

# INTRODUCTION TO STATS

7 days Session

7pm - 8pm → 8:30 pm

## ① Basics To Advance

{DATA Scientist, Data Analyst, BUSINESS  
INTELLIGENCE TOOLS}

2 days Basics

### ① DESCRIPTIVE STATS

↓

{① Measure of Central Tendency }  
{② Measure of Dispersion }

Summarizing the data.

Histograms, Pdf, Cdf, Probability,

Permutation, Mean, Median, Mode,

Variance, Standard deviation

i) Gaussian Distribution

② LogNormal Distr

③ Binomial Distr

④ Bernoulli's Distr

⑤ Poach Distr {Power law}

⑥ Standard Normal Distr → python

⑦ Transformation and Standardization

⑧ Q-Q plot

⑨

### ② Inferential Stats

⑧ Z test → python

t test → python

ANOVA → F test

CHISQUARE:

HYPOTHESIS TESTING {P values}

Confidence Intervals

Z table, t table

## What is Statistics?

Statistics is the science of collecting, organizing and analyzing data. { Better Decision Making }

## Dyn Data? 2

Data : Facts or pieces of information that can be measured

Eg : The IQ of a class

{ 98, 97, 60, 55, 75, 65 }

Ages of students of a class

{ 30, 25, 24, 23, 27, 28 } → DATA

## Types of Statistics

### ① DESCRIPTIVE STATS

It consists of organizing and summarizing data

### ② Inferential Stats

Technique where in we used the data that we have measured to form

Conclusions

### ① Classroom of Maths student (20)

Marks of the 1<sup>st</sup> Sem

84, 86, 78, 72, 75, 65, 80, 81, 92, 95, 96, 97, - - -

Eg : Descriptive Stats

① What is the average marks of the students in the class?

Inferential Stats

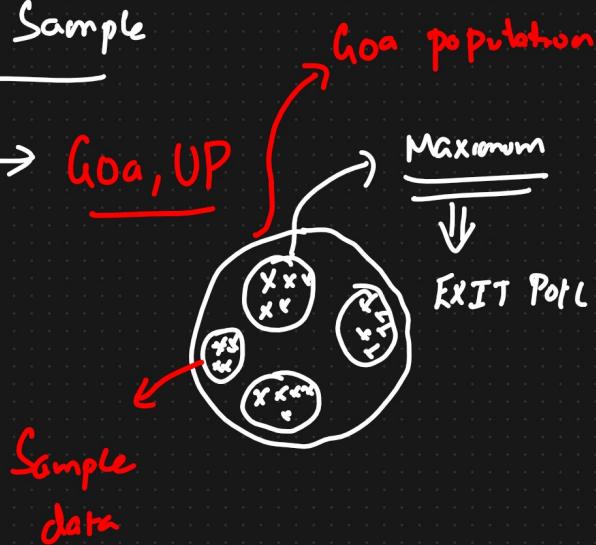
Sample 1 → Population.

Eg : Are the marks of the students of this classroom similar to the age of the Maths classroom in the college?

Population And Sample

Elections

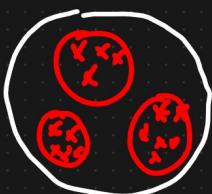
EXIT Poll



Population ( $N$ )      Sample ( $n$ )

Sampling Techniques

① Simple Random Sampling : Every member of the population ( $N$ ) has an equal chance of being selected for your sample ( $n$ )



② Stratified Sampling : Where the population ( $N$ ) is

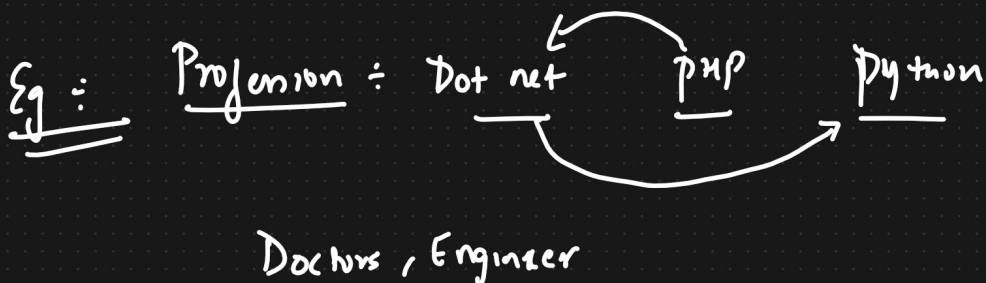
Split into non-overlapping groups (strata)

Eg: Gender → Male ✓ Survey

Female ✓

Age group

(0-10) (10-20) (20-40) (40-100)

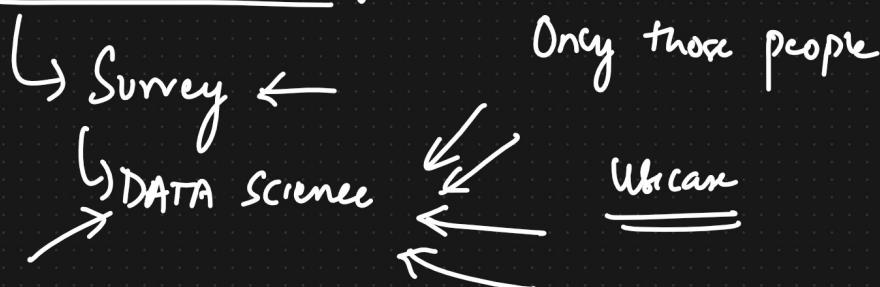


③ Systematic Sampling Thanos

( $N$ ) →  $n^{\text{th}}$  individual

Eg: Mail → Survey (Covid)  
→ 8<sup>th</sup> person → Survey

④ Convenience Sampling



Eg: EXIT POLL → Household Survey

{Random Sampling} ↓  
↓ Survey → Women

RBI → Household Survey



Survey → Women

Eg: Drug  $\rightarrow$  Tested  $\Rightarrow$   
 $\Downarrow$

What Kind of Sample ??

Variables :

A variable is a property that can take on

any value  $\{182, 178, 168, 150, 160, 170\}$

Eg: Height =  
Weight  $\{78, 99, 100, 60, 50, \dots\}$

Two Kinds of Variables  $\rightarrow$  Eg: Age  
Weight  
Height

- ① Quantitative Variable  $\rightarrow$  Measured Numerically, {Add, Subtract, multiply, divide}
- ② Qualitative / Categorical Variables

Eg: Gender  $\begin{cases} M \\ F \end{cases}$  {Based on some characteristics we can define Categorical Variables}

Eg: IQ

$\frac{0-10}{\Downarrow}$	$\frac{10-50}{\Downarrow}$	$\frac{50-100}{\Downarrow}$
Low IQ	Medium IQ	Good IQ

<u>Blood group</u>	<u>Tshirt size</u>
A+ve	L
B+	XL
O+	M
AB+	S

Quantitative

Discrete Variable

Eg: Whole number  
No. of Bank Accounts

Continuous Variables

Eg: Height = 172.5, 162.5 cm, 163.5 cm,  
Weight = 100kg, 99.5, 99.75

Eg: 2, 3, 4, 5, 6, 7,

Rainfall = 1.1, 1.25, 1.35 - - -

② Total of children in a family

Eg: 2, 3, 4, 5,

Eg: What kind of variable Gender IS? Categorical

② What " " " Marital Status? . "

③ River Length? Continuous

④ population of the state is? Discrete

⑤ Song length? Continuous

Blood pressure? Continuous.

PIN CODE ?  $\left\{ \begin{array}{l} \text{Discrete or categorical} \\ \hline \end{array} \right\}$

## ⑥ Variable Measurement Scales

4 types of Measured Variable

Colors, Gender, Type of flower

① Nominal data { Categorical data } → Classes

② Ordinal → Order of the data matters, value does not

③ Interval → Order matters, values also matter, natural zero is not present

④ Ratio.

Eg:

Students (Marks)	Rank	
100	1	
95	2	
57	4	
85	3	

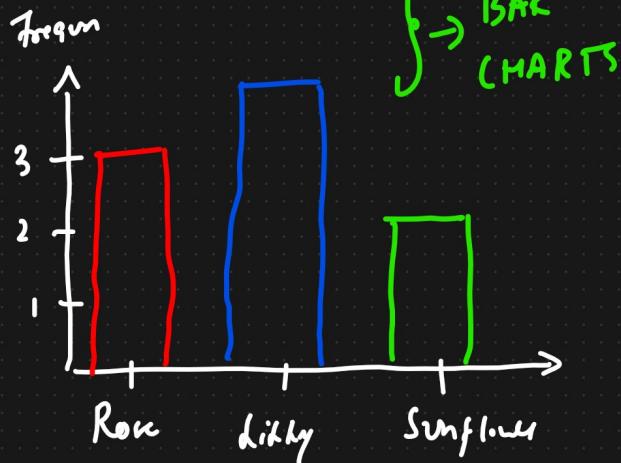


TemperaturesFahrenheit70 - 8080 - 9090 - 100100 ] Fahrenheit(F) Ratio data {Assignment}Frequency Distribution

Sample character : Rose, lilly, Sunflower, Rose, lilly, Sunflower,  
 Rose, lilly, lilly

Flower	Frequency	Cumulative Frequency
Rose	3	3
lilly	4	7
Sunflower	2	9

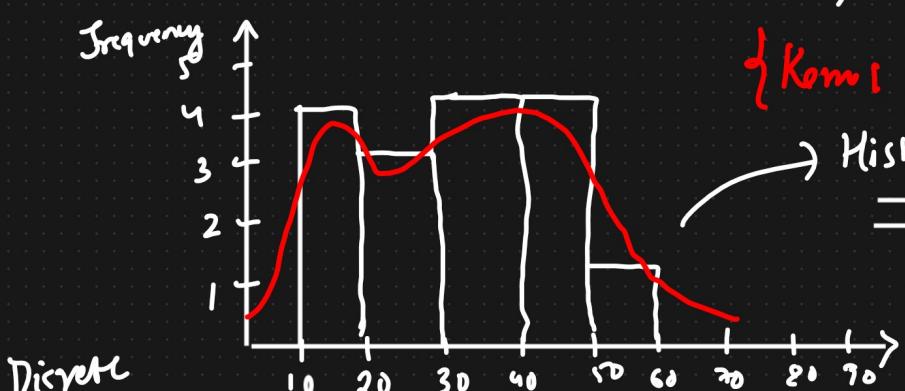
(I) BAR GRAPH



## ② Histograms $\div$ Continuous

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

Bins = 10



Discrete

pdf  $\div$  Smoothing of histograms

{ Kernel density Estimator }  
Histograms

BAR  $\uparrow$  VS Histogram  $\rightarrow$  Continuous

pdf: probability density function

# Basics To Intermediate Stats

- ① Measure of Central Tendency
- ② Measure of dispersion
- ③ Gaussian Distr
- ④ Z score
- ⑤ Standard Normal Distr

## ① Arithmetic Mean for Population & Sample

Mean (Average)

Population (N)

$$\downarrow$$

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

Sample (n)

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

$$= \underline{\underline{3.2}}$$

$$\mu = \sum_{i=1}^N \frac{x_i}{N}$$

$$= \frac{1+1+2+2+3+3+4+5+5+6}{10}$$

$$= \frac{32}{10} = 3.2$$

## ④ Central Tendency

- ① Mean ✓
- ② Median ✓
- ③ Mode ✓

Refers to the measure used to determine the centre of the distribution of data.

