

21/07/19

Central Tendency: [to represent data in a single number]

- Mean: $\bar{x} = \frac{\sum x_i}{n}$ (Q) $\bar{x} = \frac{\sum f_i x_i}{\sum f_i}$

- Median: even = $\text{avg}\left(\frac{n}{2}, \frac{n+1}{2}\right)$
odd = $n/2$

- Mode : most frequent item

- Percentile: It tells the position of observation in dataset.

$$P_x = \frac{x(n+i)}{100}$$

- Quartile - Q_1 : contains 1st 25% of data.

- Decile - divides data into 10 parts.

Measure of Variation:

- Range $(\text{max} - \text{min})$
- Inter Quartile Range $(Q_3 - Q_1)$
- Variance
- Standard Deviation

Variance: Measure of variability from the mean

$$\text{Variance} = \sigma^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n} \quad (\text{population})$$

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

Standard Deviations:

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

$$\sigma = \sqrt{\frac{\sum (x_i - u)^2}{n}}$$

Mostly $x_i - u > x_i - \bar{x}$
in order to adjust them
we divide $\frac{(x_i - u)^2}{n}$ &
 $\frac{(x_i - \bar{x})^2}{n-1}$ ie, divides by $(n-1)$

This is to adjust the balance
btw S & σ

[downward bias problem]

Skewness: is used to measure the symmetry of the distribution of data

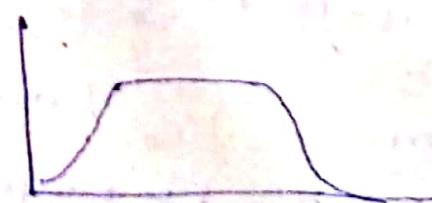
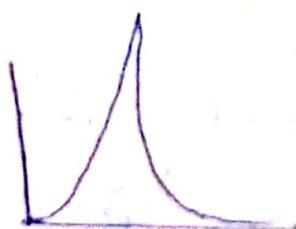
Skewness = 0 \rightarrow Symmetric

-ve \rightarrow left skewed

+ve \rightarrow Right skewed



Kurtosis: at the top of the relational graph is sharp or blunt (flat).



Using all these plots

Random Variable: Function which assigns a real number to each sample point in a sample space.

Ex: tossing 3 coins.

$$S = \{HHH, HHT, HTH, TTH, HTH, THH, HTT, TTT\}$$

$$\text{no. of tails}(X) = f(x)$$

$$\Rightarrow \{0, 1, 1, 2, 1, 2, 2, 3\}$$

Discrete Random Variables:

A Random Variable which takes only finite values.

In finite set of observations.

Continuous Random Variable:

A Random Variable which

takes infinite values in finite observations.

Probability Mass Functions:

Probability that the D.R.V 'X' takes on the value

x_i , given by $P(x_i) = P(X=x_i)$ for $i=1, 2, 3, \dots, n$

$$(i) P(x_i) \geq 0$$

$$(ii) \sum P(x_i) = 1$$

Ex: If we flip a coin 2 times.

Random Variable (X) = no. of tails

$$S = \{HH, HT, TH, TT\}$$

$$RV(X) = \{0, 1, 2, 2\}$$

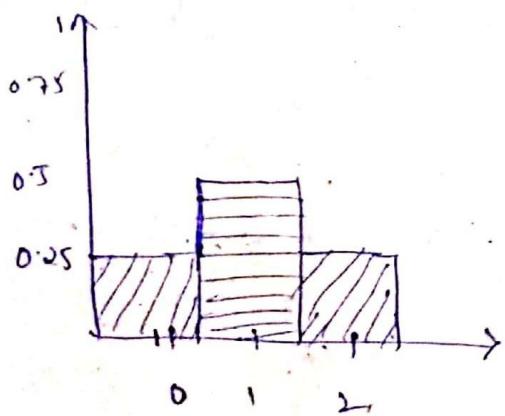
$$f(1) = P(HT, TH) = \frac{2}{4} = \frac{1}{2} \geq 0$$

$$f(0) = P(HT) = \frac{1}{4} \geq 0$$

$$f(2) = P(TT) = \frac{1}{4} \geq 0$$

$$E[f(x)] = \frac{1}{4} + \frac{1}{4} + \frac{1}{2} = \frac{1}{2}$$

x	0	1	2
f(x)	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{2}$



Probability histogram.

