

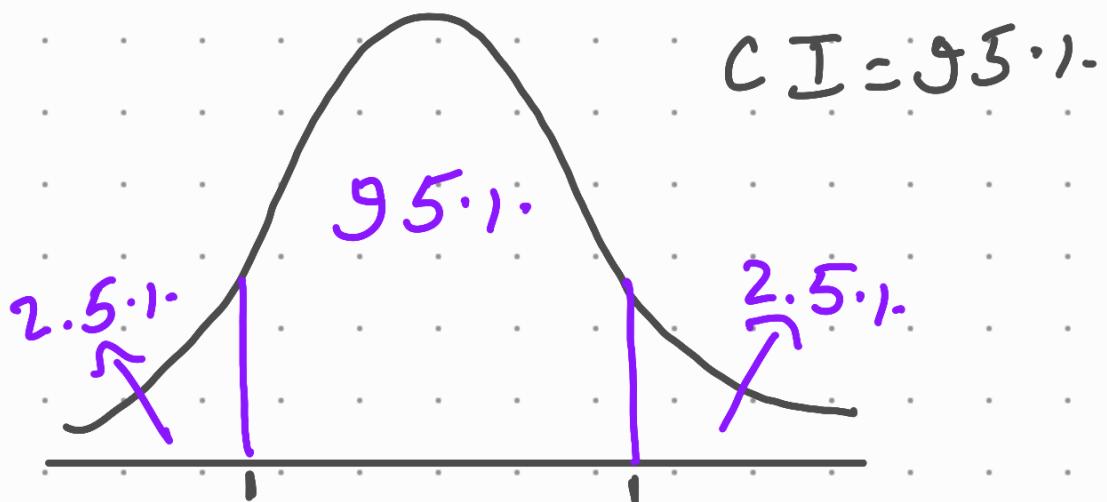
e.g.

Dataset			Output Feature
Input Features			
Size of House	No. of Rooms	Location	Price
$x_1$	$x_2$	$x_3$	$y$

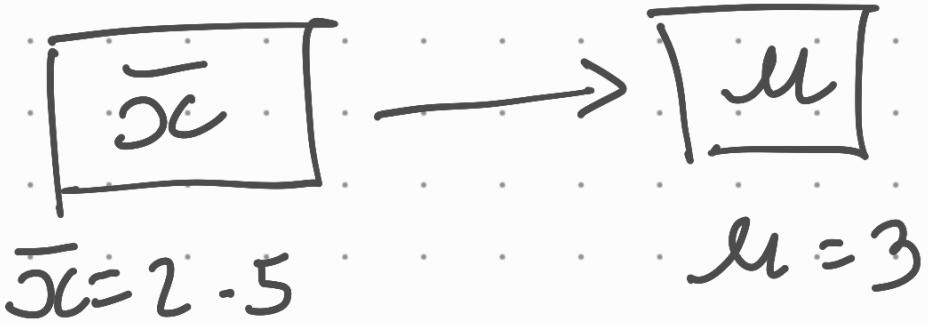
$$\begin{aligned}
 P_b(y|x_1, x_2, x_3) &= \frac{P_b(y) * P_b(x_1|x_2|x_3|y)}{P_b(x_1|x_2|x_3)} \\
 &= \frac{P_b(y) * P_b(x_1|x_2|x_3|y)}{P_b(x_1|x_2|x_3)}
 \end{aligned}$$

