23, 26, 22 }

weight → {

### Descriptive Stats

Population size often very large

Relationship?

Avg → 2.

age

Age & weight

### Inferential Stats

Avg age  
of class

2.

Avg age

University

1000 students  
(N)

sample (n)

50

boys

girls

96%

### Sampling techniques

#### ① Simple

→ Every has equal chance for going

Random

member

of sample

(n)

#### Sampling

of other population (N)

of being selected

random

N=8

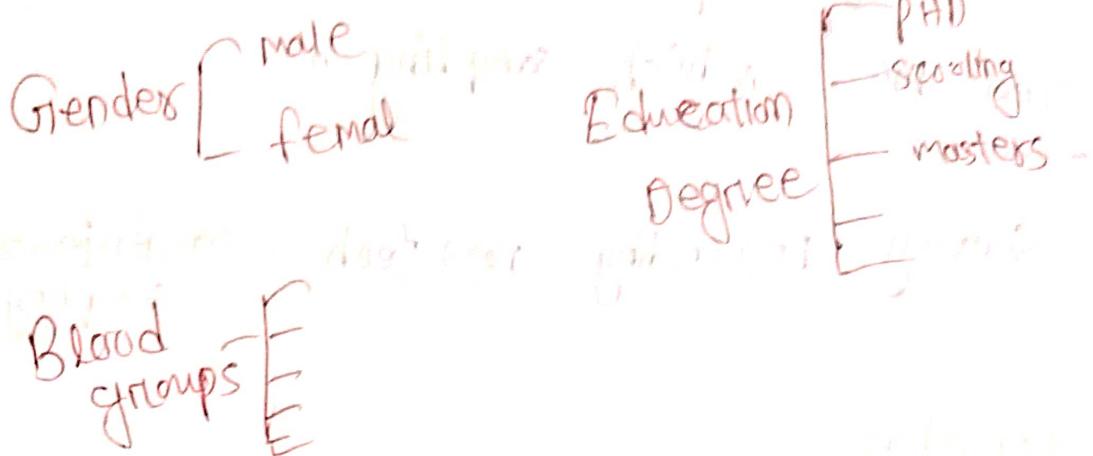
0 0 0 0  
0 0 0 0

n=3

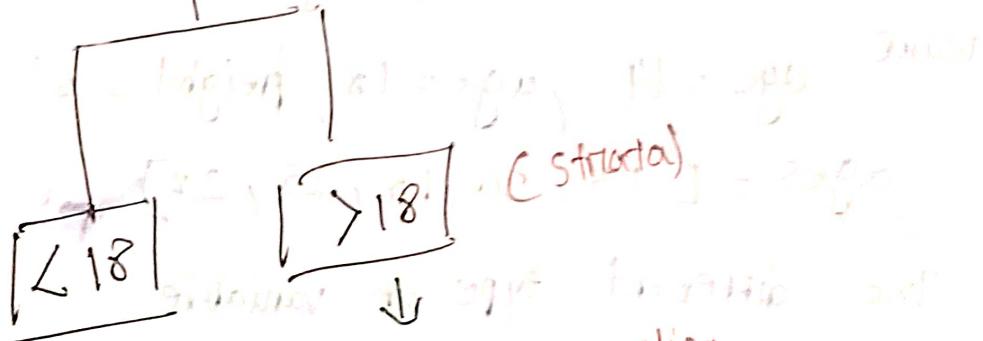


② Stratified Sampling

strata  $\rightarrow$  groups  $\rightarrow$  layers  $\rightarrow$  categories

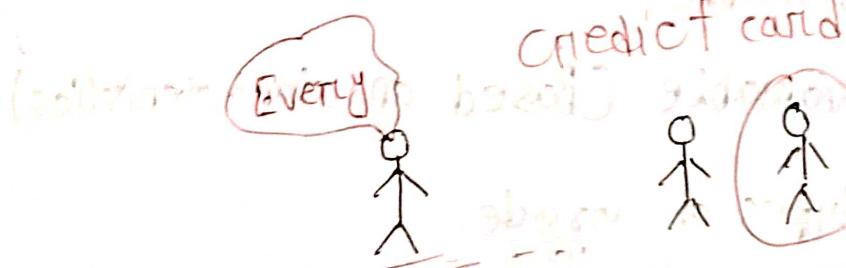


to find in population



random sampling.

③ Systematic Sampling:-



every nth person of a population

population

framing

selection

random sample

systematic sample

convenience sample

④ Convenience Sampling  
Only those who are interested in the survey will only participate.

Tell me which sampling?

→ Survey regarding new tech (convenience sampling).

① Variable  
variable is a property that can hold any value. age = 14, age = 18, height = 5' ages = [13, 20, 18, 25, 23]

Two different type of variable

① Quantitative Variable →

Mathematical variable  
e.g - age, height, weight, temperature.

② Qualitative Variable.

Categorical variable (Based on characteristics)

e.g Gender, Types of Grade.

Quantitative Variable

Discrete variable

e.g whole number

Number of children, I.Q

Continuous variables

weight, Height, speed

# STATISTICS

## Agenda

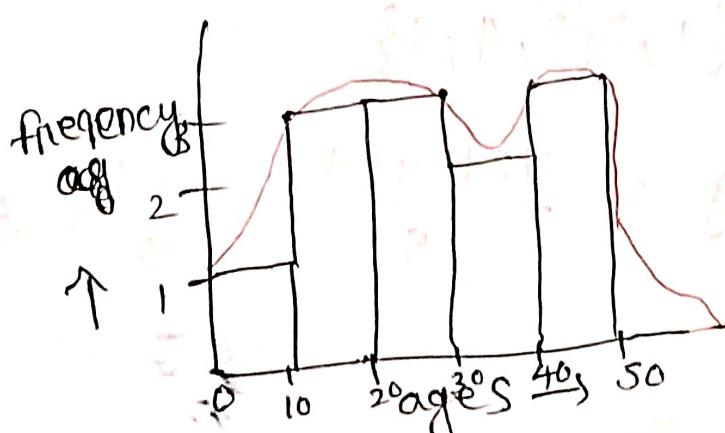
- ① Histograms
- ② Measure of central tendency
- ③ Measure of Dispersion
- ④ Percentile & Quartiles.
- ⑤ 5 Number Summary (Box plot)

