

STATISTICS FOR DATA SCIENCE DAY 5

* Inferential Statistics

(i) Hypothesis Testing

(ii) p-value

(iii) Confidence Interval

(iv) Significance Value

χ^2 -test

t-test

Chi-square test

ANOVA test (F-test)

3- Distributions :-

(i) Bernoulli

(ii) Binomial

(iii) Poisson

TRANSFORMATION

HYPOTHESIS TESTING

Steps of hypothesis testing

① Null hypothesis [the default one will be in null hypothesis]

e.g. (i) Person is not a criminal until and unless he is met person we will say he is not criminal.

(ii) Let, we have tossed a coin, now we check whether coin is fair or not

so, By default we consider "Coin is fair" null hypothesis

② Alternate hypothesis [opposite of Null hypothesis]

e.g. (i) Coin is not fair.

(ii) Person has committed crime

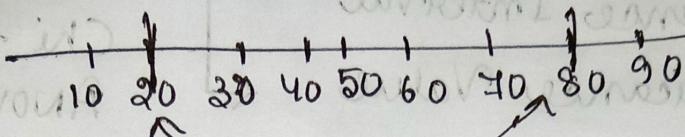
• Perform Experiments

Let, say I have tossed 100 times coin
mean = 50, Standard deviation = 10

50 times head



Fair



60 times head

Fair

70 times head

We go to

Domain expect to get a range

and that range is called confidence interval

$$C.I = [20 - 80]$$

coin is fair

If we get 75 heads then null hypothesis is accepted

If we get 10 heads then null hypothesis is rejected

and alternate hypothesis is accepted

* We reject the null hypothesis [outside C.I.]

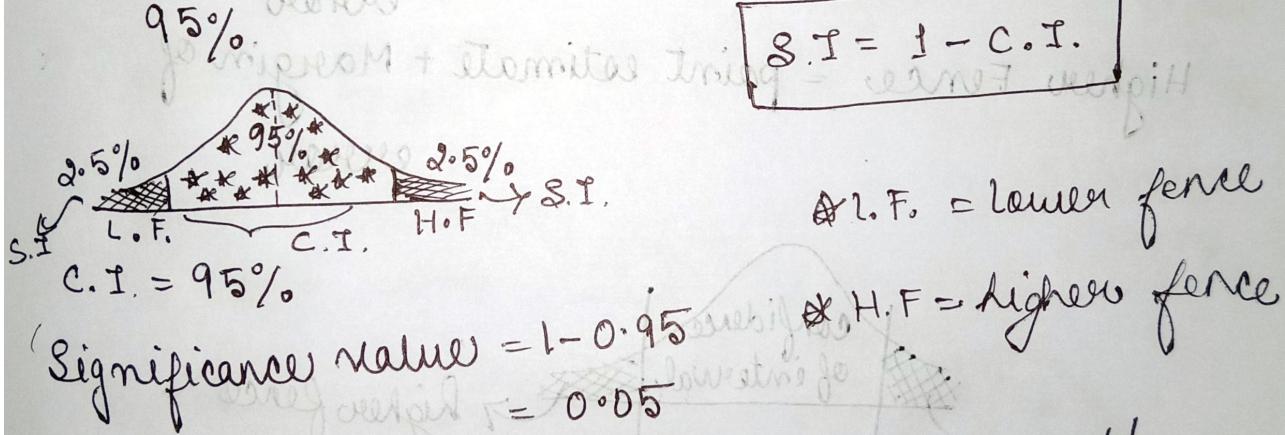
* We fail to reject the null hypothesis [within C.I.]

e.g., whether the person is criminal or not [Murder yesterday for ~~case~~ ^{new case}]

- (i) Null hypothesis : Person is not criminal
- (ii) Alternate hypothesis : Person is criminal
- (iii) Experiment / Proof : DNA, fingerprint, weapons, eye, witness, footage

Judge [take decision]
Conclusion :- Based on Judge's decision

Confidence Interval (C.I.)



