

#Statistics#

- ① ① What is statistics?
- ② Data
- ③ Types of statistics?
- ② ① Population and Sample
- ② Types of Sampling Techniques
- ③ ① Variables
- ② Kinds of Variables
- ③ Variable Measurement Scales
- ④ ① Frequency Distribution
- ② Bar Graph
- ③ Histogram
- ⑤ ① Arithmetic Mean for Population and sample.
- ② Measure of Central Tendency
- ③ Median
- ⑥ ① Mode
- ② Measure of Dispersion
- ③ Percentiles and Quantiles
Percentile
- ⑧ Five Number Summary
- ⑨ Box plot
- ⑩ Why Sample Variance is divided by $n-1$?
- ⑪ Gaussian Distribution / Normal Distribution
- ⑫ Normalization
- ⑭ ① Probability
 - ② Addition Rule ("or") Mutual Exclusive
 - ② Non Mutual Exclusive
 - ⑮ ① Non Multiplication Rule ("and") Independent
 - ② Dependent
- ⑯ Permutation and Combination
- ⑰ P-value, Hypothesis testing, Confidence Interval, Significance Value.
- ⑱ Type 1 and Type 2 Error
- ⑲ One Tail and two tail Test
- ⑳ Confidence Interval
- ㉑ One sample z-test
- ㉒ One Sample t-test
- ㉓ Chi Square Test
- ㉔ Covariance.
- ㉕ Pearson Correlation Coefficient
- ㉖ Spearman rank correlation coefficient
- ㉗ P-value and Significance value.
- ㉘ Log Normal Distribution
- Bernoulli Distribution
- Binomial Distribution
- Power Law { Pareto Distribution }
- ㉙ Central Limit Theorem
- ㉚ Anova
- ㉛ One way anova

Day 1

What is statistics?

Statistics is the science of collecting, organizing and analyzing data
For better Decision making

Data :- Facts or pieces of information that can be measured.

Ex: ① The IQ of a class { 98, 97, 60, 55, 75, 65 }
② Age of students of a class { 30, 25, 24, 23, 27, 28 }

Types of Statistics

① Descriptive statistics

- It consists of organizing and summarizing data.

② Inferential statistics

- It is a technique where we used the data that we have measured to form conclusions.

Ex: Classroom of Maths student there are around 20 students and now I want to find Marks of 1st semester. Marks with respect to percentage are { 84, 86, 78, 72, 75, 65, 80, 81, 92, 95, 96, 97, ... }

what type of question comes in Descriptive statistics?

- Average { Mean, Median, Mode }

→ what is the average marks of the students in the class

{ Standard Deviation } etc

→ what is the percentage of student passing out from class?

what type of question comes in Inferential statistics?

→ Are the marks of the students of this classroom similar to the marks of the Maths classroom in the college?

(N)

Population and Sample (n)

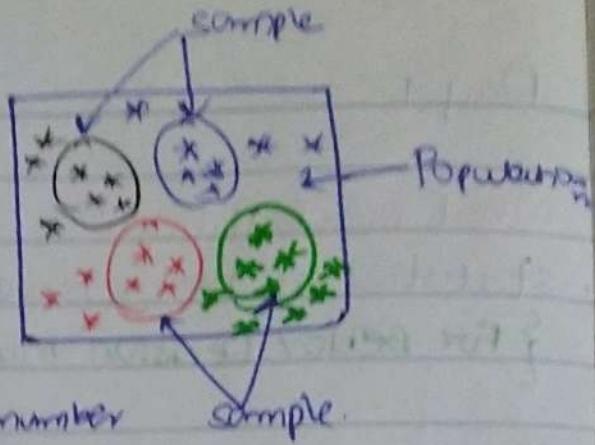
Ex:- Elections \rightarrow Exist Poll \rightarrow Exit Poll

\rightarrow In Exit poll, we cannot ask every person that

where did you vote. so we take sample of

data from different different particular region and maximum number

of vote will create exit poll.



Sampling Techniques

① Simple Random Sampling :- Pick up any people Eg:- Exit Poll

\rightarrow When performing simple Random Sampling, every member of the Population (N) has an equal chance of being selected for your sample (n).

② Stratified Sampling :- Stratified sampling is a technique where the

Population (N) is split into non-overlapping groups. (strata layers)

Ex :- Gender

- \rightarrow Male \rightarrow survey
- \rightarrow Female \rightarrow survey

• Age (0-10) (10-20) (20-30) (40-100)

• Profession Doctor / Engineer

③ Systematic Sampling :- Here from the Population (N) we pick up any n^{th} individual.

Eg:- Mail \rightarrow survey (could) \rightarrow 8th person \rightarrow survey.

④ Convenience Sampling / Voluntary Response Sampling :- Only those people who are basically interested in this will basically do it.

Eg:- Data Science Survey.

Note:- ① Exit Poll { Simple Random Sampling }

② LSI \rightarrow Household \rightarrow women { Stratified Sampling }

③ Drug \rightarrow Test { To whom to give given based upon that we will select or decide }

Variation 3 - A variable is a property that can take on many values.
 E.g.: Height = {130, 160, 170, 180} cm
 weight = {30, 40, 50, 60} kg

Kinds of variables :-

① Quantitative Variable
 - Measured numerically.

Eg:- Age, Height
 we can add, subtract etc

② Qualitative / Categorical Variable.

- Based on some characteristic we can derive categorical variable.

Eg:- Gender → Male we cannot add,
 Female subtract etc

Eg:- Blood Group, T-shirt size, IQ.

0-10	10-50	Blood IQ
Less IQ	Medium IQ	High IQ

Quantitative Variable.

Discrete Variable

Eg: Whole Number
 No. of Bank Account 2, 3, 4.
 Total Number of children in family
 2, 3, 4.

Continuous Variables

Eg:- Height = {162.3 cm, 170.1 cm}
 weights = {100kg, 92.3kg}

Variables Measurement Scales.

- ① **Nominal Data** {Categorical Data} → split in classes → Eg: Gender, Color, etc
 ② **Ordinal Data** → Order of the data matters, but values does not.

Eg:- Students (Marks) Rank }
 100 1 }
 96 2 }
 57 3 }
 85 4 }
 44 5 }

- ③ **Interval Data** → Order matters, values also matter, natural zero is not present.

