

for Joff... 29 Nov in table # 18 21045

Date _____
Page _____

if test is less than -2.045 or greater
than +2.045, reject the null hypothesis.

(b) calculate test statistics

$$f = \frac{\bar{x} - u}{S_{\bar{x}}} = \frac{140 - 100}{20/30} = \frac{40}{3.33} = 12.095$$

$$f = 12.095$$

*

Since $t = 12.095 > 2.045$ Reject the null hypothesis

Conclusion: Medication used has affected
the intelligence
Medication has increased the intelligence

ca



When to use T-test vs Z-test

Do you know the population std σ ?

Yes

No



Use t-test

{ Is the sample size above 30? }

Yes

No

{ Use t-test }

Use z-test



Type I and Type II errors

Reality null hypothesis true or null hypothesis is false

Decision null hypothesis is true or null hypothesis is false

Outcome 1 we reject the null hypothesis when in reality it is false \rightarrow Good

Outcome 2 we reject the null hypothesis when in reality it is true \rightarrow Type I error

Outcomes we retain the null hypothesis when in reality it is false \rightarrow Type II error

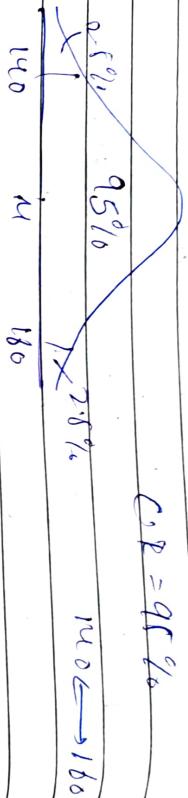
Outcomes we retain the null hypothesis when in reality it is true \rightarrow Good.

Confidence Interval and Margin of Error

Date _____

Page _____

$$\mu = 160$$



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

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

$$x = 20.45$$

$$x = 3$$

Confidence Interval

We construct a confidence interval to help estimate what the actual value of the unknown population mean is.

point estimate + margin of error

$$2.5\%$$

$$\bar{x} + z_{0.975} \frac{\sigma}{\sqrt{n}}$$

$$\alpha = 0.025$$

How we can use it - test under

$$20.05 \rightarrow 20.07$$

also if we want to conduct hypothesis testing

$$2.5\%$$

$$95\%$$

$$2.5\%$$

$$1.96$$

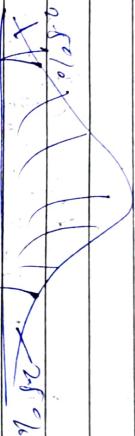
$$1.96$$



Q) On the regional section of CAT Exam, the standard deviation is known to be 100. A sample of 25 test takers has a mean of 520. construct 95% C.I about the mean?

$$\bar{X} = 520 \quad \sigma = 100 \quad n = 25 \quad t \cdot P = 0.95$$

$$\alpha = 0.05$$



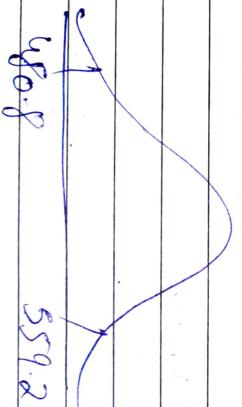
$$1 - 0.025 = 0.975$$

$$\bar{x} + Z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$$

findable

$$\text{Lower CI} = 520 - (1.96) \times \frac{100}{\sqrt{25}} = 480.8$$

$$\text{Higher CI} = 520 + (1.96) \times \frac{100}{\sqrt{25}} = 559.2$$



$$480.8$$

$$559.2$$

I am 95% confident that the mean CAT score lies between 480.8 and 559.2.

Bayesian statistics (Bayes' theorem)

Bayesian statistics is an approach to data analysis and parameter estimation based on Bayes' theorem.

