

Statistics:- the science of collecting, organizing & analyzing data.

Data: "facts or pieces of information"

2 types

Descriptive stat

Statistics

- It consists of organizing  
& summarizing of data

Inferential statistics

It consists of a technique  
to  
form some conclusions

Ex: Let's say - 20 classrooms - Aged age of a one classroom.  
- [21, 20, 18, 34, 17, 22, 24, 25, 26, 23, 22]

Random Sampling

100000

3] Systematic Sampling -  
Selecting  $n^{\text{th}}$  people from crowd.

Eg: Every 5<sup>th</sup> person come, he will  
go & ask for credit card.

4] Convenience Sampling - only those who are interested  
will be participating.

Stratified Sampling

→ Gender - male/female ✓

→ Edu. Degree - BE

master  
PHD

High school  
these  
are  
overlapping

Blood group  
✓ A+  
B+

strata  
Layers → clusters  
↓  
non overlapping

Data science  
Art  
Engg medical  
X  
✓ Interested people.

Qualitative Variable

Categorical  
Good  
Bad  
Married  
Unmarried

Histogram

Quantitative Variable

Discrete Variable

e.g.: no. of children in a family [0, 1, 2, 3]  
e.g.: whole number [1, 2, 3]

Continuous Variable

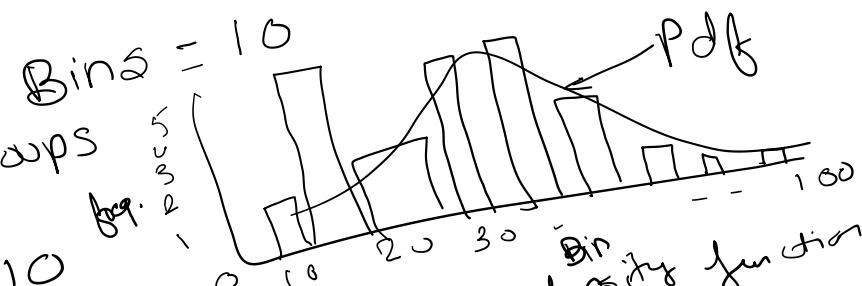
e.g. length, weight  
3.2, 5.5, 32.5

Age: {10, 12, 13, 14, 16, 17, 25, 35, 40, 65, 55, 32}

1] Sort the no. [Ascending ]  
2] Bins → {groups} or gaps → no. of groups

3] Bin size → size of  
[0 - 100]

$$\text{Bins} = 10 \quad \text{Bin}_2 = \frac{100}{10} = 10 \\ \text{Bin size} = \frac{100}{100} = 10$$



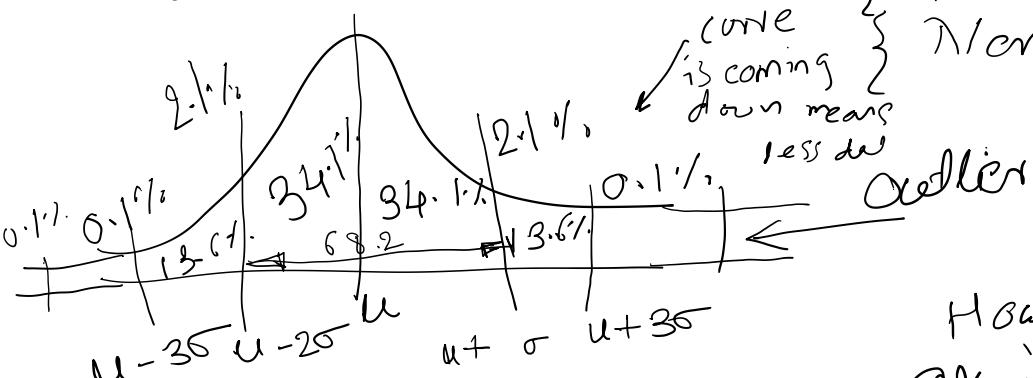
what is probability density function  
→ Smoothing histogram  
→ through kernel density estimator

$x = \{2, 4, 6, 10, 12, \dots\}$

{Empirical rule}

{Gaussian Distribution}

{Normal Distribution}



How much % of data available at center.

