

Statistics

Descriptive

(Description of the dataset)

Inferential

(Hypothesis Testing)

(Infer)

ML

algorithms

(KNN)

- ① Measures of Dispersion \leftrightarrow Data Imputation

(Mean, Median, Mode,

Range, Standard deviation)

Variance, IQR)

$$\begin{array}{l} \text{mean} = \\ \text{median} = \\ (\text{mode}) \end{array}$$

Symmetric / Normal

Gaussian

- ② Data Distribution

Skewed distribution

Right

Left

CTC

mean > median >

mode

mean <

median <

mode

$$\begin{aligned} \text{SNF} \rightarrow \mu = 0 \\ \sigma = 1 \end{aligned}$$



- ③ Central Limit Theorem

4

Z-score

$$\frac{x_i - \mu}{\sigma}$$

5

Sample vs Population

PMF

6

Random Variables

Discrete

Continuous

PDF

7

Visualization & Plots (Regplot, Catplot)

→ Scatterplot-, Barchart,

Box/Violin Plot (outliers),

Correlation heatmap,

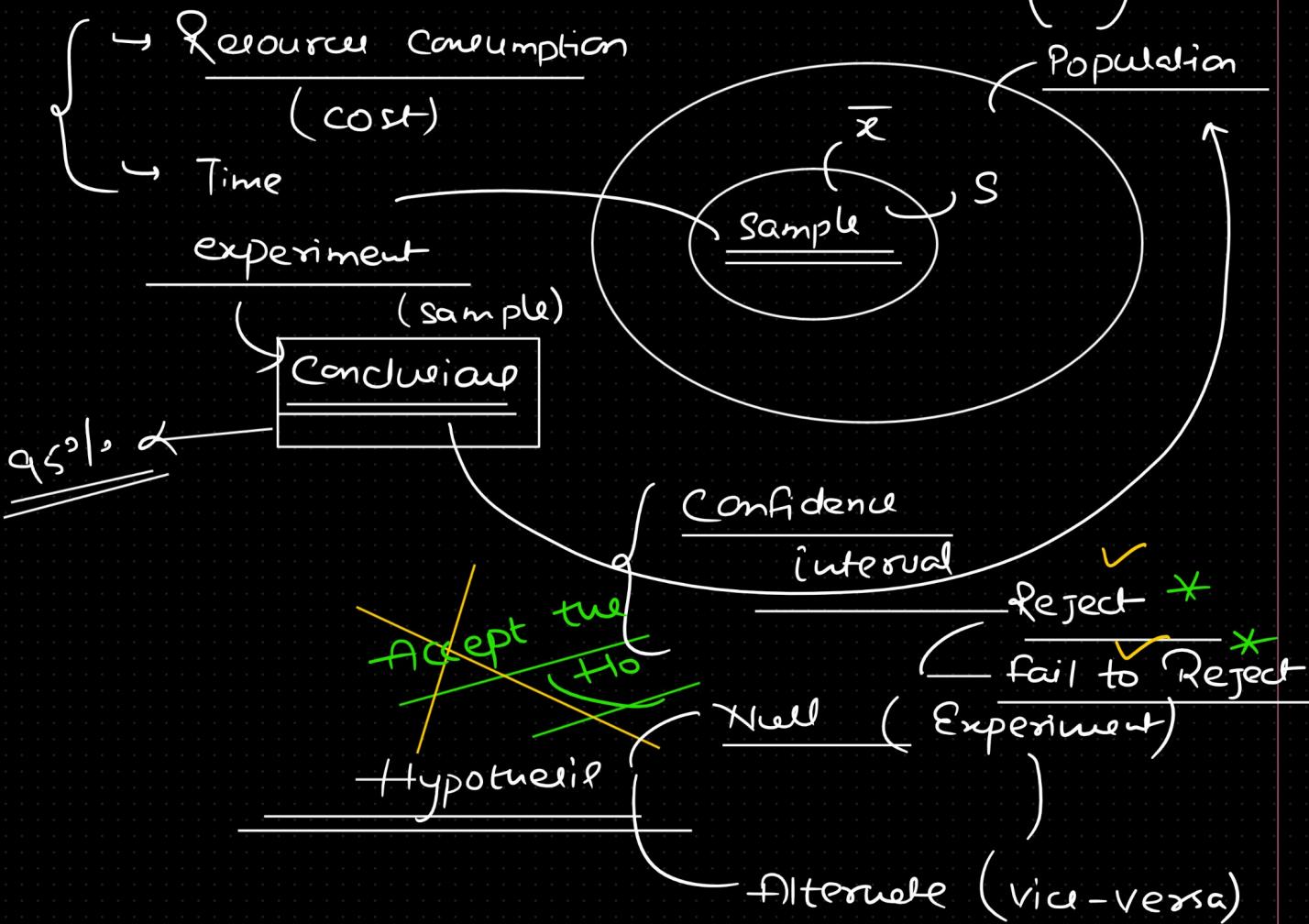
(feature selection)