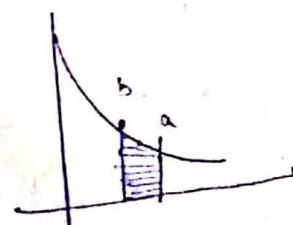
Probability Density function (continuous Rv):

If X is a continuous Random Variable the function $f(x)$ is Probability Density function

$$f(x) = P(a \leq x \leq b) = \int_a^b f(x) \cdot dx$$

if (i) $f(x) \geq 0$

$$(ii) \int_{-\infty}^{\infty} f(x) dx = 1$$



$$\text{Ex} \quad P(x) = \begin{cases} x & 0 \leq x < 1 \\ 2-x & 1 \leq x < 2 \\ 0 & \text{other} \end{cases}$$

$$\Rightarrow \int f(x) dx$$

$$\Rightarrow \int_{0.5}^1 f(x) dx + \int_{1.5}^{1.5} f(x) dx$$

$$\Rightarrow \int_{0.5}^1 x dx + \int_1^{1.5} 2-x dx$$

$$\Rightarrow \left[\frac{x^2}{2} \right]_{0.5}^{1.5} + 2[x]_1^{1.5} - \left[\frac{x^2}{2} \right]_1^{1.5}$$

$$\Rightarrow \frac{1}{2} [1.5^2 - 0.5^2] + 2[1.5 - 1] - \frac{1}{2} [(1.5)^2 - 1^2]$$

$$\Rightarrow \frac{1}{2} [0.25] + 2[0.5] - \underline{\underline{1.25}}$$

$$\Rightarrow 1 + \left(\frac{0.25 - 1.25}{2} \right)$$

$$\Rightarrow 1 + \frac{0.50}{2}$$

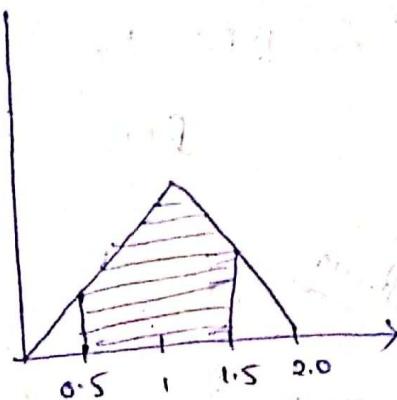
$$\Rightarrow 1 + 0.25 = \underline{\underline{1.25}}$$

Cumulative Distribution function (CDF) :-

[irrespective of continuous or discrete function, we can use Cumulative distribution function (CDF)].

C.D.F of R.V 'x' is defined as the probability that

R.V 'x' takes a value "less than or equal to x"



$$C.D.F = F_X(x) = P(X \leq x)$$

$$f(a) = P(X \leq a) = \sum_{x \leq a} f(x)$$

5/7/19

Histogram:

Pictorial representation of frequency distribution.

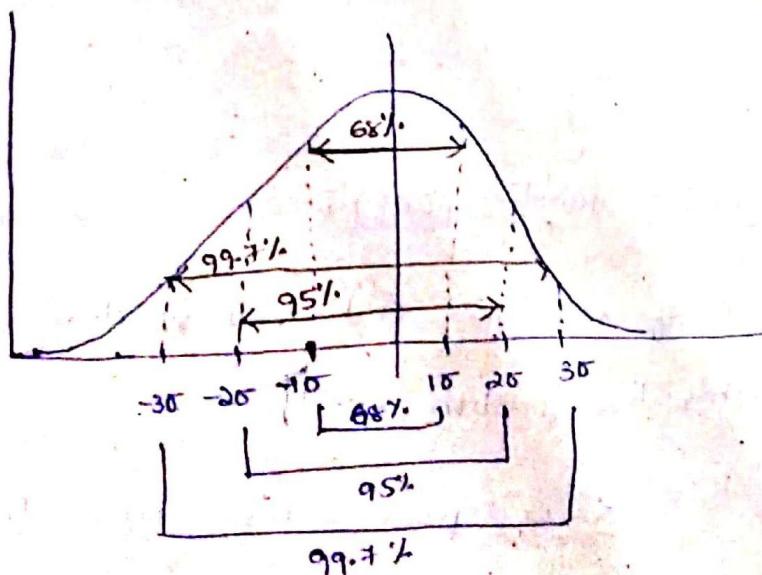
which groups the data into groups.

