

23-08-2025

Agenda: statistics



$Z\text{-stat}$ :

used → hypothesis test

→  $n > 30$  (CLT)

→ standardised how far the sample mean is from the population mean in units standard deviate.

$$Z = 1.50 \quad \alpha = 0.05$$

$$Z_{\text{critical}} = 1.645$$

$Z_{\text{score}} < Z_{\text{critical}} \rightarrow$  Fail to reject  $H_0$

→ True / False

p-value:

used → hypothesis test

pvalue  $\approx$  small value → evidence against  $H_0$

pvalue  $\approx$  large value → not enough evidence.

$p = 0.067 \rightarrow$  if null true  
 $\alpha = 0.05$  chance result are this extreme

$\rightarrow p > \alpha$  (fail to reject  $H_0$ )  
 $< \alpha$  (reject  $H_0$ )

$p = 0.01 \rightarrow$  1% prob that mean of your sample is due random.

$p \downarrow$  more  $\uparrow$

$0.01 < 0.02 < 0.03 < 0.04 < 0.05$

reject  $H_0$

Type I & Type II Error :-

$TP, TN, FP, FN \rightarrow$  Type II Error (missed detection)  
↓  
Type I error (false alarm)

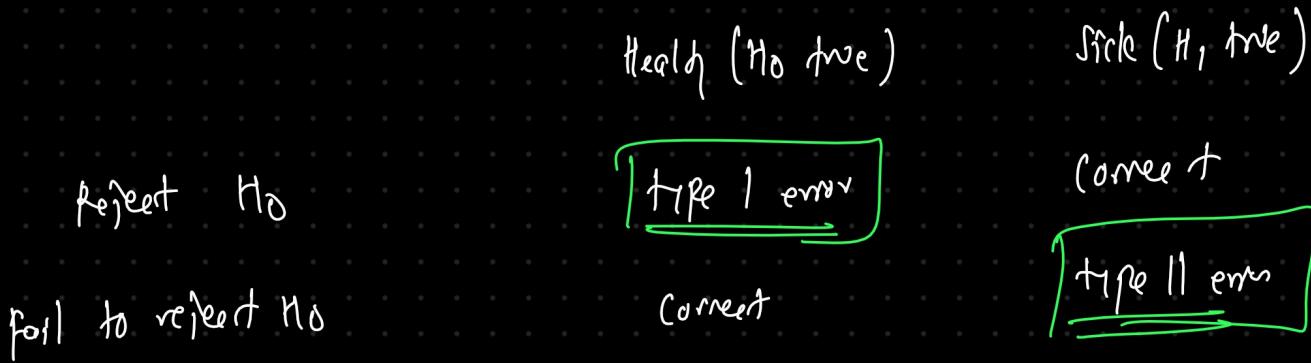
$H_0$ : The defendant is innocent.  
 $H_1$ : The defendant is guilty.

Imagine:

Real	Predicted	
Apple	Apple	→ TP
None	None	→ TN
Apple	None	→ FN
None	Apple	→ FP

$H_0$ : patient is healthy

$H_1$ : patient has disease



Z-test

- population  $\sigma$  (known)
- sample size ( $n > 30$ )
- Data  $\sim$  normal (CTL applies)

T-test

- population  $\sigma$  is unknown instead we use sample  $\sigma$
- works even with small sample size ( $n < 30$ )
- Data should be  $\sim$  normal

Degrees of freedom! —  
(df)

5 students:

avg marks of 5 students: 72

student-marks: [ A B C D ]  
fixed

$$72 = \frac{72 + 80 + 64 + 52 + (?)}{5} \text{ solve for } (?)$$

$$df = n - 1$$

---

T-test:

Pizza problem:

$H_0$ : avg delivery time  $\leq 30$  minutes :  $\mu \leq 30$

$H_1$ : takes longer than they claim :  $\mu > 30$

$$n = 25$$

$$\bar{x} = 32 \text{ minutes}$$

$$s = 4 \text{ minutes (sample } \sigma \text{)}$$

T-test formula:

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} = \frac{32 - 30}{4/\sqrt{25}} = 2.5$$

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

$\alpha = 0.05, df = 24 \rightarrow$  lookup on t-table  $\rightarrow t_{0.05, 24} \rightarrow 1.711$   
(one tailed)

$$t_{\text{calc}} = 2.5, t_{\text{critical}} = 1.711$$

$t_{\text{calc}} > t_{\text{critical}} \rightarrow \underline{\text{reject } H_0}$

Chi-squared distribution:

→ Categorical data (count, frequencies)  
→ variable →  $\chi^2$

(1) goodness of fit / Independence Test

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

$O \rightarrow \text{observed freq}$   
 $E \rightarrow \text{expected frequency}$

## (2) Test of variance (spread)

$$\chi^2 \rightarrow \frac{(n-1)s^2}{\sigma^2}$$

$s^2 \rightarrow$  sample variance  
 $\sigma^2 \rightarrow$  population variance  
 $n \rightarrow$  no. of sample

