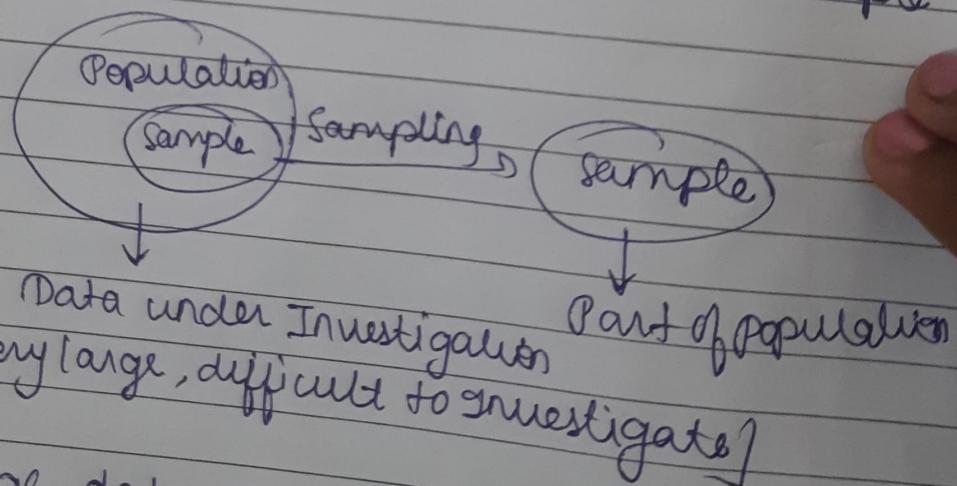


Population (N)

It includes all the elements from a set of data. The population is the whole set of values or individuals you are interested in. Population characteristics are mean (μ), standard deviation (σ), proportion (P), median, percentiles etc. The value of population is fixed. These characteristics are called population distribution.

- * Sample (n) → subset of population data and the set of values you actually use in your estimation.
This sample has some quantity computed from values eg. mean (\bar{x}), standard deviation (s), sample proportion etc. this is called sample distribution.



- * If large dataset then go for sample
of small then population

Statistics

kisi bhi data ke andar se info nikalne ke
upar mathematical analysis using particular
methods

descriptive

→ Jab ham (entire)
population data

ke upar analysis karte
hai then descriptive
is statistics :

→ measures of central
tendency (mean, median, mode)

→ measures of variability →
range, mean Absolute Deviation,
variance, Std

→ measures of shape → σ
skewness

→ percentiles → Boxplot

→ covariance and
correlation

Inferential statistics

→ sample ke upar
analysis.

sample

→ Central limit theorem

→ Hypothesis testing

→ Z-test

→ T-test

→ Chi-square-test

* Measurement of central tendency *

① Mean / Avg

$$\text{mean} = \frac{\text{Sum of all data}}{\text{Total no. of data}}$$

Import numpy as np

Import pandas as pd

$$\text{ar} = \text{np.array}[4, 5, 6, 2, 1, 8, 3, 6]$$

$$\text{np.sum(ar) / len(ar)}$$

$$\Rightarrow 4.8$$

$$\text{np.mean(ar)} \Rightarrow 4.8$$

② median → data

→ first sort

→ dataset['Age'].mean
 → sns.histplot(x="Age",
 data=dataset, bins=[
 i for i in range(0, 81, 10)])
 → plt.show()

$$\begin{aligned} &\rightarrow \text{even}, \frac{n_2 + (n+1)}{2} \\ &\rightarrow \text{odd} \rightarrow \frac{(n+1)}{2} \end{aligned}$$

$$\text{① np.median(dataset['Age'])}$$

$$\text{② dataset['Age'].median(),}$$

③ mode → most repeated

$$\text{class } \hookrightarrow \begin{array}{l} \mu = 82 \\ \text{std} = 20 \\ n = 81 \\ \bar{x} = 90 \end{array}$$