Bayes' theorem

Probability \rightarrow Independent Events \rightarrow Dependent Events

① Independent events

e.g.: Rolling a die
 $\{1, 2, 3, 4, 5, 6\}$



$$\rightarrow P(Y) = 3$$

$$P(A_1) = \frac{1}{6} \quad P(A_2) = \frac{1}{6}$$

$$P(Rain|Y) = P(R) \times P(Y|R)$$

Tossing a coin

$$P(H) = 0.5 \quad P(T) = 0.5$$

conditional probability

$$\frac{2}{5} \times \frac{3}{2} = \frac{6}{25}$$

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

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

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

$$P(A)$$

$$\Pr(A|B) = \frac{\Pr(A) \times \Pr(B|A)}{\Pr(B)}$$

A, B = events

\rightarrow already occurred

$\Pr(A|B)$ = probability of A given B is true

$\Pr(B|A)$ = " " B " " A is true

$\Pr(A)$, $\Pr(B)$ = independent probability of A and B

Dataset

\rightarrow independent

\rightarrow obj/depends

Size of House	No. of Rooms	Location	Price
x_1	x_2	x_3	y

$$\Pr(y|x_1, x_2, x_3) = \frac{\Pr(y) \times \Pr(x_1, x_2, x_3|y)}{\Pr(x_1, x_2, x_3)}$$

Baye's theorem

CHI-SQUARE TEST

The Chi-square test for goodness of fit test claims about population proportions categories. It is a non parametric test that is determined on categorical features and nominal data.

There is a population of birds who likes different colors birds

	Actual Sample	Expected Sample
Yellow Bird	43	22
Red Bird	13	17

	Actual Sample	Expected Sample
Orange Bird	13	59

→ Observed categorical distribution

Theory categorical distribution

Goodness of fit data

In a science class of 15 students, there are 9 girls. This class fit the theory the handled. Only 10 people are left handled.

of people are left handled

10/15 handled = 0.67

9/15 unhandled = 0.60

25/15 = 1.67

A

light handed 6/15

total 15 Now we need to add

it to the total



Chi Square for Goodness of fit

In 2010 census of the 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 a 2020, weight of $n = 500$ individuals were sampled. Below are the results

≤ 50	$50 - 75$	> 75
140	180	200

Using $\alpha = 0.05$, would you conclude the population different of weights has changed in the last 100 years?

2010 Expected	$\leq 50\text{ kg}$	$50 - 75$	> 75
	200	300	500

2020 $n=500$ observed	≤ 50	$50 - 75$	> 75
	140	180	200

Expected	≤ 50	$50 - 75$	> 75
	0.2×500 = 100	0.3×500 = 150	0.5×500 = 250

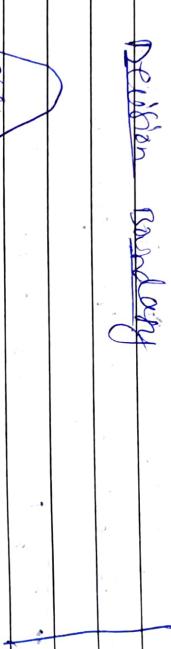
null hypothesis: H_0 : the data meets the expectation
 research hypothesis: H_1 : the data does not meet the expectation

$$\textcircled{2} \quad \alpha = 0.05 \quad CI = 95\%$$

\textcircled{3} Degree of freedom

$$df = k - 1 = 3 - 1 = 2$$

\textcircled{4} Decision Boundary



CV
critical value

α
 χ^2

If χ^2 is greater than 5.99, reject H_0 .
 Else we fail to reject the null hypothesis.

\textcircled{5} calculate Chi-Square Test statistics

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

2020	observed	50	140	160	200
n=500	expected	50	75	75	250



Date _____

Page _____



$$\Rightarrow \frac{(190 - 100)^2}{100} + \frac{(160 - 150)^2}{150} + \frac{(200 - 250)^2}{250}$$

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

$$\Rightarrow 16 + 0.66 + 10$$

$$\Rightarrow 26.66$$

$$\chi^2 = 26.66$$

If χ^2 is greater than 5.99 , reject Ho

else
we fail to reject the null hypothesis

$$26.66 > 5.99 \text{ Reject Ho}$$

The weight of 2020 population are different
than those expected in the 2010 population.