[68.2 - 95.4 - 99.7]

✓ Empirical Formula

Z-score

$$= \frac{x - \mu}{\sigma}$$

If expressed in standard normal distribution. Items = {1, 1, 2, 3, 4, 5, 5, 8, 7, 5, 9, 10, 10, 20} = 14 items

- used for Standardization

- means making all input units in one.

$$\text{of 5.1. percentile} = \frac{95}{100} (m) = \frac{95}{100} \times 15 = 14.25$$

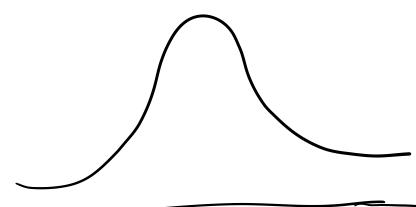
14.25 is index so we choose 20 = 20 days

Log-normal distribution

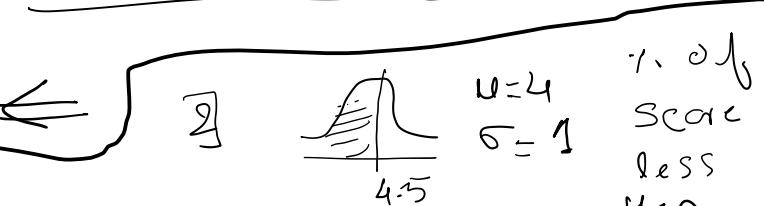
always right skewed

$$x =$$

$$\Rightarrow \log_e(u) \Rightarrow$$



$$\leftarrow e^x$$



e.g. 1



$$\mu = 4$$

$$\sigma = 1$$

Z-Score-

$$= \frac{5.5 - 4}{1} = 1.5$$



From Z table - 1.5 → 0.93319

Area of tail = 1 - 0.93319

$$AUC = 0.06681$$

Z Table -



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 lower than 85?

1] lower than 85

$$\mu = 100, \sigma = 15$$

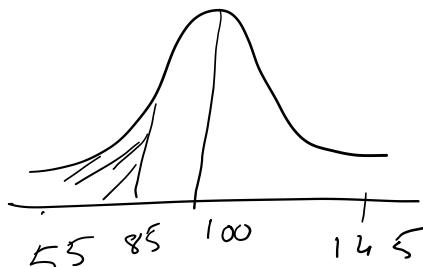
2] higher than 85

3] between 85 to 100

-1 of score  
betw 100 to 125  
 $\sim 5.15.1.$

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

$$= \frac{-15}{15} = -1$$



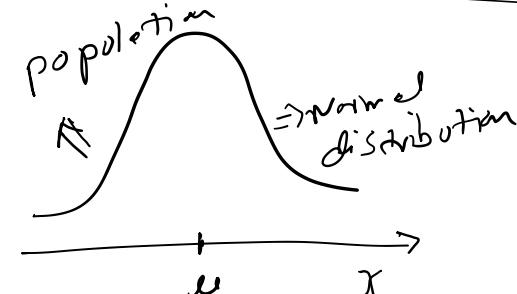
$$\text{Area of Bod} = 0.1587$$

$$= 1 - 0.1587 = 84.13\%$$



$$\begin{aligned} & 100 \quad 125 \\ & -125 - 100 = \frac{5}{15} = 0.9515 - 0.5 \\ & \sim 0.4515 \\ & \sim 45.15\% \end{aligned}$$

Central limit theorem:-



If data follows, normal or non-normal distribution,

If multiple samples i.e.  $n \geq 30$ , and if we take its sample mean, and if we plot, it follows normal or gaussian distribution.

The more the value of  $n$ , more the data follows normal distribution.

$n$  is greater than 30

$$\begin{aligned} \{n_1, \dots, n_n\} &\rightarrow \bar{x}_1 \\ \{\bar{x}_1, \dots, \bar{x}_n\} &\rightarrow \bar{x}_2 \end{aligned}$$

$$\{n_1, \dots, n_n\} \rightarrow \bar{x}$$

Probability : measure of the likelihood of an event.

