

Central Limit Theorem

The central limit theorem states that if you have a population with mean ' μ ' and standard deviation ' σ ' and take sufficiently large random samples from the population with replacement, then the distribution of the sample means, it will be approximately normally distributed

Probability:

Probability is a measure of the likelihood of an event

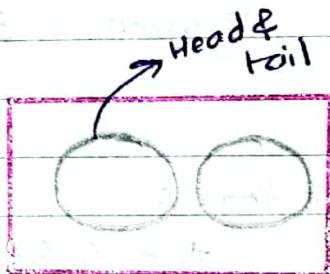
Eg.: Tossing a fair coin

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

Mutual Exclusive Event

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

- a) Tossing a coin
- b) Rolling a dice



Non Mutual Exclusive Events

Two events can occur at the same time.

Mutual Exclusive Event Problem Statement:

→ What is the probability of coin landing on heads or tails?

Addition Rule for mutual exclusive events

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

* What is the probability of getting 1 or 6 or 3 while rolling a dice?

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

Non Mutual Exclusive Event Problem Statement:

$$\begin{aligned} P(A \text{ or } B) &= P(A) + P(B) - P(A \text{ and } B) \\ &\quad \uparrow \\ \text{Addition rule for non mutual exclusive event} \end{aligned}$$

Bag of Marbles: 10 Red, 5 Green, 3 (R & G).

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

$$\begin{aligned} P(A \text{ or } B) &= P(A) + P(B) - P(A \cap B) \\ &= \frac{10}{19} + \frac{9}{19} - \frac{3}{19} \\ &= \frac{13}{19} + \frac{6}{19} = \frac{19}{19} = 1 \end{aligned}$$

Dependent Event

Two events are dependent if they affect one another

- * What is the probability of rolling a 6 sided dice to get a "5" and then a "3".

↓
Since it is independent event

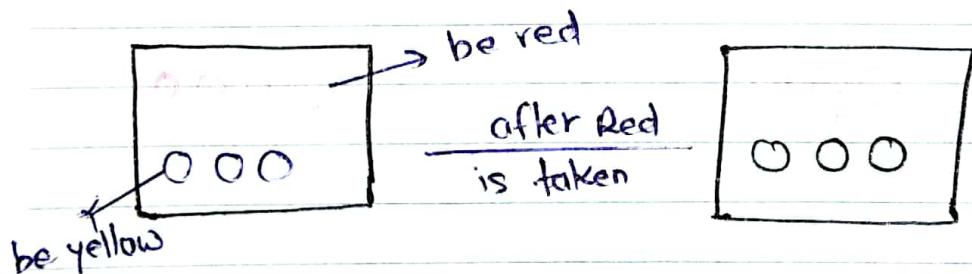
$$P(3) = \frac{1}{6} \quad \& \quad P(5) = \frac{1}{6}$$

Multiplicative rule:

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

$$\begin{aligned} &= \frac{1}{6} \times \frac{1}{6} \\ &= \underline{\frac{1}{36}} \end{aligned}$$

- * In bags of marble, 4 Red marble, 3 Yellow marble. What is the probability of drawing a "orange" and then drawing a "yellow" marble from the bag.



~~Soln~~ $P(\text{Red}) = \frac{4}{7}$

$$\therefore P(\text{Yellow} | \text{Red}) = \frac{3}{6} = \frac{1}{2}$$

conditional probability

Now, $P(\text{Red and Yellow}) = P(\text{Red}) \times p(\text{Yellow} / \text{Red})$

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

Permutation

- each of several possible ways in which a number

Q. Considered the letters 'A', 'B' and 'C': How many ways can these letters be arranged?

$$\rightarrow \underline{3} \times \underline{2} \times \underline{1} = 6 \text{ ways}$$

First place has 3 choice 2 choices for second

$$\begin{array}{ll} \rightarrow ABC & \rightarrow ACB \\ \rightarrow BAC & \rightarrow BCA \\ \rightarrow CAB & \rightarrow CBA \end{array} \quad \left. \begin{array}{l} \text{3 ways} \\ \text{3 ways} \\ \text{3 ways} \end{array} \right\} 6 \text{ option}$$

n be total no. of object

r be no. of selection,

* Orders matter in permutation

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

Combination

The number of combinations of 'n' objects taken 'r' at a time can be calculated using the combination formula.

$$C(n, r) = \frac{n!}{r!(n-r)!}$$

Covariance

- It is a measure of the joint variability of two random variables.