- * In S.I. region we reject the null hypothesis
- * In C.I. region we fail to reject the null hypothesis

Point Estimate :- The value of any statistic

that estimates the value of a parameter is called Point Estimate.

e.g. mean of sample is used to estimate the mean of parameter

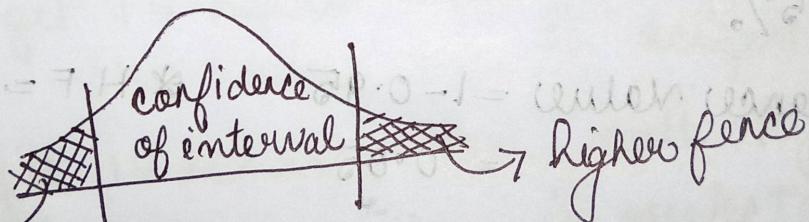
$\bar{x} \xrightarrow{\text{estimate}} \mu \xrightarrow{\text{parameter}}$
statistic

Parameter \Rightarrow Population Mean

$$\text{Point estimate} \pm \text{Margin} = \text{Parameter}$$
$$(\bar{x}) \qquad \qquad \qquad \text{of error} \approx (\mu)$$

Lower Fence = point estimate - Margin of error

Higher Fence = point estimate + Margin of error



$$\text{Margin of error} = \frac{6}{\sqrt{n}}$$

where, α = significance value

$$\frac{6}{\sqrt{n}} = \text{standard error}$$

Problem: On the quant test of CAT Exam, a sample of 25 test takers have a mean of 520 with a sample standard deviation of 100. Construct a 95% C.I. about the mean?

Sol: Given, $n = 25$, $\bar{x} = 520$, $s = 100$
 $C.I. = 95\%$, S.V. = $1 - C.I. = 1 - 0.95 = 0.05$

Lower C.I. = Point estimate - Margin of error

$$= 520 - \frac{100}{2 \cdot \sqrt{25}} \cdot 0.05$$

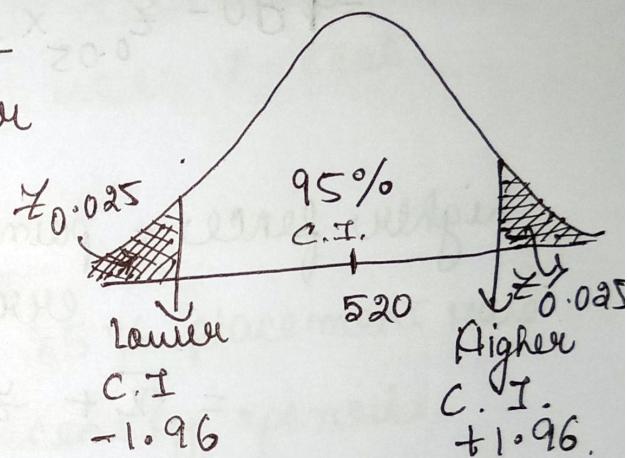
$$= 520 - \frac{20}{0.025}$$

$$= 520 - 1.96 \times 20$$

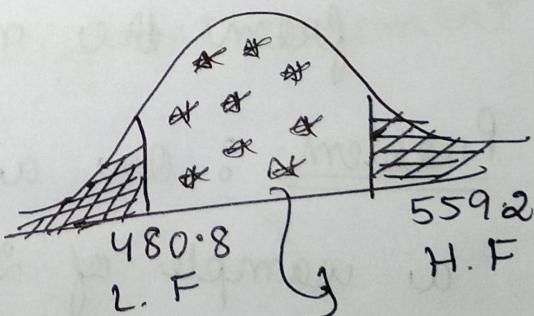
$$= 520 - 39.2 = 480.8$$

$$\text{Higher C.I.} = 520 + 39.2$$

$$= 559.2$$



$$1 - 0.025 = 0.975$$

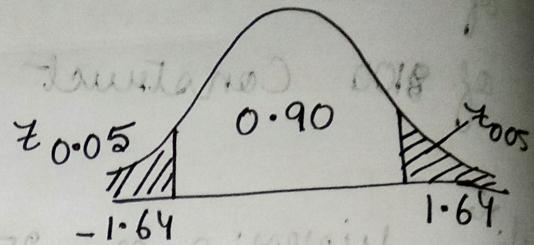


Value came under this fail to reject null hypothesis

Problem $\bar{x} = 480$, $\sigma = 85$, $n = 25$, C.I = 90%

$$S.V. = 1 - C.I = 1 - 0.90 = 0.10$$

Lower fence = \bar{x} - Margin of error



$$= \frac{480 - z_{0.05} \cdot \frac{85}{\sqrt{25}}}{480}$$

$$= 480 - z_{0.05} \times 17 = 480 - 1.64 \times 17$$

$$= 452.12$$

Higher fence = point estimate + Margin of error

$$= \bar{x} + z_{0.025} \frac{\sigma}{\sqrt{n}} = 480 + 1.64 \times 17$$

$$= 507.8$$

* here ± 1.64 is the standard deviation from the mean

Problem : On a quant test of CAT exam a sample of 25 test takers has a mean of 520, with a sample standard deviation of 80. Construct 95% C.I about the mean?