## Histograms

$$\text{Ages} = \{10, 12, 14, 18, 24, 26, 30, 35, 40, 42, 43, 50\}$$

### steps

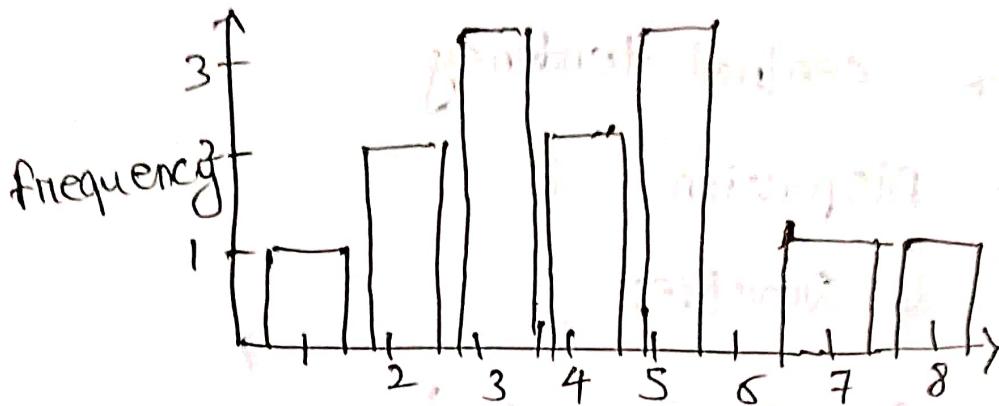
- ① sort the numbers
- ② Bins  $\rightarrow$  No. of groups
- ③ Bin size  $\rightarrow$  size of bins.  $\text{bins} = 10$   
 $\text{bin size} = \frac{\text{max-min}}{\text{bins}} = \frac{50}{10} = 5$



To smoothen this histogram we need probability density function (pdf)

## ② Discrete continuous

NO. of Insta account = [2, 3, 5, 1, 4, 15, 3, 7, 18  
3, 2, 4, 15]



To smoothen this we use probability mass function (PMF)

## ② Measures of central tendency (CT)

- ① mean
- ② Median
- ③ Mode

A measure of CT is a single value attempt to describe a set of data identifying central position.

Mean:

values = {1, 2, 3, 4, 5}

$$\text{mean/avg} = \frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$$

## Population ( $N$ )

Sample ( $Cn^+$ )

Population mean ( $\mu$ )

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

## Simple Mean ( $\bar{x}$ )

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

## Population age

$$= \{24, 25, 2, 1, 128, 27\}$$

$$\text{M} = \frac{24 + 23 + 21 + 28 + 27}{6}$$

Sample age

$$= \{24, 21, 1, 27\}$$

$$\bar{x} = \frac{24+2+1+27}{4}$$

$$= \frac{54}{4} = 13.5$$

## Practical application

Suppose you have a dataset

We can replace central tendency.

NAN with the

## ② Median

$\{1, 2, 3, 4, 5\}$   $\{1, 2, 3, 4, 5, \textcircled{100}\}$

$$\bar{x} = 3 \longrightarrow \bar{x} = 19.16$$

by adding one big value mean changed drastically.

These issues sol<sup>n</sup> is median

Steps to find out median

① sort the numbers

② Find the central number {if no. of element odd,  
central element}

if, even avg of central elements.

③ Mode [most frequent occurring elements]

e.g. Weather reports in a day over the time.

weather = [cloudy, sunny, sunny, cloudy, sunny, sunny]

we can say mostly cloudy way sunny.

### ③ Measure of Dispersion

① Variance ( $\sigma^2$ )

② Standard deviation ( $\sigma$ )

#### Variance

population

variance ( $\sigma^2$ )

sample

variance ( $s^2$ )

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

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

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

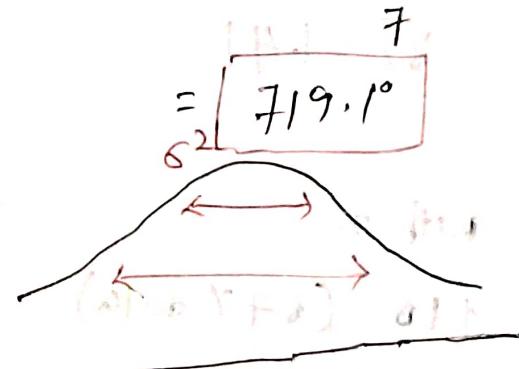
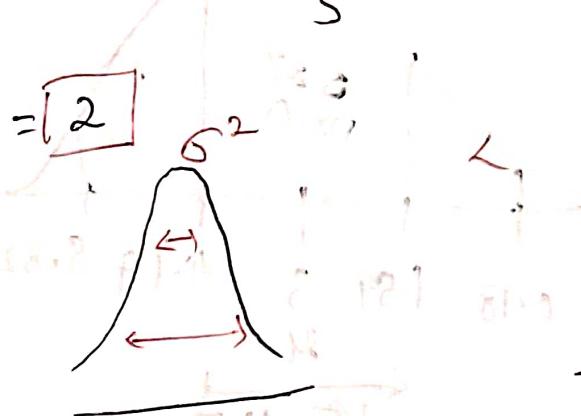
$$M = \frac{15}{5} = 3$$

$$\sigma^2 = \frac{(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2}{5}$$

$$\{1, 2, 3, 4, 5, 6, 80\}$$

$$M = \frac{101}{7} = 14.4$$

$$\sigma^2 = \frac{(1-14.4)^2 + (2-14.4)^2 + \dots + (80-14.4)^2}{7}$$



NOTE:  $(n-1)$  degree of freedom

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

sample variance is an estimator for population variance.

When we applied to sample data, the population variance is a biased estimator for sample variance. It tends to underestimate the amount of variability.

To solve this underestimation problem, statisticians found out by dividing  $(n-1)$ .