2010
2020

Specifically used for hypothesis testing
to analyze or compare the ratios of
two different groups

Date _____
Page _____

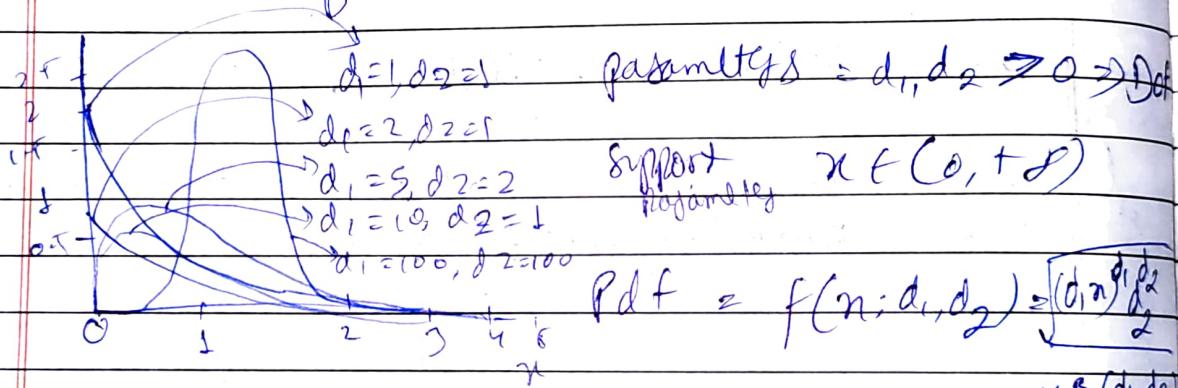
F distribution

In probability theory and statistics, the F-distribution or F-ratio, also known as Snedecor's F dist. or the Fisher-Snedecor dist. (after Ronald Fisher and George W. Snedecor) is a continuous probability distribution that arises frequently as the null dist. of a test statistic, most notably in the analysis of variance (ANOVA) and other F-tests.

used in
ANOVA

Fisher-Snedecor
prob. density func

using F-test we compare
two groups of more
groups so two D.F.



$$\pi B\left(\frac{d_1}{2}, \frac{d_2}{2}\right)$$

Samples
of two
groups
of data

Beta function

$$B(m,n) = \frac{(m-1)! (n-1)!}{(m+n-2)!}$$

$$\Rightarrow m/n / (m/n)$$

F distribution with d_1 and d_2 degrees of freedom
distribution of a random variable

$$\lambda = \frac{S_1^2}{S_2^2}, \quad S_{1,2} \Rightarrow \text{Independent variables}$$

S_1^2, S_2^2

with Chi-square dist.

$d_1, d_2 \Rightarrow$ respective degrees of freedom

f - test Evidence Ratio f-test

f - test of Variance Ratio test

① The following data shows the no. of webs produced daily at some days by 2 workers A and B

A B

40 39

30 38

38 41

41 33

38 32

35 39

40 40

$\lambda = 2.05$

Surplus
in no. of
goods
for day

34

Providence

① Null Hypothesis $H_0: \sigma_1^2 = \sigma_2^2$

$H_1: \sigma_1^2 \neq \sigma_2^2$

dist. to
varied
varied
varied
varied

(Q) calculation of variance &

$$S_1^2 = \frac{1}{n-1} \sum_{i=1}^{n-1} (x_i - \bar{x})^2$$

workers A

x_1	\bar{x}_1	$(x_1 - \bar{x}_1)^2$	x_2	\bar{x}_2	$(x_2 - \bar{x}_2)^2$
40	37	9	39	37	4
30	37	49	38	37	1
38	37	1	41	37	16
41	37	36	33	37	16
38	37	1	32	37	25
35	37	4	39	37	4
$\bar{x}_1 = 37$		$\sum (x_1 - \bar{x}_1)^2 = 80$	$\bar{x}_2 = 37$		
			$\bar{x}_2 = 37$		$\sum (x_2 - \bar{x}_2)^2 = 84$

