



STATISTICS FOR THE DATA SCIENCE

Part - 4

- STUDENT T-DISTRIBUTION
- BAYES THEOREM
- CHI - SQUARE TEST
- F - DISTRIBUTION
- ANOVA AND ITS TYPES

#Value_freeContent



@Krishan kumar

 (B) Student t-distribution → Statistical analysis using Z-score
we need population standard deviation (σ).

* How do we perform an analysis when we don't know the population standard deviation.

Ans, Student t-distribution

in Z-score,

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

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



s = Sample std.

* Degree of freedom →

$$d.f. = n - 1$$

n ⇒ Sample size

we 3 people

X □ □

► We can find area under the curve by using t-value OR dof.
for t-distribution where population std (not given)

only
Sample std (given)

Problem

In the population the average IQ is 100. A team of researchers want to test a new medication to see if it has either a +ve or -ve effects on intelligence, or no effect at all. A sample of 30 participants who have taken the medicine has a mean of 140 with a std of 20. did the medication affect intelligence? C.I. \Rightarrow 95%.

Ans $\mu = 100$, $\sigma = 20$, $n = 30$, $\bar{x} = 140$, C.I. \Rightarrow 95%, $\alpha = 0.05$

① Null hypothesis (H_0) $\rightarrow \mu = 100$

② Alternate " (H_1) $\rightarrow \mu \neq 100$ { 2 Tail test }

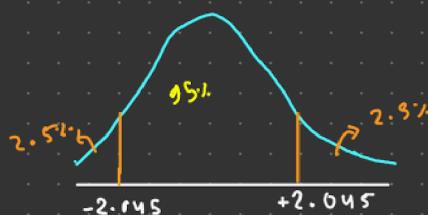
only apply when the value of a mean is either $\uparrow\downarrow$ OR $\uparrow\downarrow$

③ $\alpha = 0.05$

④ Degree of freedom :-

$$d.f \rightarrow 30 - 1 \rightarrow \underline{\underline{29}}$$

⑤ Decision boundary :-



by using t-table,

$$d.f = 29$$

$$\alpha = 0.05$$

we got \Rightarrow +2.045

if t-test is less than -2.045 and greater than 2.045 , we reject the null hypothesis.

⑥ Calculate t-test statistic,

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \Rightarrow \frac{140 - 100}{\frac{20}{\sqrt{30}}} \Rightarrow \frac{40}{3.65} \Rightarrow \underline{\underline{10.96}}$$

⑦ Conclusion

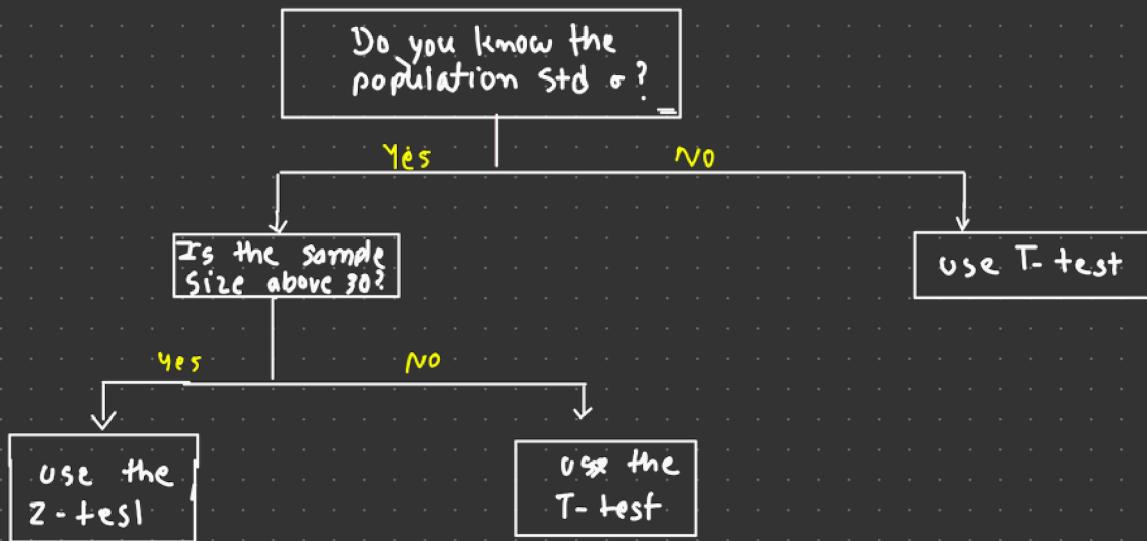
Decision Rule: if t is less than -2.045 and greater than 2.045 we reject the null hypothesis.