## ② Standard deviation ( $\sqrt{s^2}$ )

This is the square root of variance.

e.g

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

$$\bar{x} = 3, s^2 = 2$$

$$\sqrt{2} = 1.41$$

with in

$\pm 1\sigma$  (67% data)

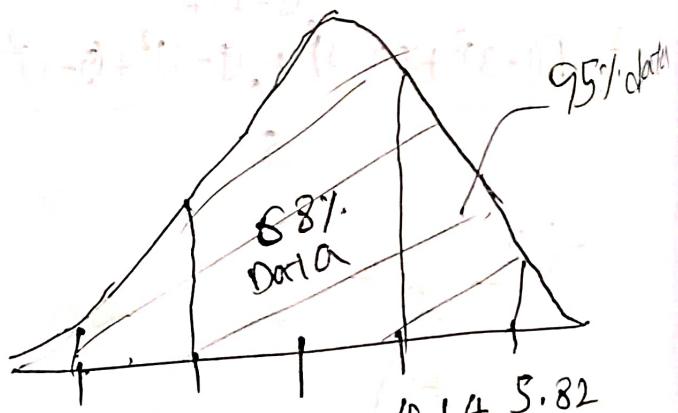
$$0.18 \quad 1.59 \quad 3$$

88%  
Data

$\pm 2\sigma$  (95% data)

$$\xleftarrow{-2\sigma} \xrightarrow{+2\sigma} -1.96 \quad +1.96$$

$\pm 3\sigma$  (99.7% data)



## ④ Percentile & Quartile

percentages = {1, 2, 3, 4, 5, 6, 7, 8} .

percentage of even numbers

$$\text{Ans. } P_0 = \frac{\text{No. of evens}}{\text{Total no. of Numbers}} = \frac{4}{8} = 0.5 = 50\%$$

### Percentile

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

99 percentile = It means the person has got better marks than 99% of the other students.

Dataset = {2, 2, 3, 3, 5, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9, 10}

(10), 11, 11, 12 }

what is the percentile rank of 10?

Percentile rank of  $x = \frac{\text{No. of values below } x}{n} \times 100$

$$= \frac{16}{20} = 80 \text{ percentile.}$$

Percentile rank of 6 =  $\frac{7}{20} = 35 \text{ percentile}$

Q. Find 25 percentile

$$\boxed{\text{Value } (i) = \frac{\text{Percentile} \times n}{100}}$$

$$= \frac{25}{100} \times 20 = 5^{\text{th}} \text{ index}$$

Ans. = 5

Q. Find 30 percentile

$$\text{value } (i) = \frac{30}{100} \times 20 = 6^{\text{th}} \text{ index}$$

Ans = 15

5.6 percentile

$$\text{value } (i) = \frac{56}{100} \times 20 = 11.2$$

Its not a whole number so

$$i = 11.2 + 1 = 12.2$$

and i = whole no. person is 12

Ans = 8

## ⑤ 5 numbers Summary

① minimum

② first Quartile (25 percentile)  
(Q1)

③ median

④ third Quartile (75 percentile)  
(Q3)

⑤ maximum.

} remove  
the  
outliers  
(BOX plot)

$$\{-15, 1, 2, 3, 3, 4, 13, 5, 5, 6, 6, 7, 8, 9, 9, 27\}$$

fence = [lower fence  $\leftrightarrow$  higher fence]

lower fence =  $Q1 - 1.5(IQR)$

IQR =  $(Q3 - Q1)$

Interquartile range

higher fence =  $Q3 + 1.5(IQR)$

$$Q1 = \frac{25}{100} \times 15 = 3, Q3 = \frac{75}{100} \times 15 = 12$$

~~Q3 = 12~~  $\cancel{Q3} \times \cancel{15}$  ~~12~~  $\cancel{12}$

$$IQR = 7 - 3 = 4$$

$$LF = Q1 - 1.5(4)$$

$$= 3 - 6 = -3$$

$$HF = Q3 + 1.5(4)$$

$$= 7 + 6 = 13$$

$$[-3, 13]$$

$$\text{outliers} = -15, 27$$

① minimum = -15

② Q1 = 3

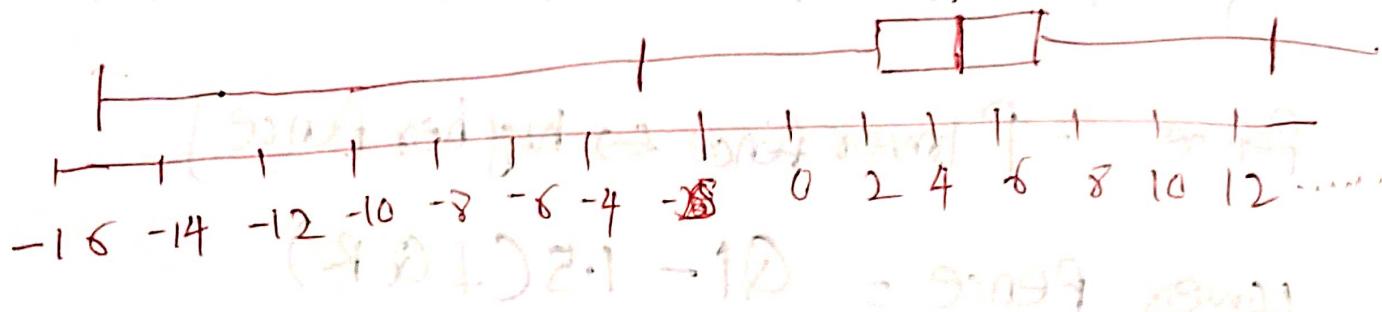
③ median = 5

④ Q3 = 7

⑤ maximum = 27

(101, 101)