Age	Weight	Shows
12	60	age ↑ → weight ↑
13	65	age ↓ → weight ↓
14	68	
17	60	

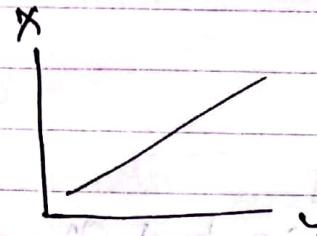
$$\text{Cov}(X, Y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

It is variance with own.

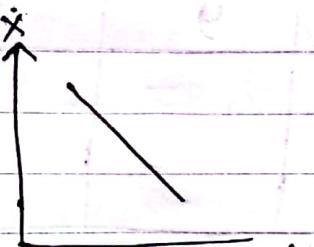
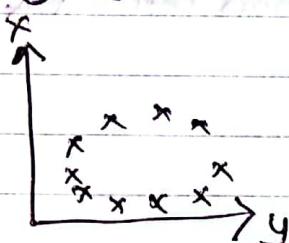
$$\text{Cov}(X, X) = \sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$$

In +ve Covariance

$$\begin{aligned} x \uparrow &\rightarrow y \uparrow \\ x \downarrow &\rightarrow y \downarrow \end{aligned}$$

In -ve Covariance

$$x \uparrow \rightarrow y \downarrow$$

In 0 covarianceSpearman Correlation

→ It is for non linear as covariance 'works well in only linear datasets.'

$$r_s = \frac{\text{Cov}(R(x), R(y))}{\sqrt{(R(x))^2} \sqrt{(R(y))^2}}$$

Rank

Rank always in ascending

X	Y	R(X)	R(Y)
10	4	4	1
8	6	3	2
7	8	2	3
6	10	1	4

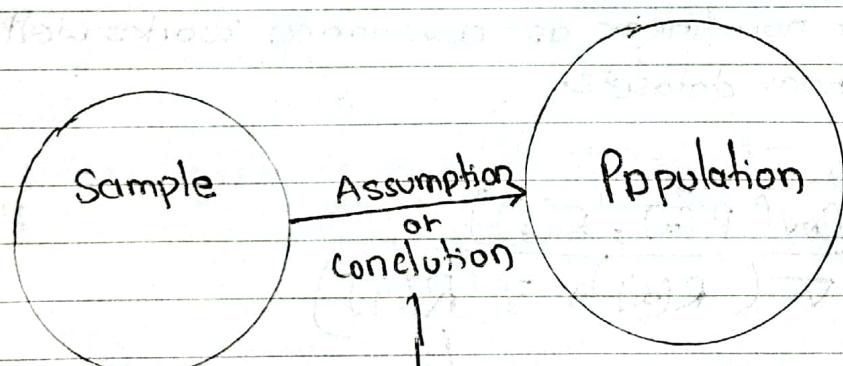
* Why do we need correlation? ~~will be~~

X	y	z	O/P

- If correlation between X & O/P is high, then it is important feature
- If there exist very +ve correlation bet' X & y, then we can drop the one

Hypothesis testing:

Before coming into it, let's cover Inferential Statistics



To validate Assumption or conclusion
we perform hypothesis testing

Steps of hypothesis testing:

① Null hypothesis: default value.

like: Person hasn't committed the crime

② Alternate Hypothesis: Opposite of null hypothesis

Eg: Person is criminal

③ Perform Experiment \rightarrow DNA, finger Print, Weapons
Eye witness, footage.