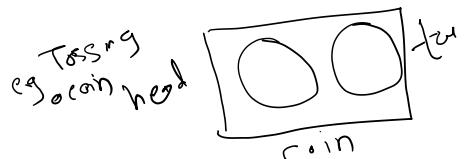
E.g. tossing a fair coin,  $p(H) = 0.5$   $p(T) = 0.5$   $Q = 1 - P$

{BERNOULLI'S DISTRIBUTION}  $\rightarrow$  TWO OUTCOMES

E.G. ROLLING A DICE? {1,2,3,4,5,6}

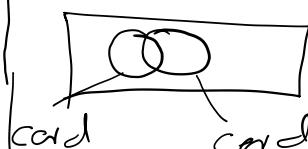
$P(1) = 1/6$   $P(2) = 1/6$

1] Mutual exclusive events: two events are mutually exclusive if they can not occur at the same time.



Rolling a dice

2] Non mutual exclusive events: Picking randomly from the deck of card, two are "hearts and king" are not mutual exclusive



e.g. Bag of marble:

Red, Green, Red & Green  
Red or green

what is the probability of coin landing on heads or tails

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

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

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

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

Addition Rule

$$= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{3}{6}$$

$$= 0.5 = \frac{1}{2}$$

\* multiplication Rule → Two events are independent if they do not affect one another

1) Independent events e.g. Tossing a coin  $P(H) = 0.5$   $P(T) = 0.5$

2) Dependent events here one event not impacting other event.

↳ Two events are dependent if they affect one another.

e.g. Bag of color marbles.  $\{ \begin{matrix} \text{o o o} \\ \times \times \end{matrix} \} \xrightarrow{\text{one removed}} 7$

$$P(O) = \frac{4}{7} \rightarrow P(X) = \frac{3}{6} \rightarrow P(O) = \frac{3}{5} \rightarrow P(X) = \frac{2}{4}$$

e.g. prob of rolling "5" & then "3"

with normal 5x5 sides die

- multiplication rule here

$$P(A \& B) = P(A) \times P(B)$$

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

Bag of marbles : 10 red, 6 green 3(red and green)

When picking randomly from a bag of marbles, what is the probability of choosing a marble that is red or green?

$$P(R) = \frac{10}{19} \quad P(G) = \frac{6}{19} \quad P(R \text{ or } G) = \frac{3}{19}$$

Addition Rule for non-mutual ex event

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

$$= \frac{10}{19} + \frac{6}{19} - \frac{3}{19} = \frac{13}{19}$$

c) orange marble  
d) yellow marble  
& 3 orange  
→ prob. of orange  
& then yellow → multiplicative rule & conditional events

$$P(O \& Y) =$$

$$= \frac{3}{12} \times \frac{3}{6} = \frac{12}{42} = \frac{2}{7}$$

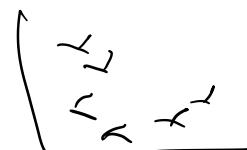
"hives Booyes"

## \* Covariance

$$\begin{Bmatrix} x_1 & y_1 \\ x_2 & y_2 \end{Bmatrix}$$



$$\begin{Bmatrix} x_1 & y_1 \\ x_2 & y_2 \end{Bmatrix}$$



Aj<	at
12	20
15	30
19	50
15	55
x_i	y_i

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

x	y
1	3
2	4
3	6

$$\bar{x} = 2$$

$$\bar{y} = 4.333$$

$$\text{Cov}(x, y) = \frac{1}{2} [(3-2)(3-4.333) + (2-2)(4-4.333) + (3-2)(6-4.333)] = 1.497$$

can be written as

$$\downarrow$$

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

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

Pearson & Spearman

Rank correlation

\* Spearman coefficient :- A Pearson coefficient

# Inferential stat.

}] P value :- Out of all

It is the probability of the null hypothesis is true?