Normal Distribution:

Normal distribution of data usually allow the use of parametric test to analyse the data

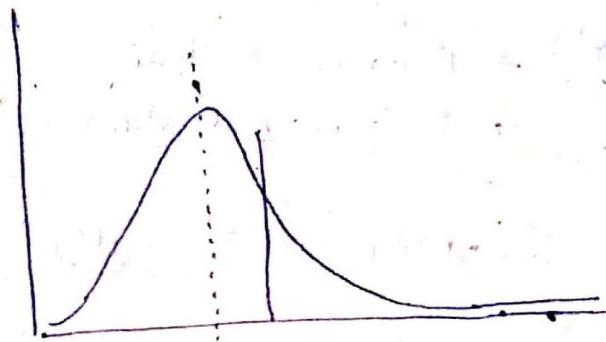
Note In normal distribution mean, median, mode are located at same place (i.e., middle)

- Parametric test has lots of tools, it is well developed system



Positively skewed data

The mean shift right side



To measure normality

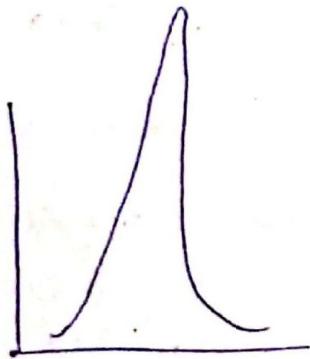
(+ve skewed data)

measure of distribution.

— Skewness [symmetry of the distribution]

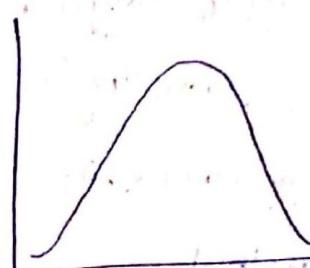
— Kurtosis [peakedness] — measure of peakedness.

for normal distribution: Kurtosis = 0



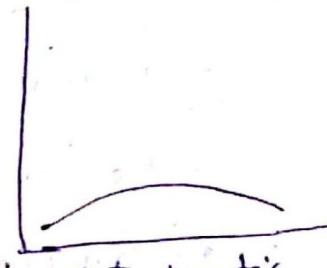
Lepto kurtic

(a)



meso kurtic

(b)



platy kurtic

(c)

Standard error = $\frac{\sigma}{\sqrt{n}} \rightarrow \text{standard deviation}$
/ $\sqrt{\text{no. of observations}}$

If $\times 3$ times of standard error. Then we can't say it is in normal distribution.

$$K < 3\sigma_x \quad (\text{standard error})$$

$$K < 3 \left[\frac{\sigma}{\sqrt{n}} \right].$$

Why

$$\text{skewness} < 3 \left(\frac{\sigma}{\sqrt{n}} \right)$$

- ~~July~~
- Remove outliers
 - increasing the data size

Sampling and Estimation:

Sampling is a process of selecting subgroup from population to make inferences about population parameters such as mean, median, proportion, standard deviation etc.

Steps used in Sampling process:

1. Identification of target population
 2. Decide the sampling frame
 3. Determine the sample size
 4. Sampling method
 - Probabilistic sampling
 - Non-probabilistic Sampling
- Probabilistic Sampling
- Random Sampling. → If population is homogeneous, it is used
 - stratified Sampling
 - clustered Sampling
 - bootstrapping (bagging)

→ Stratified Sampling:

- divide population with equal classes.
- Then the ratio of proportions should be equal in each class

→ Clustered Sampling:

→ only important classes only considered

→ neglecting the negligible class.