Distribution plot,

Histogram

Inferential Statistics

$$\mu \quad \sigma$$



$$N = 10,000$$

$$n = 100$$

$$N > n$$

Population Mean	Sample Mean
$\mu = \frac{\sum_{i=1}^N x_i}{N}$	$\bar{X} = \frac{\sum_{i=1}^n x_i}{n}$
N = number of items in the population	n = number of items in the sample

$N = 10,000$

$n = 100$

Population Standard Deviation

$$\sigma = \sqrt{\frac{\sum(X_i - \mu)^2}{N}}$$

σ = population standard deviation
 N = the size of the population
 X_i = each value from the population
 μ = the population mean

Sample Standard Deviation

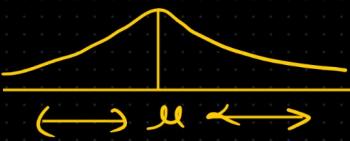
$$s = \sqrt{\frac{\sum(X_i - \bar{x})^2}{n-1}}$$

s = sample standard deviation
 n = the size of the sample
 X_i = each value from the sample
 \bar{x} = the sample mean

SD → measure of

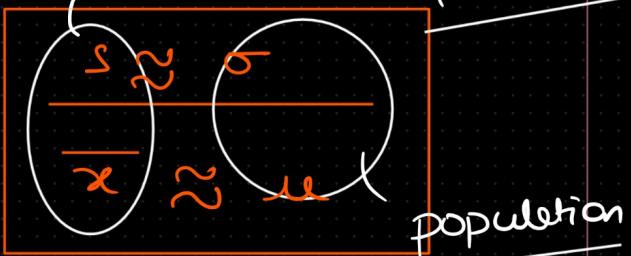
dispersion

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



$\longleftrightarrow \mu \longleftrightarrow$
 sample → valid

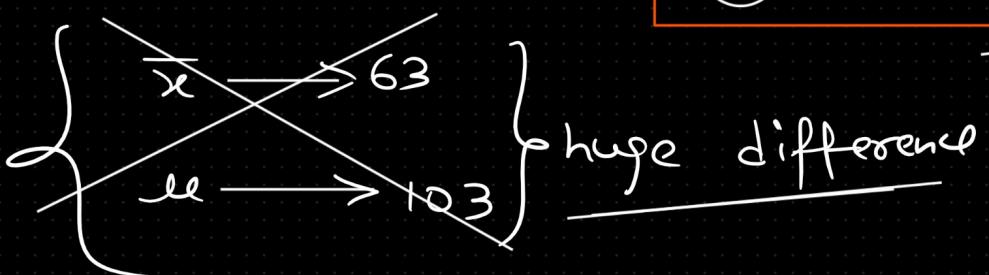
parameters



Why sample SD, division by

$n-1$??

Task 1



(ML)

Confusion Matrix

Type I and Type II Error		<u>Actual</u>
<u>Null hypothesis is ...</u>	<u>True ✓</u>	<u>False ✓</u>
<u>Rejected ✓</u>	Type I error False positive Probability = α	Correct decision True positive Probability = $1 - \alpha$
<u>Not rejected</u>	Correct decision True negative Probability = $1 - \alpha$	Type II error False negative Probability = β

Scribbr

$H_0 \Rightarrow$
Heart attack has no link with drinking coffee

False (Reality)

(Domain Experts)

$H_0 \rightarrow$ actually (True)

Reject H_0

Domain

Experiment

experts

Reality

Experiment \rightarrow fail to

(Good or
bad) reject H_0 .

(TN)

$H_0 \rightarrow$ True

Exp \rightarrow Reject H_0

Bad decision

(α) Type 1
error

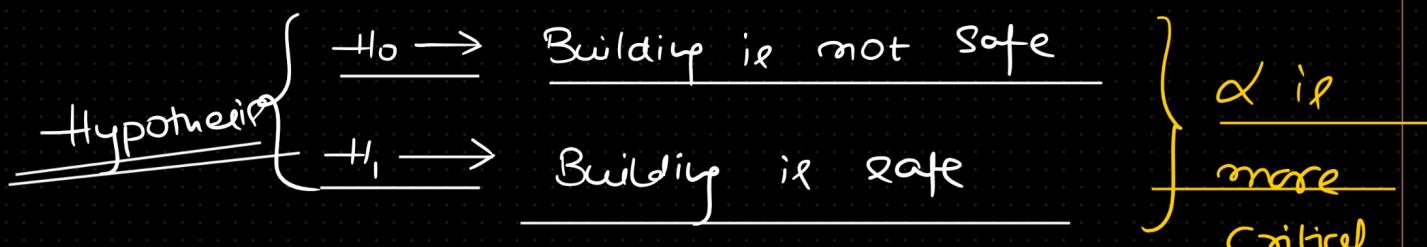
(Reality)