Example:

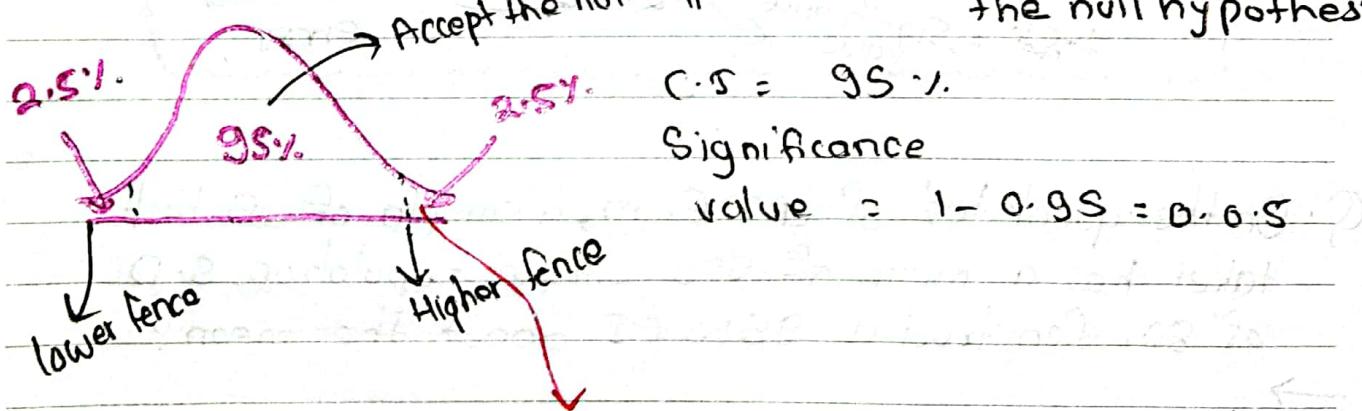
- 1) Null hypothesis: coin is fair
- 2) Alternate hypothesis: coin is not fair
- 3) Experiment:

Head : 70 { 40 { SS
Tail : 30 { 60 { 4S

Then, we go to confidence interval. Say, if C.I is [20 - 80], then our experiment make null hypothesis true.

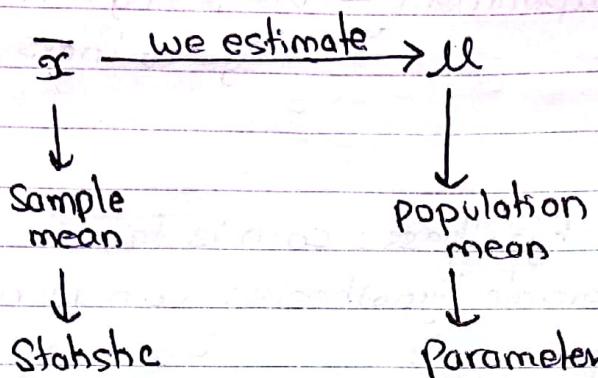
Hence Null hypothesis accepted.

What is 95%. C.I means? Accept the null hypothesis [we fail to reject the null hypothesis]



We Reject the null hypothesis

Point Estimate: The value of any statistics that estimates the value of a parameter is called point estimate.



$$\text{Parameter} = \text{Point Estimate} \pm \text{Margin of error}$$

So, for confidence interval

$$\text{Lower fence} = \text{Point Estimate} - \text{Margin of error}$$

$$\text{Higher fence} = \text{Point Estimate} + \text{Margin of error}$$

$$\text{Margin of error} = Z_{\alpha/2} \cdot \sqrt{\frac{\sigma^2}{n}}$$

standard error.

