

1 May

NOTES KASI

## 1 Population (N)

Collection of all items

e.g. New York University Students

## Sample (n)

Subset of Population

This is called Statistics.

Less costly, less time consuming

### Sample

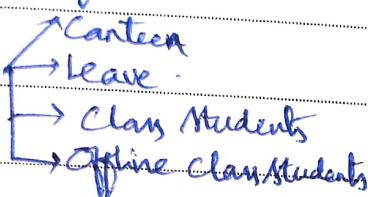
Randomness

Representativeness

Accurately represents the members of entire population

### Groups

NYU Students



Sample: 50 students

Random :-

Representativeness :- Considering all groups

## 2 Various types of Data :-

Based on type :-

1. Categorical :- Car brands, BMW, Audi, ... Yes/No Questions ; ~~Age~~ Climate Seasons
2. Numerical :-
- Discrete - No. of Children in class - Integer
  - Continuous - Weight of a Person (depends on scaling...)  $10/3 = 3.33 \dots 3.\bar{3}$  63.4789 kg

Based on level of Measurement :-

1. Qualitative
- Nominal - Car brands - No Order
  - Ordinal - Seasons - Summer  $\rightarrow$  Rainy  $\rightarrow$  Winter - Order

2.

2. Quantitative
- Interval - No True zero -  ${}^{\circ}\text{C} / {}^{\circ}\text{F}$  (No True zero)
  - Ratio - Real Scenario's - True zero -  
e.g. Two Apples 6 apples eg. Temp  ${}^{\circ}\text{K}$  (kelvin)  
2 : 6 (has True zero)  ${}^{\circ}\text{K}$

Today:  $5^{\circ}\text{C}$  or  $41^{\circ}\text{F}$

Yesterday:  $10^{\circ}\text{C}$  or  $50^{\circ}\text{F}$

(In terms of  ${}^{\circ}\text{C}$   $\rightarrow$  we feel more temp)  $\rightarrow$  double  
(In terms of  ${}^{\circ}\text{F}$   $\rightarrow$  we don't feel diff)

## NOTES

These ~~ster~~ scales are artificially created.

In terms of  ${}^{\circ}\text{K}$  (Has True Zero).

Neg. Can be colder than  $-273\text{ }^{\circ}\text{C} \equiv \underline{0\text{ }^{\circ}\text{K}} = -459\text{ }^{\circ}\text{F}$

$5\text{ }^{\circ}\text{K} : 10\text{ }^{\circ}\text{K}$   
(We can)

$\times 5\text{ }^{\circ}\text{C} : 10\text{ }^{\circ}\text{C}$   
(We don't)



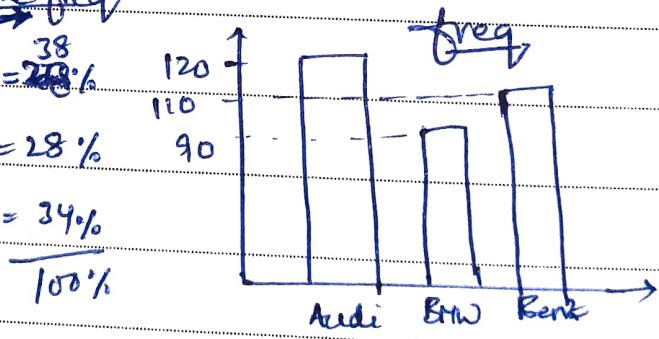
\* Visualize Data - Can be done by knowing its Type & Measurement level for Categorical Variables

1. Freq. Distribution Table
2. Bar Chart
3. Pie Chart
4. Pareto Diagram

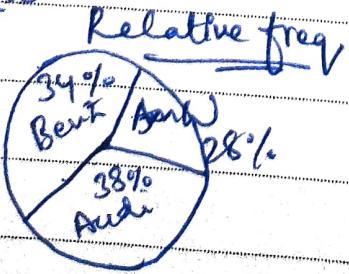
### 1. Freq. Distribution Table

Car	Freq	Relative freq
Audi	120	$\frac{120}{320} \times 100 = 38\%$
BMW	90	$\frac{90}{320} \times 100 = 28\%$
Benz	110	$\frac{110}{320} \times 100 = 34\%$
	<u>320</u>	<u>100%</u>

### 2. Bar Chart / Column Chart



### 3. Pie Chart



### 4. Pareto Diagram

Sum of Relative freq.