$\{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, \boxed{100}\}$

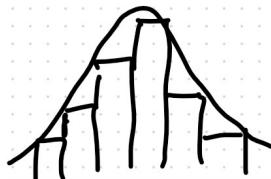
$$\text{Mean} = \frac{32 + 100}{11}$$

$$= \frac{132}{11} = 12$$

$$\begin{aligned} M &= 3.2 \\ &+ 100 \\ M &= 12 \end{aligned}$$

$\xrightarrow{\quad}$  outliers

Distribution



Median

$\{1, 1, 2, 2, 3, \boxed{4}, 5, 5, 6, \boxed{100}, 112\}$

1) Sort the numbers

Odd number = 11

$$\boxed{M=12}$$

$$\left\{ \begin{array}{l} \text{Median} = 3 \\ \hline \text{Median} = 3.5 \end{array} \right\} \text{ Avg} = \frac{3+4}{2} = 3.5$$

$$\begin{array}{l} M = 3.2 \\ \hline M = 12 \end{array}$$

$$\begin{array}{l} \text{Median} = 3 \\ \hline \text{Median} = 3.5 \end{array}$$

$\{$  Median works well with outlier  $\}$

$\equiv$

$$\textcircled{2} \quad \text{Mode} = \{1, 2, \underbrace{2, 2}_{2}, 3, 4, 5, \underbrace{6, 6, 6}_{3}, 7, 8, \underbrace{100, 100, 100, 100}_{2}\}$$

Mode = {Most frequent Element}

$\equiv$

$\boxed{\text{Mode} = 6} \rightarrow \text{Measure of Central Tendency}$

Type of flower

petal length

petal width

DATA SET  
 $\downarrow$

Mode

Rose

Lily

Sunflower

-

$\left\{ \begin{array}{l} \text{Missing} \\ \text{Value} \end{array} \right\} \rightarrow \text{Most frequent element}$

10% Missing data

Ages of Students }  
 =  
 Age  
 25  
 26  
 =  
 =  
 32  
 34  
 38

Mean? ✓  
 Median?  
 Mode ??

### (f) Measure of Dispersion

① Variance

② Standard Deviation

① Variance ✓

$$\mu \stackrel{=} \rightarrow \{1, 1, 2, 2, 4\} \quad \frac{80}{5} = 2$$

$$\mu \stackrel{=} \rightarrow \{2, 2, 2, 2, 2\} \quad \frac{10}{5} = 2$$

↓  
Dispersion  
Spread

### Population Variance

$$\sigma^2 = \frac{N}{n} \sum_{i=1}^N (x_i - \mu)^2 = \frac{10.84}{6} = 1.81$$

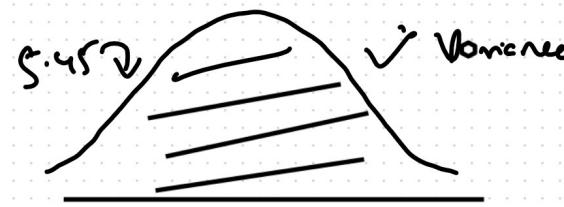
$x$	$\mu$	$x - \mu$	$(x - \mu)^2$
1	2.83	-1.83	3.34
2	2.83	-0.83	0.6889
2	2.83	+0.83	0.6889
3	2.83	0.17	0.03
4	2.83	1.17	1.37
5	2.83	2.17	4.71
	$\mu = 2.83$		$10.84$

### Sample Variance

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



Variance is more ??



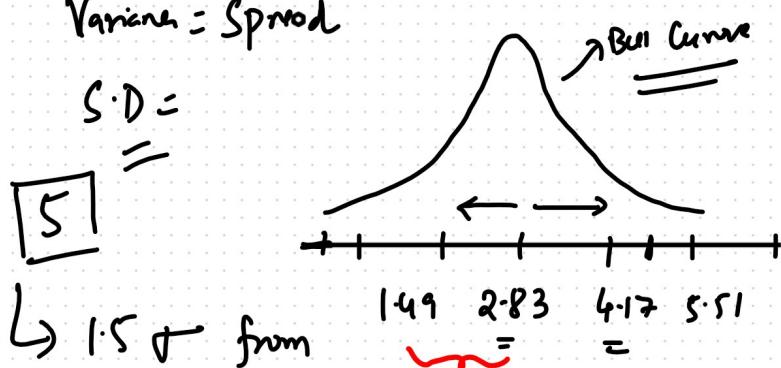
$$\sigma = \sqrt{\text{Variance}} = \sqrt{1.81} = 1.345$$

Variance = Spread

S.D. =

5

$$\begin{array}{r} 2.83 \\ 1.34 \\ \hline 4.17 \end{array} \quad \begin{array}{r} 2.83 \\ 1.34 \\ \hline 1.49 \end{array}$$



$$\begin{array}{r} 4.17 \\ 1.34 \\ \hline 5.51 \end{array}$$

the  $\mu$

## ④ Percentiles And Quartiles {Find outliers?}

Percentage : 1, 2, 3, 4, 5

% of the numbers that are odd?

% = # of numbers that are odd

$$\begin{aligned} & \text{Total Numbers} \\ &= \frac{3}{5} = 0.6 = 60\% \end{aligned}$$

Percentiles (GATE, CAT, GMAT, SAT)

↳ Dfn : A percentile is a value below which a certain percentage of observations lie.

Data set : 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12

11?

What is the percentile ranking of 10?  $n=20$

$x=10$   
Percentile Rank of  $x = \frac{\# \text{ of values below } x}{n} \times 100$

11?

$$= \frac{4}{20} \times 100 = 20\%$$

$$= \frac{17}{20} \times 100 = 85\%$$

② What value exists at percentile ranking  
of 25%?  $\boxed{75\%}$  ??

$$\text{Value} = \frac{\text{Percentile}}{100} \times (n+1)$$
$$= \frac{25}{100} \times (21) = 5.25 \rightarrow \text{Index Position} = \boxed{5} \rightarrow 25\%$$
$$= \frac{75}{100} \times (21) = 15.75 \Rightarrow \boxed{9} \text{ Answer} \rightarrow \text{index} =$$

## Five Number Summary

- ① Minimum
- ② First Quartile (Q1)
- ③ Median
- ④ Third Quartile (Q3)
- ⑤ Maximum

Removing the outliers

Outlier?? [27 > 13]

{ 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27 }

$$25\% \quad \frac{25}{100} \times (19+1)$$

~50 ✓

[Lower fence  $\longleftrightarrow$  Higher fence]

$$\frac{25}{100} \times 20 = \underline{\underline{5}} \Rightarrow \text{index} = 3 \Rightarrow$$

$$\text{Lower fence} = Q_1 - 1.5(\text{IQR})$$

$$Q_1 = \underline{\underline{3}} \checkmark$$

$$\text{Upper fence} = Q_3 + 1.5(\text{IQR}) \quad Q_3 = (75\%)$$

$$Q_3 = \underline{\underline{7}} \checkmark$$

$$\text{Interquartile Range (IQR)} = Q_3 - Q_1$$

$$= 7 - 3 = \underline{\underline{4}}$$

$$\text{Lower Fence} = Q_1 - 1.5(\text{IQR})$$

[Lower Fence  $\longleftrightarrow$  Higher Fence]

$$= 3 - 1.5(4)$$

$$[ -3 \longleftrightarrow 13 ]$$

$$= 3 - 6 = \boxed{-3} \checkmark$$

$$\text{Higher Fence} = Q_3 + 1.5(\text{IQR})$$

$$= 7 + 1.5(4)$$

$$= 7 + 6 = \underline{\underline{13}}$$

Remaining data

1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27

$$\text{Minimum} = 1$$

$$Q_1 = 3$$

$$\text{Median} = 5$$

$$Q_3 = 7$$

$$\text{Max} = 9$$

DATA

Visualisation

→ 5 Number Summary

Box plot

27

Box plot



23

23

$$\sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1} \rightarrow \begin{array}{l} \text{Bessel's correction} \\ \text{Degree of freedom} \end{array}$$

STATS  
Why sample variance is  $n-1$ ?

## Day 3 - Advance Statistics

## ① Distributions

- ↳ Normal Distr
  - ↳ Standard Normal Distr
  - ↳ Z score
  - ↳ Log Normal Distr
  - ↳ Bernoulli's Distr
  - ↳ Binomial Distr

# Practical

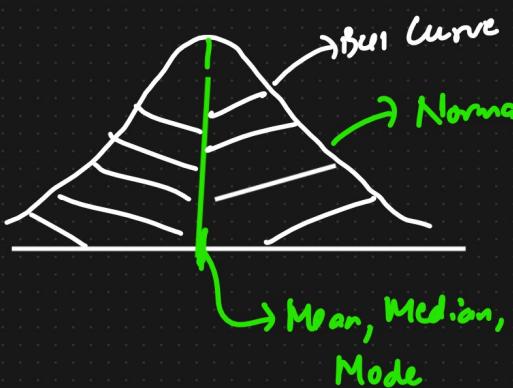
- ① Mean, Median, Mode
  - ② Variance, Standard deviation
  - ③ Histogram, pdf, Bar plot, Violin plot
  - ④ IQR
  - ⑤ Log Normal Distribution

## ① Distributions

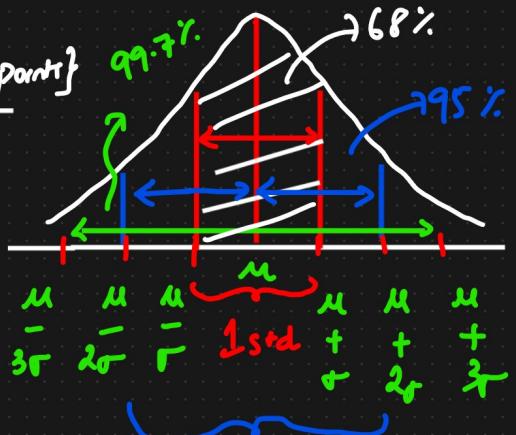
$$A_{\text{Q13}} = \{24, 26, 27, 28, 30, 32, \dots\}$$



