

collect, organize, analyze.

## Statistics for Data Science:

### Descriptive

- organize/summarize data

- Central Tendency

- Dispersion

- Distr. of data

- Histogram, pdf, pmf.

### Inferential

- using data to form conclusion

- ZD Test

- t -Test

- Chi square

- ANOVA

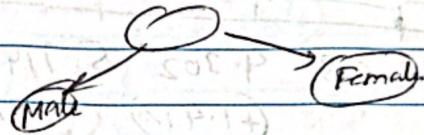
(N) Population: group we're interested in studying.

(n) Sample: subset / representative subset of sample.

## Sampling Techniques:

① Simple Random Sampling: every member of N has equal chance of being selected (n).

(Layering) ② Stratified sampling: non overlapping groups.  
- divide into



③ Systematic Sampling: Mall → outside  
(Survey) ↳ every  $n^{\text{th}}$  person.

④ Convenience Sampling (Voluntary Response Sampling)

- select a preferable grp of ppl for sampling.

## Variables

### Quantitative

→ take any value w/in a range

### Qualitative/Categorical

Gender → M  
F.

Discrete

distinct, countable values

Continuous

[infinite values]

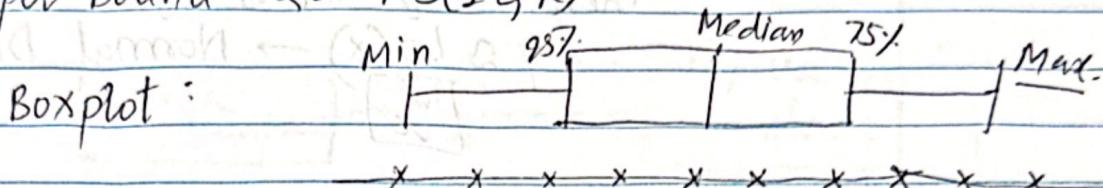
## Outliers:

### 5 Number summary:

- ① Minimum ② First quartile (25%) ③ Median } Boxplot  
 ④ Third quartile (75%) ⑤ Maximum

$$\text{Min. Low Bound} = Q_1 - 1.5(\text{IQR}) \quad \text{IQR} = Q_3 - Q_1.$$

$$\text{Max. Upper Bound} = Q_3 + 1.5(\text{IQR})$$



\* Z-score: data points more than 3 s.d. away from mean. (Normally distributed data).

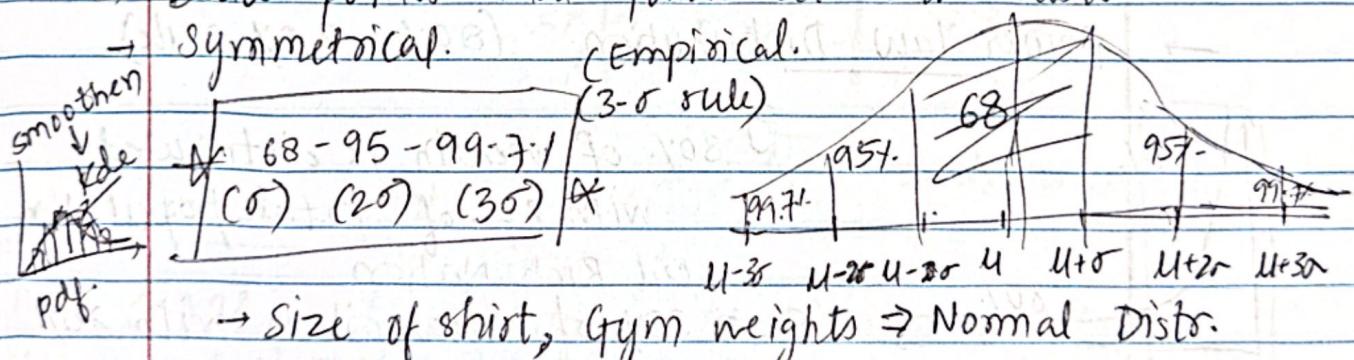
\* DBScan: points that aren't a part of Dense Clusters.