FF, P, R, B, E, S, 8, 22, F, 8, E, L, Box Plot



$$(Q1) 2.1 + 8.0 = 10.1 \text{ (approx)}$$

$$(Q3) 8.1 + 8.0 = 16.1 \text{ (approx)}$$

$$(\text{Min}) -15 + 8.0 = -7.0 \text{ (approx)}$$

$$(\text{Max}) 27 + 8.0 = 35.0 \text{ (approx)}$$

$$\text{Median} 5 + 8.0 = 13.0 \text{ (approx)}$$

$$(\text{Outlier}) 10.1 + 8.0 = 18.1 \text{ (approx)}$$

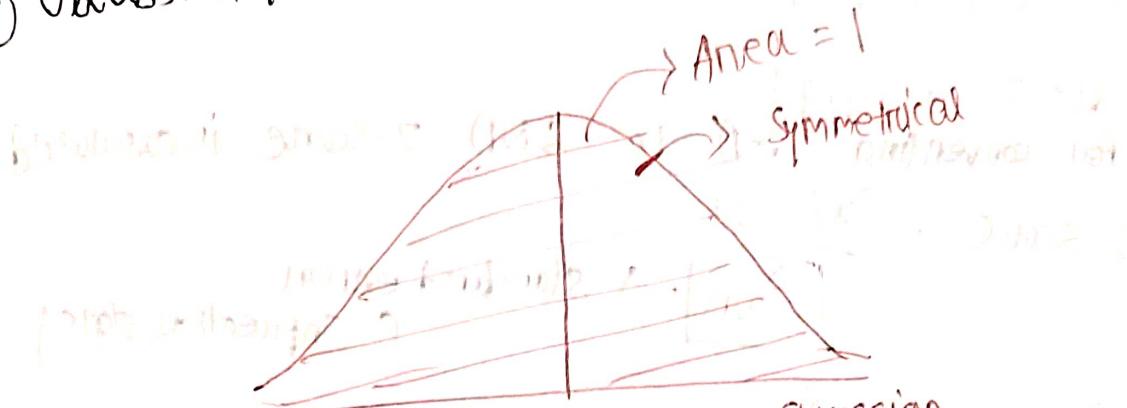
$$(\text{Outlier}) 16.1 + 8.0 = 24.1 \text{ (approx)}$$

$$(\text{Outlier}) 18.1 + 8.0 = 26.1 \text{ (approx)}$$

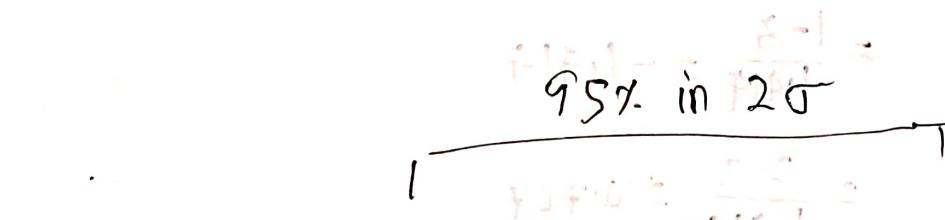
## Agenda

- ① Normal Distribution
- ② Standard Normal Distribution
- ③ Z-score
- ④ Log Normal Distribution

### ① Gaussian / Normal Distribution

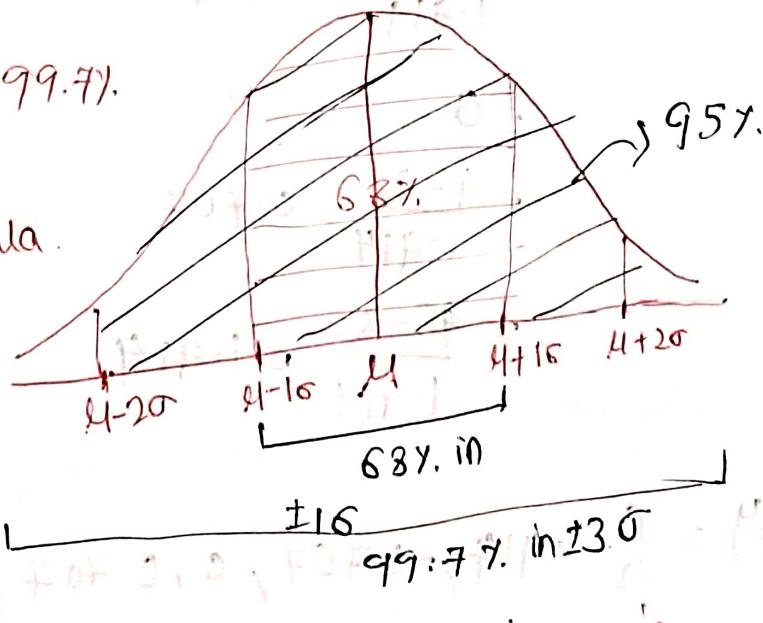


Most of the data follows a symmetrical distribution.



This (68-95-99.7)

is called  
empirical formula.



How to decide a distribution is gaussian or not?

Ans. Using QQ plot

② Standard Normal Distribution ( $\mu=0, \sigma=1$ )

$X \sim$  Gaussian Distribution ( $\mu, \sigma$ )

↓ transform

$Y \sim$  Standard Normal Distribution (SND), ( $\mu=0, \sigma=1$ )

$X = \{1, 2, 3, 4, 5\}$  (data points)

$\mu = 3, \sigma = 1.4$   
for converting  $X$  to SND z-score is calculated.

$$Z \text{ score} = \frac{x_i - \mu}{\sigma/\sqrt{n}} \rightarrow \text{standard error}$$

(Inferential stats)

$$Z \text{ score} = \frac{x_i - \mu}{\sigma} \quad (n=5, \text{ each data point is a sample})$$

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

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

$$= 0$$

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

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

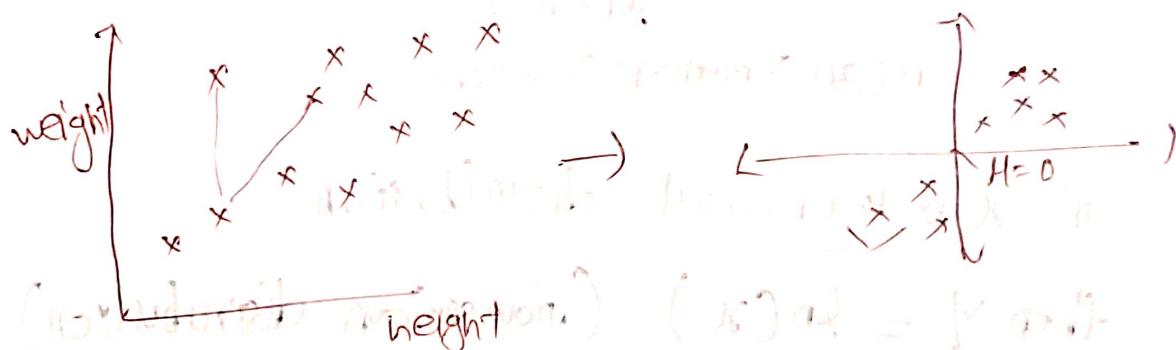
$$Y = \{-1.414, -0.707, 0, 0.707, 1.414\}$$

Why the above things required?

why?