↳ Bayes' Theorem

Confidence Interval and Margin of Error



## Point Estimate



## Confidence Interval

Point Estimate  $\pm$  Margin of Error

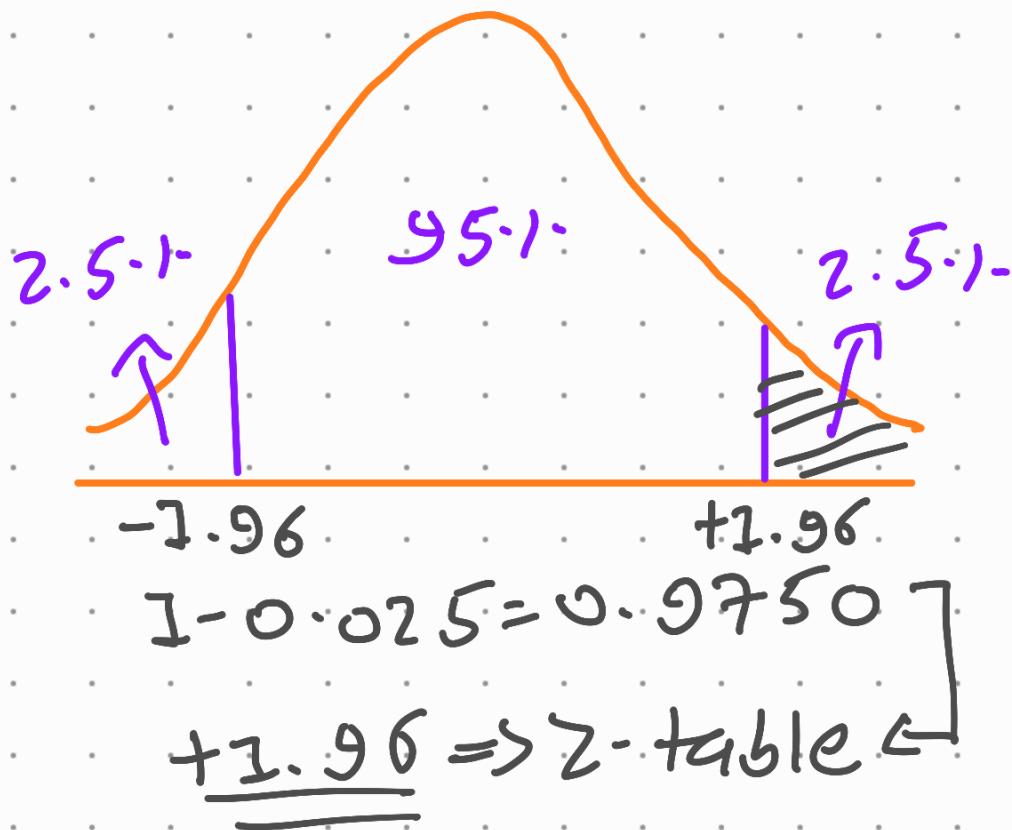
$$Z\text{-Test} \Rightarrow \bar{x} \pm Z_{\alpha/2} * \frac{\sigma}{\sqrt{n}}$$

$$T\text{-Test} \Rightarrow \bar{x} \pm t_{\alpha/2} * \frac{s}{\sqrt{n}}$$

e.g. On the verbal section of exam, the standard deviation is known to be 100. A sample of 30 test takers has a mean of 520. Construct 95-1. C.I. about the mean.

$$\alpha = 0.05$$

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

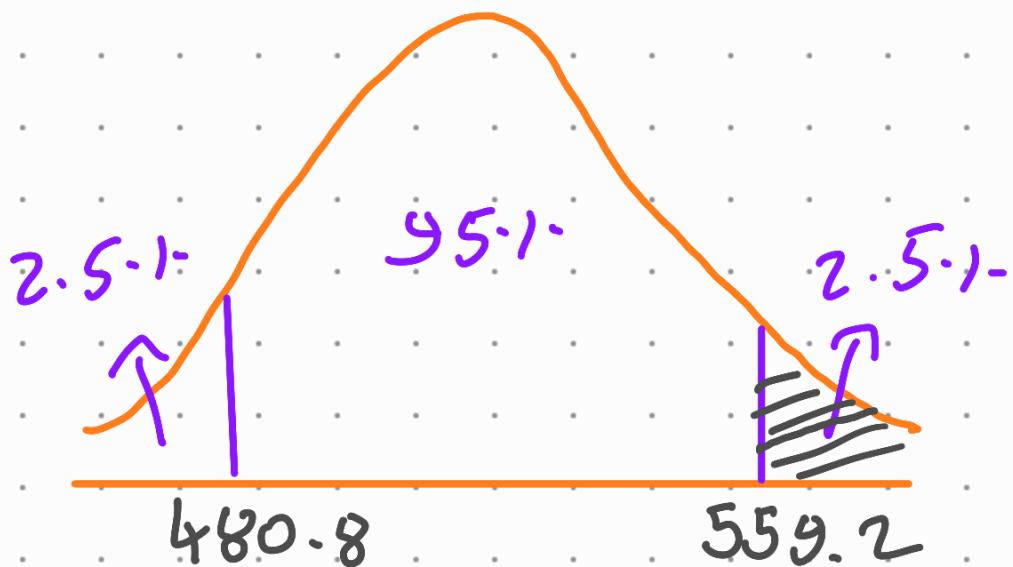


Lower CI

$$= 520 - (1.96) * \frac{100}{\sqrt{25}}$$
$$= 480.8$$

Higher CI

$$= 520 + (1.96) * \frac{100}{\sqrt{25}}$$
$$= 559.2$$



Conclusion: I am 95.1% confident about the mean score is between 480.8 and 559.2.