## Types of Distribution:

- Normal Distribution and Empirical Formula.

→ Data points that follow bell curve dist.

→ Symmetrical. (Empirical)



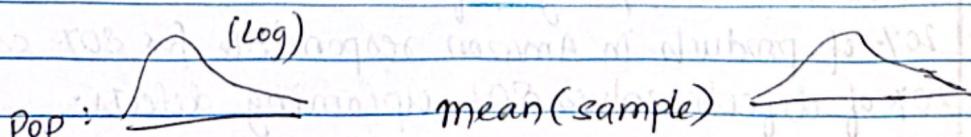
→ Size of shirt, Gym weights ⇒ Normal Distr.

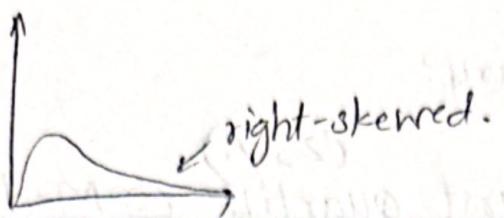
- Central Limit Theorem:

$$n \geq 30.$$

• If you take multiple samples and from a population and plot their means, the resulting distribution would be normally distributed.

• The original population need not be normally distributed





## → Log Normal Distribution:

if  $X \approx$  Log Normal Distribution.

then  $\sqrt{\text{natural log}}$

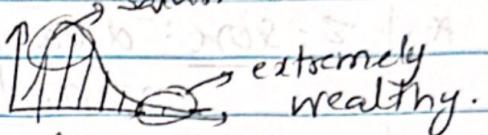
$y \approx \ln(x) \rightarrow$  Normal Distribution.

$\boxed{\log_e}$

$$X \approx \exp(y)$$

e.g. • Wealth Distribution

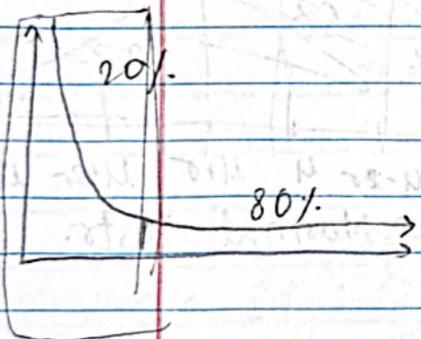
• Comments by Length on a channel.  
(Less long paragraphs).



Data Transformation: Log Dist  $\rightarrow$  Normal.

$\rightarrow$  better to train ML models.

## → Power Law Distribution: (80% - 20% rule)



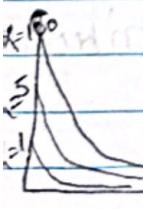
e.g. 80% of wealth is distributed with 20% of total population

Oil Rich Nation

- 80% of total oil is with 20% of nations.

Transform to Normal Distr. using Box Plot Transformation.

## → Pareto Distribution: follows Power Law Distribution.



$\alpha$  increases, height of Pareto increases.

e.g. - 20% of products in Amazon responsible for 80% sales.  
- 20% of defects solves 80% upcoming defects.

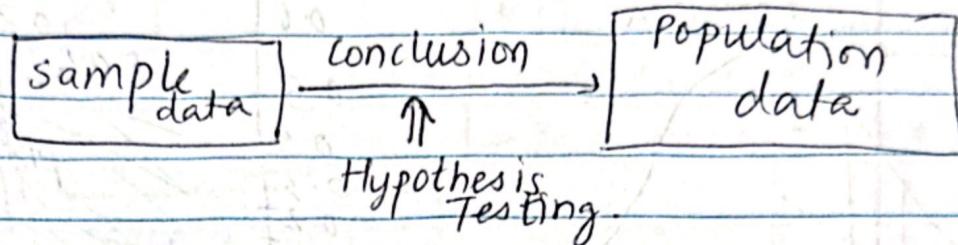
## Pareto:

A Relationship between Log Normal and Pareto?