Eg:- Temperature (F) 70-80 80-90 90-100

Distance 10-20 20-30 30-40

Ratio :-

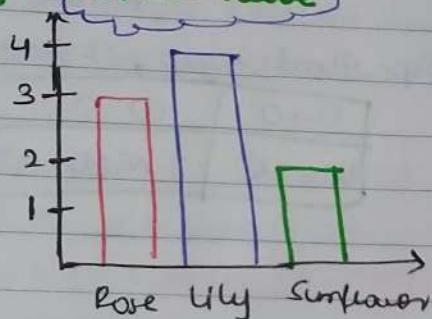
(4)

Frequency Distribution

Sample dataset = Rose, Lilly, Sunflower, Rose, Lilly, Sunflower, Rose, Lilly, Lilly

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

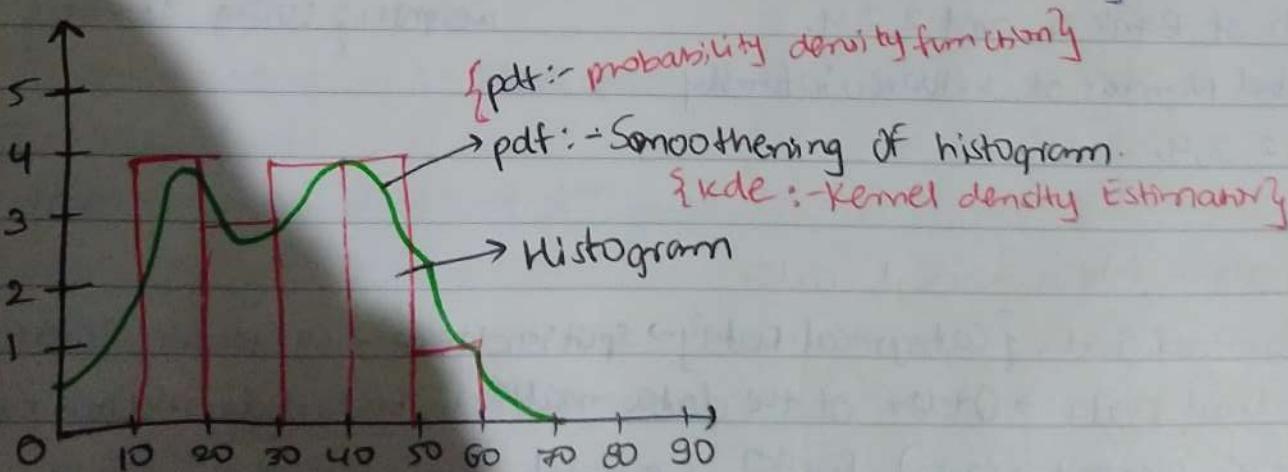
Bar Graph :- Discrete value



Summarizing the data

Histogram :- Continuous value

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



Day 2
① Arithmetic
Mean
P...
 $x = \frac{1+2+2}{3} = 1.66$

② M

- R

③ N

D

Day 2

① Arithmetic mean for Population and Sample

Mean (Average)

Population (N)

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

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

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

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

$$\mu = 3.2$$

Sample (n)

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

$$\bar{x} = 3.2$$

② Measure of Central Tendency

① Mean ② Median ③ Mode

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

③ Median

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

$$\mu = 3.2 \leftarrow \text{Mean}$$

$$\text{Data} = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100\}$$

$$\mu = 12 \leftarrow \text{Mean}$$

$$\text{Data} = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100\}$$

Median :- Step ① Sort the data.

Step ② No. of element : 11

Step ③ Selectable middle number

$$1, 1, 2, 2, 3, 4, 5, 5, 6, 100$$

$$\text{Median} = 3$$

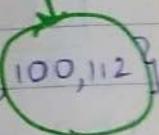
Outliers



$$\text{Data} = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100, 112\}$$

Updated

Outliers



Step ① : Sort the data

Step ② : No. of element : 12

Step ③ : Select middle 2 numbers

$$1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100, 112$$

$$\text{Median} = \frac{3+4}{2} = \frac{7}{2} = 3.5$$

Note:- Median work very well with the outliers.

HW

(1)

④ Mode :-

Most frequent Element

Data :- {1, 2, 2, 3, 4, 5, 6, 6, 6, 6, 7, 8, 100, 200}

$$\boxed{\text{Mode} = 6}$$

→ measure of central tendency

Data :- {1, 2, 2, 3, 4, 5, 6, 6, 6, 6, 7, 8, 100, 100, 100, 100}

$$\boxed{\text{Mode} = 100}$$

Note :- Huge difference between both the mode due to outlier.

Outliers

Where to use Mode? → works well with categorical data.

Measure of Dispersion :-
① Variance

② Standard Deviation

Data = {1, 1, 2, 2, 4}

$$\boxed{\mu = 2}$$

Data = {2, 2, 2, 2, 2}

$$\boxed{\mu = 2}$$

Now how we will identify that the above two data are totally different from each other. For that we use Variance.

① Variance

Population Variance (σ^2)

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

Sample Variance (s^2)

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

Data : {1, 2, 2, 3, 4, 5}

page 10 explained

Population
② Standard Deviation

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

Sample Standard Deviation

$$s = \sqrt{\text{Variance}}$$

Vari
S.D

Percent

Percent
%

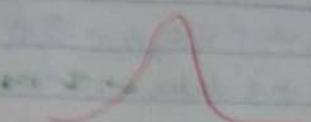
Percentile

HW

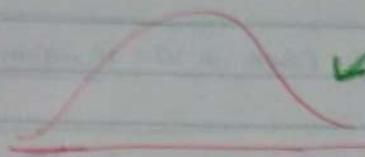
(6)

Q3

x	μ	$x - \mu$	$(x - \mu)^2$
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



Variance is more?



mode

$$\mu = 2.83$$

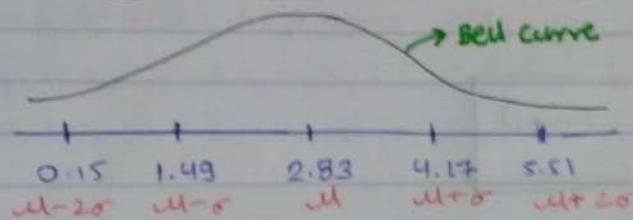
$$\Sigma = 10.84$$

$$\sigma^2 = \frac{10.84}{6} = 1.81$$

Variance

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

Standard Deviation



Variance = Spreadness of data

S.D = From mean how far the data is

totally

Percentiles and Quartiles

Percentage :- 1, 2, 3, 4, 5

% of the numbers that are odd?

$$\% = \frac{\text{Number of numbers that are odd}}{\text{Total Number}} = \frac{3}{5} = 0.6 = 60\%$$

Percentile 3 - A percentile is a value below which a certain percentage of observation lie.

Data = {2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12}

Note 1: Don't do the data ~~as it's more not relevant~~

What is the percentile ranking of 10?

$$\text{Percentile Rank of } 10 = \frac{\text{Number of value below } 10}{n} \times 100$$

$$= \frac{16}{20} \times 100$$

$$\text{Percentile Rank of } 10 = 80\%$$

What value exists at percentile ranking of 25%?

$$\text{Value} = \frac{\text{Percentile}}{100} \times (n+1)$$

$$= \frac{25}{100} \times 21$$

$$= 5.25 \rightarrow \text{Index value.}$$

Index element at 5

Index element at 5 is 5

Index element at 6

Index element at 6 is 5

$$\therefore \text{Value} = \frac{5+5}{2}$$

$$\text{Value} = 5 \rightarrow 25\%$$

Five Number Summary { Removing the outliers }

① Minimum

② First Quartile (Q_1) 25%

③ Median

④ Third Quartile (Q_3) 75%

⑤ Maximum

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

Step 1) Sort the Data. ~~order recursive.~~

[Lower fence \leftrightarrow Higher fence]