$$S_1^2 = \frac{80}{n-1} = \frac{80}{5} = 16$$

$$S_2^2 = 84$$

$$\Rightarrow S_1^2 = 16$$

Calculation of variance Ratio Test (F test)

$$F = \frac{S_1^2}{S_2^2} = \frac{16}{12} = \underline{\underline{1.33}}$$

For H_0 also follows right skewed dist.

Date _____
Page _____

Decision Rule

$$\begin{aligned} H_A &= \delta - 1 = 5 \\ d t_2 &= 8^{-1} = 2 \\ \alpha &= 0.05 \end{aligned}$$

CV

(3.9415)

Decision
 F_{test} is greater than 3.9415, reject the null hypothesis

1.3 Since 3.9415, we fail to reject the null hypothesis

reject H_0

Conclusion:

Worker B is not more stable or efficient when compared to worker A

Analysis of variance (ANOVA)

Def'nt ANOVA is a statistical method used to compare the means of 2 or more groups

- ① ANOVA → independent factors (variable)
- ② levels

e.g. factor = medicines (factor)

levels = sing long range (degrees)

mode of payment → (factor)

MAP principle (MDS) MDTI factors



Assumptions of ANOVA

1. C.I.T

Date _____
Signature _____

To
Group

① Normality of sampling distribution of mean
The distribution of sample mean is normally
distribution

Gaussian dist.

② Absence of outliers
outlying score need to be removed from the dataset. outliers need to be removed?

③ Homogeneity of variance
population variance in different levels of each independent variable are equal
 $\sigma_1^2 = \sigma_2^2 = \sigma_3^2$ If one variable do not
influence the other variables

④ Samples are independent and random
randomly selected.

6)

Types of ANOVA (3 types)

Date _____
Page _____

(3) Factor

- ① One way ANOVA if one factor with at least 2 levels, other levels are independent.

Eg Doctor wants to test a new medication to relieve headache.

They split the participants in 3 conditions (long, very short, medium).

Doctor ask the participants to rate the headache from 1 to 10.

factor

Medication \rightarrow factor.

long strong very long \rightarrow levels are independent.

5	4	20
3	4	6
-	-	-
-	-	-

- ② Repeated Measures ANOVA if one factor at least 2 levels, levels are dependent.

Running \rightarrow factor.

Ones \rightarrow Day 1 Day 2 Day 3
are dependent. $\begin{array}{ccc} 5 & 4 & 9 \\ 7 & - & - \end{array}$



(B) Factorial analysis of Two or more factors (each of which with atleast 2 levels). Levels can be independent and dependent.

Running \rightarrow fastly

to	Day 1	Day 2	Day 3	Result
levels	1	2	3	
sex	♂	♀	♀	
gender (male)	9	8	7	
gender (female)	2	4	6	
factor	7	8	3	

Hypothesis testing in ANOVA (partitioning of variance)

in the ANOVA

is based on how many samples we have

null hypothesis $H_0 : \mu_1 = \mu_2 = \dots = \mu_k$

alternative hypothesis H_1 : At least one of the sample means is not equal

$$H_0, H_1 \text{ or } H_0 + H_1$$

to start cannot write like this

if all means are not equal
one mean is not equal to other
difference is not equal to some other
diff. mean

use F-test

F = variance between samples
variance within samples

$$H_0 : \bar{X}_1 = \bar{X}_2 = \bar{X}_3$$

H_1 : At least one sample mean not equal.

variance	1	2	3	4
within samples	4	3	3	2
1	6	6	6	6
2	5	5	5	5
3	3	3	3	3
4	2	2	2	2
5	1	1	1	1
6	4	4	4	4

$$\bar{X}_1 = 3 \quad \bar{X}_2 = 4 \quad \bar{X}_3 = 4$$

variance is spread

One way ANOVA

one factor with at least 2 levels, levels are independent.