## ① Gaussian / Normal Distribution



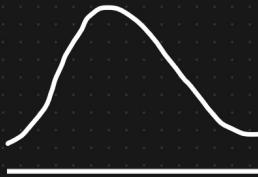
Datase t  
{ 100 datapoints }



## Empirical Formula

$$68 - 95 - 99.7\% \text{ Rule}$$

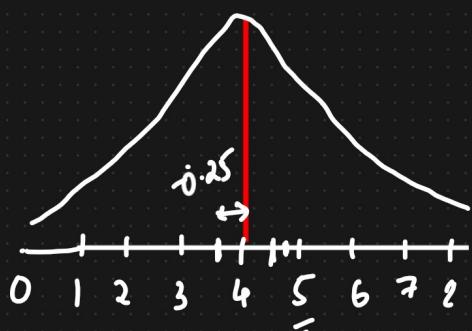
Eg: ① Height → Normally Distributed



Domain Expert → {Doctor}

② Weight    ③ IRIS DATASET

$$\text{Eg: } \mu = 4 \quad \sigma = 1$$



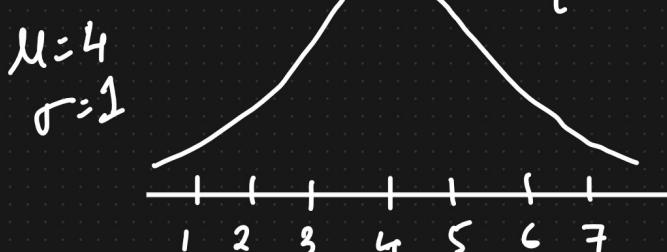
4.5 → Standard deviation

+0.5 sd

4.75 ??

$$\begin{aligned} Z\text{ score} &= \frac{x_i - \mu}{\sigma} \rightarrow \text{re} \\ &= \frac{4.75 - 4}{1} = 0.75 \text{ sd} \end{aligned}$$

$$Z\text{ score} = \frac{3.75 - 4}{1} = -0.25$$



$$\leftarrow Z\text{ score} = \frac{x_i - \mu}{\sigma} \quad \begin{matrix} \mu = 4 \\ \sigma = 1 \end{matrix}$$

$$\left\{ -3, -2, -1, 0, 1, 2, 3 \right\}$$

$$Z(1) = \frac{1-4}{1} = -3 \quad Z(3) = \frac{3-4}{1} = -1 \quad Y \sim SND(\mu=0, \sigma=1)$$

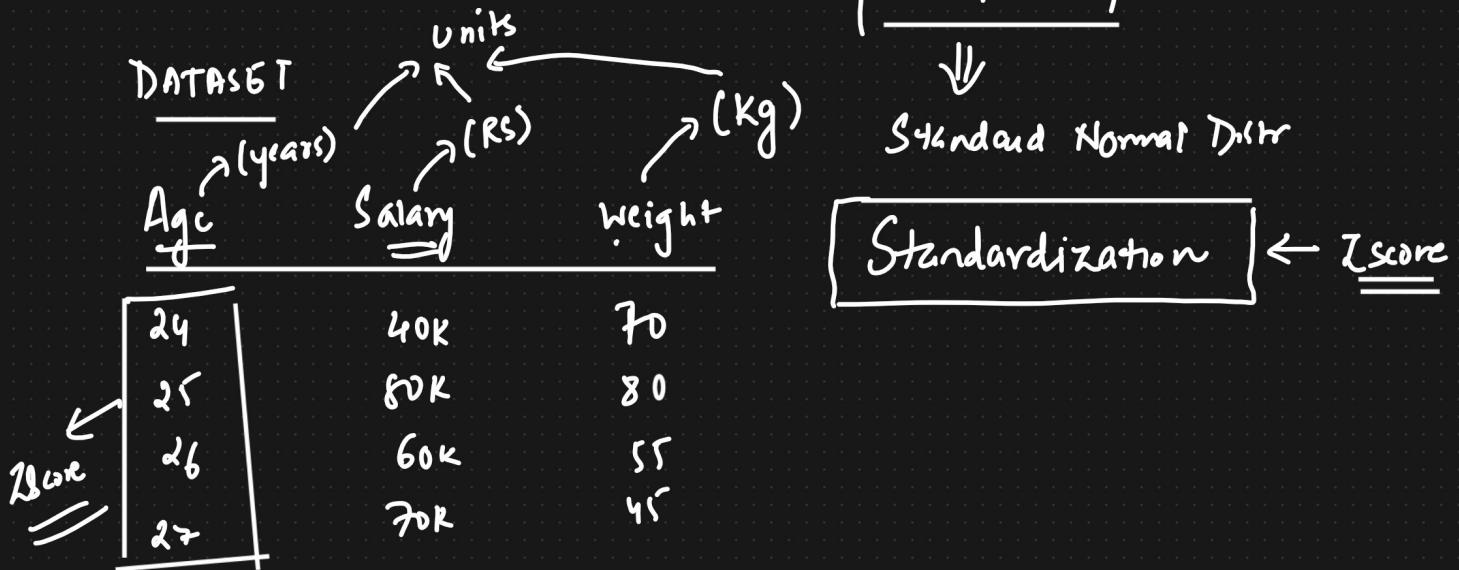
$$Z(1) = \frac{2-4}{1} = -2$$

$\left\{ 1, 2, 3, 4, 5, 6, 7 \right\} \rightarrow \text{Normal Distr}$

$\downarrow$       Standard Normal Distr  
 $\downarrow \quad (\mu=0, \sigma=1)$

$\{-3, -2, -1, 0, 1, 2, 3\}$  ↗  
Satisfying this  
property

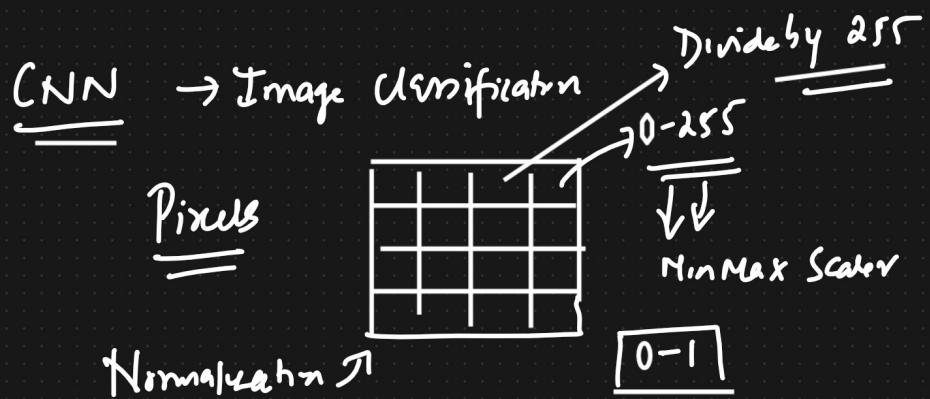
## Practical Application



Normalization  $\leftarrow \frac{\text{Value} - \text{Min}}{\text{Max} - \text{Min}}$

$(-1 \rightarrow 1) \hookrightarrow (0 \rightarrow 1)$

Min Max Scaler  $\rightarrow (0 \rightarrow 1)$



Practical Eg  $\{ \text{India vs SA} \}$

① ODI Series  $\downarrow 2021$  (CRICKET)

Scores Average  $2021 = 250$

Standard Deviation = 10

Compared to both the scores  
in which year Rishabh Pant final

Series < Team final score = 240 Score was better??

2020

Scenes Average 2020 = 260

Standard Deviation = 12

Team final score = 245

2021

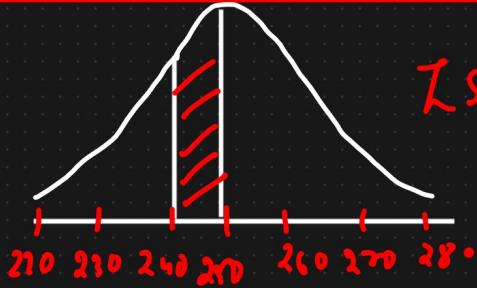
$$Z \text{ score} = \frac{x_i - \mu}{\sigma} = \frac{240 - 250}{10} = \frac{-10}{10} = -1 \}$$

2020

$$Z \text{ score} = \frac{x_i - \mu}{\sigma} = \frac{245 - 260}{12} = \frac{-15}{12} = -1.25 \}$$

In 2021 ✓

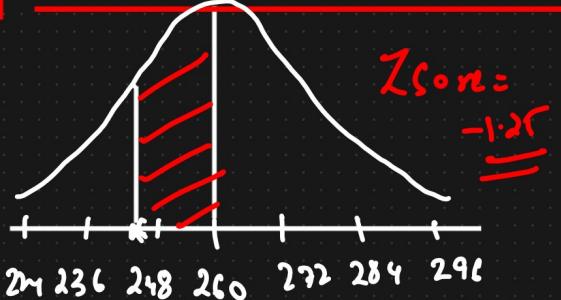
$$\boxed{\mu = 250 \quad x_i = 240 \quad \sigma = 10}$$



Final Match

In 2020 ✓

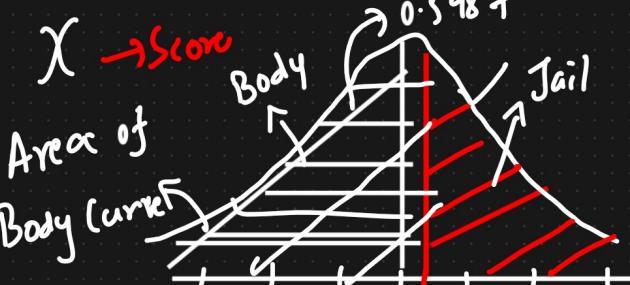
$$\boxed{\mu = 260 \quad x_i = 245 \quad \sigma = 12}$$



Assignment

Stats Interview Question

3



Question? ∵ What percentage of scores falls above 4.25?

$$\mu = 4 \quad \sigma = 1$$

$$\boxed{1 - \text{Left Area}}$$

of 1 2 3 4 0.25 6 7 } 1 - 0.5987

$$Z = \frac{x_i - \mu}{\sigma} = \frac{4.25 - 4}{1} = 0.25 = 0.4013$$

$\Downarrow 40\%$

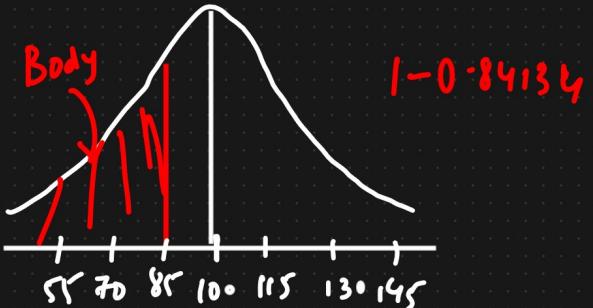


- ④ In India the average IQ is 100, with a standard deviation of 15. What percentage of the population would you expect to have an IQ lower than 85?