Inferential Stats:

### \* Hypothesis Testing:

- test an assumption regarding a population parameter.



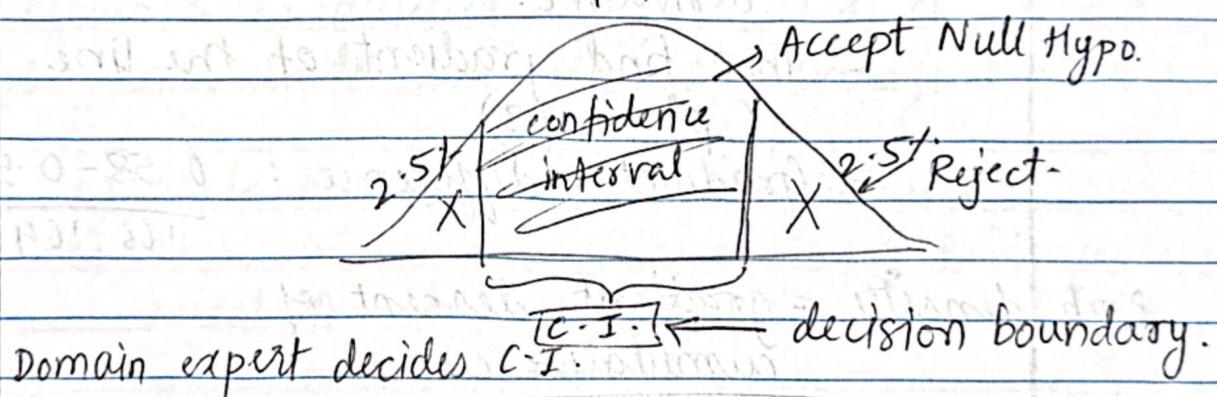
④ Null Hypothesis ( $H_0$ )

② Alternative Hypothesis ( $H_1$ )

③ Experiment:

95% C.I.

[1 - C.I.]  
5% significance  
value.



Random var:

map output of random process to a number.

dice, coin toss etc.

