

Anova

- 1) Recap
- 2) Anova → One Way Anova
- 3) Assumption for Anova
- 4) Kruskal Wallis Test
- 5) QQ Plot
- 6) Levene Test
- 7) Shapiro's Test

⇒ Hypothesis Testing →

⇒ checking if the results (sample data) we are seeing is significantly different than expected value.

→ Z test → data is normal
• chance that data lies in extreme region

→ T test → • Matis equation → changed from Pop to Samp

→ Chi-Square → Categorical Variable
→ As data is Categorical → # of Values

⇒ Exp dist
data

← →
Compare the
diff

Observed
data

$$\chi^2 = \sum \frac{(\text{Exp} - \text{Act})^2}{\text{Exp}}$$

χ^2 being large \rightarrow wrong Hypothesis
 Reject Hypothesis

$\Rightarrow \cdot \text{Exp} \rightarrow$

- 1) Equal prob
- 2) Based on Hist distn

Independent Variable (Cat)

	A	B	C	D	Total
Cat	10	20	30		
Dog	20	40	60		
Horse	30	60	90		

Anova \rightarrow T Test \rightarrow 2 Samples at a time.

Aerofit \rightarrow 3 products \rightarrow Income

KP781	10L	12L	30L	40L	...
KP781	8L	10L	11L	--	--
KPS81	7L	12L	--	--	--

\Rightarrow
KP281
Ku81

KP281
KPS81

KP781
KPS81

- \Rightarrow n dist' Cat \rightarrow ${}^n C_2$ test to perform
 \Rightarrow Higher Chances of Error \rightarrow as each test will have their own error'

Why not to do pairwise T-Test

A	62
B	65
C	68
	Red
	pred

$$\begin{aligned} 65 &\rightarrow -3 \\ 65 &\rightarrow 0 \\ 65 &\rightarrow 3 \end{aligned}$$

error is 0

I only get $1/3$ corr

\Rightarrow ① T Test \rightarrow 2 Sample

χ^2 Test \rightarrow 2 Categorical variables [Cat vs Cat]

Anova \rightarrow >2 Sample

TTest

Num vs Cat
 \hookrightarrow 2nd

\Rightarrow Gender vs Income \rightarrow T-Test
(2 value)
(Num)

\Rightarrow Gender vs Product $\rightarrow \chi^2$
(Cat)
(Cat)

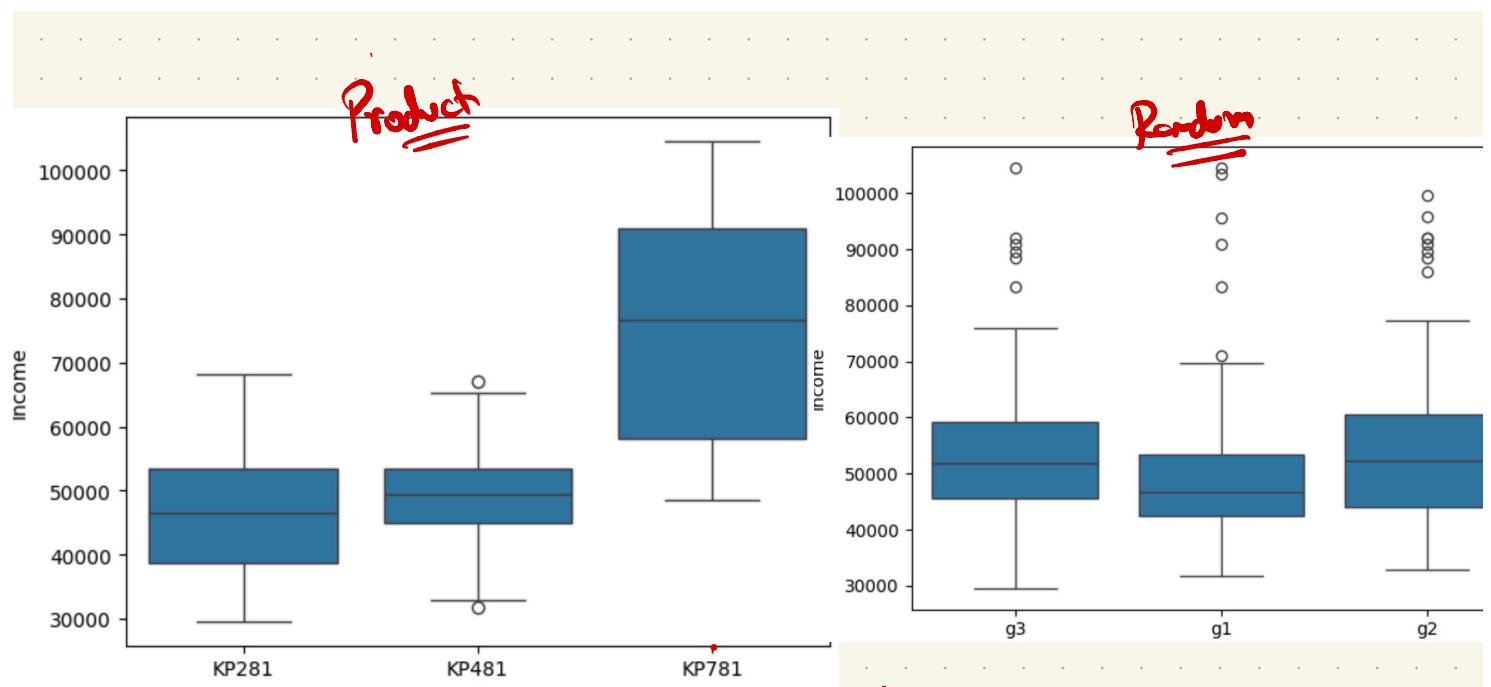
\Rightarrow Product vs Income \rightarrow ANOVA
(3 value)
(Num)