(Ans)  $Z = \frac{85 - 100}{15} = \frac{-15}{15} = -1$

and

JQ, 90 to 120

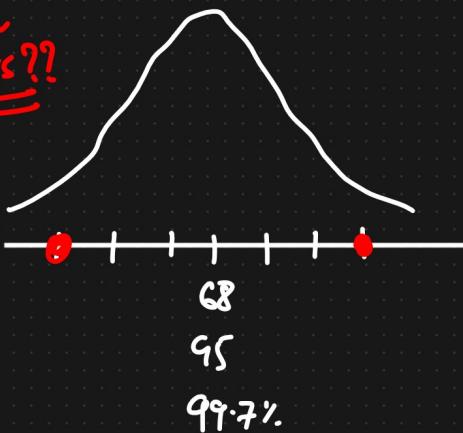


1 - 0.84134

# DAY 4 - STATS

- ① IQR - python ✓
- ② Probability ✓
- ③ Permutation And Combination ✓
- ④ Confidence Intervals ✓
- ⑤ P value ✓
- ⑥ Hypothesis Testing ✓

Aftr 3<sup>rd</sup> Sd  
outliers??



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

④ Probability : Probability is a measure of the likelihood of an Event

Eg : Roll a dice  $\{1, 2, 3, 4, 5, 6\}$

$Pr(\cdot 6) = \frac{\# \text{of way an event can occur}}{\# \text{of possible outcome}}$

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

Toss a coin  $\{H, T\}$

$$\boxed{Pr(H) = \frac{1}{2}}$$

## ② Addition Rule (Probability, "or")

### Mutual Exclusive Event

Two Events Are mutual exclusive if they cannot occur at the same time

Eg: Rolling a die  $\{1, 2, 3, 4, 5, 6\}$   
 $\{1, 2\}$

Tossing a coin  $\{\text{H}, \text{T}\}$

### Non Mutual Exclusive

Multiple events can occur at the same time

Eg: Deck of cards  $\{\text{Q}, \text{K}\}$

① If I Toss a coin, what is the probability of the coin landing on heads or tails?

Ans) Mutual Exclusive Addition Rule

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

$$\boxed{\Pr(A \text{ or } B) = 1}$$

Roll a die

$$\begin{aligned} \Pr(1 \text{ or } 3 \text{ or } 6) &= \Pr(1) + \Pr(3) + \Pr(6) \\ &= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} \\ &= \frac{3}{6} = \frac{1}{2} = 0.5 \end{aligned}$$

Non Mutual Exclusive

You are picking a card randomly from a deck.  
 What is the probability of choosing a card  
 that is Queen or a heart?  
 $\rightarrow (52)$

Ans) Non mutual Exclusive

$$P(Q) = \frac{4}{52} \quad P(\text{Heart}) = \frac{13}{52} \quad P(Q \text{ and Heart}) = \frac{1}{52}$$

Addition Rule for non mutual exclusive Events

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

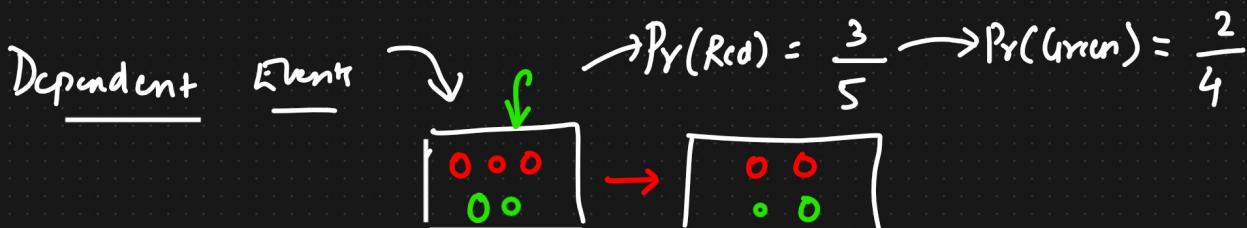
$$\begin{aligned} P(Q \text{ or Heart}) &= P(Q) + P(\text{Heart}) - P(Q \text{ and Heart}) \\ &= \frac{4}{52} + \frac{13}{52} - \frac{1}{52} \\ &= \frac{16}{52} \approx \frac{1}{3} \end{aligned}$$

### ③ Multiplication Rule

#### {Independent Events}

Eg: Rolling a dice  $\{1, 2, 3, 4, 5, 6\}$

1, 1, 2, Each & Every are independent  $\rightarrow$  Red marble •



# Naive Bayes {conditional probability}

## Independent Events

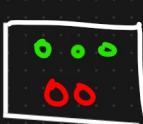
④ What is the probability of rolling a "5" and then a "4" in a dice?

Ans) Independent Event

Multiplication Rule

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

$$P(5 \text{ and } 4) = \frac{1}{6} * \frac{1}{6} = \frac{1}{36}$$



$$\frac{3}{5}$$

$$\frac{2}{4}$$

④ What is the probability of drawing a Queen and then a Ace from a deck of cards?

$$P(Q \text{ and } R) = P(Q) * P(R|Q)$$

$\underbrace{\qquad\qquad}_{\text{Sunt}} =$

Ans) Dependent

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

↑ conditional probability  
↓ Bayes theorem

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

$$= \frac{4}{52} * \frac{4}{51}$$



④

Permutation and Combination

Permutation

School trip {Chocolate factory} → Dairy, 5 star, Milky bar, Eclairs, Crem,  
 Student {Assignment} SITK

Student ↗  
 →

$$\underline{6} \times \underline{5} \times \underline{4} = \underline{\underline{120}}$$

⑤

Dairy, Crem, Milky

$$n = 6$$

Milky, Crem, Dairy

$$r = 3$$

Permutation

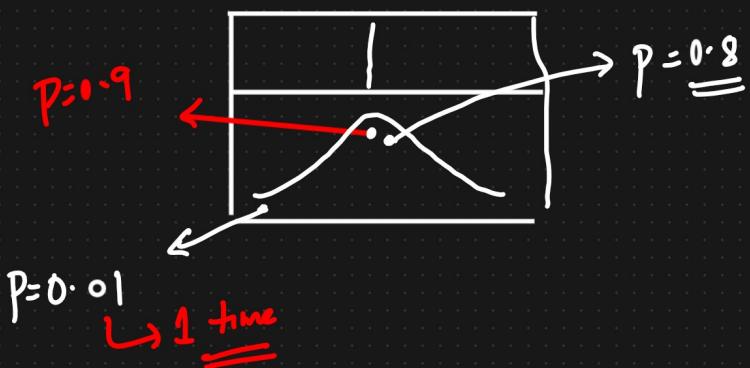
$$\begin{aligned} n_{Pr} &= \frac{n!}{(n-r)!} = \frac{6!}{(6-3)!} \\ &= \frac{6 \times r \times 4 \times 3!}{3!} \\ &= \underline{\underline{120}} \end{aligned}$$

⑥ Combination

Dairy	Crem	Eclair
—	—	—

$$\begin{aligned} n_C &= \frac{n!}{r!(n-r)!} = \frac{6!}{3!(6-3)!} \\ &= \frac{6^2 \times 5^2 \times 4^2 \times 3!}{3 \times 2 \times 1 \times 3!} \\ &= \underline{\underline{20}} \end{aligned}$$

# ① P value { Many people get's confused }



Every 100 time I touch the mouse pad 80 times I touch this specific region

## Hypothesis testing, Confidence Interval, Significance Value, - - -

Coin  $\rightarrow$  Test whether this coin is a fair coin or not by performing 100 tosses

Shady coin  $p(H) = \underline{100\%}$

$$\boxed{P(H) = 0.5 \quad P(T) = 0.5}$$

50 times Head (The coin is fair)

### Hypothesis Testing

① Null Hypothesis : Coin is fair

✓ ② Alternate Hypothesis : Coin is unfair

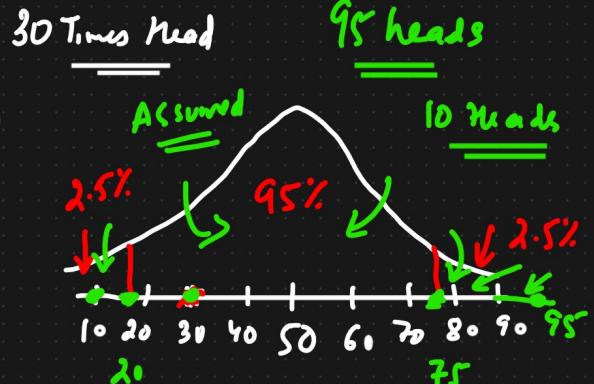
③ Experiment

④ Reject or Accept the Null Hypothesis

$$100\% - 5\% \\ \Downarrow \\ 95\% \\ \boxed{CI}$$

### Significance value

$$\alpha = 0.05 \text{ of Domain Expert}$$



1

25

$$\alpha = 0.20$$

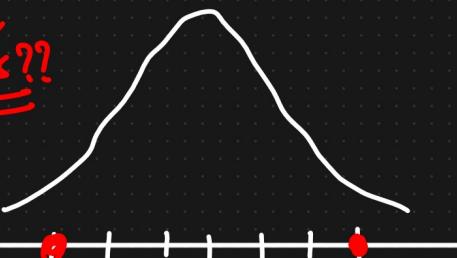
$$\alpha = 0.3$$

$$CI = \underline{70\%}$$

# DAY 4 - STATS

- ① IQR - python ✓
- ② Probability ✓
- ③ Permutation And Combination ✓
- ④ Confidence Intervals ✓
- ⑤ P value ✓
- ⑥ Hypothesis Testing ✓

*Aftr 3<sup>rd</sup> Sd outliers??*



68  
95  
99.7%

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

④ Probability : Probability is a measure of the likelihood of an Event

Eg : Roll a dice  $\{1, 2, 3, 4, 5, 6\}$

$Pr(\cdot 6) = \frac{\# \text{of way an event can occur}}{\# \text{of possible outcome}}$

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

Toss a coin  $\{H, T\}$

$$\boxed{Pr(H) = \frac{1}{2}}$$

## ② Addition Rule (Probability, "or")

### Mutual Exclusive Event

Two Events Are mutual exclusive if they cannot occur at the same time

Eg: Rolling a die  $\{1, 2, 3, 4, 5, 6\}$   
 $\{1, 2\}$

Tossing a coin  $\{\text{H}, \text{T}\}$

### Non Mutual Exclusive

Multiple events can occur at the same time

Eg: Deck of cards  $\{\text{Q}, \text{K}\}$

① If I Toss a coin, what is the probability of the coin landing on heads or tails?