probability distr. Function

(contn) Probability Density Function

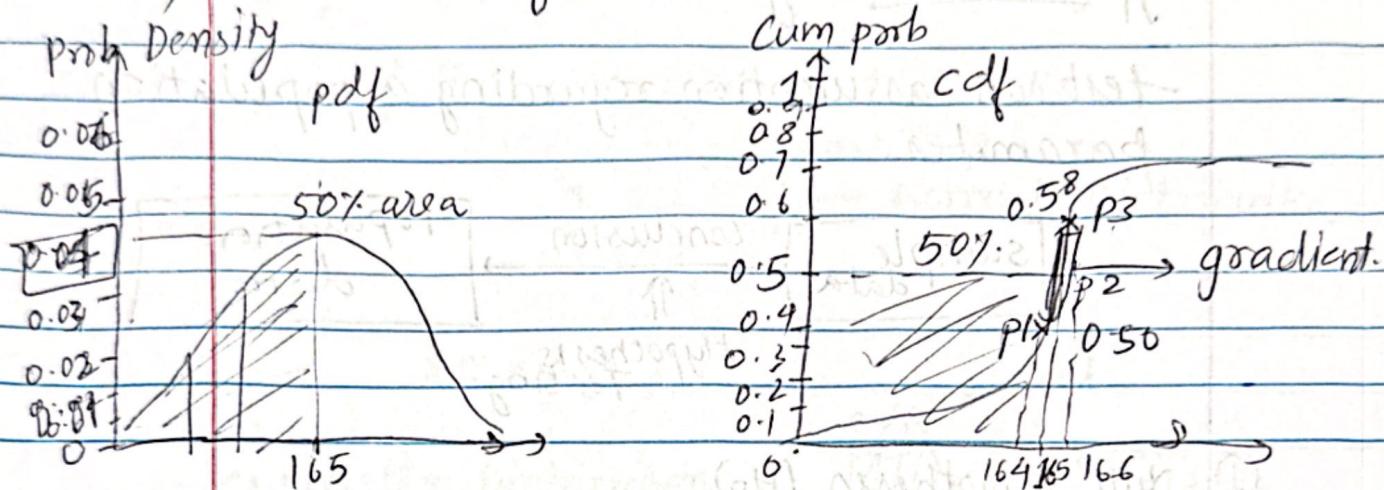
(discr.) Prob. Mass Function

Continuous Random var. ← check its distribution using pdf.

Discrete random var → use pmf to check its distribution.

## Area under curve (pdf)

i) Distribution of continuous random variable.



- How to find actual probability at 50% from cdf?

- Take 2 points' (above and below) derivative.

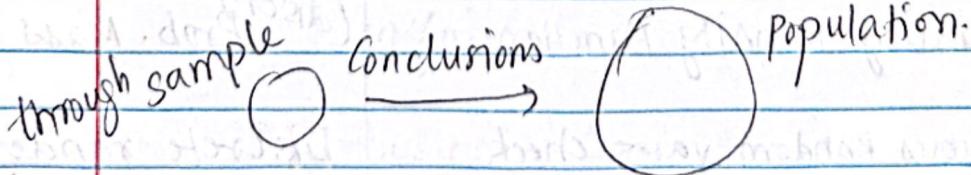
- Then find gradient of the line (p<sup>1</sup> to p<sup>3</sup>)

- Gradient : difference :  $\frac{0.58 - 0.50}{166 - 164} = 0.04$

prob density = gradient descent of cumulative curve.

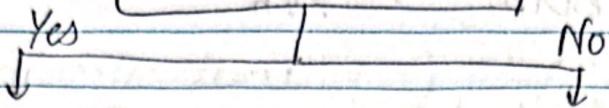
→ Hypothesis Testing and Statistical Analysis

- |                     |          |         |                                 |
|---------------------|----------|---------|---------------------------------|
| ① Z-test            | ② t-test | Average | ③ Chi-Square → Categorical Data |
| ④ ANOVA → Variance. |          |         |                                 |



→ When to use t-test vs z-test?

Do you know population std dev  $\sigma$ ?



$H_0 \neq$ : Two Tailed Test  
 $H_1 >$ : Right-Tailed

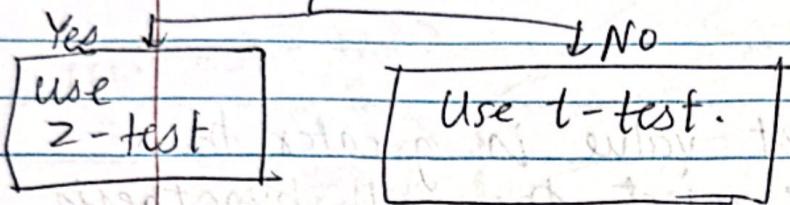
t-test

Is the sample size  $> 30$ ?

Use t-test  
(sample size = 30)

$H_1 <$ : Left Tail

\* We need the area under curve (range) of the C.I.  
→ This can be found using z-table:



Q: The average heights of all residents in a city is 168cm with a population std  $\sigma = 3.9$ . A doctor believes the mean to be different. He measured the height of 36 individuals and found the average  $\approx 169.5$  cm.

- State Null and Alternate Hypothesis
- At a 95% C.I., is there enough evidence to reject the null hypothesis?