<u>Age (years)</u>	<u>weight (kg)</u>	<u>height (cm)</u>
24	72	150
26	78	160
32	84	165
33	92	170
34	87	150
35	83	180
38	80	175

→ To bring above data to same scale (standardization)



### Normalization

$$[0, 57, 11, 6],$$

$$[3 - 9] \text{ (we provide range)}$$

e.g. min-max scales

$$x_{\text{scaled}} = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}}$$

data will be ranged

from (0, 1)

### Standardization

Data

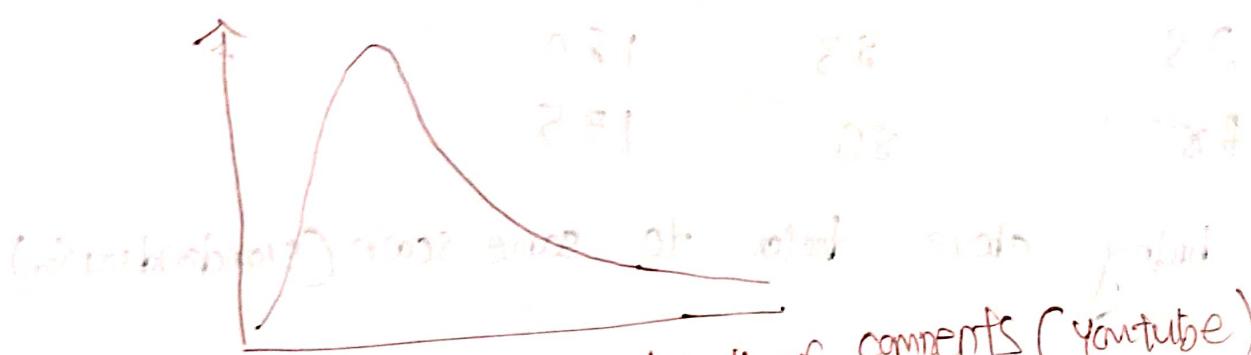
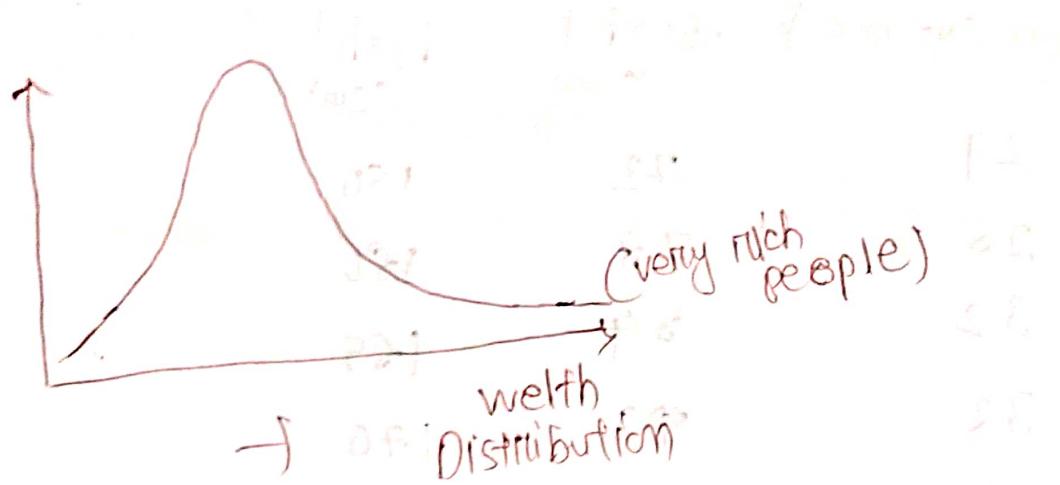
↓ standardize

standard data

$$(H=0, \sigma=1)$$

$$[-3, 3]$$

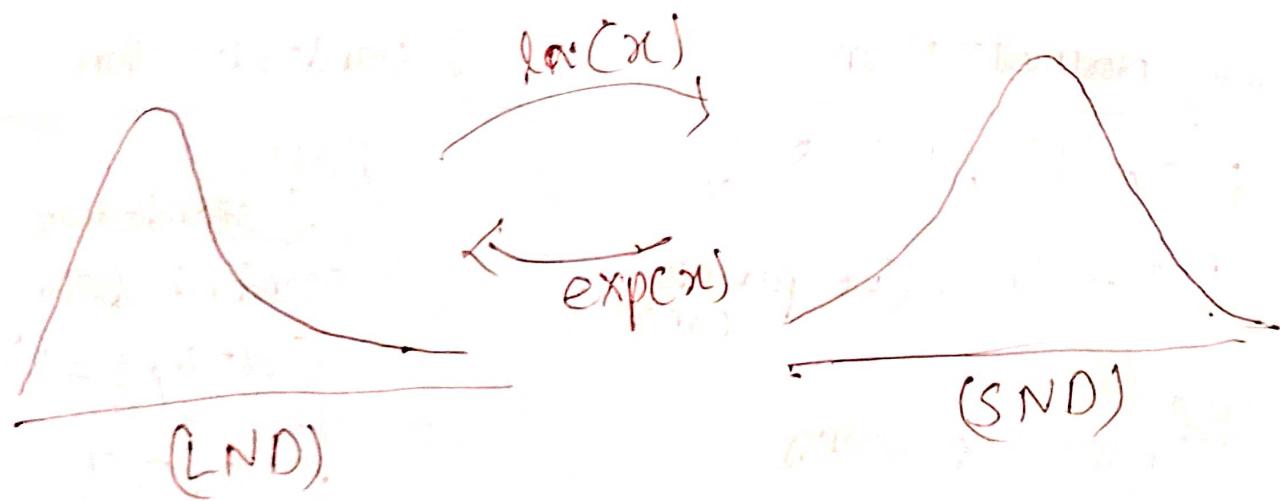
## ④ Log Normal Distribution (LND)