## Chi Square Test

→ The Chi Square Test for Goodness of Fit

- This test claims about the population proportions.
- It is a non-parametric test that is performed on categorical Data
- Categorical Data [Ordinal and Nominal] data.

e.g. (1) There is a population of mule who likes different color bikes.

Theory      sample

Yellow Bike      1/3      22

Red Bike      1/3      17

Orange Bike      1/3      59

Theory (categorical)  
Distribution

Observed (categorical)  
Distribution

## \* Goodness of Fit Test

e.g. ① In a Science class of 75 students, 11 are left handed. Does this class fit the theory that 12% of people are left handed.

→

	O	E
Left Handed	11	9
Right Handed	64	66
	75	75

② In 2010 census of the city, the weight of the individuals in a small city were found to be the following.

<50kg	50 - 75	>75kg
20%	30%	50%

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

$\leq 50\text{kg}$	$50 - 75$	$> 75\text{kg}$
Observed		
140	160	200

Using  $\alpha = 0.05$ , would you conclude the population difference of weights has changed in the last 10 years?



2020

Expected

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

2020  
Observed  
 $n=500$

$\leq 50\text{kg}$	$50 - 75$	$> 75\text{kg}$
Observed		
140	160	200



	$< 50$	$50 - 75$	$> 75$
Expected	$0.2 \times 500$ $= 100$	$0.3 \times 500$ $= 150$	$0.5 \times 500$ $= 250$

(a) Null Hypothesis  $H_0$  = The data meets the expectations  
(b) Alternate Hypothesis  $H_1$  = The data does not meet the expectations