Sol: Given, $n = 25$, $\bar{x} = 520$, $s = 80$

$$S.V. = 1 - 0.95 = 0.05 \quad C.I = 95\%$$

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

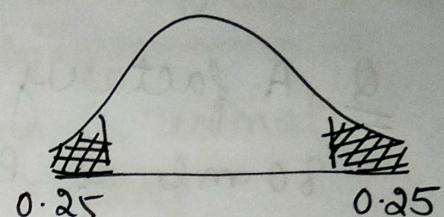
t-test

Degree of freedom

$$= n-1 = 25 - 1 = 24$$

$$\text{Lower C.I.} = 520 - t_{0.05/2} \left(\frac{80}{5} \right)$$

$$= 520 - 2.064 \times 16 \\ = 486.976$$



$$\text{Higher C.I.} = 520 + 2.064 \times 16 = 553.024$$

* Since sample standard deviation is given therefore we use t-test.

1 Tail and 2 Tail Test

* Colleges in Town A has 85% placement rate.

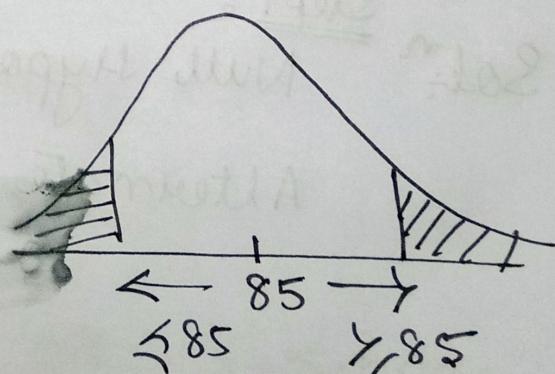
A new college was recently opened and it was found that a sample of 150 students had a placement rate of 88% with a standard deviation of 4%.

Does this college has a different placement rate with 95% C.I.?

* greater than 85%



1 tail test Right side



* lesser than 85%



1 tail test Left side

- ① χ^2 -test
- ② t-test

Hypothesis Testing Problem

Q A factory has a machine that fills 80ml of Baby medicines in a bottle. An employee believes the average amount of baby medicine is not 80ml. Using [Two tailed Test] 40 samples, he measures the average amount dispersed by the machine is to be 78ml with a standard deviation of 2.5.

State (a) Null & Hypothesis

(b) At 95% c.i., is there enough evidence to support machine is working properly or not.

Given : $\mu = 80 \text{ ml}$, $n = 40$, $\bar{x} = 78$, $s = 2.5$

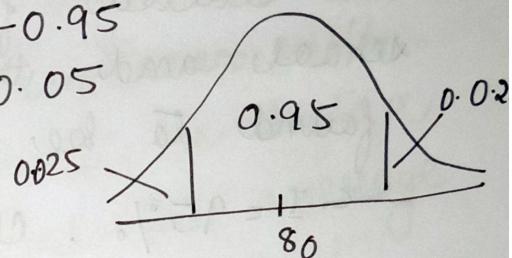
$$\text{C.I.} = 95\% = 0.95$$

~~$$\therefore \text{C.I.} = 1 - 0.95 = 0.05$$~~

Step 1 : Null Hypothesis : $\mu = 80$ i.e. machine is working

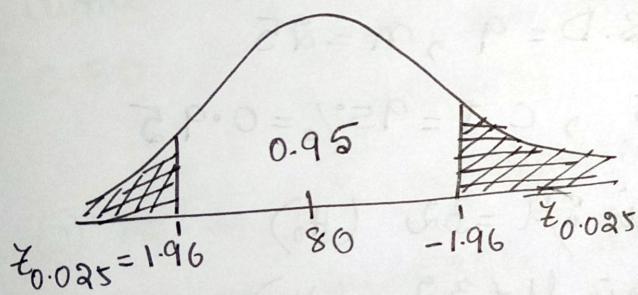
Alt. hypothesis : $\mu \neq 80$ i.e. machine is not working

Step 2 : $S. \sqrt{f(x)} = 1 - \text{C.I.} = 1 - 0.95 = 0.05$



Step 3 : Z -test

$$Z_{0.025} = 1.96$$



* When $n > 30$
sample s.d is given
use Z -test

* When $n < 30$ and
sample s.d is given
use t -test

Decision Boundary

1.96 & -1.96 is the boundary.