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

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

$$\text{IQR} = Q_3 - Q_1$$

$$Q_1 = 3$$

$$Q_3 = 7$$

$$\text{Lower fence} = Q_1 - 1.5(\text{IQR}) = 3 - 1.5(4) = 3 - 6 = -3$$

$$\text{Higher fence} = Q_3 + 1.5(\text{IQR}) = 7 + 1.5(4) = 7 + 6 = 13$$

[-3 \leftrightarrow 13]

Data = [1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9]

Minimum = 1

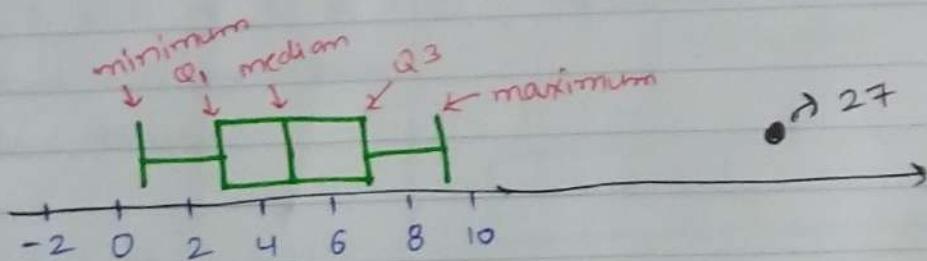
Q₁ = 3

Q₃ = 7

Median = 5

Maximum = 9

BOX PLOT



Note:- Box Plot is used to see the outlier in visualization.

why Sample Variance is divided by $n-1$?

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

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

Denominator :- Small

$$\therefore s^2 = \text{Large}$$

Denominator : Large

$$\therefore s^2 = \text{Small}$$

Because we are underestimating true variance.

$$\therefore s^2 \approx \sigma^2$$

Day 3
① Distribution
App = {24,
Gauss}

Eg:-

Eg:-

more
of sta

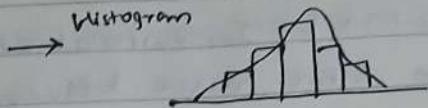
Z_{SOM}

more

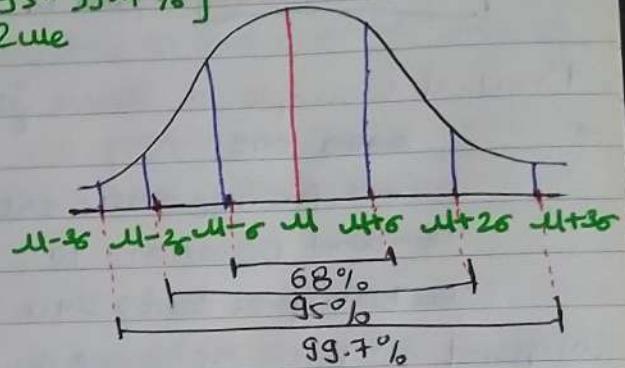
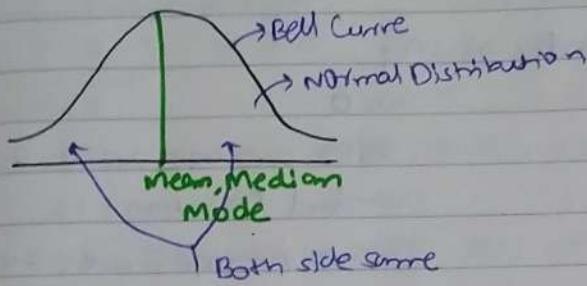
10

Day 3

- ① Distribution of Dataset
- $$\text{Age} = \{24, 26, 27, 28, 30, 32, \dots\}$$

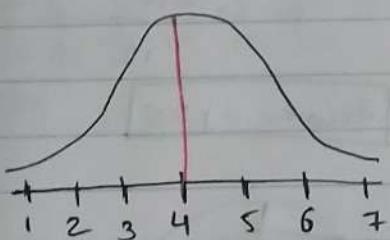


* Gaussian / Normal Distribution $\{68-95-99.7\% \}$ Rule



Eg:- Height \rightarrow Normally Distributed \rightarrow Domain Expert \rightarrow Doctor weight, Iris dataset

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



Where does 4.75 lies in terms of standard deviation?

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

$$= \frac{4.75 - 4}{1}$$

= 0.75 sd to right

Where does 3.75 lies?

$$= \frac{x_i - \mu}{\sigma} = \frac{3.75 - 4}{1} = -0.25 \text{ sd to left}$$

Now, what happen if we apply z-score to every element. \rightarrow Normal Distribution

$D_1 = \{1, 2, 3, 4, 5, 6, 7\} \rightarrow$ Initially data.

$D_2 = \{-3, -2, -1, 0, 1, 2, 3\} \rightarrow$ After applying z-score for every element
 \downarrow Standard Normal Distribution
 $(\mu=0, \sigma=1)$

$$y \sim \text{SN}(\mu=0, \sigma=1)$$

* Now, why we do this?

Age (years)	Salary (Rs)	Weight (kg)
-------------	-------------	-------------

\rightarrow Converting the data into same scale and the process is known Standardization.

(12)

Normalization :- In normalization we have an option i want to shift the entire value between 0 to 1. or (-1 to 1)

MinMaxScalar → used in Image classification

Practical Example of Zscore {India vs SA}

→ 2020 2021

Series Average 2021 = 250

Standard Deviation = 10

• Bas Team final score = 240

2020

Series average 2020 = 260

Standard Deviation = 12

Team Average score = 245

Compared to both the series in which year team perform better?

2021

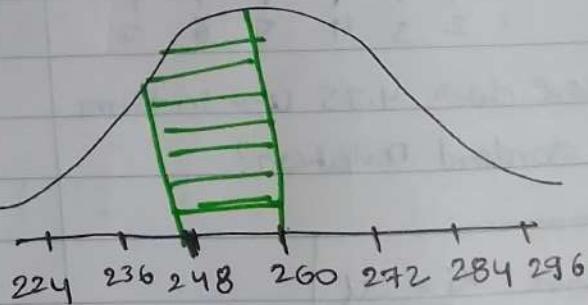
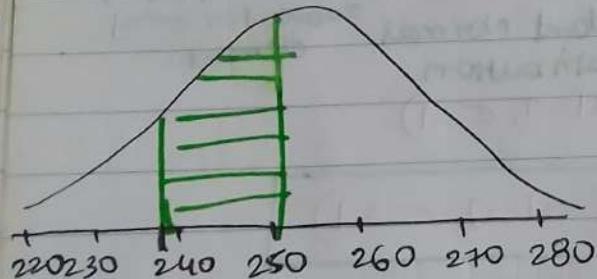
$$\begin{aligned} Z\text{score} &= \frac{x_i - \mu}{\sigma} = \frac{240 - 250}{10} \\ &= \frac{-10}{10} \end{aligned}$$

∴ $Z\text{score} = -1$

2020

$$\begin{aligned} Z\text{score} &= \frac{x_i - \mu}{\sigma} = \frac{245 - 260}{12} \\ &= \frac{-15}{12} \end{aligned}$$