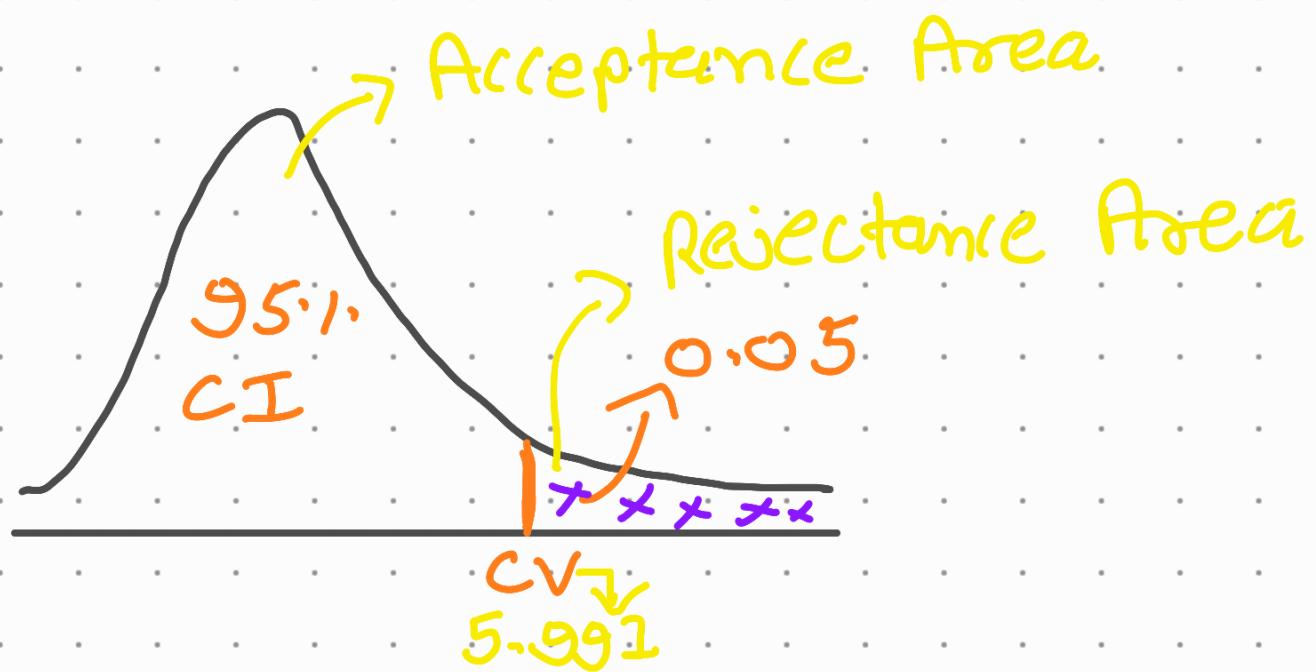
(c)  $\alpha = 0.05$  CI = 95%

(d) degree of freedom

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

$k$  = number of categories

## (e) Decision Boundary



- If  $\chi^2$  is greater than 5.991, Reject the Null Hypothesis ( $H_0$ )
- Else we fail to reject the Null Hypothesis ( $H_1$ )

## (f) Calculate $\chi^2$ -Test stats

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

$$\begin{aligned}
 &= (140 - 150)^2 + (160 - 150)^2 + (200 - 250)^2 \\
 &= \frac{1600}{100} + \frac{100}{250} + \frac{2500}{250} \\
 &= 16 + 0.66 + 12
 \end{aligned}$$

$$\boxed{\chi^2 = 26.66} > 5.99, \text{ Reject}$$

the Null Hypothesis  $H_0$ .

### Conclusion

→ The weights of 2020 populations  
are different than those  
expected in 2010.

# Analysis of Variance (ANOVA)

→ ANOVA is a statistical method used to compare the means of 2 or more groups.

→ ANOVA

- ① Factors (variable)
  - ② Levels

e.g. ① medicine (Factors)

[Dosage] 5 mg 10 mg 15 mg  
Levels

## ② Mode of Payment

[Levels] Cray, PhonPe, IMP5, NEFT

## → Assumptions In Anova

## ① Normality of Sampling Distribution of Mean

→ The distribution of sample mean is normally distributed.

## ② Absence of Outliers

→ Outlying score need to be removed from the dataset.

## ③ Homogeneity of Variance

→ Population variance in different levels of each independent variable are equal

$$[\sigma_1^2 = \sigma_2^2 = \sigma_3^2]$$

## ④ Samples are independent and random.

### Types of ANOVA (3 Types)

① One Way Anova: One factor with at least 2 levels. These levels are independent.

e.g. ① Doctor wants to test new medication to decrease headache.

→ They split the participants in 3 conditions [10mg, 20mg, 30mg]  
→ Doctor ask the participants to rate the headache between [1-10]

Medication → Factor

10 mg 20 mg 30 mg → [Levels]

5	7	2	Rating to Headaches
3	4	6	
-	-	-	
-	-	-	
-	-	-	

② Repeated Measures ANOVA: One factor with at least 2 levels. These levels are dependent.

e.g.

Running → Factor

Day-1 Day-2 Day-3 → Levels

8	5	4
7	4	9
-	-	-
-	-	-

③ Factorial ANOVA: Two or more factors each of which with at least 2 levels. These levels can be independent or dependent.

e.g.

Running → Factor

		Day-1	Day-2	Day-3 → Levels
Gender	Male	8 9 2 7	5 4 4 8	4 3 6 3
	Female			

Hypothesis Testing In ANOVA

(Partitioning of Variance In The ANOVA)

(a) Null Hypothesis  $H_0: \mu_1 = \mu_2 = \mu_3 = \dots = \mu_k$

(b) Alternate Hypothesis  $H_1$ : At least one of the sample mean is not equal

$$\times \boxed{\mu_1 \neq \mu_2 \neq \mu_3 \dots \mu_k}$$

## Test Statistics

$$F = \frac{\text{Variance between Sample}}{\text{Variance within Sample}}$$

Variance Between Sample

$x_1$	$x_2$	$x_3$
1	6	5
2	7	6
4	3	3
5	2	2
3	1	4
$\bar{x}_1 = 3$	$\bar{x}_2 = 29/5$	$\bar{x}_3 = 4$

$$H_0 = \bar{x}_1 = \bar{x}_2 = \bar{x}_3$$

$H_1$  = At least one sample mean is not equal

# One Way ANOVA

→ One factor with at least 2 levels are independent

e.g.

- ① Doctors want to test a new medication which reduces headache. They split the participant into 3 condition [15mg, 30mg, 45mg]. After on the doctor ask the patient to rate the headache between [1-10]. Are there any differences between the 3 conditions using  $\alpha = 0.05$ ?