1) $H_0 \rightarrow$ They are independent / no impact

2) Distⁿ \rightarrow to calculate Statistic
+
Tail
+
P-value

\Rightarrow Product vs Income
(Cat)

\Rightarrow Group divides data in random 3 groups
and compare in Group



A Qty which changes when data differ from exp.
Act

⇒ Anova → Analysis of Variance.

F-ratio =

Variance b/w groups
Variance within groups

Variance b/w group → is higher when there is a lot of dependence on group.

Variance within group → should be higher for random group.

\Rightarrow	<u>14</u>	A	\rightarrow	<u>18-20</u>	$ $	A \rightarrow 12-30
	<u>23</u>	B	\rightarrow	<u>20-25</u>		B \rightarrow 12-30
	<u>29</u>	C	\rightarrow	<u>27-30</u>		C \rightarrow 12-30

Case 1

Variance within group \rightarrow Low

Case 2

Variance b/w group \rightarrow Some

Some

\Rightarrow Var within is lower for dependent case,
the group

Var b/w the is higher for dependent case.
group

$$f \text{ ratio} = \frac{\text{Var b/w the group}}{\text{Var within the group}}$$

f ratio is higher when data is dependent on group.

F ratio is above a threshold \rightarrow reject H_0

High ratio \downarrow \rightarrow Low p-value \rightarrow reject H_0

Assumptions of Anova \rightarrow

- ① Data is gaussian
- ② Data is independent across each record
- ③ Equal variance in diff groups

If assumptions of Anova fails \rightarrow Kruskal Wallis Test

\Rightarrow Kruskal \rightarrow

$H_0 \rightarrow$ Population median of all groups is same.

$H_a \rightarrow$ at least 1 median is different.

One-way ANOVA (Analysis of Variance) assumes the following:

Homogeneity of variances:

This assumption states that the variances within each group being compared are roughly equal.

Normality of data:

This assumption refers to the distribution of residuals or errors being approximately normal.

It does not necessarily require the raw data to be normally distributed.

Independence of observations:

This assumption means that the observations within each group are not related or dependent on each other.

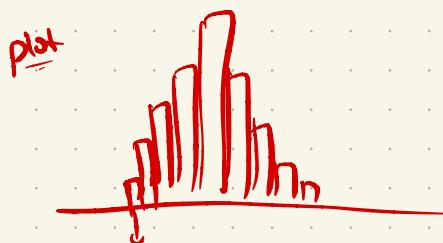
While equal sample sizes are desirable for one-way ANOVA, they are not always a strict assumption.

① Normality Test → If data is Gaussian or not.

A) QQ plot → Quantile Quantile plot.

⇒ plot a histogram of your data

+
Plot a histogram of ideal Normal Curve!