$Z\text{score} = -1.25$



Team perform well in the year 2021

IMP IMP
Interview

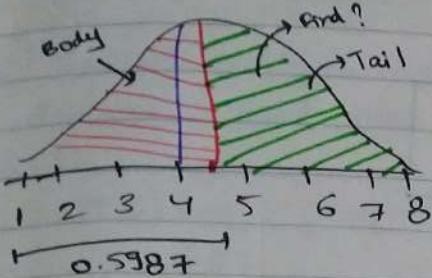
body
+ + +
1 2 3
← 0.5

a) In
per
85

find
y
+
55
← 0

(12)

want to

IMP IMP
Interviews-

Question 3 - what percentage of scores fall about

$$4.25? \mu = 4, \sigma = 1$$

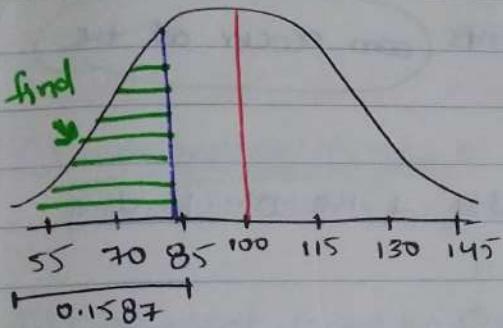
$$z = \frac{x_i - \mu}{\sigma} = \frac{4.25 - 4}{1} = 0.25$$

Z score actually help to find out the area of the body curve.

$$z = 0.25 = 0.5987 = \text{Body}$$

$$\therefore \text{Area Final} = 1 - 0.5987 = 0.4013 = 40\%$$

- a) 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?



$$z = \frac{x_i - \mu}{\sigma} = \frac{85 - 100}{15} = \frac{-15}{15} = -1$$

$$z = -1 = 0.1587$$

$$\therefore \text{Area Final} = 0.1587 = 15\%$$

Day 4

(iv)

Probability :- Probability is a measure of the likelihood of an event
Eg:- Rolling a dice = {1, 2, 3, 4, 5, 6}

$$P(6) = \frac{\text{Number of ways an event can occur}}{\text{Number of possible outcomes}} = \frac{1}{6}$$

Addition Rule (Probability, "or")

① Mutual Exclusive Event :- Two events are mutual exclusive if they cannot occur at the same time.

Eg:- Rolling a dice = {1, 2, 3, 4, 5, 6}
Tossing a coin = {H, T}

② Non Mutual Exclusive Event :- Multiple events can occur at the same time.
Eg:- Deck of card {Q, K}

question) If I toss a coin, what is the probability of the coin landing on heads or tails?

Sol:- Mutual Exclusive {Addition Rule}

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

Q2) Roll a dice $\rightarrow P(1 \text{ or } 3 \text{ or } 6)$

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

$$= \underline{\underline{3/6}}$$

$$= \underline{\underline{1/2}}$$

Non Mutual Exclusive Event {Addition Rule}

- Q) You are picking a card randomly from a deck. What is the probability of choosing a card that is Queen or a heart?

Sol:-

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

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

Multiplication Rule (Probability, "and")

- ① Independent Events :- ~~Each~~ Every one independent

Eg:- Roll a dice = {1, 2, 3, 4, 5, 6}

- ② Dependent Events :- One event affect the another event.

$$\begin{array}{c} \text{○ ○ ○} \\ \text{○ ○} \end{array} \rightarrow P(E) = \frac{3}{5} \Rightarrow \begin{array}{c} \text{○ ○} \\ \text{○ ○} \end{array} \rightarrow P(G) = \frac{2}{4}$$

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

Sol:- Independent Event {Multiplication Rule}

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

$$\therefore P(5 \text{ and } 4) = P(5) * P(4)$$

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

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

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

Dependent Event {Multiplication Rule}

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

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

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

$$= 0.0060$$

Permutation and Combination

① Permutation :- All the possible arrangement of the element in any order.

$$n_{P_r} = \frac{n!}{(n-r)!}$$

(EX :- School trip to Chocolate factory → Dairy, Sweets, Milky Bar, Eclairs, Cakes, Silk.

~~Student one~~ Task :- whichever chocolate you see and write only 3 entry

$$\text{Student} \rightarrow \underline{6} \times \underline{5} \times \underline{4} = 120 \quad \begin{cases} \text{Dairy, Cereals, Milky} \\ \text{Milky, Cereals, Dairy} \end{cases}$$

Using formula;

$n = 6$... option $n = 3$... to write

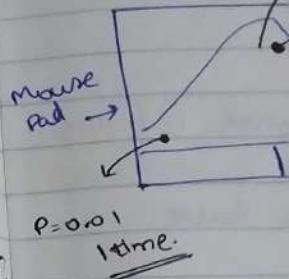
$${}^n P_r = \frac{n!}{(n-r)!} = \frac{6!}{3!} = \frac{6 \times 5 \times 4 \times 3!}{3!} = 120$$

② Combinations:- Only unique entry no swaping or rearranged allowed

{ Dairy , Ghee , Milky } → allowed
{ milky , Ghee , Dairy } → Not allowed

$${}^nC_r = \frac{n!}{r!(n-r)!} = \frac{6!}{3!3!} = \frac{6 \times 5 \times 4 \times 3!}{3! \times 3!} = 20 //$$

valeur



Hypothesis Testing

Coin → Test w

To check mete

Hypothesis Test

- i Null Hypothesis
 - ii Alternative Hypothesis
 - iii Experiment
 - iv Reject or Accept

We perform
should be
be away from

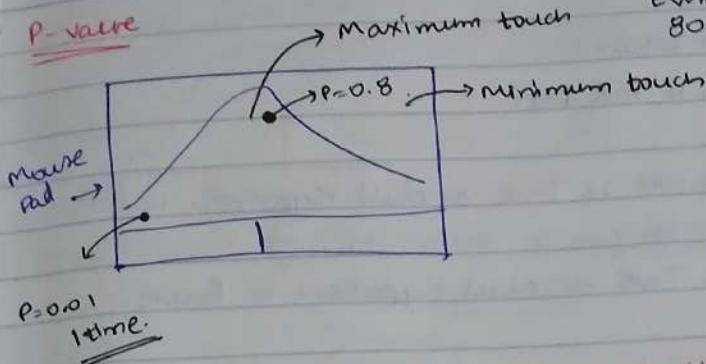
(16)

(17)

my order

Eclairs,

only 3 entry?



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

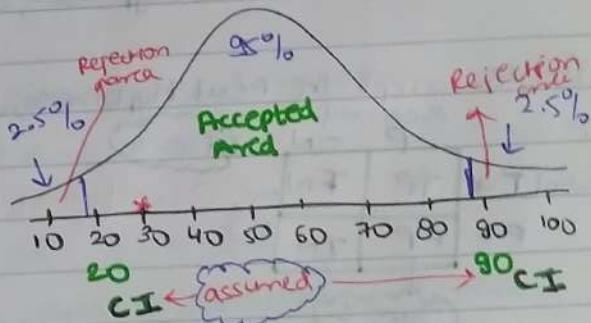
Hypothesis Testing, Confidence Interval, Significance Value,

Coin \rightarrow Test whether this coin is a fair coin or not by perform 100 tosses.