→ Bootstrapping:

- Non-probabilistic Sampling:

1. Convenience Sampling

2. Voluntary Sampling.

- Sampling Distribution:

Probability distribution of static such as mean, standard

deviation computed from random sample

Population

1	2	3	4	5	6	7	8	9	10
5	10	15	20	25	30	35	40	45	50

Sample of size 2:

Sample	5,5	5,15	10,5	10,15	10,45	20,15	45,15
--------	-----	------	------	-------	-------	-------	-------

mean	5	10	7.5	12.5	27.5	17.5	30
------	---	----	-----	------	------	------	----

	50,20	25,10
--	-------	-------

	35	17.5
--	----	------

If we take sample sufficiently large sample size from population. The mean of all samples follow approximately normal distribution with mean equal to population mean and std deviation equal to $\frac{\sigma}{\sqrt{n}}$

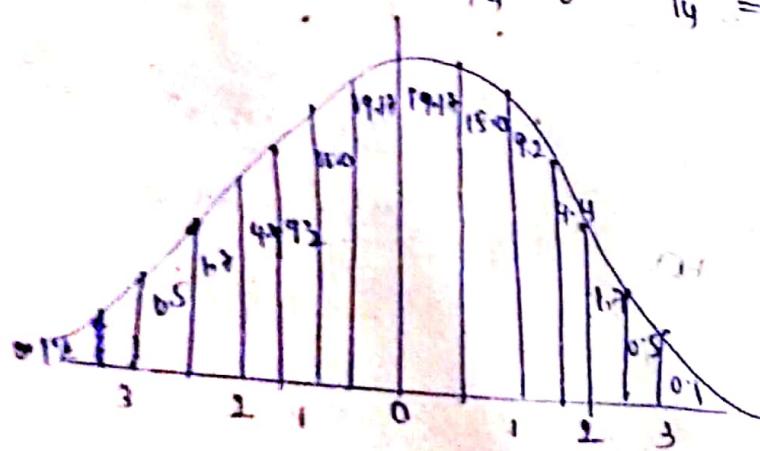
Standard Normal Distribution:

It is a normal distribution, it has mean '0' & stand. deviation '1'

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

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

$$P(X < z_2) \Rightarrow P\left(Z < \frac{z_2 - \mu}{\sigma}\right) = \frac{14}{14} = P(Z < 1)$$



Confidence Intervals:

CI is the range in which the value of population parameter is likely to be with certain probability.

$$\bar{X} \pm z \frac{\sigma}{\sqrt{n}}$$

Ques If we measure the height of 40 randomly chosen men and get a mean height 175cm. We also know that standard deviation is 20cm

Sol

$$S = 40$$

$$\sigma = 20$$

$$\bar{X} = 175$$

$$\bar{X} \pm z \frac{\sigma}{\sqrt{n}}$$

$$= 175 \pm 1.96 \times \frac{20}{\sqrt{40}}$$

$$= 175 \pm 6.20$$

$$\Rightarrow 168.8 \text{ cm to } 181.2 \text{ cm}$$

6/8/19

Hypothesis testing:

Null hypothesis - When comparing any two, there is no difference

Alternate hypothesis

Rejection of Null hypothesis

let sample - $n=10$

Hypothesis is an

$$\bar{X} = 44.5$$

assumption about

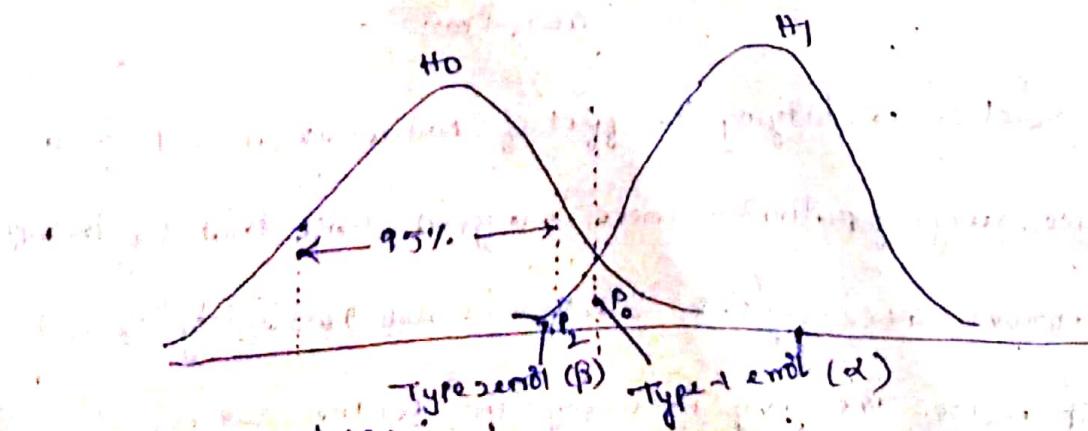
$$\sigma = 11.41$$

population

parameters

$$H_0: \mu \geq 45 \quad \text{Null hypothesis}$$

$$H_1: \mu < 45 \quad \text{Alternate hypothesis}$$



Here though (P_0) is from

H_0 , but we are considering it as not $\in H_0$, since it is

not in the 95% confidence of the H_0 , it is considered false.
i.e. type-I error (α).

The μ_2 belongs to H_1 , but it is in the range of 95% confidence interval, though it is belongs to H_1 , it is accepted as belongs to H_0 . This is of type 2 error (β)

~~flaging~~
 α - false rejection of H_0

β - false acceptance of H_0

Ex: A Researcher thinks that a knee surgery patients go to physical therapy twice a week. (instead of 3 times) their recovery period will be longer. Average recovery time for knee surgery patients is 8.2 weeks.

Sol: $\mu_0 = 8.2$ weeks

Always test on null hypothesis

$$H_0: \mu \leq 8.2$$

$$H_1: \mu > 8.2$$

but not on alternate hyp

\rightarrow Always alternate hypothesis is accepted.

Ex: A Researcher is studying the effect of Radical exercise program on knee surgery patients. There is a good chance that the therapy will improve recovery time, but there is also possibility it will make it worse. The avg. Recovery time for knee surgery patients is 8.2 weeks.

$$H_0: \mu = 8.2$$

$$H_0: \mu = 8.2$$

$$H_1: \mu \neq 8.2$$

Ex: Salary of Machinelearning expert on average is atleast \$ 10,000

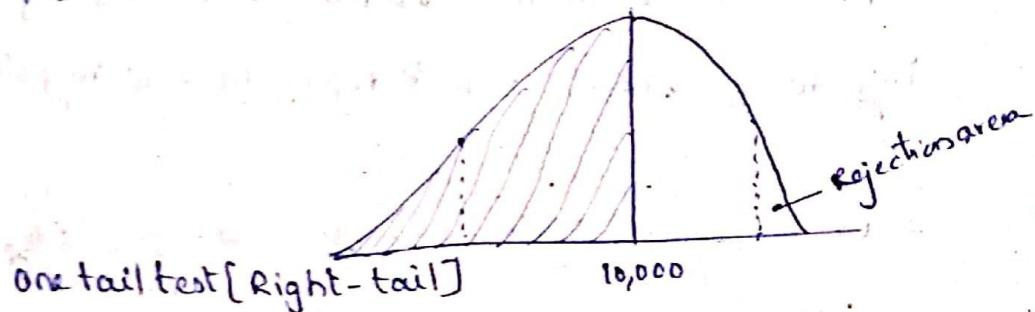
$$\$ 10,000$$

$$H_0: \mu = 10,000$$

$$H_0: \mu \leq 10,000$$

$$H_1: \mu > 10,000$$

Ex:



8/8/19

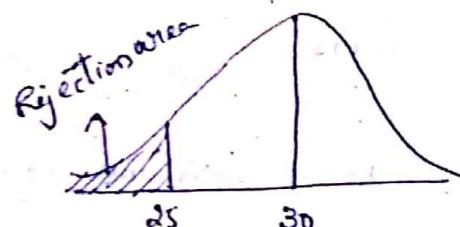
Ex: Average waiting time at London Heathrow airport security

check is less than 30 min.