Q of have to find whether it justifies or not

Here we are given pop - std, mean
③ no. of sample data > 30
hence we use z-test

$$z\text{-test} = \frac{\bar{x} - \mu}{\left(\frac{\sigma}{\sqrt{n}}\right)}$$

\bar{x} = sample mean

μ = pop mean

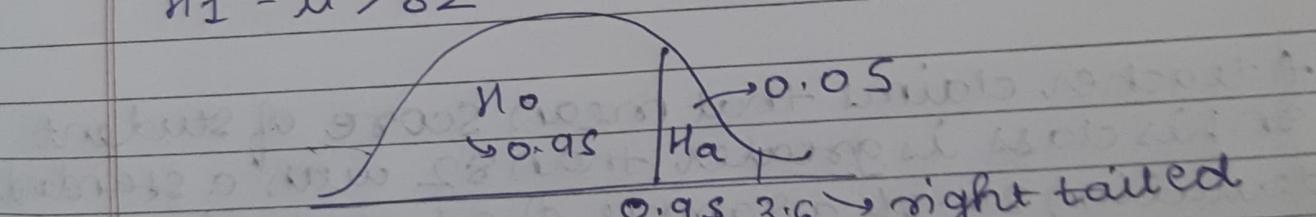
σ = pop std

n = no. of sample data

$$H_0: \mu \neq 82$$

$$H_1: \mu > 82$$

$$\alpha = 0.05$$



0.95 3.6 → right tailed

$$Z\text{ value} = 1.6408 \rightarrow Z_{\text{cal}} = 8.6$$

$$Z_{\text{cal}} = \frac{90 - 82}{\left(\frac{20}{\sqrt{81}}\right)} = \frac{8}{\left(\frac{20}{9}\right)} = \frac{72}{20} = 3.6$$

Z_{cal} is in H_a area hence teachers claim that $\mu > 82$ is true

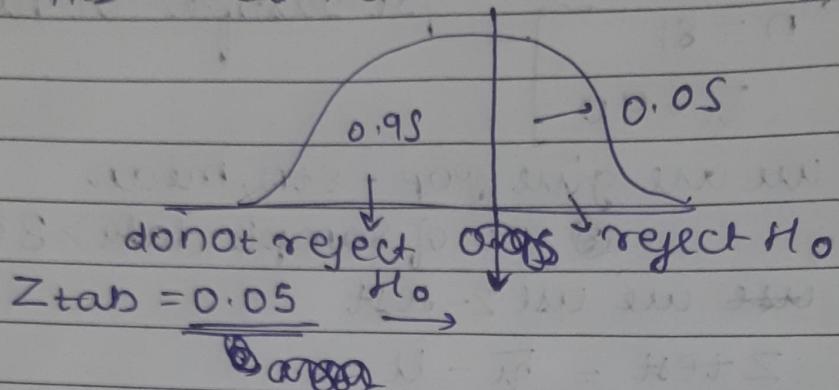
The teacher believe true mean > 82

$$H_0 = \text{mean} \neq 82$$

$$H_1 = \text{mean} > 82$$

default level of significance = 0.05

mean what is believed to be correct



→ you look for z-value where area to the left is 95% (0.95)

$$\text{area} = 0.95 \rightarrow Z_{\text{tab}} \rightarrow \underline{\underline{1.64}} \quad \begin{matrix} 1.60 \\ 0.40 \\ 1.64 \end{matrix}$$

$$Z_{\text{cal}} = 3.6$$

which means we reject H₀.
hence it is true

- Q A teacher claims that mean score of student in his class is greater than 82 with a standard deviation 20. If a sample of 81 students was selected with a mean score of 90.

Q2) Imagine you work for an e-commerce company, and your team is responsible for analyzing customer purchase data. You want to find out whether a new website design has led to a significant increase in the average purchase amount compared to the old design.

Below is data of random 30 samples of purchase they made on website

$$\text{old design data} = [45.2, 42.8, 38.9, \dots, 4]$$

$$\rightarrow \text{new design data} = [48.5, 49.1, 50.2, \dots]$$

$$\text{Population StD} = 2.5$$

→ $H_0: \bar{W}_{\text{new}} = \bar{W}_{\text{old}}$

$H_{01}: \bar{W}_{\text{new}} > \bar{W}_{\text{old}}$

$$\bar{W}_{\text{new}} = \bar{W}_{\text{old}}$$

$$\bar{X}_{\text{new}}, \bar{X}_{\text{old}}$$

$$\sigma = 2.5$$

$$n = 30$$

$$Z_{\text{test}} = \frac{\bar{X} - \mu}{\left(\frac{\sigma}{\sqrt{n}}\right)} \quad n \geq 30$$

$$Z_{\text{test}} = \frac{\left(\bar{W}_{\text{new}} - \bar{W}_{\text{old}}\right) - (\mu_{\text{new}} - \mu_{\text{old}})}{\left(\frac{\sigma}{\sqrt{n}}\right)}$$
$$= \frac{\left(\bar{W}_{\text{new}} - \bar{W}_{\text{old}}\right)}{\left(\frac{\sigma}{\sqrt{n}}\right)}$$