Margin of error = $Z_{\alpha/2} \cdot \sqrt{\frac{\sigma^2}{n}}$

$1-\alpha$ = Significance value

Q. On the quant test of CAT Exam, a sample of 25 test taker has a mean of 520 with a population S.D of 80. Construct a 95% CI about the mean?



$$n = 25$$

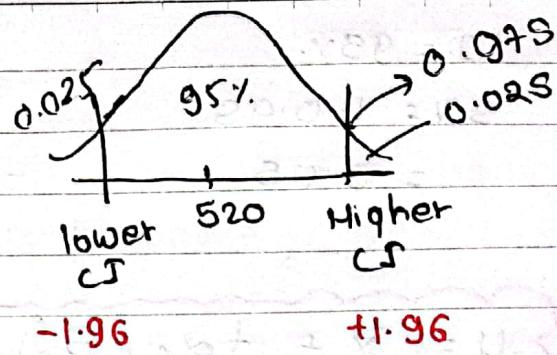
$$\bar{x} = 520$$

$$\sigma = 100$$

$$\text{C.I.} = \bar{x} \pm Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \Rightarrow \sigma = 1 - 0.95 \\ = 0.05$$

We know,

Point Estimate \pm marginal error
parameter.



Also,

Lower CI = Point Estimate - Marginal error

$$= 520 \pm Z_{0.05} \cdot \frac{80}{\sqrt{25}}$$

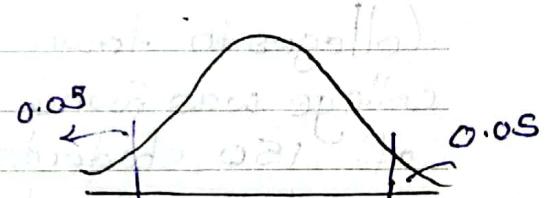
$$= [501, 551.36]$$

Q. $\bar{X} = 480$, $s = 88$, CI = 90%, $n = 25$

$$\alpha = 1 - 0.9 = 0.1$$

$$\text{CI} = 480 \pm Z_{0.05} \cdot \frac{88}{\sqrt{25}}$$

$$= 480 \pm 1.64 \times 17$$



Q. On the quant test of CAT exam, a sample of 25 test takers has a mean of 520 with a sample standard deviation of 80. Construct 85% CI about the mean?



$$\bar{X} = 520$$

$$S = 80$$

$$CT = 95\%$$

$$S.V = 1 - 0.95 \\ = 0.05$$

When sample S.D is given, then t-test is performed.

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

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

$$= 520 \pm 2.064 \times \frac{80}{\sqrt{5}}$$

$$= [486.976, 553.024]$$

One tail and Two Tail

Colleges in town A has 85% placement rate. A new college was found open and it found that a sample of 150 students had a placement rate of 88% with a standard deviation of 4%. Does this college has a different placement rate with 95% CT?



When it say 'different placement rate', it can be greater than 85 or less than 85. It is two tail problem.

Hypothesis Testing:

A factory has a machine that fills 80ml of baby medicine in a bottle. An employee believes the average amount of baby medicine is ^{not} 80 ml. Using 40 Samples, he measures the average amount dispersed by the machine to be 78ml. with S.D of 2.5.

- State null and alternate hypothesis
- At 95%, CS, is there enough evidence that support machine is working properly or not?

Ans:

- Null hypothesis: (H_0) : $\mu = 80$
Alternate hypothesis (H_1) : $\mu \neq 80$

(b) →

$$\mu = 80$$

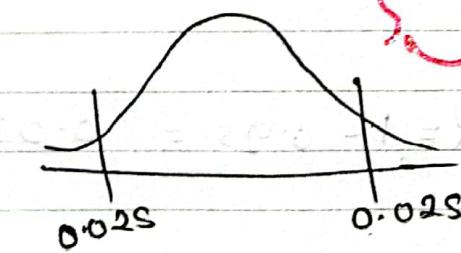
$$CS = 0.95, \alpha = 1 - 0.95 = 0.05$$

$$n = 40$$

$$S \bar{v} = 2.5$$

Z test

If $n \geq 30$ or population SD \rightarrow Z test
Else $n \leq 30$ & sample std \rightarrow t-test



$$Z = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \Rightarrow \frac{78 - 80}{\frac{2.5}{\sqrt{40}}} = -5.05$$

Conclusion:

Decision rule: If $Z = -SOS$ is less than -1.96 or greater than 1.96 , reject null hypothesis.