Ans) Mutual Exclusive Addition Rule

$$\Pr(A \text{ or } B) = \Pr(A) + \Pr(B)$$

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

$$\boxed{\Pr(A \text{ or } B) = 1}$$

Roll a Die

$$\begin{aligned}\Pr(1 \text{ or } 3 \text{ or } 6) &= \Pr(1) + \Pr(3) + \Pr(6) \\ &= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} \\ &= \frac{3}{6} = \frac{1}{2} = 0.5\end{aligned}$$

Non Mutual Exclusive

You are picking a card randomly from a deck.  
 What is the probability of choosing a card  
 that is Queen or a heart?  
 $\rightarrow (52)$

Ans) Non mutual Exclusive

$$P(Q) = \frac{4}{52} \quad P(\text{Heart}) = \frac{13}{52} \quad P(Q \text{ and Heart}) = \frac{1}{52}$$

Addition Rule for non mutual exclusive Events

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

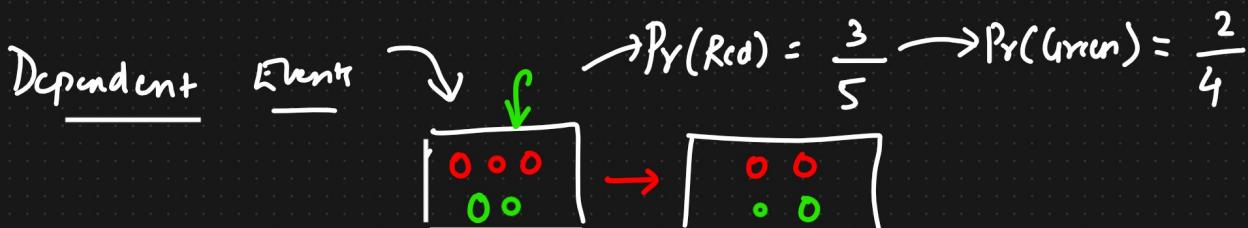
$$\begin{aligned} P(Q \text{ or Heart}) &= P(Q) + P(\text{Heart}) - P(Q \text{ and Heart}) \\ &= \frac{4}{52} + \frac{13}{52} - \frac{1}{52} \\ &= \frac{16}{52} \approx \frac{1}{3} \end{aligned}$$

### ③ Multiplication Rule

#### {Independent Events}

Eg: Rolling a dice  $\{1, 2, 3, 4, 5, 6\}$

1, 1, 2, Each & Every are independent  $\rightarrow$  Red marble •



# Naive Bayes {conditional probability}

## Independent Events

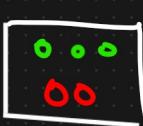
④ What is the probability of rolling a "5" and then a "4" in a dice?

Ans) Independent Event

Multiplication Rule

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

$$P(5 \text{ and } 4) = \frac{1}{6} * \frac{1}{6} = \frac{1}{36}$$



$$\frac{3}{5}$$

$$\frac{2}{4}$$

④ What is the probability of drawing a Queen and then a Ace from a deck of cards?

$$P(Q \text{ and } R) = P(Q) * P(R|Q)$$

Sunt =

Ans) Dependent

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

↗ conditional probability  
↓  
 Bayes theorem

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

$$= \frac{4}{52} * \frac{4}{51}$$



Permutation and Combination

Permutation

School trip {Chocolate factory} → Dairy, 5 star, Milky bar, Eclairs, Crem,  
 Student {Assignment} SITK

Student ↗  
 →

$$\underline{6} \times \underline{5} \times \underline{4} = \underline{\underline{120}}$$

⑤

Dairy, Crem, Milky

$$n = 6$$

Milky, Crem, Dairy

$$r = 3$$

Permutation

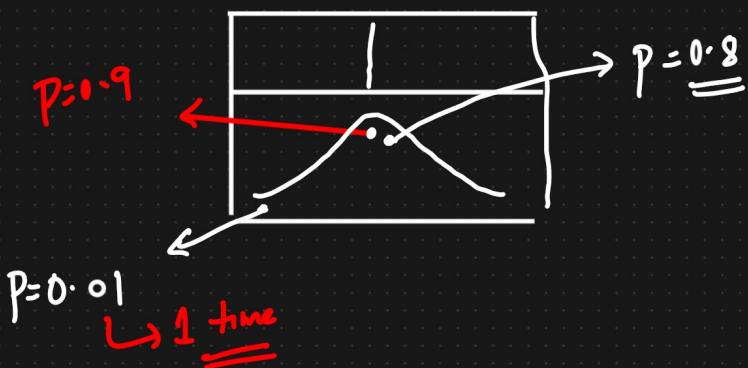
$$\begin{aligned} n_{Pr} &= \frac{n!}{(n-r)!} = \frac{6!}{(6-3)!} \\ &= \frac{6 \times r \times 4 \times 3!}{3!} \\ &= \underline{\underline{120}} \end{aligned}$$

⑥ Combination

Dairy	Crem	Eclair
—	—	—

$$\begin{aligned} n_C &= \frac{n!}{r!(n-r)!} = \frac{6!}{3!(6-3)!} \\ &= \frac{6^2 \times 5^2 \times 4^2 \times 3!}{3 \times 2 \times 1 \times 3!} \\ &= \underline{\underline{20}} \end{aligned}$$

# ① P value { Many people get's confused }



Every 100 time I touch the mouse pad 80 times I touch this specific region

## Hypothesis testing, Confidence Interval, Significance Value, - - -

Coin → Test whether this coin is a fair coin or not by performing 100 tosses

Shady coin  $p(H) = \underline{100\%}$

$$\boxed{P(H) = 0.5 \quad P(T) = 0.5}$$

50 times Head (The coin is fair)

### Hypothesis Testing

① Null Hypothesis : Coin is fair

✓ ② Alternate Hypothesis : Coin is unfair

③ Experiment

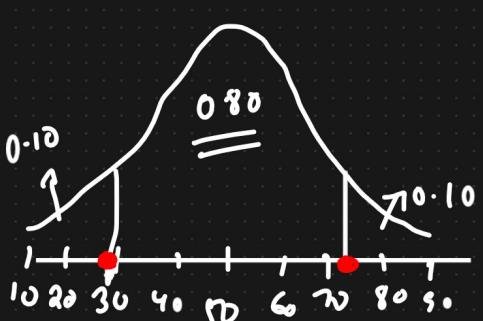
④ Reject or Accept the Null Hypothesis

$$100\% - 5\% \\ \Downarrow \\ 95\% \\ \boxed{CI}$$

### Significance value

$$\alpha = 0.05 \text{ of Domain Expert}$$

$$\alpha = 0.20 \quad \alpha = 0.3 \quad CI = \underline{70\%}$$



1

25

## Today Topics

- ① Type 1 and Type 2 Error ✓
  - ② One Tailed and 2 Tailed Test ✓
  - ③ Confidence Interval ✓
  - ④ Z-test, t-test, Chi-Square Test
- ① Type 1 and Type 2 Error

Null Hypothesis ( $H_0$ ) = Coin is fair

Alternate Hypothesis ( $H_1$ ) = Coin is not fair

## Reality check

Null Hypothesis is True or Null Hypothesis is False

## Decision

Null Hypothesis is True or Null Hypothesis is False

## Outcome 1 :

We reject the Null Hypothesis, when in reality it is false → Yes = Morris

Outcome 2 : We reject the Null Hypothesis, Person - Death Sentence

When in reality it is true → = Type 1 Error

Outcome 3 : We retain the Null Hypothesis

Accept

When in reality it is false → Type 2 Error

when in reality it is true

Outcome 4 : We Accept the Null Hypothesis when in reality it is true → Good

	P	N
T	TP	TN
F	FP	FN

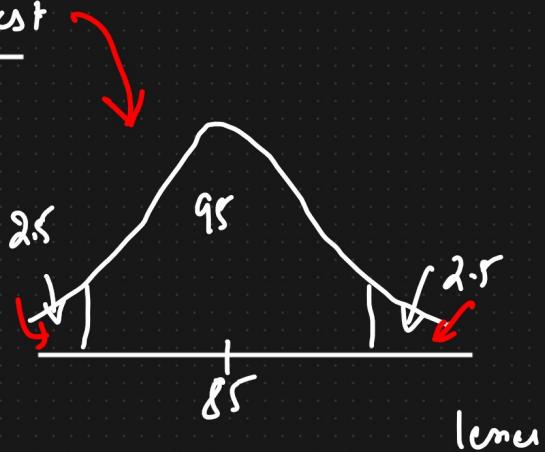
Type 2

Type 1

## ② 1 Tail and 2 tail Test

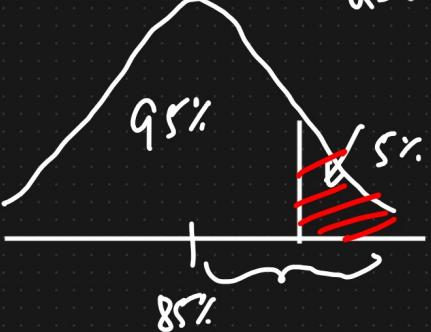
Eg.: Colleges in Karnataka have an 85% placement rate. A new college was recently opened and it was found that a sample of 150 students had a placement rate of 88%. With a standard deviation 4%. Does this college has a different placement rate?  $\alpha = 0.05$  85%

### 2 tailed Test



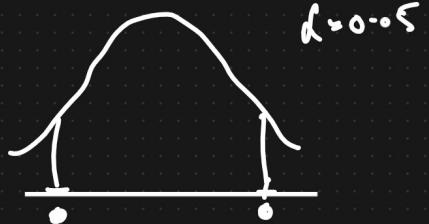
Does this college have a placement rate greater than 85%?

$$\alpha = 0.05$$



③

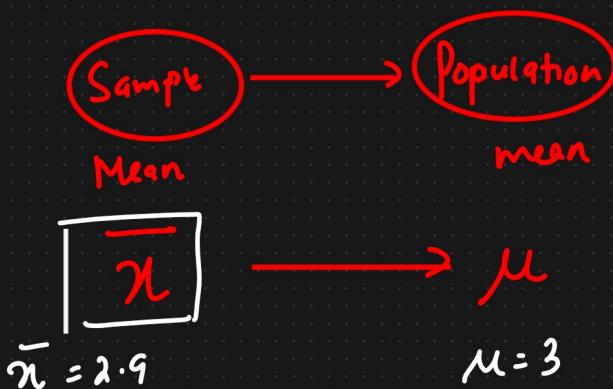
## Confidence Intervals



### Point Estimate

The value of any statistic that estimates the  
Value of a parameter ✓

### Inferential Stats



### Confidence Intervals

Point Estimate  $\pm$  Margin of Error

Q) On the Quant test of CAT Exam, the standard deviation is known to be 100. A sample of 25 test takers has a mean of 520 score. Construct a 95% CI about the mean?