Descending Order of freq.

Bar Chart Data + Pie Chart data can be seen.

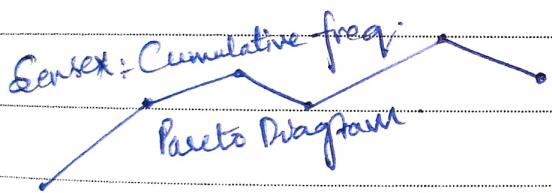
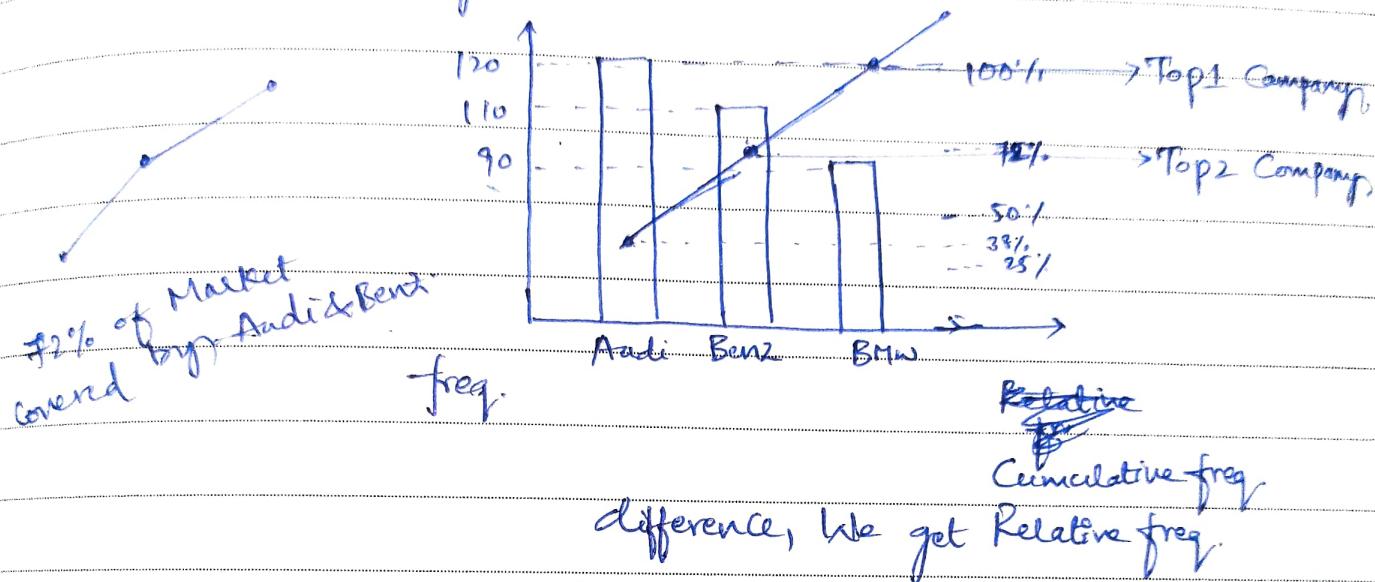
NOTES

Recording: Descending Order of freq.

Car Freq. Relatifefreq. Cumulativefreq.