$$H_0: \mu \geq 30$$

$$H_1: \mu < 30$$

one tail test [left-tail]

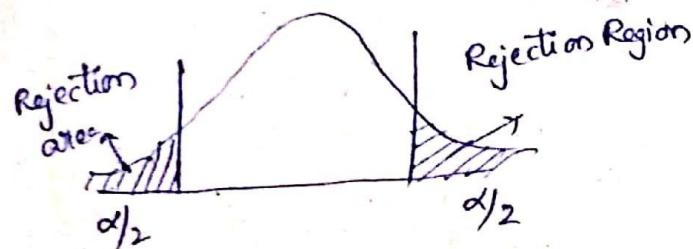


Ex: Avg salary of male and female mba students at graduation is different.

$$H_0: \mu_m = \mu_f$$

$$H_1: \mu_m \neq \mu_f$$

two tailed test



Steps for hypothesis testing:

- State the null hypothesis
- State the alternate hypothesis
- Determine the value of α
- find Z -score associated with alpha level
- Find the test statistic using formulae
- If the value of the static is less than Z -score of alpha level [p-value is less than α -value] reject the null hypothesis otherwise don't reject the null hypothesis.

Ex $\bar{x} = 44.5$; $S = 11.41$, $n = 10$ ($\alpha = 95\%$, $\alpha = 5\%$, type I error)

$$\text{H}_0: \mu \geq 45$$
$$\text{H}_1: \mu < 45$$
$$44.5 \pm 1.96 \cdot \frac{11.41}{\sqrt{10}}$$
$$= 44.5 \pm \frac{11.41}{3.162}$$
$$= 44.5 \pm \frac{11.41}{3.162} \Rightarrow 37.42 \text{ to } 51.57$$

if < 37.42 & > 51.57 are rejection areas.

If $\mu_0 = 45$, then $\gamma 37.42$ & < 51.57 there is no statistical difference.

14/8/19

25, 60, 43, 56, 32, 43, 47, 59, 39, 41

$$n = 10$$

$$s = 11.41$$

$$\bar{x} = 44.5$$

$$\mu_0 = 45$$

Performance initiatives taken by the HR Team are not upto mark [one tail test].

$$H_0 : \mu \geq 45 \quad (\text{null hypothesis})$$

$$H_1 : \mu < 45 \quad (\text{Alternative})$$

commonly $\alpha = 5\%$. it may be 1%, 2% or any thing.

so, $C.I \rightarrow \text{confidence interval} : 95\%$.

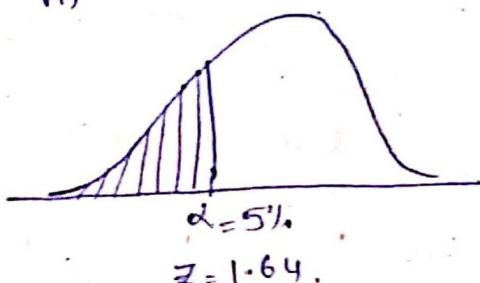
$$\text{confidence interval} \quad \bar{x} \pm z \frac{\sigma}{\sqrt{n}}$$

$$\sqrt{10} = 3.162$$

$$\bar{x} \pm z \frac{\sigma}{\sqrt{n}}$$

$$\Rightarrow 44.5 \pm (1.64) \left(\frac{11.41}{\sqrt{10}} \right)$$

$$\Rightarrow \mu = 38.5$$



allowed region, > 37.42 & < 51.57

$\therefore \mu = 38.5$ is in the allowed region.

so, H_0 accepted, H_1 Rejected.

1) $H_0: \mu < 45$

$H_a: \mu \geq 45$

$\alpha = 5\%$

$Z = 1.64$. [One tail test]

$$Z \Rightarrow \frac{\bar{x} - \mu_0}{\left(\frac{\sigma}{\sqrt{n}}\right)} = \frac{44.5 - 56}{\left(\frac{11.41}{\sqrt{10}}\right)} \approx \frac{-8.5}{\left(\frac{11.41}{3.162}\right)} \Rightarrow -3.12$$

Pvalue = -3.12.