To check whether coin is fair the probability should be $P(H) = 0.5$ $P(T) = 0.5$
i.e. 50 times head ^{and} 50 times tail.

Hypothesis Testing

- i) Null Hypothesis :- Coin is fair
- ii) Alternate Hypothesis :- Coin is unfair
- iii) Experiment ... z test, t test ...
- iv) Reject or accept the Null Hypothesis



We perform 100 tosses we get 30 times head. Now we over experiment should be near to the mean. How do we defined that how far it can be away from the mean? Significance Value. $\alpha = 0.05 = 5\%$ {domain expert}
 $\therefore 100\% - 5\% \rightarrow 95\% \rightarrow$ Confidence Interval

Day 5.

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.

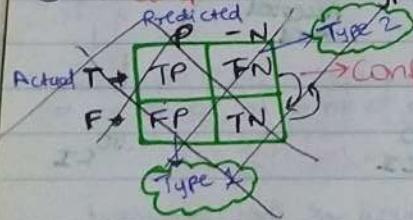
Outcome :-

① we reject the Null Hypothesis, when in reality it is false \rightarrow Yes

② we reject the Null Hypothesis, when in reality it is true \rightarrow Type 1 Error

③ we accept the Null Hypothesis, when in reality it is false \rightarrow Type 2 Error

④ we accept the Null Hypothesis, when in reality it is true \rightarrow Yes

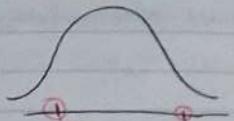


1 Tail and 2 tail Test

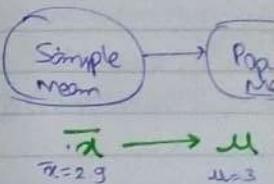
Q] Colleges in Karnataka have an 85% placement rate. A new college was recently opened and it was found that a sample of 150 students had a placement rate of 88% with a standard deviation 4%. Does this college has a different placement rate? $\alpha = 0.05$

\rightarrow Two tail test because percentage can increase or decrease of placement rate.

Confidence Interval



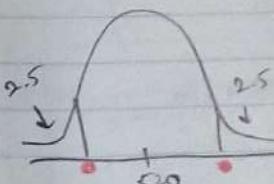
Point Estimate :- The



Confidence Interval :-

a) On the world test of to be 100. A sample of 95% CI about the

Sol:- $\sigma = 100$, $n=25$,

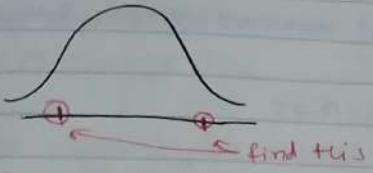


$$\therefore CI = \bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

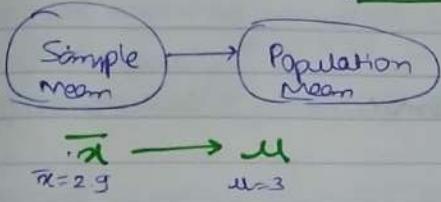
$$= 520 \pm \frac{z_{0.05}}{\sqrt{25}} \left(\frac{100}{n} \right)$$

$$= 520 \pm 20.025$$

(18) Confidence Interval



Point Estimate :- The value of any statistic that estimates the value of a parameter.

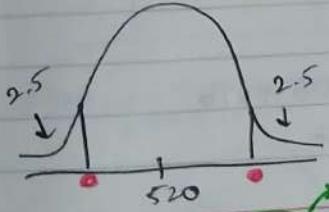


Confidence Interval = Point Estimate ± Margin Error

- a) On the quant test of CAT exam the population 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?

$$\text{Sol: } \sigma = 100, n = 25, \alpha = 0.05, \bar{x} = 520$$

Note:- Population standard deviation is given then we z-test, $n \geq 30$



$$\therefore \text{CI} = \bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$$= 520 \pm \frac{z_{0.05}}{2} \left(\frac{100}{\sqrt{25}} \right)$$

$$= 520 \pm z_{0.025} \left(\frac{100}{5} \right)$$

$$= 520 + z_{0.025}(20) : = 520 - z_{0.025}(20)$$

$$\text{Now, } 1 - 0.025 = 0.975$$

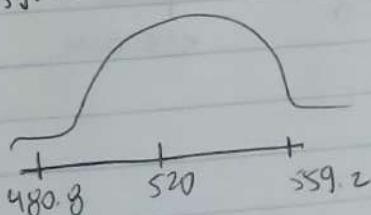
$$\text{Now } 0.975 = 1.96$$

$$\therefore \text{CI} = 520 + 1.96(20)$$

$$= 559.2$$

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

$$= 480.8$$



(20)

Q) On the aptitude 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?

$$\text{Sol: } n = 25, \bar{x} = 520, s = 80, \alpha = 0.05$$

Note:- Here population standard deviation is not given then we will use t-test.

CI = Point Estimate ± Margin of Error

$$CI = \bar{x} \pm t_{\alpha/2} \left(\frac{s}{\sqrt{n}} \right)$$

standard Error

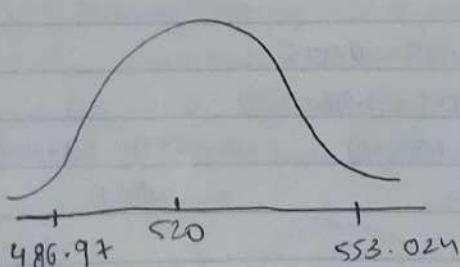
$$\text{Upper Bound} = \bar{x} + t_{0.05/2} \left(\frac{s}{\sqrt{n}} \right)$$

Imp Degree of freedom = 25 - 1 = 24

$$\therefore \text{Upper Bound} = 520 + t_{0.025} \left(\frac{80}{\sqrt{25}} \right)$$

$$= 520 + 2.064$$

$$= 553.024$$



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

$$\text{Lower Bound} = 520 - t_{0.025} \left(\frac{80}{\sqrt{25}} \right)$$

$$= 520 - 2.064 \times 16$$

$$= 486.97$$

One Sample z-test

- ① Population standard deviation σ
- ② $n > 30 \rightarrow$ Sample Size

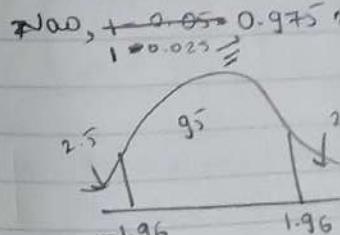
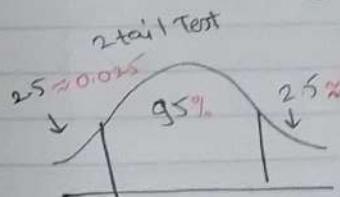
a) In the population, the average researchers wants to test effect on intelligence, or no who have taken the medication the intelligence? $\alpha = 0.05$,

Sol) ① Define Null Hypothesis (H_0)

② Alternative Hypothesis (H_1) = ...