$$t = 10.96 > 2.045 \quad \{ \text{we Reject the null hypothesis} \}$$

Final → Medication has increased the intelligence.

★ Interview Questions

When to use t-test Vs Z-test



* Types 1 and Type 2 Errors

Reality: Null hypothesis is true OR H_0 is false

Decision: H_0 is True OR H_0 is False.

Conclusion

* outcome : 1

we reject the H_0 (Null hypothesis) when in reality it is false \rightarrow Good

* outcome : 2

we reject the null hypothesis when in reality it is True \rightarrow True \rightarrow Type 1 Error

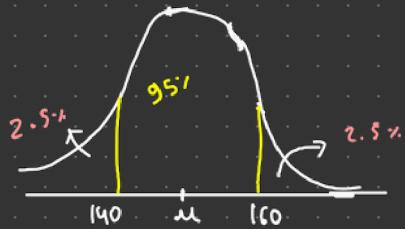
* outcome : 3

we retain (accept) the H_0 , when in reality it is false \rightarrow Type 2 Error

* outcome 4:

we retain the H_0 , when in reality it is True
 \downarrow
Good

Confidence Interval and Margin of Error:



$$\mu \rightarrow 160$$

$$C.I. \rightarrow 95\%$$

$$140 \leftrightarrow 160$$

Point estimate: A value of any statistics that estimate the value of an unknown population parameter is called "point estimate".

$$\bar{x} \rightarrow \mu$$

$$\bar{x} = 2.95 \quad , \mu = 3$$

Confidence interval: We construct a C.I. to help the estimate what the actual value of the unknown population mean is:-

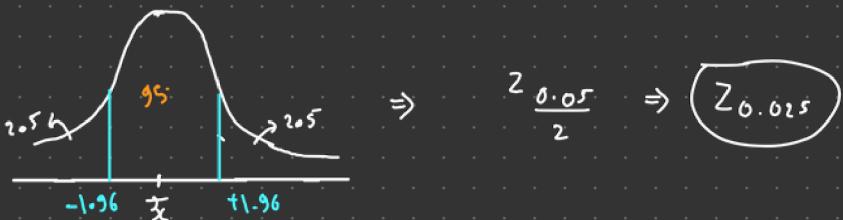
Point estimate \pm Margin Error

Z-test

$$\bar{x} \pm z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$$

$\{\alpha = \text{Significant value who decided}$
 How much C.I. is

$$\alpha = 0.05$$



Problem: On a verbal section of CAT Exam, the standard deviation is known to be 100. A sample of 25 test taken has a mean of 520. Construct a 95% C.I. about the mean.

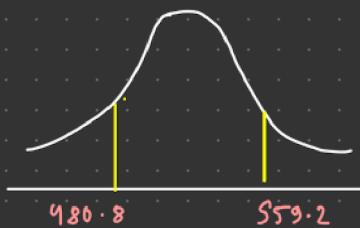
Ans $\rightarrow \bar{x} = 520, \sigma = 100, n = 25, \text{C.I. } \Rightarrow 0.95, \alpha = 0.05$



$$\bar{x} \pm z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}$$

$$\rightarrow \text{Lower C.I.} = 520 - (1.96) \times \frac{100}{\sqrt{25}} = 480.8$$

$$\rightarrow \text{Higher C.I.} = 520 + (1.96) \times \frac{100}{\sqrt{75}} = 559.2$$



\rightarrow Conclusion: I am 95% sure (confident) that the mean CAT exam score lies between 480.8 and 559.2