16/819

bulbs.
A Manufacturer purchase that are suppose to burn for a mean life time of atleast 3000 hrs. with a standard deviation of 500 hrs. If a sample of 100 bulbs is taken with mean $\bar{x} = 2800$ hrs

$n = 100$

$\bar{x} = 2800$ hrs.

$H_0: \mu \geq 3000$ hrs.

$H_1: \mu < 3000$ hrs.

- decide $\alpha = 5\%$.

confidence : 95%.

then $Z = 1.64$.



3000
 -1.645

so, single tail test

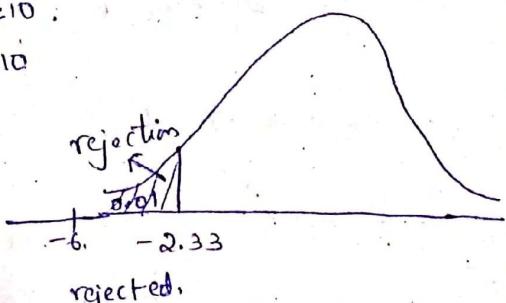
$$Z = \frac{\bar{X} - \mu}{\left(\frac{\sigma}{\sqrt{n}}\right)} \Rightarrow \frac{2800 - 3000}{\frac{500}{\sqrt{100}}} \\ \Rightarrow \frac{-200}{50} \\ \Rightarrow -4$$

Z is in rejection area so, it is rejected.

reject null hypothesis.

Q30 A company claims that its weight reducing drug will cause atleast 10 pounds within 1 month. A random sample of 64 subjects is taken and the avg. weight loss is 7 pounds. with standard deviation $s=4$ pounds & $\alpha=0.01$

$$\rightarrow n = 64 \quad H_0: \mu \geq 10 \\ \bar{x} = 7 \quad H_1: \mu < 10 \\ s = 4 \\ \alpha = 0.01$$



$$Z = \frac{\bar{X} - \mu}{\left(\frac{\sigma}{\sqrt{n}}\right)}$$

$$\Rightarrow \frac{7 - 10}{\frac{4}{\sqrt{64}}} \Rightarrow \frac{-3}{\frac{1}{4}} = -12$$

rejected.

reject null hypothesis

so, given alternate hypothesis is true.

A researcher claims that 10 years old watch 6.6 hrs of TV daily. You try to verify this with following sample data.

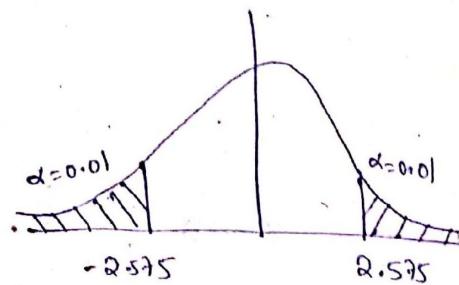
$$n = 100$$

$$\bar{x} = 6.1 \text{ hrs}$$

$$s = 2.5 \text{ hrs}$$

$$\alpha = 0.01$$

$$H_0: \mu = 6.6 \text{ hrs}$$



$$H_1: \mu \neq 6.6 \text{ hrs}$$

$$Z = \frac{\bar{x} - \mu}{\left(\frac{\sigma}{\sqrt{n}}\right)} \Rightarrow \frac{6.1 - 6.6}{\frac{2.5}{\sqrt{100}}} \Rightarrow \frac{-0.5}{0.25} \Rightarrow \frac{-5}{25} \times \frac{10}{10} \Rightarrow -\frac{1}{5} \times \frac{10}{10} \Rightarrow 0.2 \times 2 = \underline{\underline{0.2}}$$

it is ~~notin~~ ^{not} -2.575 i.e., in rejection region.