③ State Alpha $\alpha = 0.05$

④ State Decision Rule { Two T Test }



⑤ Calculate Test Statistic

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

(20)

(21)

has a mean
confidence

One Sample z-test

- ① Population standard deviation is given
- ② $n \geq 30 \rightarrow$ Sample Size

a) In the population, the average IQ is 100 with a standard deviation of 15. Researchers wants to test a new medication to see if there is any effect on intelligence, or no effect at all. A sample of 30 participants who have taken the medication has a mean of 140. Did the medication affect the intelligence? $\alpha = 0.05$, CI = 95%

Sol) ① Define Null Hypothesis (H_0) = $\mu = 100$

$$z = \frac{140 - 100}{15/\sqrt{30}}$$

② Alternative Hypothesis (H_1) = $\mu \neq 100$

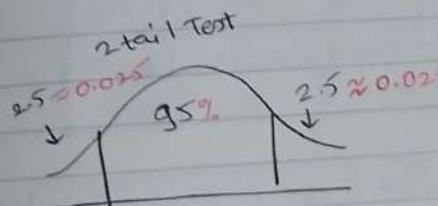
$$= \frac{40}{15/\sqrt{30}}$$

③ State Alpha $\alpha = 0.05$

$$= \frac{40}{15} \times \sqrt{30}$$

④ State Decision Rule {Two Tails Test}

$$z = 14.60$$



⑥ State Our Decision

$$14.60 > 1.96$$

\therefore If z is less than -1.96 or greater than 1.96 , reject the Null Hypothesis.

\therefore Medication increases IQ.

⑤ Calculate Test statistic.

$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \rightarrow \text{standard Error}$$

One Sample t test

$\rightarrow \mu = 100, n = 30, \bar{x} = 140, s = 20, \text{ Did medication affect intelligence?}$
 $\alpha = 0.05$

Soln -

$$\textcircled{1} H_0 = \mu = 100$$

$$\textcircled{2} H_1 = \mu \neq 100$$

\textcircled{3} Degree of freedom

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

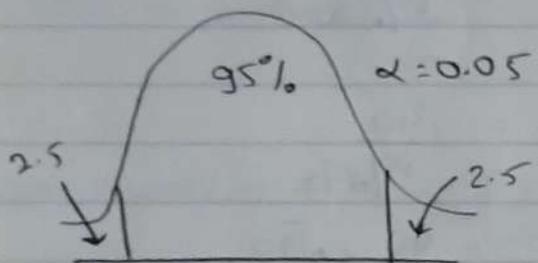
$$\therefore t = \frac{40}{20} \times \sqrt{30}$$

$$= \frac{20 \times \sqrt{30}}{2} \\ = 2 \times \sqrt{30}$$

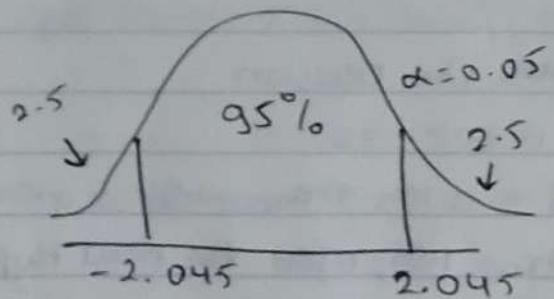
$$\boxed{t = 10.95}$$

\textcircled{4} State Decision Rule

\textcircled{6} We Reject Null Hypothesis.



$$\therefore t_{0.05} = 2.045, \text{ DDF} = 29$$



\textcircled{5} T test

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$$

$$= \frac{140 - 100}{20/\sqrt{30}}$$

Day - 6.

Chi Square Test : claims about population proportions. It is a non parametric test that is performed on categorical (nominal or ordinal) data.

Eg:- 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 $n=500$ individual were sampled. Below are the results.

< 18	18 - 35	> 35
121	288	91

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

Population 2000		
< 18	18 - 35	> 35
20%	30%	50%

$n=500$		
Observed		
$500 \times 20\%$	$500 \times 30\%$	$500 \times 50\%$
100	150	250

→ Expected

Age Category ($n=3$)		
Observation		
Expected		
121	288	91
100	150	250

Using Chi Square

Table df = 2

and 0.05 CI

is equal to

5.991

① H_0 = The data meets the distribution 2000 censit.

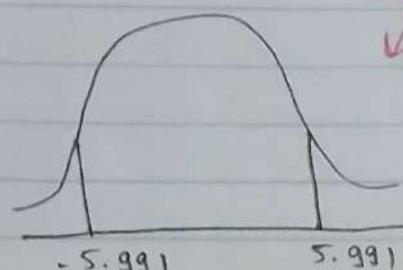
H_1 = The data does not meet the distribution 2000 censit.

② $\alpha = 0.05$ (95% CI)

③ Degree of freedom = $n-1 = 3-1 = 2$

do

④ Decision Boundary



P.T.O. →

If χ^2 is greater than 5.991 reject (H_0) Null Hypothesis

⑤ calculate test statistics.

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

$f_o \rightarrow$ Observed

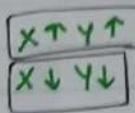
$f_e \rightarrow$ Expected

$$\begin{aligned} \chi^2 &= \sum \frac{(121-100)^2}{100} + \frac{(288-150)^2}{150} + \frac{(91-250)^2}{250} \\ &= 23.4 \\ \chi^2 &= 23.44 \end{aligned}$$

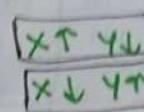
$$\chi^2 = 23.44 > 5.99 \quad \left\{ \text{Reject Null Hypothesis} \right\}$$

Covariance

weight(x)	height(y)
50	160
60	170
70	180
75	181



NO. of hours study	Play
2	6
3	4
4	3



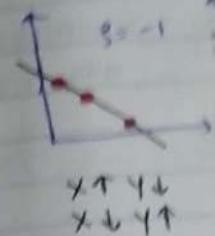
$$\text{Cov}(x,y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1} = \text{+ve or -ve or } 0 \rightarrow \text{No relation}$$

Disadvantage of Covariance

How much negative correlated or how much positive correlated.

Positive Correlation
The more towards
The more towards

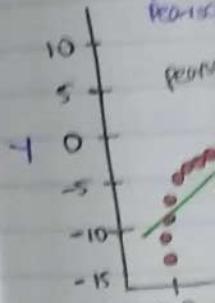
$$S(x,y) = \dots$$



Negative Correlation

Spectrum

peaks



Eg:- Height (Y)

140

160

150

145

180

165

(24)

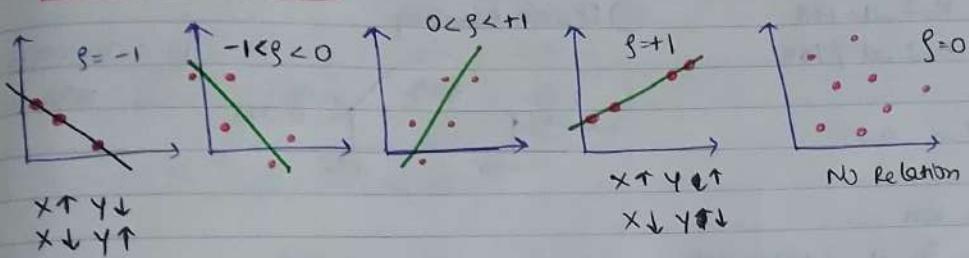
Pearson Correlation Coefficient $\{-1 \text{ to } 1\}$

The more towards +1 more positively correlation.