$Z/T \rightarrow$  is my mean different from a claim

$\chi^2 \rightarrow$  is my spread (variance) different

$H_0$

$$\text{population } \sigma^2 = 16 \text{ minutes}^2 \quad \text{mean } \pm 2 \cdot \text{STD}$$

$$\text{popul } \sigma = 4 \text{ minutes}$$

$$0.05 \rightarrow 95\% \rightarrow \underline{\text{from the mean}}$$

95% delivery should fall within  $\pm 8$  minutes of the average delivery time

$H_1$

(customers) are saying deliveries are more inconsistent.

$$H_a: \sigma^2 > 16$$

$n = 36$
$s^2 = 28$

$$\chi^2 = \frac{(36 - 1)20}{16} = 43.75$$

$$df = 36 - 1 = 35 \quad \leftarrow \quad \chi^2_{0.05, 35} = \frac{43.773 + 55.758}{2} \\ \leftarrow 0.05 = 49.765 \\ = \underline{\underline{49.8}}$$

Chi score with Chi score critical : —

$43.75 < 49.8 \leftarrow$  fail to reject  $H_0$  / accept  $H_0$

↓  
not enough  
evidence.

### Goodness of fit / Independence Test

→ examine whether two categorical variable are associated

	Apple	Banana	Orange	Total
Male	30	20	10	60
Female	20	20	10	50
Total	50	40	20	110

$H_0$  : fruit preference and gender are independent.

$H_1$  : there is an association

step: calc. expected count.

step: compute the chi-square statistic.

step: degree of freedom

step: p-value

Step 1: Expected-count

$$\rightarrow \left[ \text{Row count} \times (\text{column total}) \right] / (\text{Grand total})$$

$$* M-A \rightarrow (60 \times 50) / 110 \rightarrow 27.27$$

$$* M-B \rightarrow (60 \times 40) / 110 \rightarrow 21.82$$

$$* M-O \rightarrow (60 \times 20) / 110 \rightarrow 10.91$$

$$* F-A \rightarrow (50 \times 50) / 110 \rightarrow 22.73$$

$$* F-B \rightarrow (50 \times 40) / 110 \rightarrow 18.18$$

$$* F-O \rightarrow (50 \times 20) / 110 \rightarrow 9.09$$

Step 2: Compute the chi-square statistic

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

$$\begin{aligned}
 * M-A &\rightarrow \frac{(30 - 27.27)^2}{27.27} \rightarrow 0.273 \\
 * M-B &\rightarrow \frac{(20 - 21.82)^2}{21.82} \rightarrow 0.152 \\
 + M-O &\rightarrow \frac{(10 - 10.91)^2}{10.91} \rightarrow 0.076 \\
 * F-A &\rightarrow \frac{(20 - 22.73)^2}{22.73} \rightarrow 0.328 \\
 * F-B &\rightarrow \frac{(20 - 18.18)^2}{18.18} \rightarrow 0.182 \\
 * F-O &\rightarrow \frac{(10 - 9.09)^2}{9.09} \rightarrow 0.091
 \end{aligned}$$

$$\chi^2 = 1.102$$

Step 3: DOF

$$(df) \rightarrow (\text{number of rows} - 1) \times (\text{number of cols} - 1)$$

$$\rightarrow (2-1) \times (3-1) \rightarrow 1 \times 2 = 2$$

$$\underline{\text{Step 4}}: \quad \chi^2, \quad df, \quad p \rightarrow$$

$$p \rightarrow 0.576$$

$p > \alpha \rightarrow$  fail to reject  $H_0$

$\downarrow$

insufficient evidence.

## Confidence interval (CI)

- range of values derived from a sample statistic, likely to obtain from parameter ( $\mu$ ) .
- It is always associated with confidence level, which quantifies the level of conf that the parameter lies within the interval.
- point estimates are imperfect, we don't get range from this.

formula:

$$CI = \bar{x} \pm (\text{critical value}) \times \left( \frac{\sigma}{\sqrt{n}} \right)$$

$\bar{x} \rightarrow$  Sample mean

$\sigma \rightarrow$  Population std deviation

$n \rightarrow$  Sample size

conf level  $\rightarrow$   $\text{critical } Z \rightarrow 1.96$

$$\bar{X} = 31 \text{ minutes}$$

$$\sigma = 4 \text{ minutes}$$

sample size - 36 deliveries

$$2 \text{ critical} = 1.96$$

$$\underline{\text{MOE}}: \quad \sum_{\text{critical}} X \text{ SE} \rightarrow 1.96 \times \frac{\sigma}{\sqrt{n}} = 1.96 \times \frac{4}{\sqrt{36}} \approx 1.307$$

Construct the CI:

$$\begin{aligned} C_I &= \bar{X} \pm \text{MOE} \\ &= 31 \pm 1.307 \end{aligned}$$

$$\text{Lower bound} = 29.693$$

$$\text{Upper bound} = 32.307$$

Final: The 95% CI for the true average delivery time

$$(29.693, 32.307) \text{ minutes,}$$



python → Tools → project → start module → EDA → ML

ML:

Data gathering

process

EDA

Feature engineer

Train / Test split

Algorithm

Testing

Export

Use