- ① Null Hypothesis - coin is fair - ( $H_0$ )
- ② Alternative Hypothesis - coin is not fair - ( $H_1$ )
- ③ Perform Experiments

Significant value - 0.05  $\Rightarrow$  95% CI

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



Point estimate  $\pm$  Margin of error = Parameter  
Confidence Interval

$$\begin{aligned} \text{Lower fence} &= \text{Point estimate} - \text{margin of error} \\ \text{Higher fence} &= \text{Point Estimate} + \text{margin of error} \end{aligned} \quad \left. \begin{array}{l} \text{C.I} \\ \text{C.I} \end{array} \right\}$$

$$\text{Lower Fence} = \bar{x} - 3\sigma_{1/2} \frac{\sigma}{\sqrt{n}} \quad \text{Higher Fence} = \bar{x} + 3\sigma_{1/2} \frac{\sigma}{\sqrt{n}}$$

t test :-

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

$$\text{Lower fence} =$$

$\pm$  One Tail & 2 Tail

$$\text{Higher fence} =$$

# 1) ANOVA

$$1) H_0: \mu = 80 \quad \text{null Hypothesis}$$

The problem statement : The factory has a machine that fills 80ml of baby medicine in

Alternative  $H_1$  पर्याप्ति

a bottle . An employee believes the average amount of baby medicine is not 80ml

Z test

$$H_1: \mu \neq 80$$

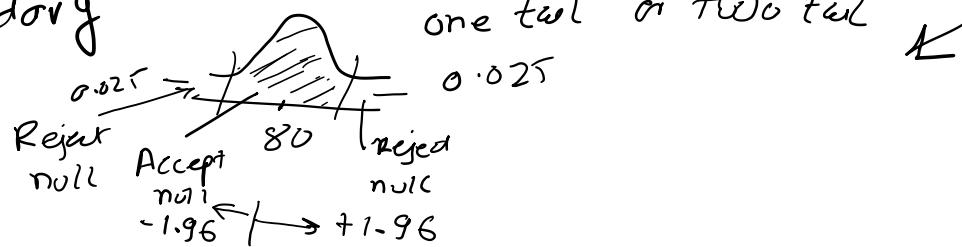
using 40 samples, he measures the average amount dispensed by the machine to be 78ml with a std. dev of 2.5

$$SV = 1 - CI$$

a] state null and alt hypothesis b] at 95% C.I, is there enough evidence to support whether machine is working properly or not.]

$$2) \text{Step} \quad \alpha = 0.05 \quad CI = 95\% \quad L = 1 - 0.95 = 0.05$$

3) Decision Boundary



$$\begin{aligned} n &= 40 \\ \bar{x} &= 78 \\ s &= 2.5 \\ Z_{\text{test}} &\approx t_{\text{test}} \\ \mu &= 80 \text{ ml} \end{aligned}$$

from Z table

4) calculate  
Test statistics

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

$$n = 40$$

sample size

$$\begin{aligned} \text{std. error} &\rightarrow S/\sqrt{n} \\ \frac{78 - 80}{2.5/\sqrt{40}} &= \frac{2 \times \sqrt{40}}{2.5} = -5.05 \end{aligned}$$

Reject

5) Decision Rule

Conclusion : If  $Z = -5.05$  is less than  $-1.96$ , we reject null hypothesis.

then we reject null hypothesis with 95% CI.

Fault in mc or  
mc is not working  
properly.

A complain was registered, the boys in a government school are underfed. average weight of the boys of age 10 is 32 kgs with s.d = 9 kg. A sample of 25 boys were selected from the government school and the average weight was found to be 29.5 kgs? with CI=95%. Choose whether it is true or false?

#### Conditions for Z-TEST

1] WE KNOW THE POPULATION S.D OR

2] WE DO NOT KNOW THE POPULATION S.D BUT OUR SAMPLE IS LARGE  $N \geq 30$ .