\Rightarrow then, H_1 is rejected.

$$\text{Confidence interval: } \bar{x} \pm \frac{s}{\sqrt{n}}$$

$$\Rightarrow 6.1 \pm 2.575 \times \frac{2.5}{\sqrt{100}}$$

$$\frac{14.1612}{2.575 \times 0.25} \\ \underline{\underline{0.64375}}$$

$$= 6.1 \pm (2.575 \times 0.25)$$

$$\frac{6.10}{0.64} \\ \underline{\underline{5.36}}$$

$$\Rightarrow 6.1 \pm 0.64 \underline{\underline{0.64}}$$

$$\Rightarrow 6.1 \pm 0.64$$

$$\Rightarrow \underline{\underline{5.36 \text{ to } 6.74}}$$

A manufacturer produces bolts with a thickness of exactly 1 inch. A customer takes a random sample of 100 bolts, and finds $\bar{x} = 1.2$ inches, $s = 0.40$ inches. Should the manufacturer's claim that the bolts be rejected?

$$n = 100$$

$$\bar{x} = 1.2 \text{ inches}$$

$$s = 0.40 \text{ inches}$$

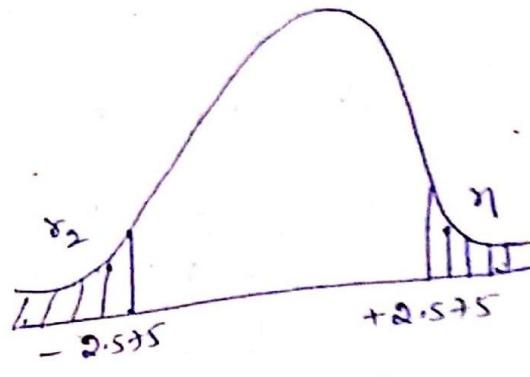
$$\alpha = 0.01$$

$$z = 2.575$$

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

$$\Rightarrow \frac{1.2 - 1}{\frac{0.40}{\sqrt{100}}} \Rightarrow \left(\frac{0.2}{0.40} \right) \Rightarrow \frac{0.2 \times 10}{0.40} \Rightarrow \frac{2}{0.40} \Rightarrow \left(\frac{2}{100} \right) \Rightarrow \frac{200}{400} \Rightarrow 5$$

$$\Rightarrow +5$$



$+s$ is in τ_f region so, it is rejected.

using confidence intervals.

$$\bar{x} \pm z \frac{s}{\sqrt{n}}$$

$$\frac{2.575^2 \times 0.40}{\sqrt{100}}$$

$$\rightarrow 1.2 \pm 2.575 \times \frac{0.40}{\sqrt{100}}$$

$$= 1.2 \pm \left(2.575 \times \frac{0.40}{10} \right)$$

$$= 1.07 \text{ to } 1.303$$

A company claims that their battery have a life of atleast 100 hrs. The following data:

$$n = 121$$

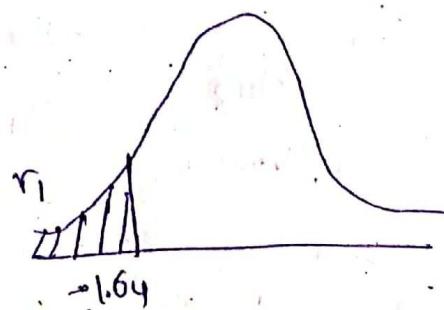
$$s = 3 \text{ hrs}$$

$$\bar{x} = 97 \text{ hrs}$$

$$\alpha = 0.05 \Leftrightarrow 5\%$$

$$H_0: \mu \geq 100 \text{ hrs}$$

$$H_1: \mu < 100 \text{ hrs}$$



$$Z = \frac{\bar{x} - \mu}{\left(\frac{s}{\sqrt{n}}\right)} \Rightarrow \frac{97 - 100}{\frac{3}{\sqrt{121}}} \Rightarrow \frac{-3}{\left(\frac{3}{11}\right)} \Rightarrow \underline{\underline{-11}}$$

hence it is in rejection region. so, reject null hypothesis

- Sample size should be ≥ 30 .
- It should follow standard normal distribution.