Audi	120	88%	$\rightarrow 38\%$
Benz	110	34%	$\leftarrow \rightarrow 72\%$
BMW	90	28%	$\leftarrow \rightarrow 100\%$

### Pareto Diagram



useful in  
Real Time

1. Histogram
2. Cross Tables
3. Scatter Plots

### freq. Distribution Table

Dataset - freq

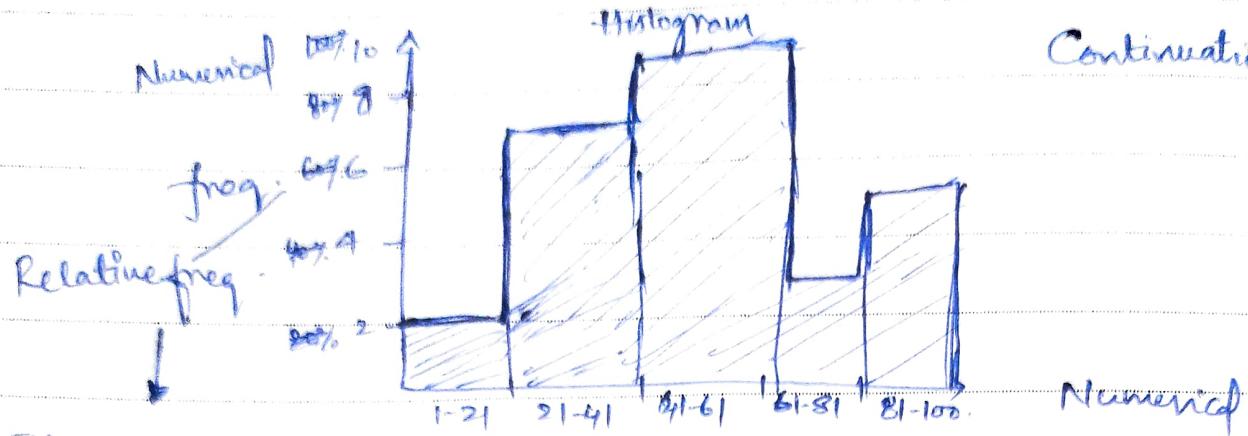
Desired intervals = 5

Interval width =  $\frac{\text{Largest No} - \text{Smallest No.}}{\text{Desired Intervals}}$

$$= \frac{100 - 1}{5} \approx 20 \quad (\text{But have to use interval width})$$

Interval Start	Interval End	Freq.	Relative freq.
1 - 20 (Include)	- 2	2	$\frac{2}{26} \times 100 \approx 7.69\%$
(Don't include) 1 - 40	- 41	7	26.9%
41 - 60	- 61	9	34.6%
61 - 80	- 81	3	11.5%
81 - 100	- 101	5	19.2%
		$\sum 26$	100%

## NOTES



88

Graph looks same, b/c if we use Relative freq.

Histogram can be done with Cinequal Intervals

Age	diff Interval
18-35	7
26-30	5
31-35	5
60+	-



\*\*\* Relationship b/w 2 Variables:

1. Categorical -  $\rightarrow$  Cross Tables

2. Numerical -  $\rightarrow$  Scatter plots (Both Variables are Numerical)

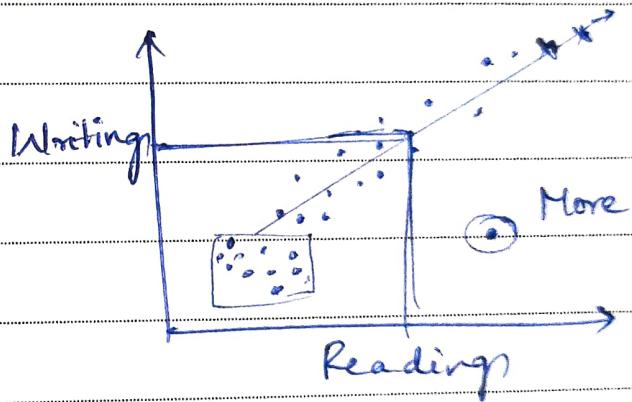
Cross Tables      Cross Table

Type of Investment	Investor A	Inv. B	Inv. C	Total
Stocks	96	185	39	320
Bonds	181	3	29	213
Real Estate	88	152	142	382
Total	365	340	210	915

Cross Tables - Side by Side Bar ChartScatter Plots

Student ID      Reading, Writing  
Scores

1	273	216
2		
3		
:		
10	390	206



from graph

~~Outliers~~  
Up trend:

Conclusion:

1. Most Students getting good score in Writing, Scoring good in Reading.
2. Average Most are concentrated here.
3. High Reading, High Writing score.
4. Bad Scores people - Performance trends in deviation while performing different tasks
5. Good at Reading, Bad @ Writing  $\Rightarrow$  Outliers.