$H_0 \rightarrow$ False

Exp \rightarrow fail to Reject H_0

Bad decision (β)

Type 2 error



Type 1 Error \rightarrow $\left. \begin{array}{l} H_0 \rightarrow \text{True (Reality)} \\ \text{Exp} \rightarrow \text{Reject the } H_0 \end{array} \right\} \text{to handle}$

$\downarrow (\alpha)$

Type I error is committed if we reject "Building is not safe" when it is not safe.

Type 2 Error \rightarrow $\left. \begin{array}{l} H_0 \rightarrow \text{false (Reality)} \\ \text{Exp} \rightarrow \text{fail to reject } H_0 \end{array} \right\}$

$\downarrow (\beta)$

Type II error is committed if we fail to reject "Building is not safe" when it is safe.

To detect whether the patient is diabetic or not.

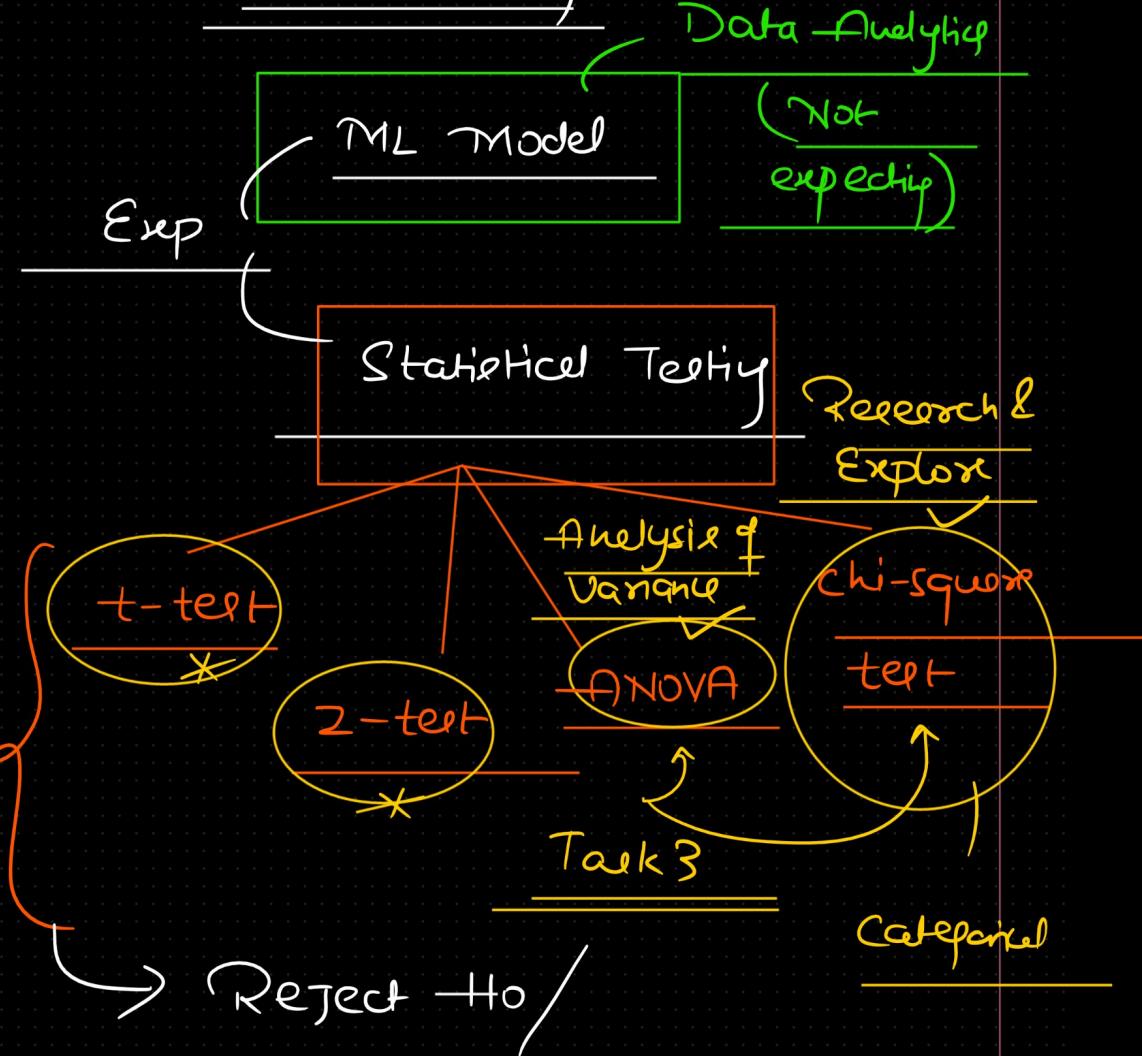
$H_0 \rightarrow$ Patient is not diabetic (0)
 $H_1 \rightarrow$ Patient is diabetic (1)