Since $Z = -SOS < -1.96$

↓
Reject null hypothesis.

$E = M \pm 80$

A complaint was registered that boys in a government school are underfed. Average weight of boys of age 10 is 32 kgs with $S.D = 9\text{kg}$. A sample of 25 boys were selected from the government school & the average weight was found to be 29.5 kgs. with $C.I = 95\%$. Check if it is true or false.

Soln

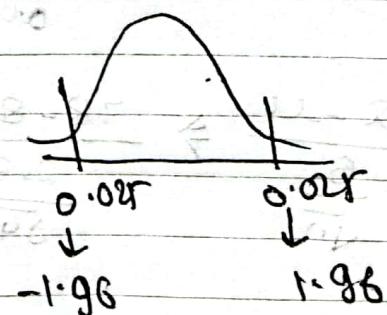
$$S.D = 9\text{kg}$$

$$M = 32\text{kg}$$

$$n = 25$$

$$\bar{x} = 29.5$$

$$C.I = 95 \therefore \alpha = 1 - 0.95 = 0.05$$



Null hypothesis (H_0) $\Rightarrow \mu = 32$
 Alternate hypothesis (H_1) $\Rightarrow \mu \neq 32$

$$Z = \frac{29.5 - 32}{\frac{9}{\sqrt{2}}} = -1.39$$

Conclusion: $-1.39 > -1.96$, So we accept the null hypothesis
 "We fail to reject null hypothesis"

Q. A factory manufactures cars with a warranty of 5 years or more on the engine and transmission. An engineer believes that the engine or transmission will malfunction in less than 5 years. He tests sample of 40 cars and finds the average time to be 4.8 years with a standard deviation of 0.50.

- (a) State the null and alternate hypothesis.
- (b) At a 2% significance level, is there enough evidence to support the idea that the warranty should be revised?

Soln,

Step 1: $H_0: \mu \geq 5$

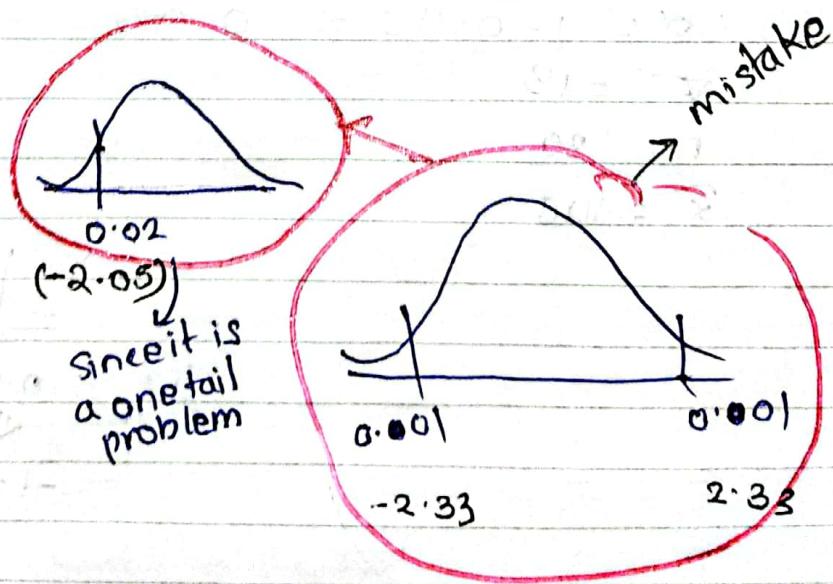
$H_1: \mu < 5$

Step 2: $C:I = 98$

$n = 40$

$\bar{x} = 4.8$

$s = 0.5$



Since it is
a one tail
problem

$$Z = \frac{4.8 - 5}{\frac{0.5}{\sqrt{40}}} = -2.529 \approx -2.05$$

Hence, we reject the null hypothesis.