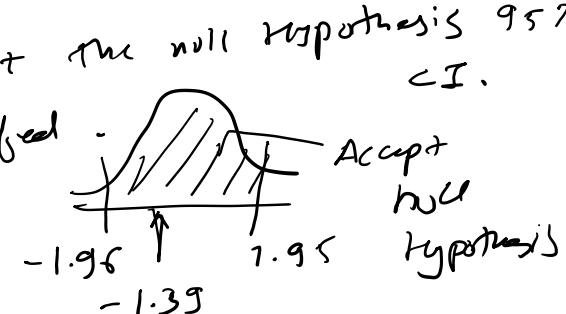
#### CONDITIONS FOR T-TEST

1] WE DO NOT KNOW THE POPULATION VARIANCE OR S.D

2] OUR SAMPLE SIZES IS SMALL I.E  $N \leq 30$

$$Z\text{-score} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{29.5 - 32}{9/\sqrt{25}} = -1.39$$

Conclusion -  $-1.39 > -1.96$  so Accept the null hypothesis 95% CI.  
Students are not underfed.



A FACTORY MFG. CARS WITH A WARRENTY OF 5 YEARS ON THE ENGINE AND TRANSMISSION. AN ENGINEER BELIEVES THAT THE ENGINE OR TRANSMISSION WILL MALFUNCTION IN LESS THAN 5 YEARS. HE TESTS A SAMPLE OF 40 CARS AND FINDS THE AVERAGE TIME TO BE 4.8 YEARS WITH A STD. DEV OF 0.50.

1] STATE THE NULL AND ALTERNATIVE HYPOTHESIS

2] AT 2% SIGNIFICANT LEVEL. IS THERE ENOUGH EVIDENCE TO SUPPORT THE IDEA THAT THE WARRENTY SHOULD BE REVISED.

$$H_0: \mu \geq 5$$

$$H_1: \mu < 5$$

From z-table =  $-2.05$   
z-score is

$$n = 40$$

$$\bar{x} = 4.8 \text{ yrs}$$

$$\text{s.d} = s = 0.50$$

$$\alpha = 0.02 \quad CI = 98\%$$

$$z\text{-score} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{4.8 - 5}{0.5/\sqrt{40}} = -2.52$$

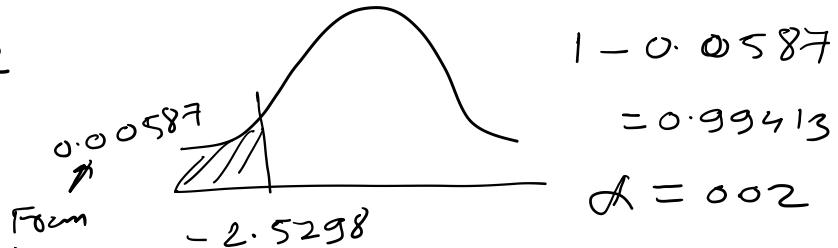
$$\alpha = 0.02 \rightarrow \text{Reject } H_0$$



So,  $-2.52 < -2.05$   
 Conclusion  $\rightarrow$  Reject null  
 $\rightarrow$  warrenty should be revised.

So, company may consider warrenty of 4 years, as most of the part will fail within 4.8 years.

using  
P-value - from table - Z-score = -2.52



EX; THE AVERAGE WEIGHT OF ALL RESIDENT IN A TOWN XYZ IS 168 POUNDS . A NUTRITIONIST BELIEVES THE MEAN TO BE DIFFERENT. SHE MEASURED THE WEIGHT OF 36 INDIVIDUALS AND FOUND THE MEAN TO BE 169.5 POUNDS WITH A STANDARD DEVIATION OF 3.6.

A) NULL AND ALTERNATE HYPOTHESIS

2] 95% IS THERE ENOUGH EVIDENCE TO DISCARD THE NULL HYPOTHESIS?