- Q) Doctors want to test a new medication with reduced headache. They splits the participant into 3 condition (15mg, 30mg, 45mg). Later on the doctor ask the patient to rate the headache b/w 1-10. Are there any differences b/w the 3 conditions using alpha = 0.052

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

- Q) Define Null and Alternative hypothesis?

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

$$H_a: \text{not all } \mu \text{ are equal}$$

- ② Significance $\alpha = 0.05$ CI $= 0.95$

- ③ Calculate Degree of freedom

$$N = 21$$



Total sample size

$$g = 3$$



No. of categories

$$n = 7$$



Sample size

df of blue samples

$$\begin{aligned} df_{\text{between}} &= a - 1 = 3 - 1 = 2 \\ df_{\text{within}} &= N - a = 21 - 3 = 18 \end{aligned} \quad \left. \begin{array}{l} df_1, df_2 \\ (2, 18) \end{array} \right\}$$

①

F table

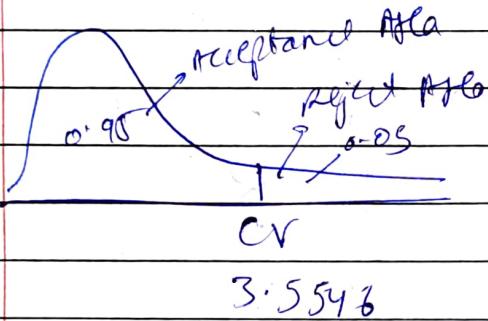
$$df_{\text{total}} = N - 1 = 20$$

$$\alpha = 0.05$$

②

To find critical value

(c) Decision Boundary ↗



for $df_1 = 2$

$$df_2 = 18$$

The value in F table
is 3.5546

②

Decision Rule

If F is greater than 3.5546, reject the null hypothesis

(d) calculate F test statistic

F = ratiannie blue sample
ratiannie within sample

SS	df	MS	F
----	----	----	---

Between 98.37

within 10.29

Total 108.96



$$\textcircled{1} \quad S.S \text{ between} = \frac{\sum (\bar{x}_i)^2}{n} - \frac{\bar{x}^2}{N}$$

$$15mg = 4+8+7+8+8+9+8 = 57$$

$$30mg = 7+8+8+7+8+7+6 = 47$$

$$47mg = 4+3+2+3+4+3+2 = 21$$

$$\Rightarrow \frac{57^2 + 47^2 + 21^2}{7} - \frac{[57^2 + 47^2 + 21^2]}{21}$$

$$= 98.07$$

$$\textcircled{2} \quad S.S_{within} = \sum y_i^2 = \frac{\sum (\bar{x}_i)^2}{n}$$

$$\sum y_i^2 = 9^2 + 8^2 + 7^2 + 8^2 + 8^2 + \dots$$

$$= 853$$

$$\Rightarrow 853 - \frac{[57^2 + 47^2 + 21^2]}{7}$$

$$\Rightarrow [10.29]$$

$S.S \rightarrow$ sum of squares

$\partial R \rightarrow$

$M.S \rightarrow$ mean square,

$$MS = \frac{SS}{df}$$

Date _____
Page _____

SS df MJ F

between 98.67 2 49.34

within 10.29 18 0.54

Total 108.96 20

F test = $\frac{MS_{\text{between}}}{MS_{\text{within}}}$

F = Variance between samples
Variance within samples

$$F = \frac{49.34}{0.54} = 89.56$$

If F is greater than 3.5548, reject the H₀.

89.56 > 3.5548

Reject the H₀.

∴ there are differences.