- ② In the population, the average IQ is 100 with a standard deviation of 15. A team of scientists wants to test a new medication to see if it has a +ve or -ve effect, or no effect at all. A sample of 30 participants who has taken the medicine has a mean of 110. Did the medication affect intelligence?

$$\text{Prop } CS = 95\%$$

~~So, the null hypothesis is that the mean is not affected by the medication.~~

$$H_0 \Rightarrow \mu = 100 \quad [\text{No effect}]$$

$$H_1 \Rightarrow \mu \neq 100 \quad [\text{Effect}]$$

Two tail problem

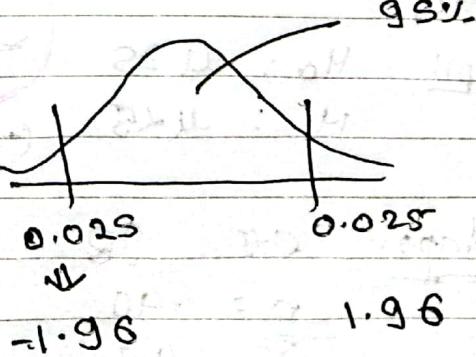
$$CS = 95\%$$

$$\alpha = 1 - 0.95 \Rightarrow 0.05$$

$$\sigma = 15$$

$$n = 30$$

$$\bar{x} = 110$$



$$Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{140 - 100}{\frac{15}{\sqrt{80}}} = 14.60$$

Since: $14.60 > 1.96$, we reject the null hypothesis, the medication has effect.

Q. The average weight of all resident in a town XYZ is 168 pounds. A nutritionist believes the true mean to be different. She measured the weight of 36 individuals and found the mean to be 169.5 pounds with the standard deviation of 3.9?

(a) Null & Alternate hypothesis

(b) 95%, is there enough evidence to discard the null hypothesis?

So,

$$H_0: \mu = 168$$

$$H_1: \mu \neq 168 \rightarrow \text{Two tail test}$$

Given,

$$\alpha = 0.05$$

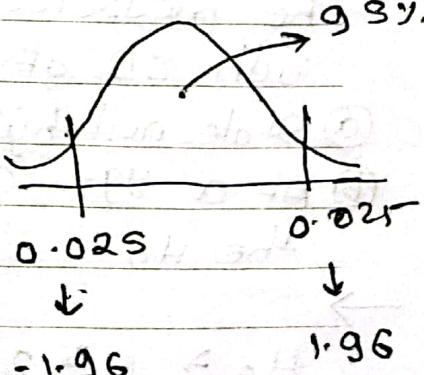
$$\bar{X} = 169.5$$

$$S = 3.9$$

$$n = 36$$

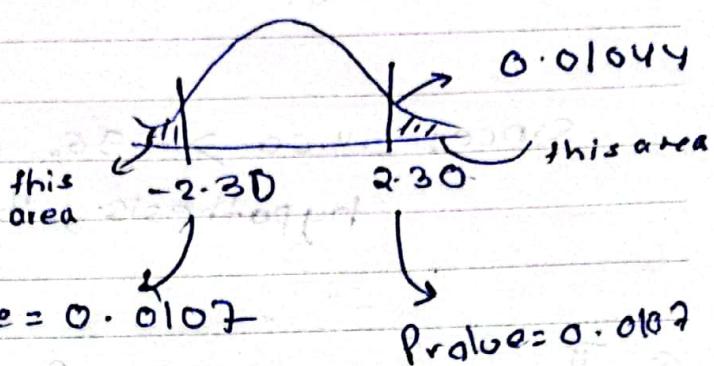
$$Z = \frac{\bar{X} - \mu}{\frac{S}{\sqrt{n}}} = \frac{169.5 - 168}{\frac{3.9}{\sqrt{36}}} = 2.3076$$