mean > median > mode

if  $x \approx$  log normal distribution

then  $y = \ln(x)$  (Gaussian distribution)



if  $x \sim LN$  then  $\ln(x) \sim N$

if  $x \sim N$  then  $e^x \sim LN$

lognormal vs normal

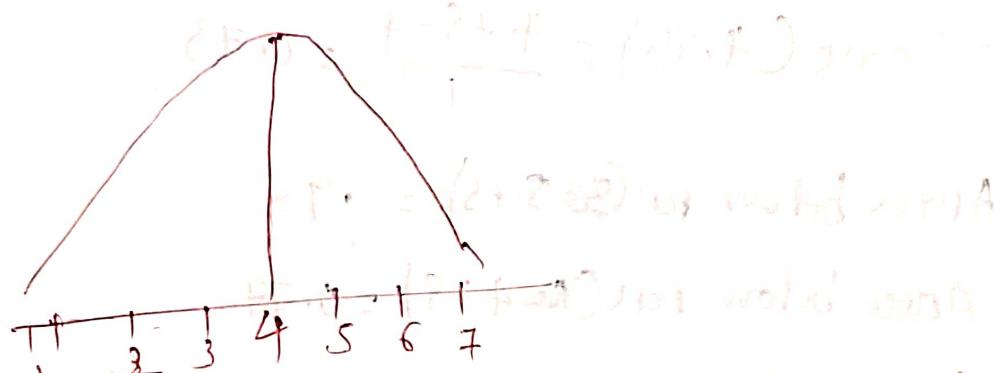
Finding % of area using Z score & Z table

Let

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

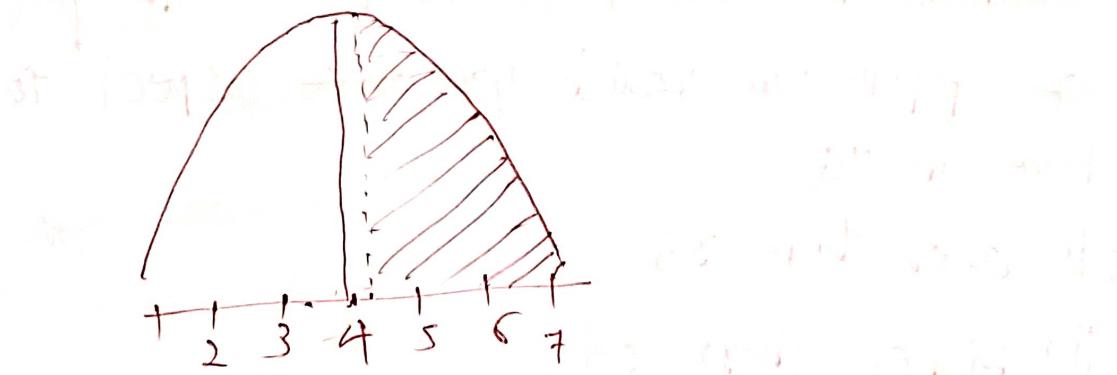
$$\text{let } \mu = 4$$

$$\sigma = \sqrt{\frac{(x_1 - \mu)^2 + (x_2 - \mu)^2 + \dots + (x_n - \mu)^2}{n}}$$



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

find % Area at 0.25 that falls above 4.25.



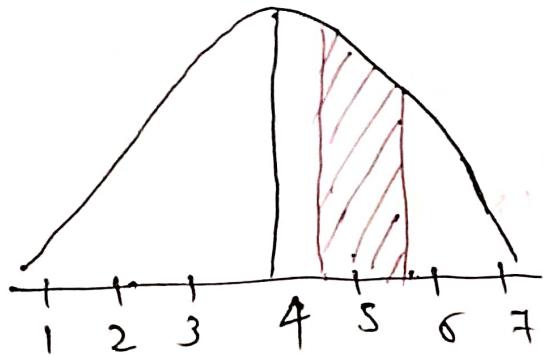
$$\text{z-score} = \frac{4.25 - 4}{1} = 0.25$$