Compare each bar → normal bar

each Quantile → normal Quantile

5

0-5%
5-10%
10-20%
⋮

divide in 100 parts →

$$\Rightarrow (\mu, \sigma) \rightarrow f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2}$$

⇒ Data → get μ, σ → Plot a Normal curve for this (μ, σ)

Shapiro Test →

$H_0 \rightarrow$ Data is Gaussian

$H_a \rightarrow$ Data is not Gaussian.

⇒ To check normality → QQplot, Shapiro Test

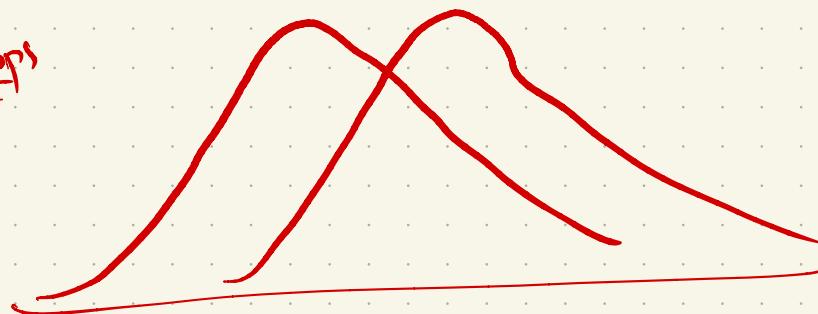
⇒ Equal Variance →

Levene Test → If there is equal variance in diff groups.

$H_0 \rightarrow$ Have same variance

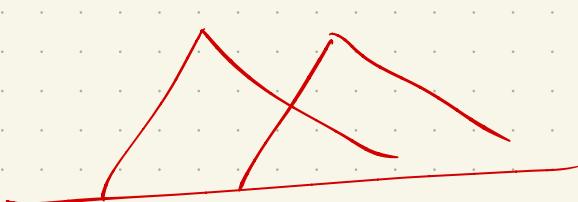
$H_a \rightarrow$ Have diff variance.

2groups



Normal vs observed

group1 vs group2
quantiles



\Rightarrow Anova \rightarrow More than 2 Cat and
You want to compare impact of Cat on a numerical
Column.

\rightarrow F-Stats \Rightarrow $\frac{\text{Var b/w group}}{\text{Var within group}}$

F \uparrow p-value \downarrow \rightarrow data is dependant on group

\rightarrow Assumption \rightarrow $\begin{cases} \rightarrow \text{Normal} \\ \rightarrow \text{equal variances} \\ \rightarrow \text{independent} \end{cases}$

\rightarrow If Assumption fail \rightarrow KU Test

\rightarrow Normality $\xrightarrow{\text{QQ plot}}$ Shapiro \rightarrow Shapiro Wilks

\rightarrow Equal var \rightarrow Levene Test

Doubts

① \rightarrow Z Test \rightarrow 2600 problem

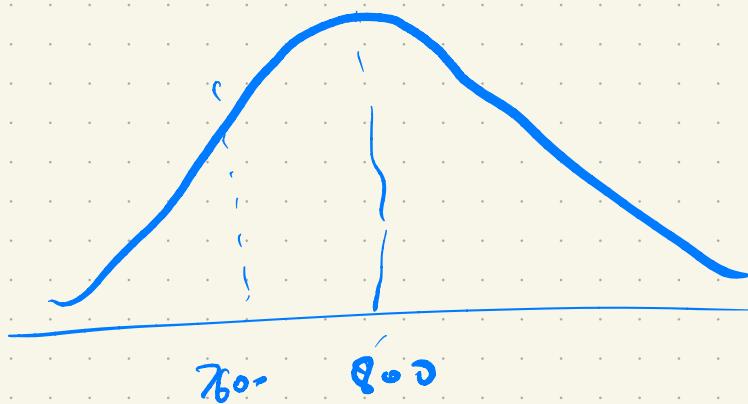
$$\Rightarrow H_0 = \text{Steps} = 8000$$