Scatter Plot data Analysis)

## NOTES

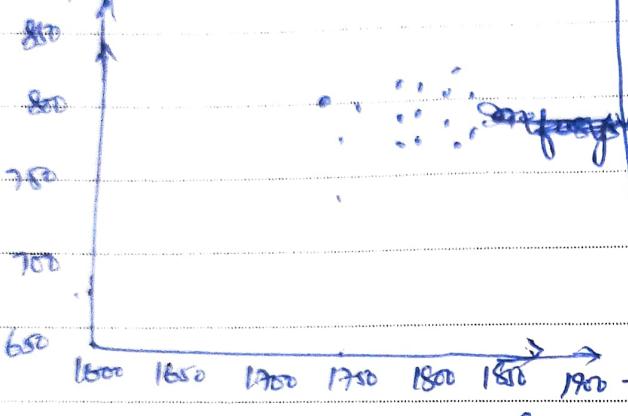
Each Share Value

+ 1800 -

+ 400 -

Lipin Pharma Share + 800 -

Lipin Pharma



78-

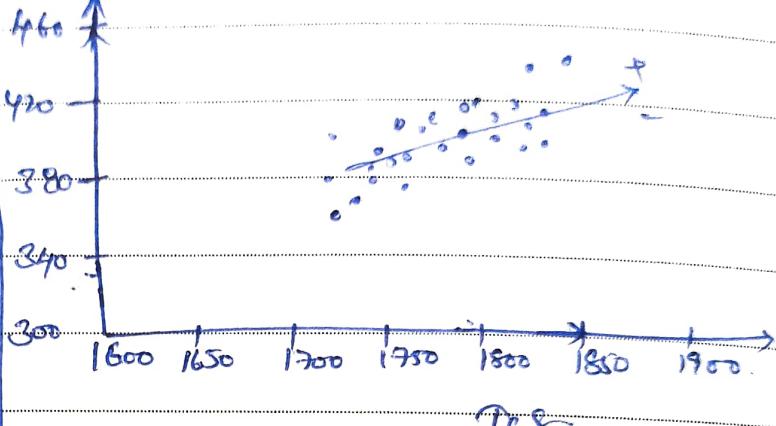
Independent

No Correlation b/w

When one is high, the other not.

Vertical somewhat.

Different



7cs

Dependent

Highly Correlated b/w.

When one is higher, the other is &  
vice versa

- \* Measures of Central Tendency: Mean, Median, Mode
- \* Measuring Skewness
- \* ~~Measures of Variability~~
  - (i) Standard Deviation & (ii) Coefficient of Variation
  - Covariance
- \* Correlation Coefficient

NOTES

### \* 3. Measures of Central Tendency:-

Mean:  $\mu$  (population)       $\bar{x}$  (sample)

Fairly affected by outliers

Student	Marks
1	10
2	20
3	30
4	40
5	50
$\rightarrow$	6
	0

$$\text{Mean} = \frac{\sum_{i=1}^N x_i}{N} = \frac{x_1 + x_2 + \dots + x_N}{N}$$

$$\text{Mean}(5) = \frac{10+20+30+40+50}{5} = \text{Marks } 30$$

$$\text{Mean}(6) = \frac{10+20+30+40+50+0}{6} = \text{Marks } 25$$

Median: Middle No. in Ordered Dataset.

No. at position  $\frac{(n+1)}{2}$  in Ordered List

$$\text{Median}_5 = \frac{5+1}{2} = 3 \rightarrow \text{Marks } 30$$

$$\text{Median}_6 = \frac{6+1}{2} = 3.5 \rightarrow \frac{30+40}{2} = 35$$

Mode: Used for both Categorical & Numerical.

10, 20, 10, 30, 40, 50.       $\rightarrow \text{Mode} = 10$

"Repeated"

Cats have  $\rightarrow$  single / Multiple modes  $\rightarrow$  no mode.