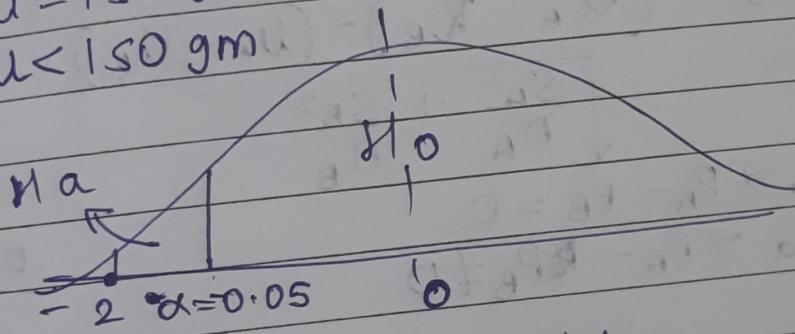
T-test

Q1) A manufacturer claims that the average weight of a bag of potato chips is 150 gm. A sample of 25 bags is taken, and the average weight is found to be 148 gm, with a standard deviation of 5 gm. Test the manufacturer's claim using a one-tailed t-test with a significance level of 0.05.

$$\mu = 150 \quad \bar{x} = 148 \quad n = 25 \quad s = 5$$

$$H_0: \mu = 150 \text{ gm}$$

$$H_a: \mu < 150 \text{ gm}$$



$$df = n - 1 = 24$$

$$t = -1.711 \quad (\text{cal}) \quad (\text{tab}) \quad (\text{as on left side})$$

$$t_{\text{cal}} = \left(\frac{148 - 150}{\frac{5}{\sqrt{25}}} \right) = -2 = (-2)$$

Hence t_{cal} is in H_a means $\mu < 150$.
Hence company's claim is wrong

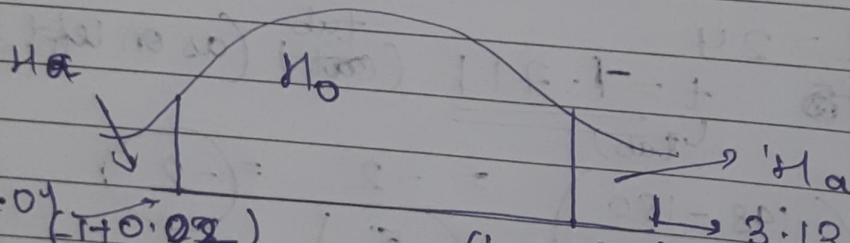
Q2) A company wants to test whether there is difference in productivity between 2 teams. They randomly select 20 employees from each team and record their productivity scores. The mean productivity score for Team A is 80 with SD of 5, while mean productivity score for Team B is 75 with a standard deviation of 6. Test at 5% level of significance whether there is a difference in productivity between 2 teams.

$$t = \frac{(\bar{x}_A - \bar{x}_B) - (\mu_A - \mu_B)}{\sqrt{\frac{s_A^2}{n_A} + \frac{s_B^2}{n_B}}}$$

$$H_0 \rightarrow \mu_A - \mu_B = 0$$

$$H_a \rightarrow \mu_A - \mu_B \neq 0$$

$$\frac{0.05}{2} = 0.025$$



$$t = -2.09 \quad \{T+0.02\}$$

$$df = 20 + 20 - 2 = 38 \quad t = 2.024$$

$$t_{cal} = \frac{(80 - 75) - 0}{\sqrt{\frac{25}{20} + \frac{36}{20}}} = \frac{5}{\sqrt{61}} = 3.13$$

$\therefore H_a$ is right; H_0 is rejected.

Here we have taken degree of freedom (df) as $n-1$ because typing speed of same sample was recorded in previous question samples were diff.

Q3 A company wants to test whether the new training program improve the typing speed of all employees. The typing speed of 20 employees was recorded before and after the training program. The data is given below. Test at 5% level of significance whether the training program has an effect on the typing speed of the employees.

• Before: 50, 60, 45, 65, 55, 70, 40, 75, ... 65

• After: 60, 70, 55, 75, 80, 50, ...

$$H_0: t_B - t_A = 0$$

$$H_a: t_B - t_A \neq 0$$

$$(1 - 0.025)$$

$$\Leftrightarrow 2.09$$

$$(0.025)$$

$$\Rightarrow 2.09$$

$$df = n - 1 = 19$$

$$t = \bar{u}_A - \bar{u}_B$$

$$\sqrt{\frac{s_A^2}{n} + \frac{s_B^2}{n}}$$

$$t_{\text{table}} = St \cdot t \cdot PPF(1 - 0.025, 19)$$

t_{table}

$$\Rightarrow 2.09302$$

$$\text{Before} = np \cdot \text{am}([50, 60, 45, \dots, 65])$$

$$\text{After} = np \cdot \text{am}([60, 70, 55, 75, 80, \dots])$$

$$\text{std_a} = np \cdot \text{STD}(\text{After})$$

$$\text{std_b} = np \cdot \text{STD}(\text{Before})$$

$$\text{mean_a} = np \cdot \text{mean}(\text{After})$$

$$\text{mean_b} = np \cdot \text{mean}(\text{Before})$$

$$t_{\text{cal}} = (\text{mean_a} - \text{mean_b}) / (\text{np} \cdot \text{sqrt}((\text{std_a}^2 + \text{std_b}^2) / (n(n-1))))$$