$P \text{ value} = 0.00587$

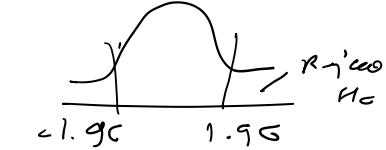
Reject null Hypothesis  $\leftarrow P \text{ value} < \alpha$

$$\bar{x} = 169.5 \quad s = 3.9 \quad n = 36 \quad \mu = 168$$

1]  $H_0: \mu = 168$       2]  $\alpha = 0.05$       3) Decision Boundary

$H_1: \mu \neq 168$

4]  $Z \text{ score} = \frac{168 - 169.5}{3.9 / \sqrt{36}} = \underline{\underline{-2.31}}$

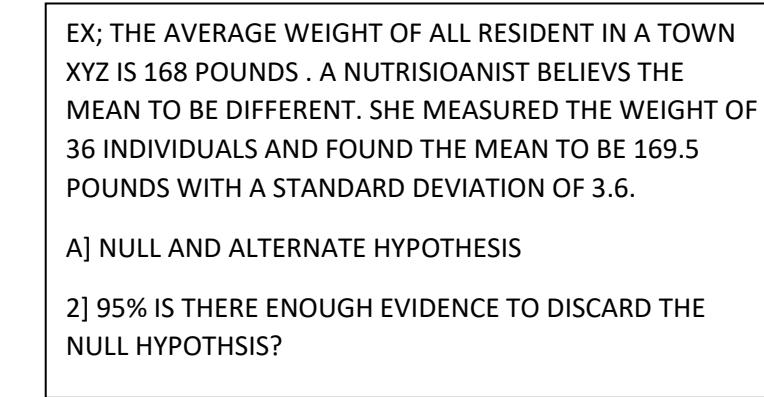


In  
2 tail test

$$\rightarrow P \text{ value} = 0.01044 + 0.01044 = 0.02088$$

$$P \text{ value} = 1 - 0.02088 = 0.97952$$

$\alpha = 0.05 \rightarrow 0.02 < 0.05 - \text{Reject null Hypothesis}$



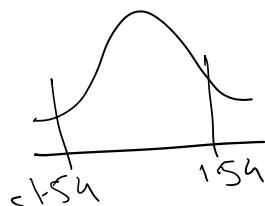
## Z test with proportions

- ① A tech company believes that the percentage of residents in town XYZ that owns a cell phone is 70%. A marketing manager believes that this value to be different. he conducts a survey of 200 individuals and found that 130 responded yes to

Owning a cell phone

(a) State the Null and Alternative Hypothesis?

(b) At a 95% CI, is there enough evidence to reject the Null Hypothesis?

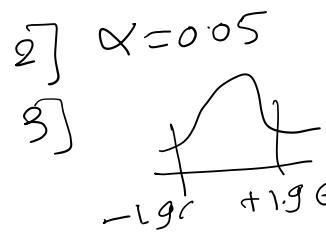


$$\begin{aligned} \text{p-value} &= 0.06178 + 0.06178 \\ &\approx 0.12356 \\ &\approx 0.12356 > 0.05 \quad \text{- accept null} \end{aligned}$$

## → Z test with proportions

- null →  $H_0 : p_0 = 70\%$   
 $H_a : p_0 \neq 70\% \quad q_0 = 1 - p_0$

② Z test with proportion



$$Z \text{ test} = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0 q_0}{n}}}$$

$-1.54 > -1.96$

- accept null

### Chi square test:

It is a non parametrical test that is performed when data is categorical and ordinal data.

q] In the 2000 U.S census the ages of individual in a small town were found to be the following.

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

In 2010, ages of n=500 individuals were sampled. below are the results.

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

Ans  $H_0$ : The data meet expected distribution  $\rightarrow$  Observed

$H_1$ : The data does not  $\rightarrow$   $\rightarrow$

Expected		
< 18	18 - 35	> 35
121	288	91

Expected		
< 18	18 - 35	> 35
100	150	250

2]  $\alpha = 0.05$  C.I = 95%.  $n = 500$

chi-sq Table

df /  $\alpha = 0.05$

3] DDF = df =  $1 < - 1 = 3 - 1 = 2$

2]  $\frac{5.991}{5.991}$

4] Decision Boundary 5.991

{  $\chi^2 < 5.991$  - Accept null

5] Chi square Test Statistics

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

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

This is greater than  
5.991

$$\chi^2 = 232.494$$

so we reject null Hypothesis.

QUE] 500 ELEMENTARY SCHOOL BOYS AND GIRLS ARE ASKED WHICH IS THEIR FAVORITE COLOR,: BLUE , GREEN OR PINK . RESULT ARE SHOWN BELOW.