Since $Z = 2.3076 > 1.96 \Rightarrow \text{Rejected.}$



Solving with P-value:

Seeing the value



$$\text{Paggregate} = 0.0107 + 0.107$$

$$\text{Paggregate} = 0.214$$

If $\text{Paggregate} \rightarrow \text{Significance level}$
accept null hypothesis

$\text{Pagg} < \text{CS}$

rejecting

A company manufactures bikes batteries with an average life span of 2 years or more years. An engineer believes this value to be less. Using 10 sample, he measures the average life span to be 1.8 years, with SD of 0.15.

⑥ State null hypothesis

⑦ At a 99% CI, is there enough evidence to discard the H_0 ?



$$H_0: \mu \geq 2$$

$$H_1: \mu < 2 \rightarrow \text{1 tail problem}$$

degree of freedom $\geq \max^{th}$ no. of choice from sample.

No. _____

Date _____

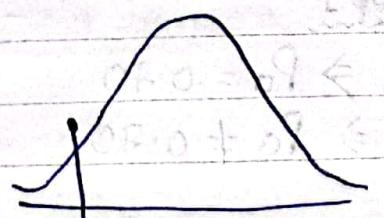
Given,

$$CI = 99\% \Rightarrow \alpha = 1 - 0.99 \\ = 0.01$$

$$n = 10 \\ S = 0.15 \quad \left. \begin{array}{l} \text{Suggest t-test} \\ \text{take negative since it is in} \\ \text{left side} \end{array} \right\}$$

$$\bar{x} = 1.8$$

$$\left. \begin{array}{l} \text{degree of freedom} = n - 1 \\ = 9 \end{array} \right\}$$



0.01

↓
2.821

take negative since it is in
left side
↓
-2.821

Now,

$$t_{\text{score}} = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} = \frac{1.8 - 2}{\frac{0.15}{\sqrt{10}}} = -4.2163$$

Since: $-4.2163 < -2.821$, we reject the null hypothesis.

The average life of battery is less than 2 years.

Z test with Proportions:

A tech company believes that the percentage of residents in town XYZ that own a cell phone is 70%. A marketing manager believes that this value to be different. He conducts the survey of 200 individuals and found that 130 responded 'Yes' owning a cell phone?

- state null and alternate hypothesis.
- At a 95%, CI, is there enough evidence to reject the null hypothesis?

→ Step 1:

$$H_0 \Rightarrow P_0 = 0.70$$

$$H_1 \Rightarrow P_0 \neq 0.70$$

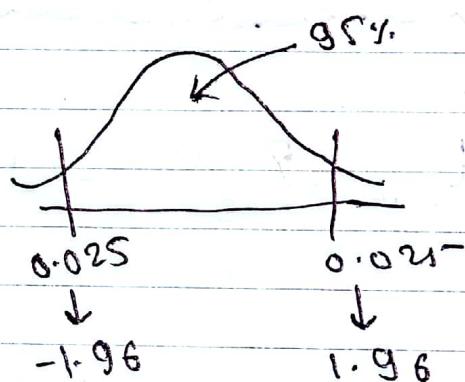
$$Q_0 = 1 - 0.70 = 0.30$$

Step 2:

$$C.S = 0.98 \Rightarrow \alpha = 0.02$$

Step 3:

$$\hat{P} = \frac{130}{200} = 0.65$$



Step 4:

$$Z \text{ test} = \frac{\hat{P} - P_0}{\sqrt{\frac{P_0 Q_0}{n}}}$$

$$\sqrt{\frac{0.70 \times 0.30}{200}}$$

$$= \frac{0.65 - 0.70}{\sqrt{\frac{0.70 \times 0.30}{200}}}$$

$$= \frac{-0.05}{\sqrt{\frac{0.21}{200}}} = -1.54$$

Since: $-1.54 > -1.96$, so we fail to reject null hypothesis.