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

#### (a) Null and Alternative Hypothesis

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

$H_1$ : Not all means are equal

(b)  $\alpha = 0.05$  CI = 0.95

#### (c) Degree of Freedom

$$N = 21 \quad n = 7 \quad a = 3$$

$$dof_{\text{between}} = a - 1 = 3 - 1 = 2$$

$$dof_{\text{within}} = N - a = 21 - 3 = 18$$

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

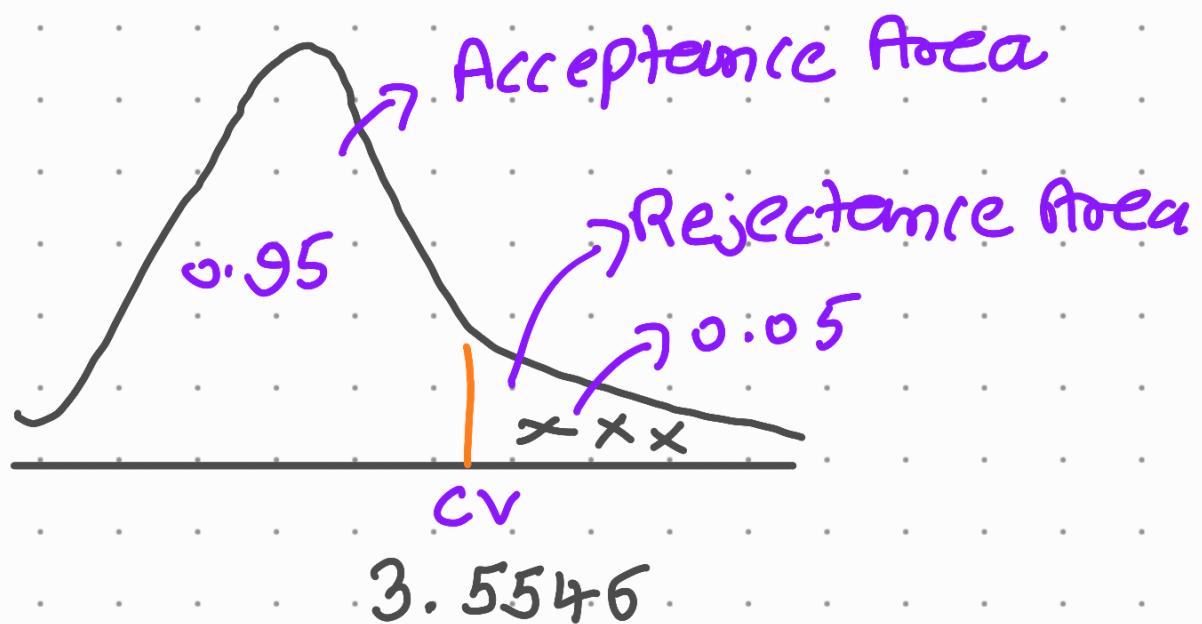
Range of value will be between

(2, 18)  $\Rightarrow$  F-table

$$\begin{matrix} \downarrow & \downarrow \\ df_1 & df_2 \end{matrix} \quad \alpha = 0.05$$

$\Downarrow$   
Critical Value

## (d) Decision Boundary



## Decision Rule

→ If the F value is greater than 3.5546, reject the null hypothesis  $H_0$ .

## (e) F-Test Stats

$$F = \frac{\text{Variance between samples}}{\text{Variance within sample}}$$

	SS	df	MS (SS/df)	
Between	98.67	2	49.34	F
Within	10.29	18	0.54	
Total	108.96	20	49.88	

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

$$IS \bar{x} = 9+8+7+8+9+8 = 57$$

$$30 \bar{x} = 7+6+6+7+8+7+6 = 47$$

$$45 \bar{x} = 4+3+2+3+4+3+2 = 27$$

$$= \frac{(57)^2 + (47)^2 + (27)^2 - [57^2 + 47^2 + 27^2]}{21}$$

$$= 98.67$$

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

$$\begin{aligned}\sum y^2 &= 9^2 + 8^2 + 7^2 + 8^2 + 8^2 + \dots \\ &= 853\end{aligned}$$

$$= 853 - \frac{[57^2 + 47^2 + 27^2]}{7}$$

$$= 10.29$$

$$F\text{-Test} = \frac{MS_{\text{between}}}{MS_{\text{within}}}$$

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

→ If F is greater than 3.5546,  
Reject the H<sub>0</sub>

→ 86.56 > 3.5546, Reject the null Hypothesis.