Solution:

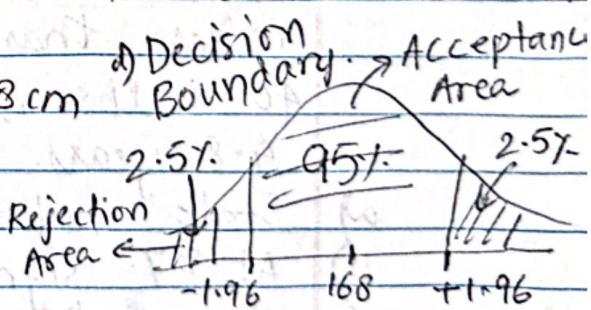
$$\mu = 168 \text{ cm} \quad \sigma = 3.9 \quad n = 36 \quad \bar{x} = 169.5 \text{ cm}$$

- Null Hypothesis ( $H_0$ ):  $\mu = 168 \text{ cm}$
- Alternate Hypothesis ( $H_1$ ):  $\mu \neq 168 \text{ cm}$

c) C.I. = 0.95 or 95%

$$\alpha = 1 - \text{C.I.} = 1 - 0.95 = 0.05$$

→ Decision Boundary:



Through decision boundary,

if z-test value falls between  $-1.96$  to  $+1.96$ , then we

fail to reject Null hypothesis.

Through z-table, we find  $[1.96]$

$$z\text{-score} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

n=1

standard error  $\rightarrow \sigma/\sqrt{n}$

d) statistical analysis

$$z\text{-test} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{169.5 - 168}{3.9/\sqrt{36}} = 2.31.$$

e) Conclusion:

- As z-test value is greater than  $+1.96$ , we reject the null hypothesis and accept the ACCEPT the alternate hypothesis.

~~2.31 > 1.96~~  
- The doctors' claim is right, the population mean is 169.5.

Question 2: Practice sum:

A factory manufactures bulbs w average warranty of 5 years with standard deviation of 0.50. A worker believes that the bulb will manufacture in less than 5 years. He tests a sample of 40 bulbs and finds the average time to be 4.8 years.

a) State null and alternate hypothesis.

b) At 2% significance, is there enough evidence to support the idea that the warranty should be revised?

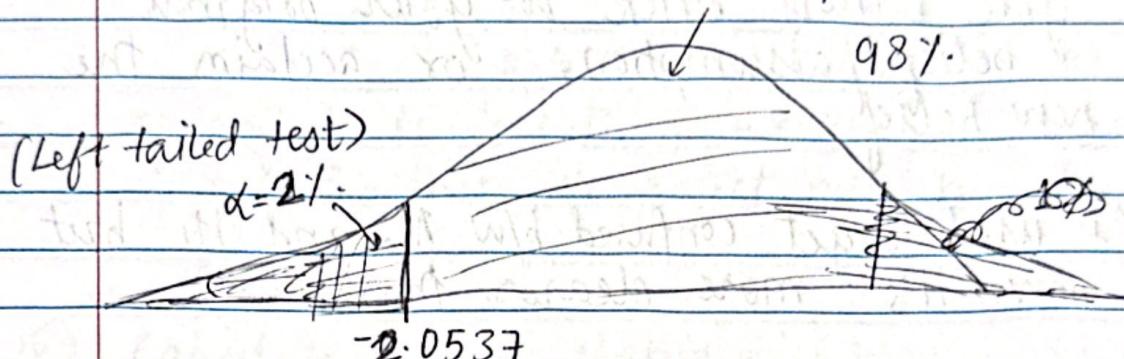
$\rightarrow \sigma = 0.50, n = 40, \mu = 5, \bar{x} = 4.8$

a)  $H_0: \mu = 5$

b)  $H_a: \mu < 5$

c) Significance level ( $\alpha$ ) = 2% or 0.02.  
C.I. =  $1 - 0.02 = 0.98$  or 98%

d) Decision Boundary: Accept Null Hypothesis



\* z-value: Standardized

critical value:  $(1 - 0.02 = 0.9800)$

Critical value = -2.0537

(if z-score value is below -2.0537, you reject the null hypothesis)

e) Calculate z-score:

$$z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{4.8 - 5}{0.5 / \sqrt{40}}$$

f) Conclusion:

$$\therefore -2.53 > -2.0537$$

(z-value) (critical value)

We reject the null hypothesis.