$$t_{\text{cal}}$$

$$\Rightarrow 2.0612005$$

$\chi^2_{\text{cal}} < \chi^2_{\text{tab}}$

H_0 is true

Chi-square test

→ goodness

→ relationship Independence

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

$\chi^2 \rightarrow$ chi-squared

O_i = observed value

E_i = Expected value.

(Q1)

A fair die is rolled 120 times and the following results are obtained:

face 1 : 22 times

face 2 : 17 times

face 3 : 20 times

face 4 : 26 times

face 5 : 22 times

face 6 : 13 times

Test at a ~~level~~ 5% level of significance whether the die is fair

$H_0 \rightarrow$ die is fair

$H_a \rightarrow$ die is not fair

$$\rightarrow df = n-1(6-1) = 5$$

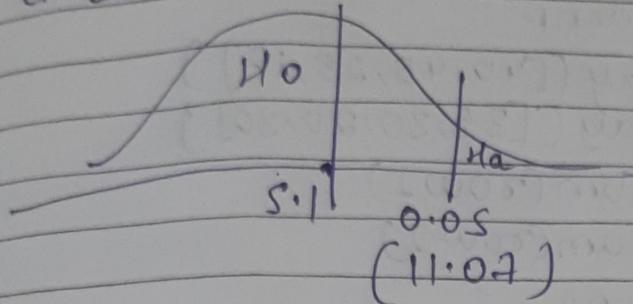
$$\chi^2_{\text{tab}} = 11.07$$

$$\chi^2_{\text{cal}} = \sum (O - E)^2$$

```

import numpy as np
ob = np.array([22, 17, 20, 26, 22, 13])
ex = np.array([20, 20, 20, 20, 20, 20])
np.sum(np.square(ob-ex))/ex
= 5.000000000000001
chical = 5.1

```



$\therefore H_0 \rightarrow$ die is fair

Q A study was conducted to investigate whether there is a relationship between gender and the preferred genre of music. A sample of 235 people was selected, and data collected is shown below. Test at 5% level of significance whether there is a significant association between gender and music preference.

	Pop	Hip hop	Classical	Rock
Male	40	45	25	10
Female	35	30	20	30

H_0 : No association

H_a : association

$$\alpha = 0.05$$

$$\begin{aligned}
df &= (row - 1)(col - 1) \\
&= (2 - 1)(4 - 1) \\
&= 3
\end{aligned}$$

$$\text{Chi}_{tab} = 7.815$$

$$\text{Expected} = \frac{\text{sum C1} \times \text{sum R1}}{\text{no. of sample}}$$

Male	40	45	25	10	→ sum R1
Fem	35	30	20	30	→ sum R2
					↓ sum C4
					↓ sum C3
					↓ sum C2
					↓ sum C1

import numpy as np

row1 = np.array([40, 45, 25, 10])

row2 = np.array([35, 30, 20, 30])

sum_r1 = np.sum(row1)

sum_r2 = np.sum(row2)

sum_r1, sum_r2

// (120, 115)

sum_cat = row1 + row2

sum_cat

array([[75, 75, 45, 40]])

sum_row = np.array([sum_r1, sum_r2])

[120, 115]

exp = []

for i in sum_row:

 for j in sum_cat:

 print(i, j)

 value = (i * j) / 235

 exp.append(value)

exp

[38.297 ,
 38.2978 ,
 ,

 19.8244]

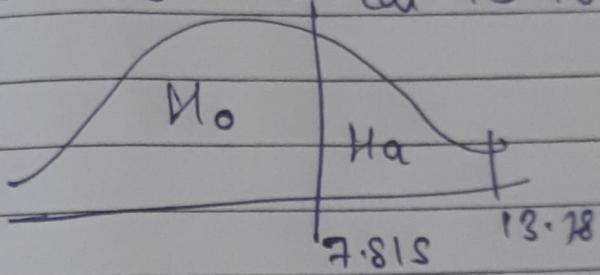
$\text{obj} = \text{np.array}([40, 45, 25, 10, 35, 30, 20, 30])$

$\text{np.sum}((\text{np.square}(\text{obj} - \text{exp})) / \text{exp})$

113.78874798

$\chi_{\text{tab}} = 7.815$

$\chi_{\text{cal}} = 13.78$



$\therefore H_0$ is wrong