Ans)  $\sigma = 100$   $n = 25$   $d = 0.05$   $\bar{x} = 520$

$$d = 1 - 0.95 = 0.05$$



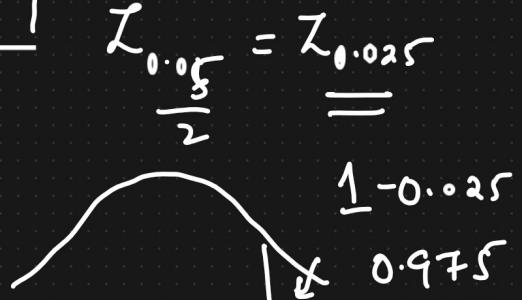
$\left\{ \begin{array}{l} \text{① Population std is given} \\ \text{② } n > 30 \end{array} \right\} \rightarrow Z \text{ test}$

Point Estimate  $\pm$  Margin of Error

$$\boxed{\bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}} \rightarrow \text{Standard Error}$$

$$Z_{0.05} = Z_{0.025}$$

$$\text{Upper bound} = \bar{x} + Z_{\frac{0.05}{2}} \frac{100}{\sqrt{25}}$$

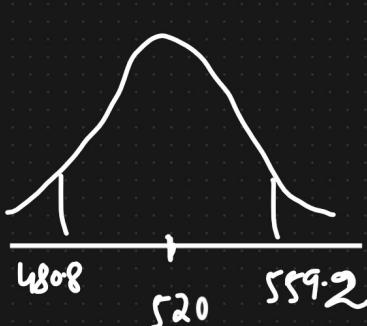


$$\text{Lower bound} = \bar{x} - Z_{\frac{0.05}{2}} \frac{100}{\sqrt{25}}$$

$$\boxed{1.96}$$

$$\text{Upper} = 520 + 1.96(20) = 559.2$$

$$\text{Lower} = 520 - 1.96(20) = 480.8$$



Stats

Find the average size of  
the shark throughout the world?

Q) On the quant test of CAT exam, a sample of 25 test takers has a mean of 520 with a standard deviation of 80. Construct 95% confidence interval about the mean?

Ans) Condition  $n=25 \quad \bar{x}=520 \quad S=80$   
 $\alpha=0.05$

Here population std is  
not given  $\rightarrow t\text{-test}$

Point Estimate  $\pm$  Margin of Error

$$\bar{x} \pm t_{d/2} \left( \frac{s}{\sqrt{n}} \right) \rightarrow \text{Standard Error } t_{0.05/2} = 2.064$$

$$\text{Upper bound} = \bar{x} + t_{0.05/2} \left( \frac{s}{\sqrt{n}} \right)$$

$$\underline{\text{Degree of freedom}} = n - 1 = 25 - 1 = 24$$

$$= 520 + 2.064 \left( \frac{80}{\sqrt{24}} \right)$$

$$= 553.024$$

$$\underline{\text{Lower bound}} = \bar{x} - t_{0.05/2} \left( \frac{s}{\sqrt{n}} \right)$$

$$= 520 - 2.064 \left( \frac{80}{\sqrt{24}} \right)$$

$$= 486.97$$

$$[486.97 \longleftrightarrow 553.024]$$

## ① One Sample Z-Test

① Population SD is given

② Sample size  $n > 30$

\*) In the population, the average IQ is 100 with a SD of 15. Researchers wants to test a new medication to see if there is positive or negative effect on intelligence, or no effect at all. A sample of 30 participants who have taken the medication has a mean of 110. Did the medication affect the intelligence?

$$\alpha = 0.05 \quad (\cdot I = 95\%)$$

$$\rightarrow \boxed{110} \checkmark$$

An) 1) Define Null Hypothesis

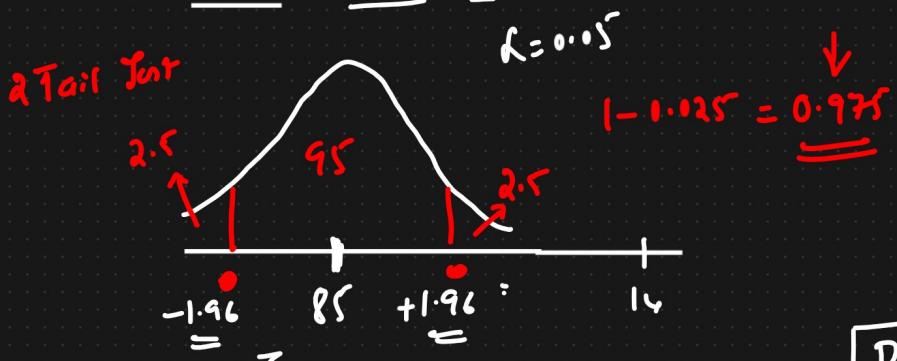
$$H_0: \mu = 100$$

2) Alternative Hypothesis  $H_1: \mu \neq 100$

③ State Alpha

$$\alpha = 0.05$$

④ State Decision Rule  $Z_{\text{table}}$



$$P \leq 0.05$$

⑤ Calculate Test Statistics

$$Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

Sample data

$\Rightarrow$  Standard Error

$$= \frac{140 - 100}{\sqrt{30}} = \frac{40}{\sqrt{30}} \times \sqrt{30} = 14.60$$

State our Decision  $\{Z = 14.60\}$

$$14.60 > 1.96 \quad Z = 14.60$$

If  $Z$  is less than  $-1.96$  or greater than  $1.96$ , reject the null hypothesis

Medication Improve the intelligence

or decrease  $? \square$  Improve  
the intelligence

## ② One Sample t-test

Z-test  $\Rightarrow$  population std

t-test  $\Rightarrow$  unknown population std

① Population the average IQ = 100

$$n = 30 \quad \bar{x} = 140 \quad s = 20$$

Did the medication affect intelligence?

$$\alpha = 0.05$$

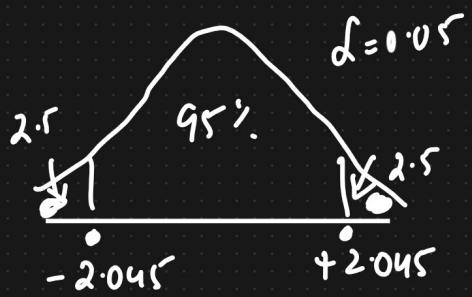
Ans) ①  $H_0: \mu = 100$

②  $H_1: \mu \neq 100$

③ Calculate the degree of freedom

$$n - 1 = 30 - 1 = 29$$

④ State Decision Rule



$$t = \frac{10.96}{\sqrt{29}} > 2.045$$

Reject Null Hypothesis

$P \leq$  significance value



⑤ T Test

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \quad \left. \begin{array}{l} \bar{x} = 140 \\ \mu = 100 \\ s = 20 \\ n = 30 \end{array} \right\} \text{Increased the } \begin{cases} \text{Intelligence.} \end{cases}$$

$$= \frac{10.96}{\sqrt{29}}$$

Reject the Null Hypothesis

## Real World Problem



# 6<sup>th</sup> DAY LIVE SESSION

- ① CHI SQUARE ✓
- ② Covariance ✓
- ③ Pearson Correlation Coefficient ✓
- ④ Spearman Rank Correlation ✓
- ⑤ Practical Implementation
  - Z-test, t-test, chi square test
- ⑥ F Test (ANOVA)

## CHI SQUARE TEST

- ① Chi Square Test claims about population proportions

It is a non parametric test that is performed on Categorical (nominal or ordinal) data.

- Q) In the 2000 Indian Census, the age of the individual in a small town were found to be the following:

Less than 18	18-35	>35
20%	30%	50%

In 2010, age of  $n=500$  individuals were sampled. Below are the results

<18	18-35	>35
121	288	91

Using  $\alpha=0.05$ , would you conclude the population distribution of ages has changed in the last 10 years?

Ans)

$<18$	$18-35$	$>35$
20%	30%	50%

{ Population } 2000

Expected

$<18$	$18-35$	$>35$
121	288	91
100	$500 \times 0.3$	$500 \times 0.5$

| n=500 |

Observed

$<18$	$18-35$	$>35$
121	288	91
100	150	250

Observation

Expected

{ Chi Square table }

- ①  $H_0$  = The data meets the distribution 2000 census       $df = 2$ ,  $\alpha = 0.05$
- ②  $H_1$  = The data does not meet " " " "
- ③ Degrees of freedom =  $n - 1 = 3 - 1 = 2$
- ④ Decision Boundary



If  $\chi^2$  is greater than 5.99 reject  $H_0$

5) Calculate Test Statistics

$$\chi^2 = \sum \frac{(f_0 - f_e)^2}{f_e}$$

$$= \frac{(121 - 100)^2}{100} + \frac{(288 - 150)^2}{150} + \frac{(91 - 250)^2}{250} \\ \approx 132.94$$

$$\chi^2 = 132.94 > 5.99 \quad \left\{ \begin{array}{l} \text{Reject the Null} \\ \text{Hypothesis} \end{array} \right.$$

$$0.11 > 0.05 \quad \boxed{0.11 \times 0.05} \quad \alpha = 0.05 \quad 0.002 < 0.05 \quad \left\{ \begin{array}{l} \text{Domain} \\ \text{Reject the null hypothesis} \end{array} \right.$$

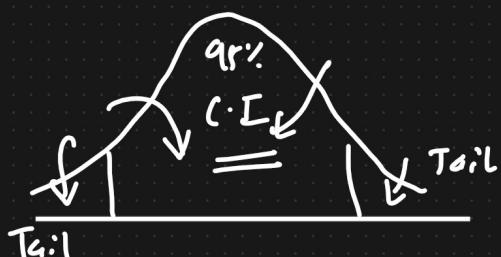
Accept the Null  
Reject the hypothesis



I Was Correct Here

$\left\{ \begin{array}{l} P\text{-value} < \text{Significance value} \\ \Downarrow \\ \text{Reject the Null Hypothesis.} \end{array} \right.$