from Z table

$$\text{Area} = 0.59 \text{ (for } z = 0.25 \text{ below)}$$

$$\text{for above} = 1 - 0.59 = 0.41 = 41\%$$

Find area of curve between 4.75 and 5.75



$$Z\text{ score}(5.75) = \frac{5.75 - 4}{1} = 1.75$$

$$Z\text{ score}(4.75) = \frac{4.75 - 4}{1} = 0.75$$

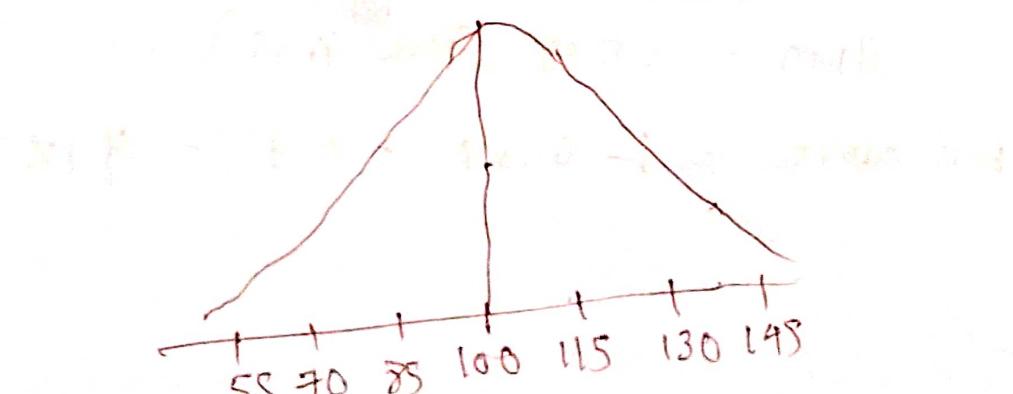
$$\text{Area below } Z(5.75) = 0.96$$

$$\text{Area below } Z(4.75) = 0.77$$

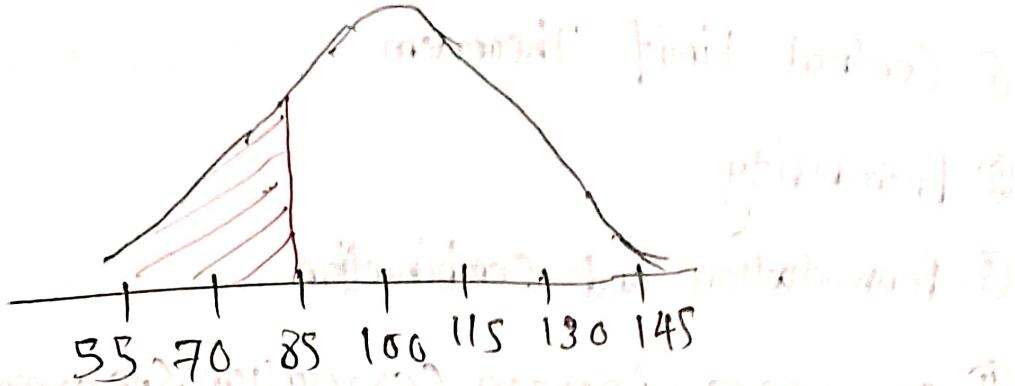
$$\text{So between } 4.75 \& 5.75 = 0.96 - 0.77 = 0.19 \\ \Rightarrow 19\%$$

Q. 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

- (i) Lower than 85
- (ii) Higher than 85
- (iii) Between 85 and 100



c)



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

$$\text{Area} = 0.1587 \Rightarrow 15.87\%$$

(ii) Higher than 85  
=  $1 - 0.1587$   
=  $0.8413 \Rightarrow 84.13\%$

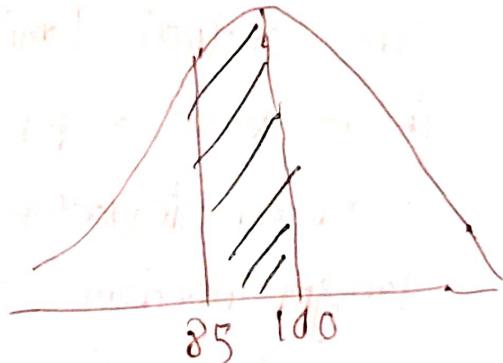
(iii) Between 85 and 100

$$Z\text{ score}(85) = -1$$

$$\text{Area} = 0.1587$$

$$Z\text{ score}(100) = 0$$

$$\text{Area} = 0.5$$



$$\text{Area between } 85 \text{ & } 100 = 0.5 - 0.1587 \\ \Rightarrow 0.3413 \Rightarrow 34.13\%$$

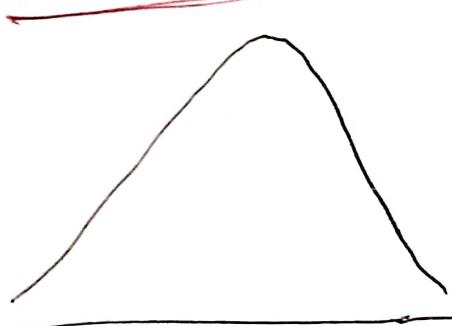
## Agenda

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

## ① Central Limit Theorem

The central limit theorem states that if you have a population with mean  $\mu$  and standard deviation  $\sigma$  and take sufficiently large random samples ( $n \geq 30$ ) from the population with replacement, then the distribution of the means will be approximately Normally Distributed.

Distributed



normal distribution



log Normal distribution

↓ any distribution with mean  $\mu$ , stand. dev  $\sigma$

$$S_1 = \{x_1, x_2, x_3, \dots, x_n\}$$

$$S_2 = \{x_1, x_2, x_3, \dots, x_n\}$$

$$S_3 = \{x_1, x_2, x_3, \dots, x_n\}$$

$$\vdots$$
  
$$S_m = \{x_1, x_2, x_3, \dots, x_n\}$$

The mean of this  $m$  sample will follow gaussian/normal distribution, if total sample size  $> 30$ .

## ② Probability

Probability is the measure of likelihood of an event.

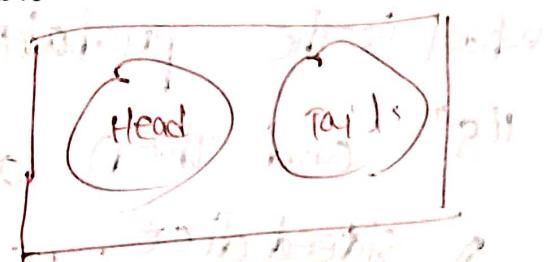
Eg. Tossing a fair coin  $P(H) = 0.5, P(T) = 0.5$

Rolling a dice  $P(1), P(2), \dots, P(6) = \frac{1}{6}$

### (i) Mutual Exclusive event

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

Eg) Tossing a coin



### (ii) Non-mutual exclusive event

Two events occur at the same time.

Eg. A bag of marbles containing red, green, yellow marbles pick two marble randomly.

## Deck of cards

→ what is the probability of choosing heart ♡ on Queen Card ?

$$P(C \cap Q) = P(C) + P(Q) - P(C \text{ and } Q)$$

$$\text{Winning probability} = \frac{13}{52} + \frac{4}{52} - \frac{1}{52} = \frac{18}{52}$$

\* Multiplication Rule

→ Dependent Events: Two events are dependent if they effect one another.