★ Bayes theorem



probability $\left\{ \begin{array}{l} \text{Independent Event} \\ \text{Dependent Event.} \end{array} \right.$

* Independent Event

ex. Rolling a dice

$$\{1, 2, 3, 4, 5, 6\}$$

$P(1) \rightarrow \frac{1}{6}$ ex, Tossing a coin

$$\begin{matrix} \vdots & P(H) \rightarrow \frac{1}{2} \\ P(6) \rightarrow \frac{1}{6} & P(T) \rightarrow \frac{1}{2} \end{matrix}$$

* Dependent Event

$$\begin{array}{ll} \text{ex. Blue balls} & \rightarrow P(R) \rightarrow \frac{2}{5} \\ \boxed{\begin{array}{c} O \ O \\ O \ O \end{array}} & \rightarrow P(B) \rightarrow \frac{3}{4} \end{array}$$

Here one event effect another events also.

ex,

$$\boxed{P(R \text{ and } B) = P(R) \times P(B|R)}$$

$$\frac{2}{5} \times \frac{3}{4} \rightarrow \frac{6}{20}$$

$$P(A \text{ and } B) = P(B \text{ and } A)$$

$$P(A) \times P(B/A) = P(B) \times P(A/B)$$

$$P(B/A) = \frac{P(B) \times P(A/B)}{P(A)}$$

Bayes theorem

$P(A/B) \rightarrow$ Conditional Probability

OR

$$P(A/B) = \frac{P(A) \times P(B/A)}{P(B)}$$

A, B → Events

$P(A/B)$ → Probability of A given B is already happened

$P(B/A)$ → Pro. of B given A is true.

$P(A), P(B)$ → The independent Pro. of A and B

* Uses of Bayes theorem in ML : (ex)

Dataset / Predict the price $\rightarrow (y=?)$

Size of the House	No. of Room	Location	Price
x_1	x_2	x_3	y

By Bayes theorem →

$$P\left(\frac{y}{x_1, x_2, x_3}\right) = \frac{P(y) \times P\left(\frac{x_1, x_2, x_3}{y}\right)}{P(x_1, x_2, x_3)}$$

Analysis by using bayes theorem → Bayesian Statistics

(() CHI SQUARE-TEST :-

It is non parametric test that

performed categorical {ordinal, nominal} data.

→ CHI SQUARE TEST for Goodness of fit test claims about population proportions [categorical variance].

ex ①

	Theory	Sample	Applied Goodness Fit
yellow bike	1/3	22	
orange bike	1/3	17	
Red bike	1/3	59	
	→ Observed categorical distribution		
	Theory categorical distribution		

ex ② Goodness of fit test :

In a student's class of 100 Student where 30 are right hand. Does this class fit the theory 12% of people are right handed.

→

	\underline{O}	$\underline{\underline{E}}$	$O \rightarrow$ Observed value
	\underline{O}	$\underline{\underline{E}}$	$E \rightarrow$ Expected value
Right handed	30	12	
left handed	$\frac{700}{100}$	$\frac{88}{100}$	$\frac{12}{100} \times 100$
	Theory categorical distribution		
	Sample info		

Problem: In 2010 Census of city, the weight of the individuals in a small city were found to be the following.

$\leq 50\text{kg}$	50 - 75	> 75
20%	30%	50%

In 2020, age of $n = 500$ individuals were sampled. Below are the results.

<50	50-75	>75
140.0	160	200

using $\alpha = 0.05$ would you conclude the population difference of weights has changed in the last 10 years.

Ans:

Expected Value	<50	50-75	>75	(2010)
	20%	30%	50%	

$$n = 500$$

Observed Value	<50	50-75	>75	(2020)
	140	160	200	

<50	50-75	>75
0.2×500	0.3×500	0.5×500
$\Rightarrow 100$	$\Rightarrow 150$	$\Rightarrow 250$

$$20\% \text{ of } 500 \rightarrow \frac{20 \times 500}{100}$$

$$30\% \text{ of } 500 \rightarrow 0.3 \times 500$$

$$50\% \text{ of } 500 \rightarrow 0.5 \times 500$$

- Null hypothesis: H_0 : The data meet the expectation
- Alternate Hypothesis: H_1 : The data does not meet the expectation.