- Q. A car company believes that the percentage of resident in city ABC that owns a vehicle is 60%. or less. A sales manager disagree with this and conduct a hypothesis testing surveying 250 residents and found that 170 responded 'Yes' to owning a vehicle.
- State the null hypothesis & alternate hypothesis
 - At 10% significance level, is there enough evidence to support the idea that vehicle ownership in city ABC is 60%.

Soln

Step 1:

$$H_0 \Rightarrow P_0 \leq 0.6$$

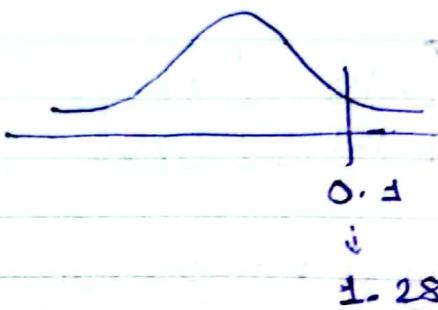
$$H_1 \Rightarrow P_0 > 0.6 \rightarrow \text{1 tail problem}$$

Step 2:

$$C.I = 90\%, \alpha = 0.1$$

Step 3:

$$\hat{P} = \frac{170}{250} = 0.68$$



Step 4:

$$z \text{ test} = \frac{0.68 - 0.6}{\sqrt{\frac{0.6 \times 0.4}{250}}}$$

$$= 2.581$$

Since: $2.581 > 1.28$, we reject the null hypothesis

Chi Square test

→ χ^2 test is a statistical test used to determine whether there is a significant association between two categorical variables.

→ It claims about population proportion.

→ It is a non-parametric test that is performed on categorical data.

In the 2000 US census, the age of individuals in a small town found to be the following:

<18	18-35	≥ 35
20%	30%	50%

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

<18	18-35	≥ 35
121	286	91

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



	<18	18-35	≥ 35
Expected	20%	30%	50%

when $n = 500$ is sampled:

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

Step 1:

Null hypothesis (H_0): The data meets the expected distribution

(H_1): The data doesn't meet the expected distribution

Step 2:

$$\alpha = 0.05, CI = 95\%$$

Step 3:

$$d.f = \text{class} - 1 \Rightarrow 3 - 1 = 2$$

Step 4:

$$\chi^2_{\alpha} = 5.991$$

decision boundary:

If $\chi^2 > \chi^2_{\alpha}$ Reject the null hypothesis

Step 5:

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

$$\text{Rejecting } \leftarrow = 232.494$$

Distribution

① Bernoulli Distribution

→ It models a random experiments with only two possible outcomes, 0 or 1.

→

$$\left\{ \begin{array}{l} P = 1 - q \\ & \& \\ q = 1 - P \end{array} \right.$$

→ 0 or 1

→ pass or fail

↳ algebra of domain of $f(x)$ is $(0, 1)$ and range is $\{0, 1\}$

Mean → Expected Value = $1 - p$

Variance ⇒ $\text{Var}(X) = p(1-p)$

or

$$P \cdot Q$$

② Binomial Distribution

→ discrete probability distribution that describes the number of successes in fixed number of independent and identical bernoulli trials.

$n \rightarrow$ number of trials

$p \rightarrow$ probability of success at each trials

↳ algebra of domain of $f(x)$

$$P_X = \binom{n}{x} p^x q^{n-x}$$

number
of success
in n trial

Mean (expected value) $\Rightarrow E(x) = n * p$

Variance $\Rightarrow V(x) = n * p * (1-p)$.

④ Power Law

→ The power law distribution is a continuous positive-only univariate distribution that describes a quantity whose probability decreases as a power of its magnitude.