$$H_a = \text{Steps} \geq 8000$$

Right tailed test

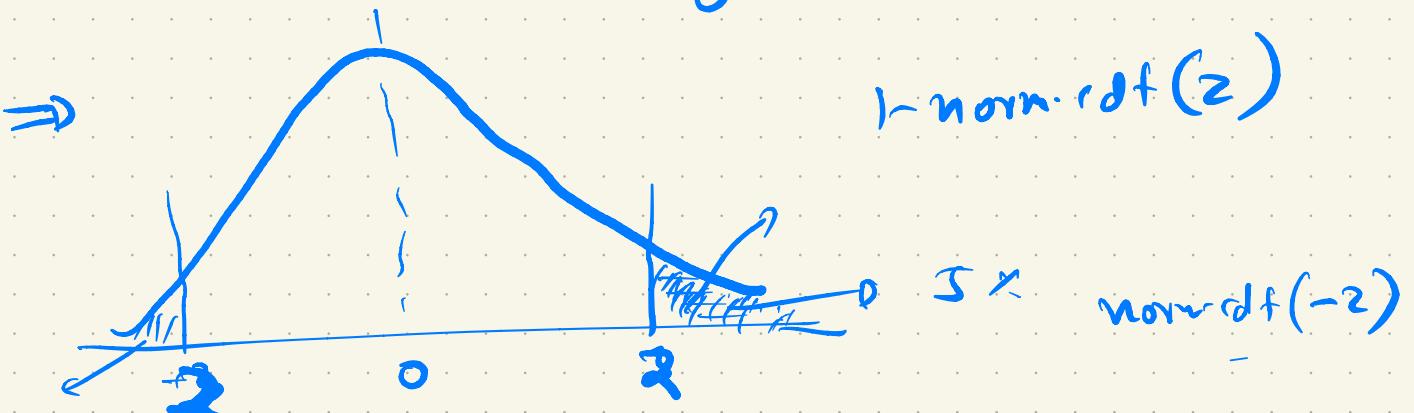
\Rightarrow Defining our dist \rightarrow Normal dist $\rightarrow (M, \sigma)$
 \rightarrow population
 \rightarrow SM \bar{y}
 CLT

$\Rightarrow H_0$ be true $\rightarrow N(8000, \frac{\sigma}{\sqrt{n}})$



p-value
 \hookrightarrow depending on
 right side
 of
 observed

\Rightarrow Define a dist with assuming H_0 is tr.-p.



$$p_{\text{value}} = 2 \times (1 - \text{norm.pdf}(\text{height_women}))$$

p-value \neq norm.cdf(1.71) - norm.cdf(-1.71)

$$p\text{-val.} = \left(1 - \text{norm.cdf}(1.71)\right) + \left(\text{norm.cdf}(-1.71)\right)$$

60% \rightarrow 1 car

28% \rightarrow 2 car

12% \rightarrow 2+ car

73 own \rightarrow 1 car \leftarrow

38 own \rightarrow 2 car \leftarrow

18 own \rightarrow 2+ car \leftarrow

129

60

28

12

73/129

38/129

18/129

60/129

28/129

12/129

73

38

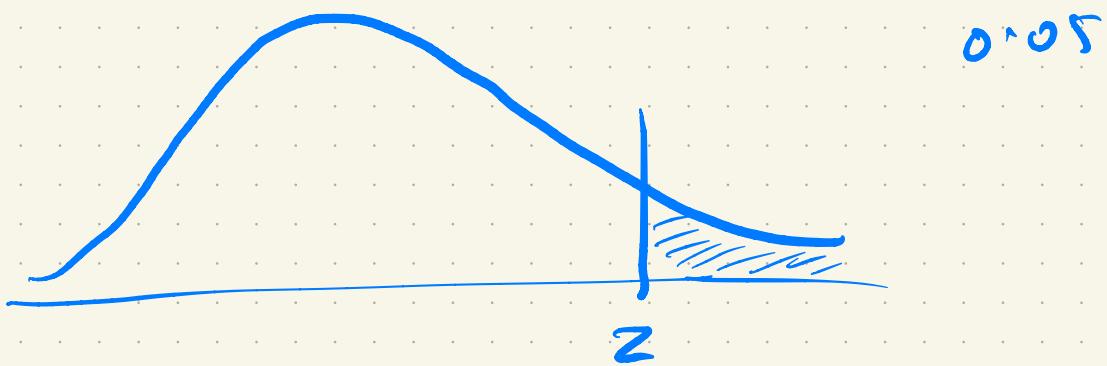
18

\Rightarrow Critical value \rightarrow
 p-value Comparing it against α .

Values associated with $p\text{-value} = \underline{\alpha}$.

\rightarrow SM data \rightarrow Z Score \rightarrow p-value

$p\text{-value} = \alpha \rightarrow$ Critical Z Score
 \downarrow
 Critical SM.