OR  
 $\left\{ \text{Accept the Null Hypothesis} \right\}$



$$P = 0.11 > 0.05$$

Accept

$$P = \boxed{0.002} < 0.05$$

Reject the Null Hypothesis

## ② Covariance

<u>X</u>	<u>Y</u>
Weight	Height
50	160
60	170
70	180
75	181

No. of hour Study	play
2	6
3	4
4	3

Quantity relationship between  $X \& Y$

Covariance

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

+ve

positive correlation

<u>X↑</u>	<u>Y↑</u>
<u>X↓</u>	<u>Y↓</u>

-ve

negative correlation

= +ve or -ve

$x \uparrow \& y$

<u>X↓</u>	<u>Y↑</u>
<u>X↑</u>	<u>Y↓</u>

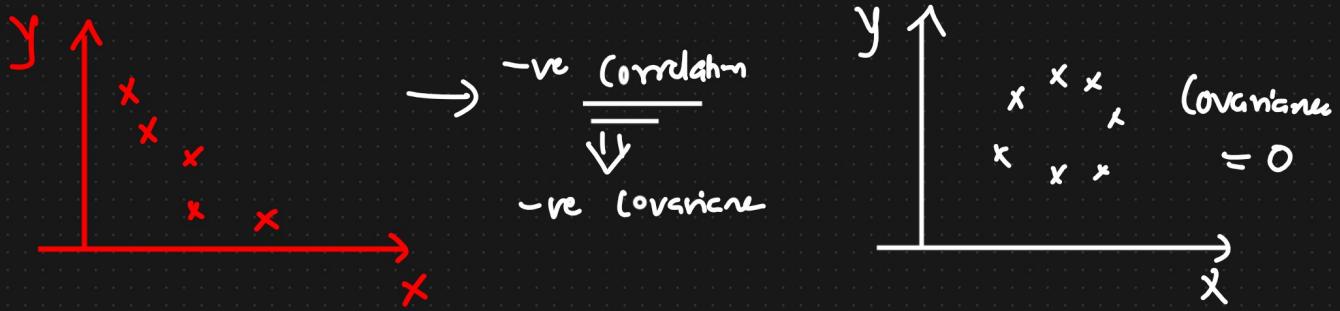
O



→ +ve Correlation

↔

↔



Disadvantage of Covariance

① Positive OR Negative ✓

$$\begin{array}{r}
 +100 \\
 -200 \\
 \hline
 -2000
 \end{array}
 \quad
 \begin{array}{r}
 +1000 \\
 -200 \\
 \hline
 = 0
 \end{array}
 \quad
 \text{f Direction}$$

② Pearson Correlation Coefficient

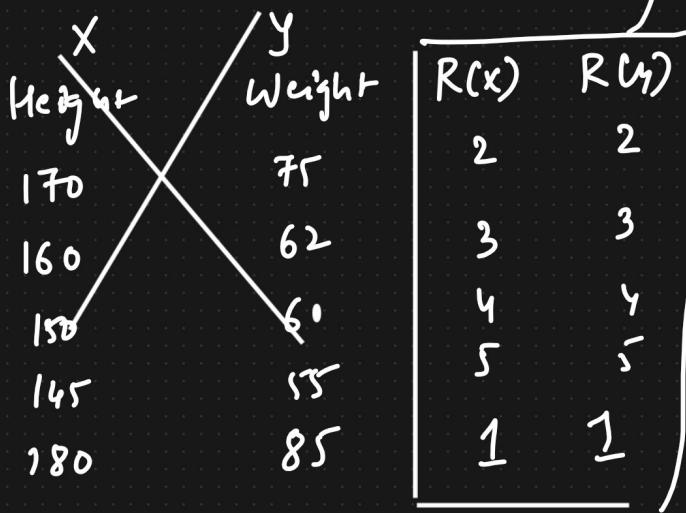
(-1 to 1)

The more towards +1 more positively correlated

The more towards -1 more negatively correlated

$$\rho_{(x,y)} = \frac{\text{Cov}(x,y)}{\sigma_x \sigma_y} = \left\{ \begin{matrix} -1 & 1 \end{matrix} \right\}$$

$$\text{Spear}(x,y) = \frac{\text{Cov}(R(x), R(y))}{R_{fx} \times R_{fy}}$$



## Non linear properties

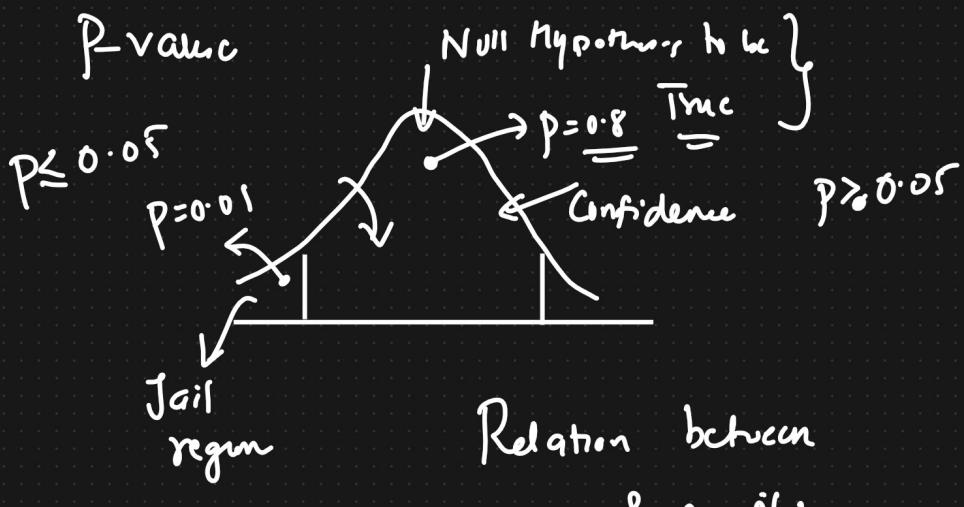
Coding

$\checkmark P \leq 0.05 \rightarrow \text{Reject the Null Hypothesis}$

$\downarrow$

probability }  
 $\alpha = 0.05$       5% probability the null hypothesis is correct

$P \geq 0.05 \rightarrow \text{Accept the Null Hypothesis}$



$P\text{-value} \quad < \quad \boxed{\text{Significance}}$

$\hookrightarrow$  Reject the Null Hypothesis

$P\text{-value} \quad \geq \quad \hookrightarrow$  Accept the Null

① P value and Significance value

# 7 day

- ① P value & Significance value
- ② Distribution
- ③ Central Limit Theorem
- ④ Bernoulli's Distr  $\xrightarrow{\text{as } n \rightarrow \infty}$  Normal Distr
- ⑤ Binomial Distr
- ⑥ Poisson's Distr {Poisson Law}
- ⓫ F Test (ANOVA)  $\xrightarrow{\text{1 hr}}$  → upload a separate video.

7080  
↓  
KRISH10 → 10%

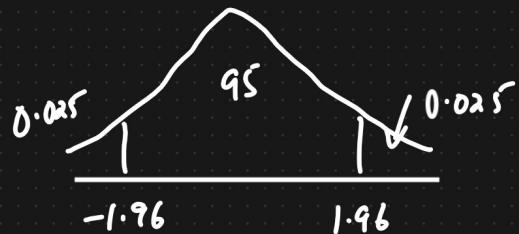
## ① P Value & Significance

↳ Define the p value

Q) The average weight of all residents in Bangalore city is 168 pounds. With a standard deviation 3.9. We take a sample of 36 individuals and the mean is 169.5 pounds. ( $\underline{I} = 95\%$ )

Ans)  $\mu = 168 \quad \sigma = 3.9 \quad \bar{X} = 169.5 \quad n = 36 \quad \alpha = 0.05$

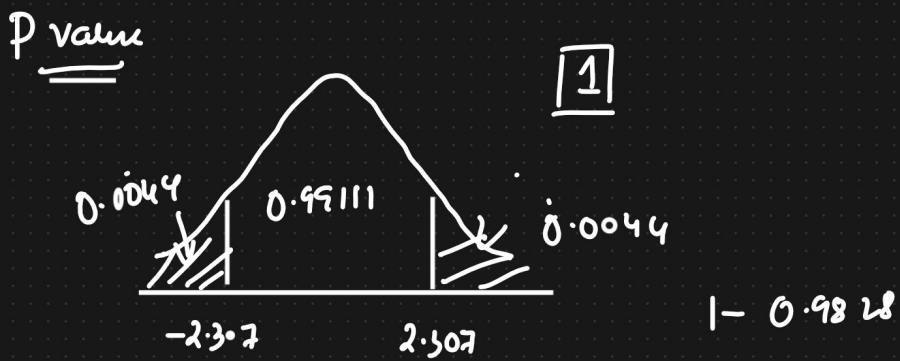
- ①  $H_0: \mu = 168$       ③ Decision boundary  $1 - 0.025 = 0.9750$   $\xrightarrow{\text{as } n \rightarrow \infty}$
- $H_1: \mu \neq 168$
- ②  $\alpha = 0.05$



⑥  $\chi^2$  Test

$$\begin{aligned} Z &= \frac{x - \mu_L}{\frac{\sigma}{\sqrt{n}}} = \frac{169.5 - 168}{\frac{3.9}{\sqrt{36}}} \\ &= \frac{1.5}{3.9} \times 6 \\ &= 2.307 \end{aligned}$$

④  $Z = 2.307 > 1.96$  Reject the Null Hypothesis



$$1 - 0.99111 =$$

$$\begin{aligned} P \text{ value} &= 0.0044 + 0.0044 \\ &= 0.0088 \end{aligned}$$

$P \text{ value} < 0.05$

$0.0088 < 0.05$  → Reject the Null Hypothesis

②

P Value  $\leq$  Significance value

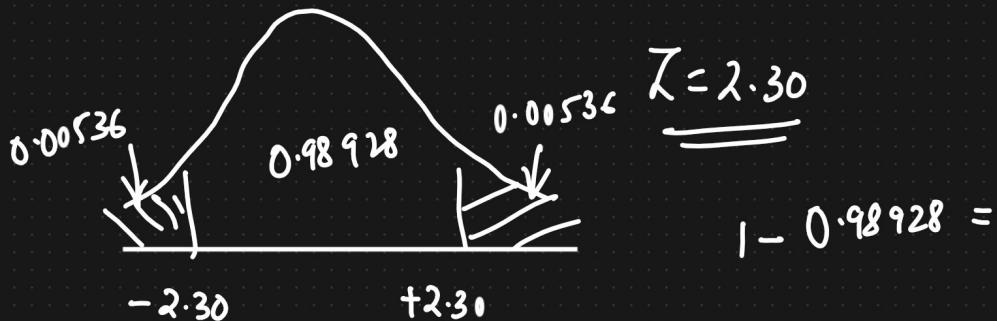


Reject the Null hypothesis

P Value  $>$  Significance value



Fail to Reject the Null Hypothesis



$2.30 > 1.96 \quad \{ \text{Reject the Null Hypothesis} \}$

$$\begin{aligned} \underline{\underline{P \text{ value}}} &= 0.00536 + 0.00536 \\ &= \underline{\underline{\quad}} \leq \alpha \Rightarrow \text{Reject} \\ &\quad \text{the Hypothesis} \end{aligned}$$

②

Average age of a college is 24 years with a standard deviation 1.5.