The more towards -1 more negatively correlation.

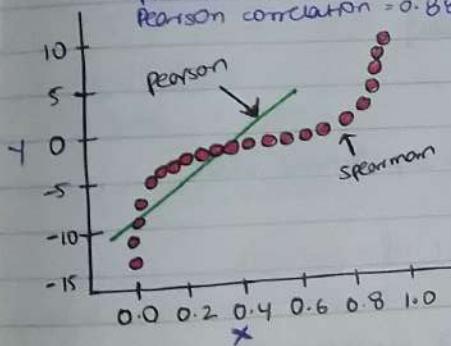
$$\rho_{(x,y)} = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y}$$

Note:- It capture only linear property.



Spearman's rank correlation coefficient

Spearman correlation = 1
Pearson correlation = 0.88



$$\text{Spear}(\pi, y) = \frac{\text{Cov}(R(x), R(y))}{R \sigma_x R \sigma_y}$$

Note:- It capture both linear and non linear property.

Eg:-

Height (x)	Weight (y)	R(x)	R(y)
170	75	2	2
160	62	3	3
150	60	4	4
145	55	5	5
180	85	1	1

Ignore

will be used in Spearman Rank

Day 7

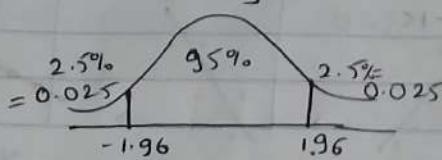
P-value and significance value

→ 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. C.I. = 95%

$$\rightarrow \mu = 168, \sigma = 3.9, n = 36, \bar{x} = 169.5 \quad \text{C.I.} = 95\% \quad \text{no} \quad [C.I. = 168 \pm 0.05 = 167.95 \text{ to } 170.05]$$

$$\begin{aligned} \textcircled{1} \quad H_0: \mu = 168 \\ H_1: \mu \neq 168 \end{aligned}$$

$$\textcircled{2} \quad \alpha = 0.05$$

Decision Boundary

$$1 - 0.025 = 0.975$$

$$z_{\text{crit}} = -1.96$$

$$\textcircled{4} \quad z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{169.5 - 168}{3.9/\sqrt{36}} = 2.307$$

$$\textcircled{5} \quad z = 2.307$$

P-value

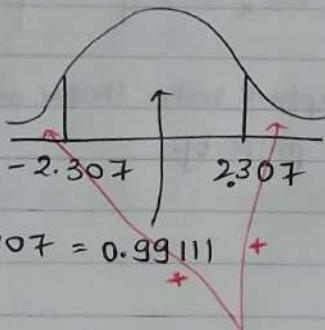
$$z = 2.307$$

$$1 - 0.99111 = 0.00889$$

$$\text{Now, } 0.00889 = 0.0044$$

$$P\text{-value} = 0.0044$$

$$= 0.0000$$



$$\begin{aligned} z &= 2.307 = 0.99111 \\ \text{table} & \quad + \\ & \quad + \end{aligned}$$

$$\text{Now, } 1 - 0.99111 = 0.00889$$

$$\begin{aligned} \text{Now, } \frac{0.00889}{2} &= 0.0044 \\ \text{Both sides!} & \quad ? \end{aligned}$$

Note

p-value < significance value.
↓

Reject Null Hypothesis

p-value > significance value
↓

Fail to Reject the Null Hypothesis

$$\begin{aligned} \text{Now, pvalue} &= 0.0044 + 0.0044 \\ &= 0.0088 \end{aligned}$$

p value < 0.05

0.0088 < 0.05 {Reject Null Hypothesis}

(27)

→ Average age of a college is 24 years with a standard deviation 1.5. Sample of 36 students mean is 25 years. with $\alpha = 0.05$, C.I = 95%, do the age vary?

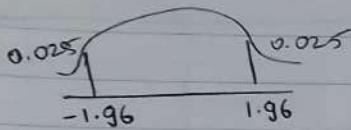
$$\text{Sol: } \mu = 24, \sigma = 1.5, n = 36, \bar{x} = 25, \alpha = 0.05$$

$$H_0: \mu = 24$$

$$H_1: \mu \neq 24$$

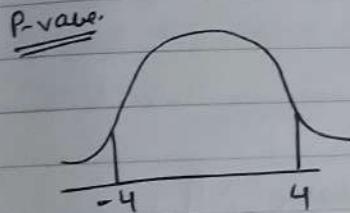
$$\alpha = 0.05$$

(3) Decision Rule.



$$(4) z\text{-score} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{25 - 24}{1.5/\sqrt{36}} = \frac{1}{1.5/\sqrt{36}} = \underline{\underline{z=4}}$$

(5) $z=4 > 1.96$ {Reject Null Hypothesis}



$$z=4 = 0.9997$$

$$1 - 0.9997 = 0.00003$$

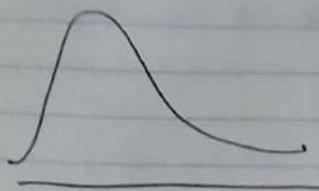
$$\text{Now, } \frac{0.00003}{2} = 0.000015$$

$$p\text{-value} = 0.000015 + 0.000015$$

$$= 0.00003$$

$0.00003 < 0.05$ {Reject Null Hypothesis}

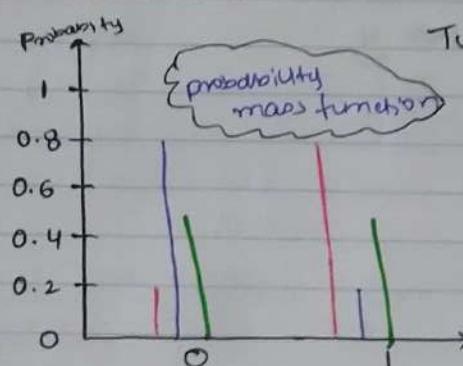
Log Normal Distribution



Eg:- ① Wealth distribution ② People writing comments.

$y \approx \text{Log Normal Distribution} \rightarrow \log(y) \rightarrow \text{Normal Distribution}$

Bernoulli Distribution



Two Outcomes {0 or 1}

$$p = 0.5$$

$$q = 1 - p = 1 - 0.5 = 0.5$$

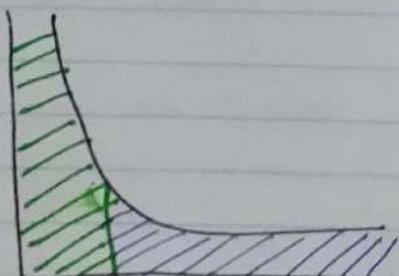
NOTE:- probability density function is for continuous variable (pdf)
probability mass function is for categorical variable (pmf)

Binomial Distribution :- Combination of Bernoulli Distribution.

$$B(n, p)$$

n - number of trials, p = probability.

Pareto Distribution {Power Law} {80-20}



Eg:- 80% of wealth is distributed with 20% of the people.