$\therefore$  The average height life of bulbs is 4.8 years!

How do you decide null / alternate hypothesis?

- The claim or belief that you already have → NULL Hypothesis.
- The claim / belief that challenges your original belief is the alternate hypothesis.

With decision boundary, critical values, you either stick to your original belief / assumption, or acclaim the new belief.

→ (I used to get confused b/w  $H_0$  and  $H_1$  but now it's more clearer than ever).

# T-Test:

→ Question:

In the population, avg IQ is 100. A team of researchers want to test a new medication to see if it has either a positive or negative effect on intelligence, or no effect at all.

A sample of 30 participants who have taken the medication has a mean of 140 with a std deviation of 20. Did the medication affect intelligence? CI = 95%.

Given:

$$\mu = 100, n = 30, \bar{x} = 140, s = 20$$

$$C.I. = 95\% \quad d = 1 - 0.95 = 0.05$$

① Null Hypothesis ( $H_0$ ) :  $\mu = 100$

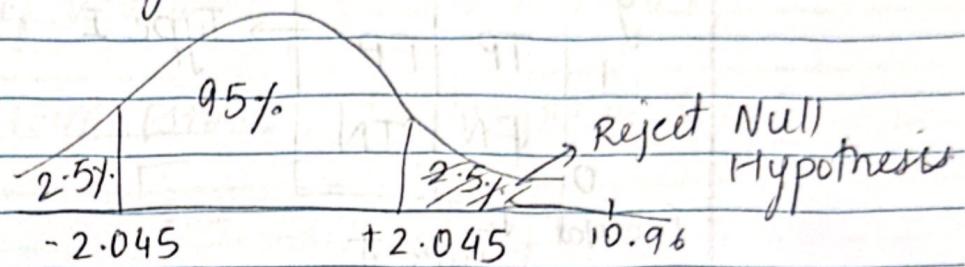
Alternate Hypothesis ( $H_a$ ) :  $\mu \neq 100$  {2-Tailed}  
[either positive or negative effect]

② Significance value :  $\alpha = 0.05$

$$\begin{aligned} \text{③ Degree of Freedom} &= n - 1 \\ &= 30 - 1 \\ &= 29. \end{aligned}$$

Degree of Freedom: Number of choices.

④ Decision Boundary:



If t-test is less than  $-2.045$  or greater than  $2.045$ , then we reject null hypothesis.

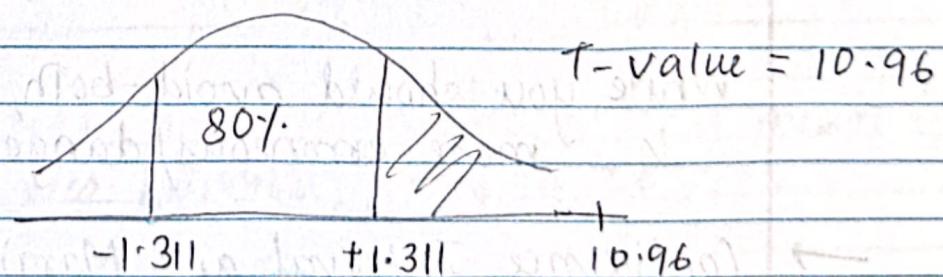
⑤ Calculate t-test statistics:

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

⑥ Conclusion: since  $10.96 > 2.045$ ,  
we reject the null hypothesis.

∴ The medication has a positive effect.

~~#Practice~~  
Scenario 2: C.I. = 80%  $\alpha = 1 - 0.80 = 0.20$ .