$$② \lambda = 0.05, C.I. \rightarrow 95\%$$

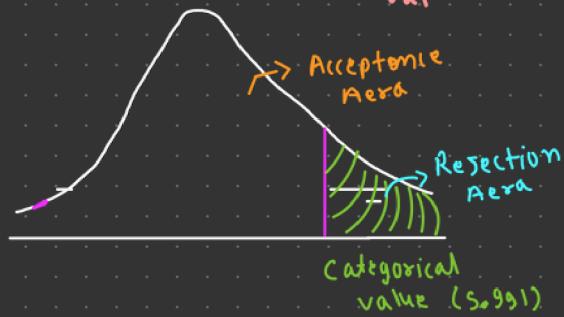
③ Degree of freedom

$$d.f = k - 1$$

$k \rightarrow$ No. of category

$$3 - 1 \rightarrow 2$$

④ Decision Boundary : We used CHI-Square table for the categories value.



denote:

$\chi^2 \rightarrow$ CHI Square (Sumbal)

If χ^2 is greater than 5.991, Reject H_0

else:

We fail to Reject the H_0

⑤ Calculate CHI-SQUARE Test-Statistic :

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

$O \rightarrow$ observed, $E \rightarrow$ expected value

E value

<50	50-75	>75
140	160	200

$$\Rightarrow \frac{140}{100}^2 + \frac{10}{150}^2 + \frac{(-50)}{250}^2$$

$$\Rightarrow \frac{1600}{100} + \frac{100}{150} + \frac{2500}{250}$$

$$\Rightarrow 16 + 0.66 + 10 \Rightarrow 26.66$$

O value

<50	50-75	>75
100	150	250

Final: If χ^2 greater than 5.99 Reject the H_0 , else Fail to reject the H_0 .

So, $\chi^2 = 26.66 > 5.99$ { we Reject the Null Hypothesis }

Cons: The weights of 2020 population are different than those expected on the 2010 population.

F-distribution:-

In probability theory and statistic the F-distribution or F-ratio also known as snedecor's F-dis. or the Fisher - snedecor distribution is continuous probability distribution that arises frequently as the null dis. of a test test statistic, most notably in the analysis of variance (ANOVA) and other F-test.

* F-test used to compare variance of mean between two group. F-test also called 'variance ratio test'

parameters:

$d_1, d_2 > 0$ (degree of freedom)

Support: $x \in (0, +\infty)$

F-distribution with d_1 and d_2 degree of freedom is the dis. of

$$\chi^2 = \frac{s_1^2/d_1}{s_2^2/d_2} \quad \begin{cases} s_1 \rightarrow \text{Independent random variable} \\ s_2 \rightarrow \dots \\ d_1 \rightarrow \text{Degree of freedom} \\ d_2 \rightarrow \text{Degree of freedom} \end{cases} \quad \left\{ \begin{array}{l} \text{CHI Square} \\ \text{distribution} \end{array} \right\}$$

F-test:

Problem: The following data show the no. of bulbs produced daily for same day by 2 worker A and B.

A	B
40	39
30	38
38	41
41	33
38	32
35	39
40	
34	

Can be consider based on the data
worker B is more stable and
efficient

$$d = 0.05$$

Ans → ① Null Hypothesis: $H_0: \sigma_1^2 = \sigma_2^2$
 $H_1: \sigma_1^2 \neq \sigma_2^2$

② calculation of variance:

$$\sigma^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

A		
x _i	\bar{x}	$(x_i - \bar{x})^2$
40	37	9
30	37	49
38	37	1
41	37	16
38	37	1
35	37	4
$\bar{x} = 37$		
$\sum (x_i - \bar{x})^2$		
$\Rightarrow 80$		

$$S_1^2 \Rightarrow \frac{(x_1 - \bar{x})^2}{n-1}$$

$$S_1^2 \Rightarrow \frac{80}{6-1} \Rightarrow \frac{80}{5} \Rightarrow \underline{\underline{16}}$$