Eg:- 80% of the company project are done by 20% of the people in a team

Cont. not
what was
focus n

S₁

S₂

S₃

S₄

S₅

Anova
compar
Medicat

Garder

Type

① One

② Repeate

Depen

③ Fact

Ver

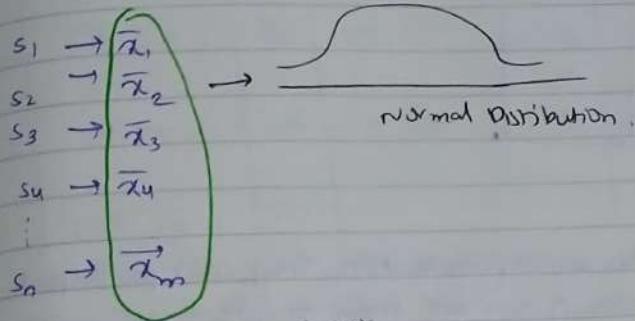
28

Central Limit Theorem

whatever distribution it may be, the sample of the distribution will follow Normal Distribution $n \geq 30$

gements.

y) → Normal Distribution



29

Anova {Analysis of Variance} → Anova is a statistical method used to compare the means of 2 or more groups. ① Factors (variable) ② levels → IMP

Medicine.

Dosage:	0mg	50mg	100mg
9	7	4	
8	6	3	
7	6	2	
8	7	3	
8	8	2	

factors:- Dosage

levels :- 0mg, 50mg, 100mg

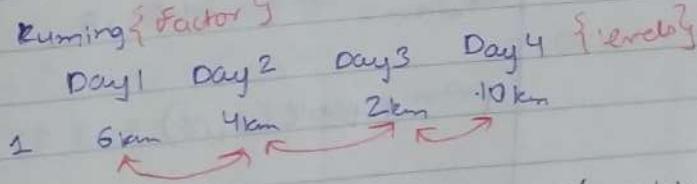
Gender :- Factors, Male, Female :- levels.

Types of Anova

① One way Anova :- One factor with atleast 2 levels, levels are independent
Eq:- Dosage.

② Repeated Measure Anova :- One factor with atleast 2 levels, but the levels are done by

dependent running {Factor}



③ Factorial Anova :- Two or more factors each of which with atleast 2 levels, levels can be either independent, dependent or both (mixed).

Eq:-

	Day 1	Day 2	Day 3
Men	9	7	4
	8	6	3
	7	5	2
Women	8	7	3
	8	8	4
	9	7	3

Factors :- Days & Gender

Level :- D1, D2, D3, Men, Women

(2)

⑤ Calculate F test

$$\text{Sum of Squares} = \sum S^2$$

$$\text{Between} = 98.67$$

$$\text{Within} = 10.29$$

$$\text{Total} = 108.96$$

$$SS_{\text{between}} = \sum (Z_i - Z_m)^2$$

$$= \sum (Z_i - Z_m)^2 = (9+8+7+8+6+6+7+8+7+8+7+6+7+6+7+7+7+8+7+7+7) = 57^2 + 1$$

$$T^2 = [57^2 + 1] = 3241$$

$$\therefore SS_{\text{between}} = 57^2 + 1 = 3241$$

$$\therefore SS_{\text{within}} = 9^2 + 8^2 + 7^2 + 6^2 + 4^2 + 1^2 = 95$$

$$SS_{\text{within}} = \sum y^2 = 9^2 + 8^2 + 7^2 + 6^2 + 4^2 + 1^2 = 285$$

$$\sum y^2 = (9^2 + 8^2 + 7^2 + 6^2 + 4^2 + 1^2) = 285$$

$$\sum y^2 = 853$$

$$\therefore SS_{\text{within}} = 853$$

$$SS_{\text{within}} = 10$$

One way Anova (F test).

Researchers want to test a new anxiety medication. They split participants into 3 conditions (0mg, 50mg, 100mg) then ask them to rate their anxiety level on scale of 1-10. Are there any difference between the 3 conditions using $\alpha = 0.05$?

0mg	50mg	100mg
9	7	4
8	6	3
7	6	2
8	7	3
8	8	4
9	7	3
8	6	2

$$\text{So 1:- } H_0 : \mu_{0\text{mg}} = \mu_{50\text{mg}} = \mu_{100\text{mg}}$$

$$H_1 : \mu_{0\text{mg}} \neq \mu_{50\text{mg}} \neq \mu_{100\text{mg}}$$

③ State the α and CI

$$\alpha = 0.05 \quad \text{CI} = 95\%$$

④ Calculate Degree of freedom

$$N = 21, n = 7$$

$\alpha = \text{level}$

$$df_{\text{between}} = a - 1 = 3 - 1 = 2$$

$$df_{\text{within}} = N - a = 21 - 3 = 18$$

$$df_{\text{Total}} = N - 1 = 21 - 1 = 20$$

⑤ State Decision Rule

$$df_{\text{between}} = a - 1 = 3 - 1 = 2 \quad (2, 18)$$

$$df_{\text{within}} = N - a = 21 - 3 = 18$$

$$\text{From F table } (2, 18) = 3.5546 \quad \{\alpha = 0.05\}$$

If f test is greater than 3.5546, reject the Null Hypothesis.

(20) ③ Calculate F test statistics.

	Sum of Square SS	df	Mean Square	F test
Between	98.67	2	= 49.34	86.56
Within	10.29	18	= 0.57	
Total	108.96	20		

$$SS_{\text{between}} = \frac{\sum (\sum a_i)^2}{n} - \frac{T^2}{N}$$

$$\begin{aligned} \sum (\sum a_i)^2 &= (9+8+7+8+8+9+8)^2 + \\ &\quad (7+6+6+7+8+7+6)^2 + \\ &\quad (4+3+2+3+4+3+2)^2 \end{aligned}$$

$$T^2 = [57^2 + 47^2 + 21^2] = 125^2$$

$$\therefore SS_{\text{between}} = \frac{57^2 + 47^2 + 21^2}{7} - \frac{125^2}{21}$$

$$\therefore SS_{\text{between}} = 98.67$$

$$SS_{\text{within}} = \sum y^2 - \frac{\sum (\sum a_i)^2}{n}$$

$$\begin{aligned} \sum y^2 &= (9^2 + 8^2 + 7^2 + 8^2 + 8^2 + 9^2 + 8^2) + \\ &\quad (7^2 + 6^2 + 6^2 + 7^2 + 8^2 + 7^2 + 6^2) + \\ &\quad (4^2 + 3^2 + 2^2 + 3^2 + 4^2 + 3^2 + 2^2) \end{aligned}$$

$$\sum y^2 = 853$$

$$\therefore SS_{\text{within}} = 853 - \left[\frac{57^2 + 47^2 + 21^2}{7} \right]$$

$$SS_{\text{within}} = 10.29$$

$$F_{\text{test}} = \frac{MS_{\text{between}}}{MS_{\text{within}}} = \frac{49.34}{0.57}$$

$$\therefore F_{\text{test}} = 86.56$$

Final Conclusion:-

$86.56 > 3.5546$, so we reject the Null Hypothesis.

under
men, women

participants
anxiety
conditions.

100mg

0mg

)

$\alpha = 0.05$

Reject