USING APLHA = 0.05, WOULD YOU CONCLUDE THAT THERE IS A RELATIONSHIP BETWEEN GENDER AND FAVORITE COLOR?

	Blue	Green	Pink	
Boys	100	150	20	270
Girls	20	30	180	230
	120	180	200	500

$$\begin{aligned} \text{DOF} &= (\text{rows}-1) \\ &\quad (\text{cols}-1) \\ &= (2-1)(3-1) \\ &= 1 \times 2 = 2 \end{aligned}$$

NULL HYPTOTHESIS  $H_0$ : - GENDER AND FAVORITE COLOR ARE RELATED

$H_1$ : GENDER AND FAVORITE COLOR ARE NOT RELATED

\* chi sq . Test Stat

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

$$(100 - 64.8)^2 + (150 - 97.2)^2$$

$$= \frac{1}{64.8} + \frac{97.2}{108} + \frac{(20-55.2)^2}{55.2} + \frac{(30-82.8)^2}{82.8} + \frac{(180-200)^2}{92}$$

$$\chi^2 = 259.7$$

It's greater than 5.991.

→ Reject  $H_0$ -null hypo →

$$\alpha = 0.05$$

From Table - 5.991

$$\text{DOF} = 2$$

$$\begin{array}{ccc} \cancel{\text{DOF}} & \cancel{\alpha} & = 0.05 \\ 1 & 2 & 5.991 \end{array}$$

→ Expected

B  
G

	Blue	Green	Pink	
B	64.8	97.2	108	270
G	55.2	82.8	92	230
	120	180	200	500

$$\frac{(30-82.8)^2}{82.8} + \frac{(180-200)^2}{92}$$

So they are not related.

## INTERVIEW QUESTIONS

Asked  
in

Amazon

1] Size of all the sharks

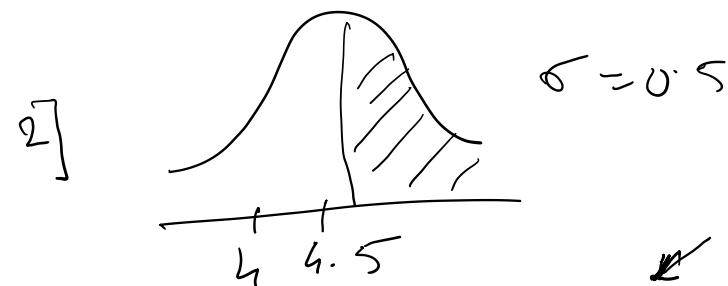
$$n=40 \quad \bar{x}=5 \text{ m} \quad s=0.5 \text{ m}$$

$$\alpha=0.05 \quad n=30 \quad -Z \text{ test}$$

confidence  $\bar{x} \pm$  margin of error

$$\text{Interval} \rightarrow \bar{x} \pm Z_{\alpha/2} \left( \frac{s}{\sqrt{n}} \right)$$

To calculate size of sharks.



non linear  
data

Diff b/w Spearman Rank

& Pearson correlation  
in linear data

\* oil Distribution in a place follows which distribution.

Ans- Pareto distribution



Power Laws

Distributions

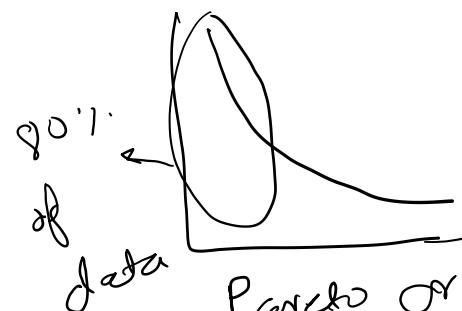
may be asked

log normal  $\xrightarrow{\text{log}} \text{normal}$

Pareto  $\xrightarrow{\text{Power law}}$  log normal

What is Z-score

Pareto Distribution



Pareto or  
Power law



Log normal

\* Disadvantage of replacing  
NaN value by mean.

\* Variance Yes - most  
values will be at