B		
x _i	\bar{x}	$(x_i - \bar{x})^2$
39	37	9
38	37	1
41	37	16
33	37	16
32	37	25
39	37	4
40	37	9
34	37	9
$\bar{x} = 37$		
$\sum (x_i - \bar{x})^2$		
$\Rightarrow 84$		

$$S_2^2 = \frac{(x_2 - \bar{x})^2}{n-1}$$

$$S_2^2 = \frac{84}{8-1} \Rightarrow \frac{84}{7} \Rightarrow 12$$

→ Calculation of variance Ratio (F-test):

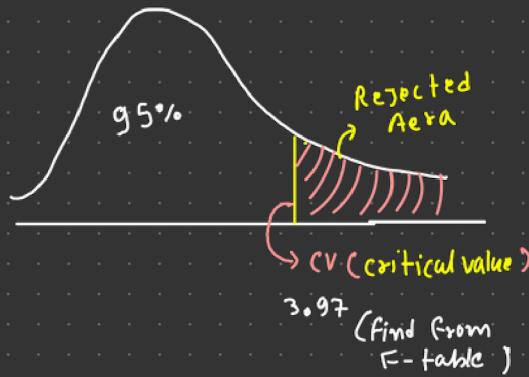
$$F = \frac{s_1^2}{s_2^2} \Rightarrow \frac{13}{12} \Rightarrow 1.33$$

③ Decision Rule:

$$df_1 \rightarrow 6-1 \rightarrow 5$$

$$df_2 \rightarrow 8-1 \rightarrow 7$$

$$\alpha \rightarrow 0.05$$



If F-test greater than 3.97,
Reject the Null Hypothesis

$1.33 < 3.97$, {so, we Reject the H_0 }

Final conclusion :-

Worker B is not efficient when compare to worker A.



(ANOVA) Analysis of variance:

Anova is a statistical method to compared the mean of 2 or more group.

ANOVA:

① factors (variable)

② levels

Ex ① Factors = Medicine

level = 5mg 10 mg 15mg [Dosage]



⇒ Assumptions in ANOVA:

- ① Normality of Sampling Distribution of mean The dis. of Sample mean is normally distributed.
- ② Absence of outliers. Outlier score need to be Remove from dataset.
- ③ Homogeneity of variance Each one of the population has same variance $\sigma_1^2 = \sigma_2^2 = \sigma_3^2$ Population variance in different levels of each independent variable are equal.
- ④ Sample are independent and Random.

⇒ Types of ANOVA:

- ① One way ANOVA
- ② Repeated Measured ANOVA
- ③ Factorial ANOVA

- ① One way ANOVA: One factor with at least 2 levels, these levels are independent.

② Repeated Measured anova:

One factor with at least 2 levels, levels are dependent.

③ Factorial Anova: Two or more factor (each of which with after 2 level). levels can be either dependent or independent.

⇒ Hypothesis testing in ANOVA:

Null Hypothesis: $H_0: \mu_1 = \mu_2 = \mu_3 = \dots = \mu_n$

$H_1:$ At least one of the mean is not equal

$$X \boxed{\mu_1 = \mu_2 = \mu_3 = \dots = \mu_n}$$

Test statistic:

$$F = \frac{\text{Variation between sample}}{\text{Variation within sample}}$$

$X_1 \quad X_2 \quad X_3$

1 6 5

2 7 6

4 3 3

5 2 2

3 1 4

$$\sum x_1 = 15$$

$$\bar{x} = \frac{15}{3} \Rightarrow 5$$

$$\sum x_2 = 14$$

$$\bar{x}_2 = \frac{14}{5}$$

$$\sum x_3 = 20$$

$$\bar{x}_3 = \frac{20}{5} \Rightarrow 4$$

$$H_0: x_1 = x_2 = x_3$$

$$H_1: \text{at least one sample mean is not equal}$$

* One way anova: One factor at least 2 levels, levels are independent

Problem: Doctor want to test a new medication which reduce headache. They split the participant into 3 condition [15mg, 30mg, 45mg]. later on the doctor ask the patient to rate the headache between [1 → 10] are there any difference between 3 condition using Alpha (α) $\rightarrow 0.05$?