If Right side
 $\alpha = 0.05 \rightarrow 1 - \text{norm.cdf}(z) = 0.05$
 $z = \text{norm.ppt}(0.95)$

Critical $\overset{\circ}{z}$

$$z = \frac{x - \mu}{\sigma} = x \rightarrow \text{critical SM}$$

$$z = \frac{x - \mu_{\text{pop}}}{\sigma_{\text{pop}}/\sqrt{n}}$$

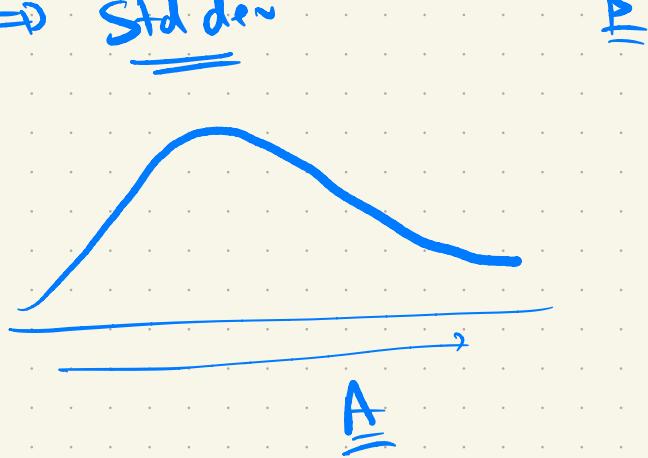
⇒ power of Test

$\beta \uparrow$ power ↓

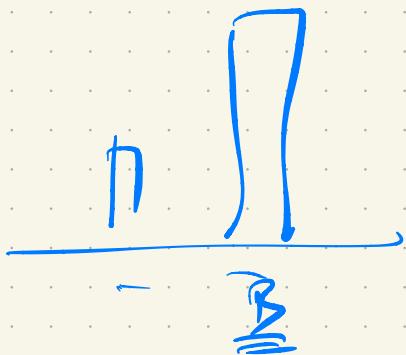
$$1 - \frac{\beta}{\alpha}$$

↳ false negative error.

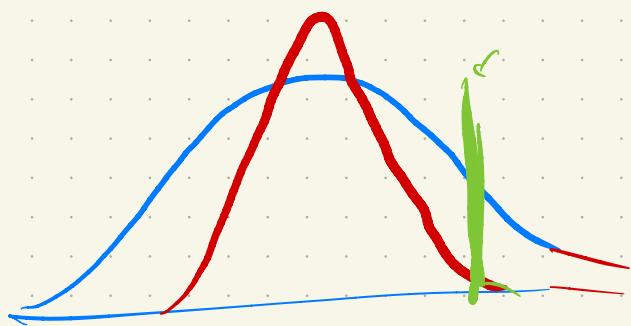
⇒ Std dev



B



P_{II}



STD ↓ $\beta \downarrow$

power ↑

If everyone
is innocent

$$\alpha = 0.02 \Rightarrow$$

98 people →

$$\beta = 0 \\ \alpha =$$

2 criminal →

98	0
2	0

If everyone
is minimizing

98 imm
2 crimi

0	98
0	2

$$\beta = 0.98$$

$$\beta = 1 - \alpha$$

\Rightarrow

92	3
2	3

CLT \rightarrow Std dev

of

$$SM_{dist} = \frac{\sigma_{pop}}{\sqrt{n}}$$

mean
of

$$SM_{dist} = \underline{M_{pop}}$$

$$\rightarrow A \rightarrow 30 \rightarrow M_A = 100$$

$$\rightarrow B \rightarrow 30 \rightarrow M_B = 120$$

$$\begin{array}{r} 180 \\ 170 \\ \hline 350 \end{array}$$

$$\Rightarrow A \rightarrow [1, 30, 40, \dots] \leftarrow 300 \rightarrow M_A = 95$$

$$B \rightarrow [5, 27, 28, \dots] \leftarrow 300 \rightarrow M_B = 48$$

$$A \rightarrow 0 = 1000$$

$$B \rightarrow 0 = 100$$

$$A \rightarrow 44-46 \leftarrow$$

$$B \rightarrow 47-49 \leftarrow$$

$$\frac{45}{48}$$

<https://colab.research.google.com/drive/11syQehH05QZwzsEkuXOx2af8WuXzPHkC?usp=sharing>

Collar link : <https://colab.research.google.com/drive/11syQehH05QZwzsEkuXOx2af8WuXzPHkC?usp=sharing>