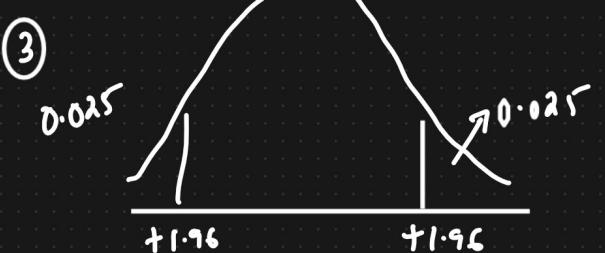
Sample of 36 student students mean is 25 years. With  $\alpha=0.05$

(i)  $H_0: \mu = 24$ , do the age vary?

$$\text{Ans) } H_0: \mu = 24 \quad \sigma = 1.5 \quad n = 36 \quad \bar{x} = 25 \quad \alpha = 0.05$$

$$H_1: \mu \neq 24$$

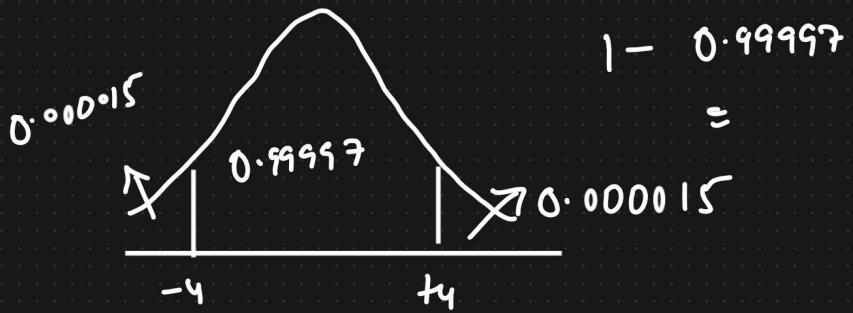
$$\underline{\underline{\alpha}} = 0.05$$



④

$$\begin{aligned} Z\text{-Score} &= \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{25 - 24}{\frac{1.5}{\sqrt{6}}} \\ &= \frac{1 \times 6}{1.5} \\ &= 4 \end{aligned}$$

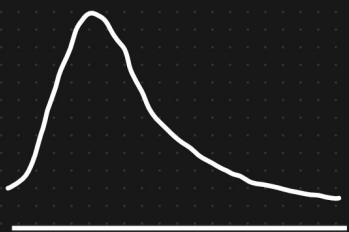
$4 > 1.96$  Reject Null Hypothesis



$$\begin{aligned} P\text{ value} &= 0.000015 + 0.000015 \\ &= 0.00003 \end{aligned}$$

Pvalue  $<$  Significance  $\left\{ \begin{array}{l} \text{Reject the} \\ \text{Null Hypothesis} \end{array} \right\}$

## ② Log Normal Distribution



Eg: ① Wealth Distribution

② People writing big comments

$\{ Y \sim \text{Log Normal Dist} \}$

$\log(y) \rightarrow$  Normal Distribution

#### ④ Bernoulli's Distribution

2 Outcomes

0 or 1

Singe Trial

$$P = 0.5$$

Tossing a coin

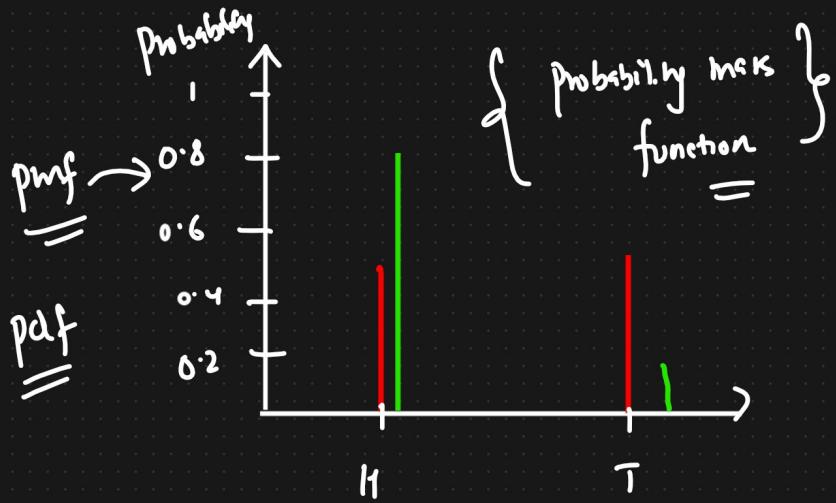
$$P(H) = 0.5 = P$$

$$q = 1 - P = 0.5$$

Do not have a fair coin

$$P(H) = 0.3 = P$$

$$P(T) = q = 1 - P = 1 - 0.3 = 0.7$$



#### ⑤ Binomial distribution

Every → Bernoulli distribution

Multiple Trial



$$p(H) = 0.5$$

$$p(H) = 0.6 \quad - \quad - \quad - \quad - \quad - \quad -$$

$$p(T) = 0.1$$

$$p(T) = 0.4$$

Power  
distribution

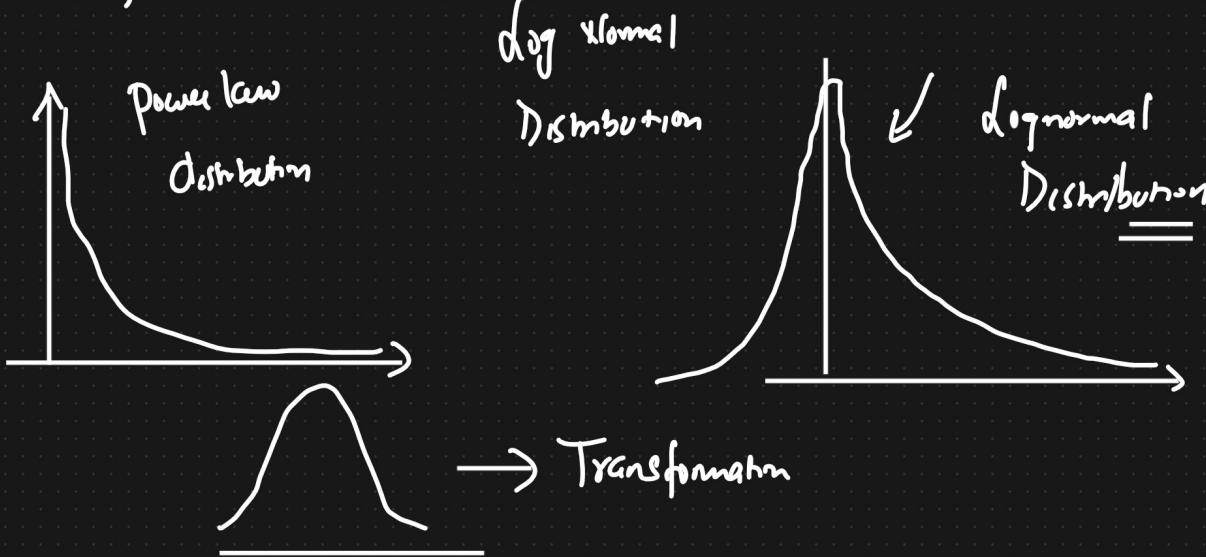


Eg: 80% of the wealth is distributed with 20% of the people  
② 80% of the company project by 20% of the people in a team

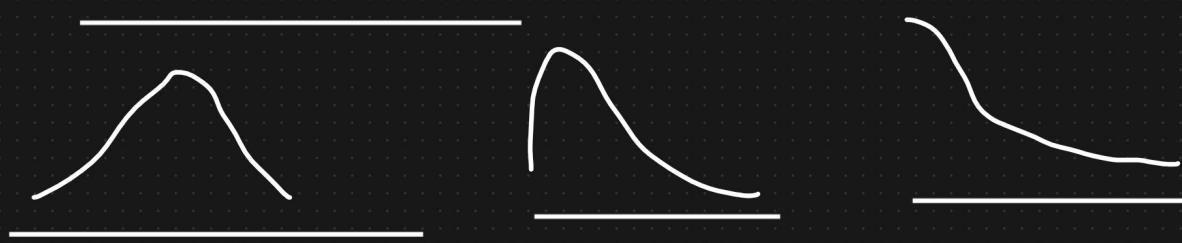
④ 80% of sales is done by the 20%

famous product.

④ 80% of the match is won by 20% of the team.



④ Central Limit Theorem



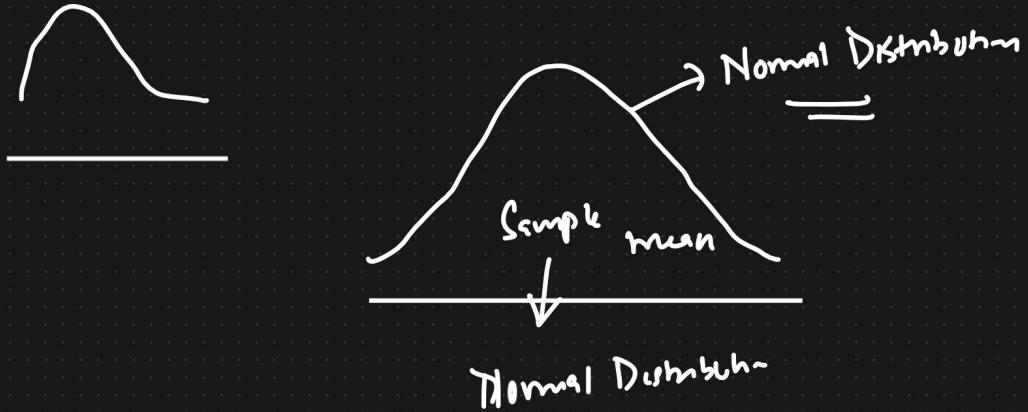
Sample size

$$S_1 \cap S_2 \rightarrow \bar{x}_1 \\ S_2 \rightarrow \bar{x}_2$$

$$\downarrow \rightarrow [n > 30]$$

$\bar{x}_1, \bar{x}_2$  }  $\rightarrow$  Sample mean

$$m \sim \{ \begin{matrix} 3 \\ S_4 \\ S_6 \\ \vdots \\ S_m \end{matrix} \} \rightarrow \begin{matrix} \overline{x}_3 \\ \overline{x}_4 \\ \vdots \\ \overline{x}_m \end{matrix}$$



## Poisson Distribution

- ① Machine Learning Algorithms ✓ → 2 Algorithms
- ② Deep Learning Algo ✓
- ③ FLASK & DJANGO ✓
- ④ MongoDB, SQL ✓
- ⑤ Blockchain Session ✓