		<u>Reality</u>	
		1	0
<u>H_0</u>	<u>Reject</u>	<u>Correct</u>	<u>Type 2 Error (Type 1)</u>
	<u>fail to Reject</u>	<u>Error (Type 2)</u>	<u>Correct</u>
<u>H_1</u>	<u>Reject</u>	<u>(Type 1)</u>	<u>Task 2</u>
	<u>fail to Reject</u>		

Type I and Type II Error		
Null hypothesis is ...	True	False
Rejected	Type I error False positive Probability = α	Correct decision True positive Probability = $1 - \beta$
Not rejected	Correct decision True negative Probability = $1 - \alpha$	Type II error False negative Probability = β

Scribbr

Statistical Testing



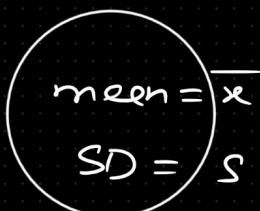
→ Reject H_0

→ Fail to reject H_0

Statistical Testing

Population

Sample



$$\text{mean} = \mu$$

$$SD = \sigma$$



$$\mu \approx \bar{x}$$

$$\sigma \approx s$$

Random

Sampling

Stratified Sampling

t-test

$\sigma = ??$

Yes

No

$$n > 30$$

Homoscedasticity

Yes

No

z-test

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

t-test

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

$$\alpha = 0.05$$

$$0.025 - 1.685$$



$$1 - 0.05$$

$$0.95$$

$$+ 1.685$$

$$0.025$$

Numerical Example: The average height of adults in a certain country is known to be 168 cm. A scientist believes that the average height of adult residents in a certain city in that country is different. The scientist measures the heights of 40 randomly selected adult residents of the city and finds an average of 171 cm with a standard deviation of 7 cm. Test the scientist's hypothesis at the 0.10 significance level.

Null $H_0 \rightarrow \mu = 168 \text{ cm}$

t-test

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

$$\bar{x} = 171 \text{ cm}$$

$$s = 7 \text{ cm}$$

$$n = 40$$

$$t = \frac{171 - 168}{\frac{7}{\sqrt{40}}} \quad \left\{ \begin{array}{l} \alpha = 0.10 \text{ (Risk)} \\ CI \Rightarrow 1 - \alpha = 1 - 0.10 \\ = 0.90 \end{array} \right.$$

$$t = 2.41$$

Alternative

Critical value ± 1.685
t-table Internet

$$dof = n - 1 = 39$$

$$\frac{2.41}{t\text{-test}} > \text{Critical value}$$

Reject H_0

↓
t-test & Critical value

The height is
different from 168 cm

Z-test, t-test \Rightarrow value

Critical value

(Internet-
t-table)

exp

Value & critical value

\rightarrow Fail to reject Ho

exp val > critical value

\hookrightarrow Reject Ho

two-tailed test \rightarrow 

one-tailed test

$$\alpha = 0.10$$

$$1 - \alpha = CI$$

A sample of 50 students using the new teaching method had an average score of 78 with a standard deviation of 10. The average score for students using the traditional method is known to be 75. We want to test the hypothesis at a 5% significance level ($\alpha = 0.05$)

$$\alpha =$$

$$\mu = 75 - H_0$$

$$\mu \neq 75 - H_1$$

$$t = \bar{x} - \mu$$

$$\bar{x} = 78$$

$$S / \sqrt{n}$$

$$S = 10$$

$$= 78 - 75$$

$$\alpha = 0.05$$

$$10 / \sqrt{50}$$

$$CI = 1 - \alpha$$

$$\approx 0.95$$

$$= 0.95$$

Refer the table

evaluated value

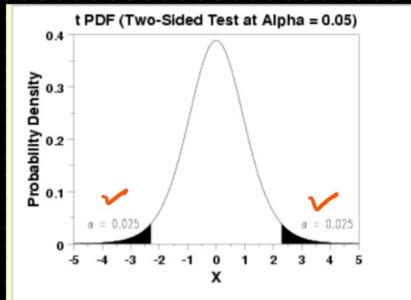
$$2.12 \geq 2.010$$

critical value

Reject H₀

one -sided

two - sided (by default)



$$\alpha = 0.05$$

$$1 - 0.025 \\ = 0.975$$

$$49, 0.975 \rightarrow 2.010$$

Diabetic