Calculate test statistics (Z -test)

$$Z = \frac{\bar{x} - \mu}{S/\sqrt{n}} = \frac{78 - 80}{2.5/\sqrt{40}} = \frac{-2}{2.5/\sqrt{40}} = -5.0596$$

Conclusion : Decision Rule → If $Z = -5.05$ is less than -1.96 or greater than $+1.96$ we reject null hypoth.

Reject the Null Hypotheses with 95% C.I
i.e. there is some faults in the machine

Q 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 kgs. A sample of 25 boys were selected from the government school and the average weight was found to be 29.5 kgs? with C.I = 95%. Check it is True or False

Sol: χ -test $\Rightarrow \mu, 30 \text{ or } \mu \leq 30 \text{ or popn s.d.}$

Given, $\mu = 32$, S.D = 9, n = 25

$$\bar{x} = 29.5, C.I = 95\% = 0.95$$

Step 1: Null hypothesis: $\mu = 32$ (H_0)
Alt. hypothesis $\mu \neq 32$ (H_1)

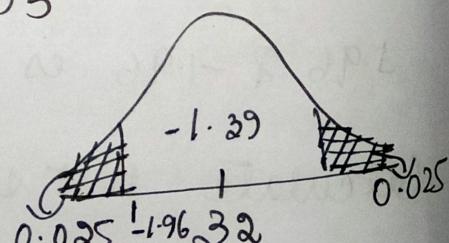
Step 2: S.V(χ) = $1 - C.I = 1 - 0.95$
 $= 0.05$

Step 3: χ -test

$$\chi_{0.025} = \frac{\bar{x} - \mu}{S.D / \sqrt{n}}$$

$$= \frac{29.5 - 32}{9 / \sqrt{25}} = \frac{-2.5}{1.8} = -1.39$$

$\chi^2(\chi_{0.025})$ by χ -test table - 1.96



Conclusion :- $1.39 < 1.96$

\therefore Accept the null hypothesis with 95% of C.I and we fail to reject null hypothesis i.e. the boys are fed well.

- 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 less than 5 years. He tests a sample of 40 cars and finds the average time to be 4.8 years with a S.D of 0.50
- ① State the null and alternate hypothesis

- ② At a 2% significance level, is there enough evidence to support the idea that the warning should be received

Sol:- Given, $\mu_0 = 5$, $n = 40$, $\bar{x} = 4.8$,
 $s = 0.50$.

- ① Step 1 : H_0 (Null hypothesis) : $\mu \geq 5$
 H_1 (Alt. hypothesis) : $\mu < 5$

- ② Step 2 : S.L. = 0.02, C.I = $1 - 0.02$
= 0.98

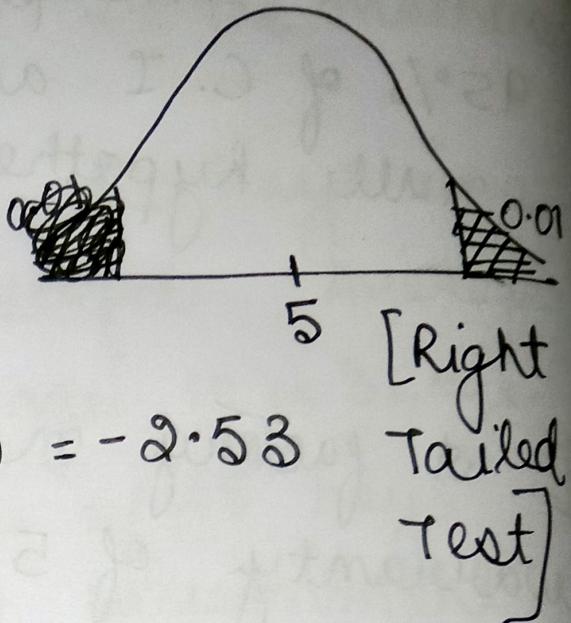
Step 3: Compute test-statistic

$$Z_{0.01} = \frac{\bar{x} - \mu}{s/\sqrt{n}}$$

$$= \frac{4.8 - 5}{0.50/\sqrt{40}} = \frac{-0.2}{0.0790}$$

$$= -2.529 = -2.53$$

$$Z_x = -2.33$$



$$\therefore -2.53 < -2.33$$

\therefore we accept the null hypothesis with 98% of C.I. and reject the alternate hypothesis i.e. yes the warning is received.