Ans:

① Define null Hypothesis:

$$H_0: \mu_{15} = \mu_{30} = \mu_{45}$$

H_1 : not all μ 's are equal

15mg	30mg	45mg
9	7	4
8	6	3
7	6	2
8	7	3
8	8	4
9	7	3
8	6	2

② State significant value:

$$\alpha = 0.05$$

$$C.I. \rightarrow 0.95$$

③ Calculate degree of freedom:

$$N = 21, \alpha = 3, n = 7$$

($a \rightarrow$ Number of Sample)

$$df_{\text{between}}: a-1 \rightarrow 3-1 \Rightarrow 2$$

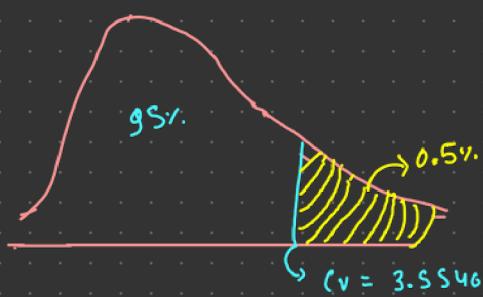
(2, 18)

$$df_{\text{within}}: N-a \rightarrow 21-3 \Rightarrow 18$$

↳ super useful to finding critical value in table

$$df_{\text{total}}: N-1 \rightarrow 21-1 \Rightarrow 20$$

F - Table



Critical value $\rightarrow 3.5546$

④ State decision rule:

If F is greater than 3.5546 , we reject the H_0 .

⑤ Calculate test Statistic:

SS = Sum of Square

SS_{between} SS_{within} SS_{total}

15 mg	30 mg	45 mg
9	7	4
8	6	3
7	6	2
8	7	3
8	8	4
9	7	3
8	6	2

$$\textcircled{1} \quad SS_{\text{between}} = \frac{\sum (\sum a_i)^2}{n} - \frac{-2}{N}$$

$$15 \text{ mg} \rightarrow 9+8+7+8+8+9+8$$

$$30 \text{ mg} \rightarrow 7+6+6+7+8+7+6$$

$$45 \text{ mg} \rightarrow 4+3+2+3+4+3+2$$

$$SS_{\text{between}} = \frac{57^2 + 47^2 + 21^2}{7} - \frac{[57 + 47 + 21]}{21}$$

$$\Rightarrow 98.17$$

$$\textcircled{2} \quad SS_{\text{within}} = \sum y^2 - \frac{\sum (a_i)^2}{n}$$

$\left\{ \begin{array}{l} y^2 \rightarrow \text{Given table value } \vec{A} \\ \vec{A} \text{ element का square करना} \end{array} \right.$

$$\Rightarrow \sum y^2 - \left[\frac{s_7^2 + u_7^2 + z_7^2}{7} \right]$$

$$\sum y^2 \Rightarrow 9^2 + 8^2 + 7^2 + 8^2 + \dots + 3^2 + 2^2 \Rightarrow 853$$

$$\Rightarrow 853 - \left[\frac{s_7^2 + u_7^2 + z_7^2}{7} \right]$$

$$\Rightarrow 10.29$$

$$\textcircled{3} \quad SS_{\text{Total}} : \quad \sum y^2 - \frac{T^2}{N}$$

$$\Rightarrow 853 - \frac{125^2}{21}$$

$$\Rightarrow 108.95$$

	SS	df	MS	F
Between	98.67	2	49.34	
Within	10.29	18	0.54	
Total	108.95	20	4.9.88	$\{ MS = \frac{SS}{df} \}$

$\Rightarrow \frac{MS_{\text{between}}}{MS_{\text{within}}} \text{ is equal to } \Rightarrow F = \frac{\text{variation between sample}}{\text{variation within sample}}$

$$\Rightarrow \frac{49.34}{0.54} \Rightarrow 86.56 \Rightarrow F = 86.56$$

Final:

If F is greater than 3.5546, we Reject the H₀

86.56 > 3.5546, so we Reject H₀.