∴ We reject the Null Hypothesis and conclude that the medication has positive effect.

## Type I and Type II Errors:

Confusion Matrix:

		1	0	Actual (y)
$\hat{y}$	1	TP	FP	→ Type I
	0	FN	TN	
Predicted				Type II

TP: Actual value is True, pred val is True.

TN: Actual value is False, pred val is False.

(Type I)

FP: Actual value is False, pred val is True.

(Type II)

FN: Actual value is False, pred val is False.

N

Scenario: Predicting whether a patient has cancer.

(FP) → Type I error: Will classify someone who doesn't have cancer as them having cancer. This is bad, but still better.

(FN) → Type II error: Classify person having cancer as not having cancer.  
- Detrimental and harmful in long run.

While you should avoid both, Type II (FN) is more erroneous/dangerous.

→ Confidence Interval and Margin of Error:

Point Estimate: A value of any statistics that estimates value of an unknown population is called Point Estimate.

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

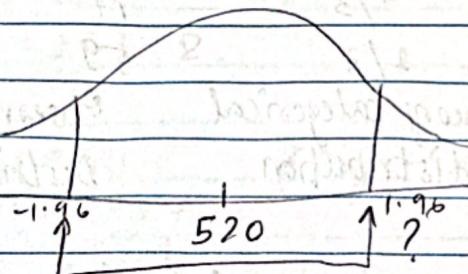
C.I.: we construct a C.I. to help estimate what actual value of unknown population mean is?

Point Estimate  $\pm$  Margin of Error

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

- Q.) On the verbal section of CAT Exam, the std deviation is known to be 100. A sample of 25 test takers has a mean of 520. Construct the 95% C.I. around the mean?

$$\rightarrow \bar{x} = 520 \quad \sigma = 100 \quad n = 25 \quad C.I. = 0.95 \quad \alpha = 0.05$$



$$\text{lower C.I.} = \bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} = 520 - (1.96) \frac{100}{\sqrt{25}} = 480.8$$

$$\text{Higher C.I.} = \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} = 520 + 1.96 \frac{100}{\sqrt{25}} = 559.2$$

Conclusion: I am 95% confident that mean CAT score lies between 480.8 to 559.2.

→ Biomedical researchers often handle categorical or nominal data.

\* Chi Square Test → Inferential Stats.

→ The chi-square test for goodness of fit test claims about population proportion of categorical variables  
→ [ordinal, nominal]

It's a non parametric test that's performed on categorical data.

Male (a pop. categorical)	Theory		Sample
	Yellow Bike	1/3	22
	Orange Bike	1/3	17
	Red Bike	1/3	59
Theory categorical distribution		Observed Categorical Distribution	

→ Chi square Fitness of Good:

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

Expected	2010			2020		
	< 50 kg	50 - 75 kg	> 75	< 50	50 - 75	> 75
	20%	30%	50%	140	160	200

→ In 2020, ages of  $n=500$  individuals were sampled.

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

(n=500 2020)

→ In the sample, the expected table would be:

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

5)

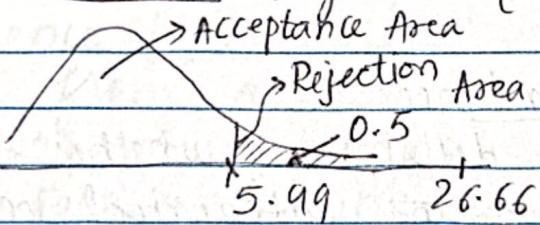
H<sub>0</sub>: Null Hypothesis: The data meets the expectation.

H<sub>1</sub>: Alternate Hyp: The data doesn't meet the expectation.

II)  $\alpha = 0.05$ , C.I. = 95% ✓ no. of categories

II) Degree of Freedom:  $df = k - 1$   
 $df = 3 - 1 = 2$

III) Decision Boundary: (Right-skewed Chi-square)



IV) If  $\chi^2$  is greater than 5.99, Reject H<sub>0</sub>. Else, we fail to reject H<sub>0</sub>.