Bag of marshals { o o o }

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

initial

NEXT

2nd effects due to event E1

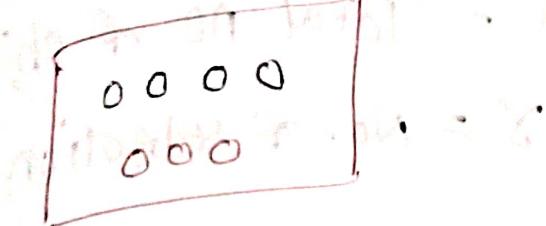
Q. what is the probability of rolling a "5" and then a "3" with a normal 6 sided dice.  $P(3) = \frac{1}{6}$ ,  $P(5) = \frac{1}{6}$ .

Multiplication rule for Independent events

$$DCA \text{ and } B) = PCA * PCB)$$

$$= \frac{1}{6} * \frac{1}{6} = \frac{1}{36}$$

② conditional probability (Naive Bayes is derived)



$$P(O) = \frac{4}{7} \rightarrow P(Y|O)$$

$$\boxed{\begin{array}{c} 000 \\ 000 \end{array}} \rightarrow \frac{3}{6} = \frac{1}{2}$$

$$P(O \text{ and } Y) = P(O) \times P(Y|O)$$

$$= \frac{4}{7} \times \frac{1}{2} = \frac{2}{7}$$

Permutation

for a school children

5 Dairy milk, Kitkat, Milky bar, Snickers,  
5 Star chocolates distributed in a  
order.

$\frac{5}{5} \times \frac{4}{4} \times \frac{3}{3}$  we can distribute in 60 ways

= 60 ways. order matters.

with permutation, order matters.

S DM, KK MB 3 } Both considered  
S DM, MB 3 } All combinations.

S KK, DM, MB 3 }

Formula

$${}^n P_r = \frac{n!}{(n-r)!} \quad n = \text{Total no. of objects}$$

$$r = \text{No. of selection}$$

$$= \frac{5!}{(5-3)!} = \frac{5 \times 4 \times 3 \times 2!}{2!} \\ = 60$$

~~E.g. 2~~

Combination: Repetition will not occur.

$$\{DM, DK, MB\} \\ \times \{DK, DM, MB\} \leftarrow$$

Formula

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

④ Covariance [feature selection]

Age	weight	Age↑ - weight↑	Age↓ - weight↓
12	40		
13	45		
14	48		
15	51		
16	54		
17	50		
18	52		

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

Hence

$$\bar{X} = 15, \bar{Y} = 51$$

$$\begin{aligned}\text{cov}(X, Y) &= \frac{(12-15)(40-51)+(13-15)(48-51)}{5-1} \\ &\quad + (15-15)(48-51)+(17-15)(60-51) \\ &\quad + (18-15)(62-51)\end{aligned}$$

$$= \frac{33+12+18+33}{4}$$

$$= \frac{96}{4} = \boxed{24}$$

(+ve co-variance)  
XT YT

Co-variance

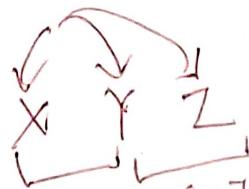
if } -ve covariance

XT YT

if } @ zero covariance NO correlation

④ Pearson co-relation coefficient. (-1, 1)

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



→ Restrict limit between

(-1, 1)

→ Tells how strengthen two features are

more

more co-related

→ Disadvantage: Useful for linear data.

## Spearman's Correlation

$$\rho_s = \frac{\text{cov}(R(x), R(y))}{\sigma(R(x)) \cdot \sigma(R(y))}$$

X	Y	R(x)	R(y)
10	4	4	1
8	6	3	2
7	8	2	3
6	10	1	4
		$\bar{x} = 2.5$	$\bar{y} = 2.5$

$$\text{cov}(R(x), R(y)) = \frac{(4-2.5)(1-2.5) + (3-2.5)(2-2.5) + (2-2.5)(3-2.5) + (1-2.5)(4-2.5)}{4-1}$$

$$\text{cov}(R(x), R(y)) = \frac{-2.55 - 0.55 + 0.25 + 2.55}{3}$$

$$\sigma(R(x)) = -1.87$$

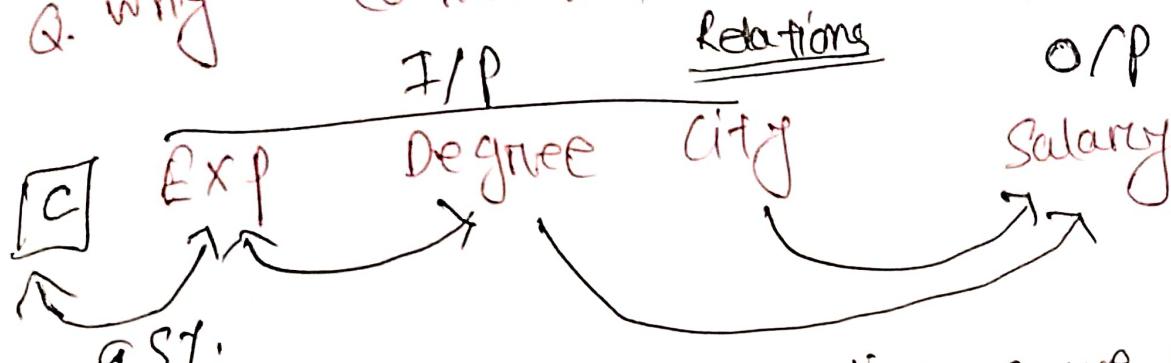
$$\sigma(R(x)) = \sqrt{\frac{(4-2.5)^2 + (3-2.5)^2 + (2-2.5)^2 + (1-2.5)^2}{4-1}}$$

$$\sigma(R(x)) = \sqrt{\frac{2.55 + 0.25 + 0.25 + 2.55}{3}}$$

$$\sigma(R(y)) = 1.87$$

$$\delta_S = \frac{-1.87}{1.87 \times 1.87} = -0.53$$

Q. why co-relation is used?



QST.

columns almost repeating so we can delete that column.