V) Calculate  $\chi^2$  Test stats:

$$\chi^2 = \sum \frac{(O-E)^2}{E} = \frac{(140-100)^2}{100} + \frac{(160-150)^2}{150} + \frac{(200-250)^2}{250}$$

$$\chi^2 = 26.66$$

Since,  $26.66 > 5.99$ , we reject the null hypothesis. There's a population difference and it has increased when compared to 2010.

Source:

UCLA

(categorical data)

→ chi square Test: (3 Types)

I) goodness-of-fit test: (2 or more categories)

- whether a set of data comes from claimed distribution or not?

e.g. - the frequencies I OBSERVE are consistent with my theoretical EXPECTED values?

II) test of homogeneity:

- determine whether 2 or more subgroup  $\in$  Pqs. share the same distribution for a single categorical variable?

- e.g.) Do ppl of diff races have same proportion of smokers and non smokers?

III) test of independence:

- way of determining whether two or more subgroups categorical variables are associated with one another in population like race and smoking, education and political affiliation.

Diff between II and III?

- pretty similar, but diff. in matter of design.

Test of Independence

- Observational Unit:

Data is collected randomly from pop., usng 2 cat. variables are observed from each unit.

- ~~random~~ random observation me compare karo based on categorical.

Test of Homogeneity

- Data is collected by randomly sampling from each sub group separately.

- Intentionally proportionally samples ~~uth~~ from sub catgry.

## ANOVA:

- statistical analysis or method used to compare the means of 2 or more groups.

Parametric [z-test, t-test] - Mean of 2 groups.

Tests : [chi-square test] → Popn proportion.

- Factors (variable) → Levels

Eg: Factor = Medicines

Levels : 5mg    10mg    15mg [Dosage]

## Assumption in ANOVA:

- 1) Normality of Sample Distribution: Mean - The distribution of sample mean is normally distributed.
- 2) Absence of outliers: Outliers needs to be removed from dataset.
- 3) Homogeneity of VARIANCE: Each one of the population has same variance.

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

- 4) Population variance in diff. levels of each independent variable are equal.

- 5) Samples are independent and random.

## Types:

- ① 1 Way ANOVA: 1 Factor, atleast 2 levels  
(levels are independent)

Eg. Dr. wants to test a new medication to decrease headache. They split participants in 3 dosages [10, 20, 30] mg. Dr. asks participant to rate the headache. [1-10]

## ANOVA - F Distribution.

~~Z~~  $\chi^2$  distribution.

### 2) Repeated Levels ANOVA:

1 Factor ; atleast 2 Levels.

Levels are dependent.

Levels	Running → Factor		
	Day 1	Day 2	Day 3
-	-	-	

If you're learning a new skill / practicing with incr. in the levels (duration), you'll improve a skill. So levels are dependent.

### 3) Factorial ANOVA: Two or more factors

(each of which w atleast 2 levels), levels can be independent and dependent.

Running   Factor		
Day 1	Day 2	Day 3   Level.

gender  
Male  
Female  
Factor  
level

### HYPOTHESIS TESTING IN ANOVA:

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

Alternate Hypothesis:  $H_1: \text{Atleast one of the mean is not equal.}$

$$\mu_1 \neq \mu_2 \neq \mu_3 \neq \dots \neq \mu_k$$

### Test statistics

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

$X_1$	$X_2$	$X_3$
a	c	e
b	d	f
$\sum X_1 = 15$	$\sum X_2 = 19$	$\sum X_3 = 22$
$\bar{x}_1 = 3$	$\bar{x}_2 = 19/5$	$\bar{x}_3 = 22/5$

$H_0: \bar{x}_1 = \bar{x}_2 = \bar{x}_3$   
 $H_A: \text{Atleast one meant not equal.}$