\* Which is best? Depends on situation.

For e.g.: Income of 1<sup>st</sup> year & 2<sup>nd</sup> year joiners in TCS  $\rightarrow$  Mean is best

Income of all employees (including seniors)  $\rightarrow$  Outliers  
 $\rightarrow$  Median is best

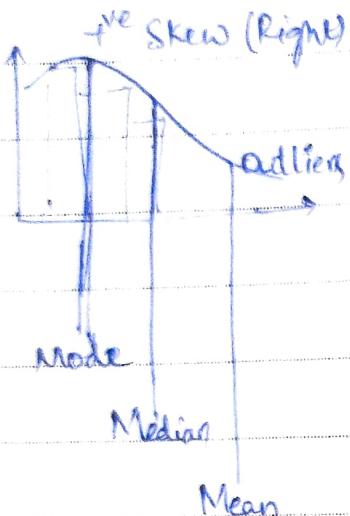
## NOTES

### To Know Asymmetry of Distribution

Skewness:

e.g. Mean Median Mode  
2.79 2 2

Mean > Median



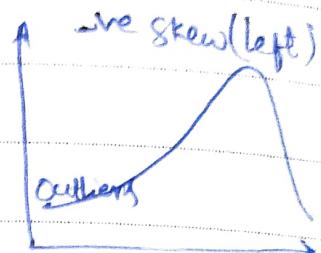
Mean Median Mode  
4 4 4

Mean = Median



Mean Median Mode  
4.9 5 5

Mean < Median



Distribution is Symmetrical

Outliers are situated in left

How skewness is used?

It gives where most of the data is situated.

$$\text{Skewness} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{\sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}}^3$$

3 Measures of Variability (i) Variance (ii) Std. Deviation.

(iii) Coef. of Variation.

(i) Variance ( $\sigma^2$ )

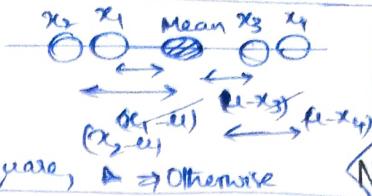
Population (100)      Sample 1 - 10  
                          Sample 10 - 10  
                          100

All 10 samples gives different measures.



Variance measures the dispersion of a set of data points around their mean. (Spread)

Mean ( $\mu$ )



NOTES

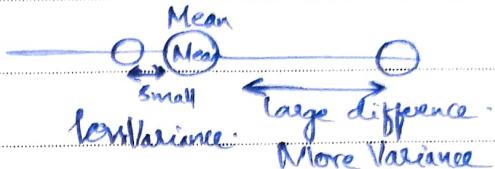
$$\text{population Variance } \sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$$

$$\text{Sample Variance } s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

\* Squaring: Error sum + 0.2 & -0.2

May Cancel out

To avoid, We square.



\* ~~Error amplified by squaring.~~

e.g.

Population:

1  
2  
3  
4  
5

$$\text{Mean} = \frac{1+2+3+4+5}{5} = 3$$

$$\text{pop. Variance } \sigma^2 = \frac{(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2}{5} = 2$$

$$\text{Consider as Sample, Sample Variance } s^2 = \frac{(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2}{5-1} = 2.5$$

Why ~~s~~  $s^2 > \sigma^2$  ?

Imaginary population

1  
2  
3  
4  
5  
5  
5

$$\text{Mean} = 3.2$$

$$\sigma^2 = 2.96$$

sample popul

1  
2  
3  
4  
5

e.g.: Salaries of 10 Team Members in our project.

{ 30000/-  
20000/-  
1,20,000/-  
45,000/-  
32,000/- }

Variance ( $s^2$ )

$$16368000000$$

3000  
3100  
3200  
3300  
3400

$$s^2 = 25000$$

Standard Deviation,  $s(\sigma)$

$$\sqrt{s^2}$$

$$\sqrt{s^2}$$

## NOTES

(iii) Coeff. of Variation = / Relative Std. Deviation.

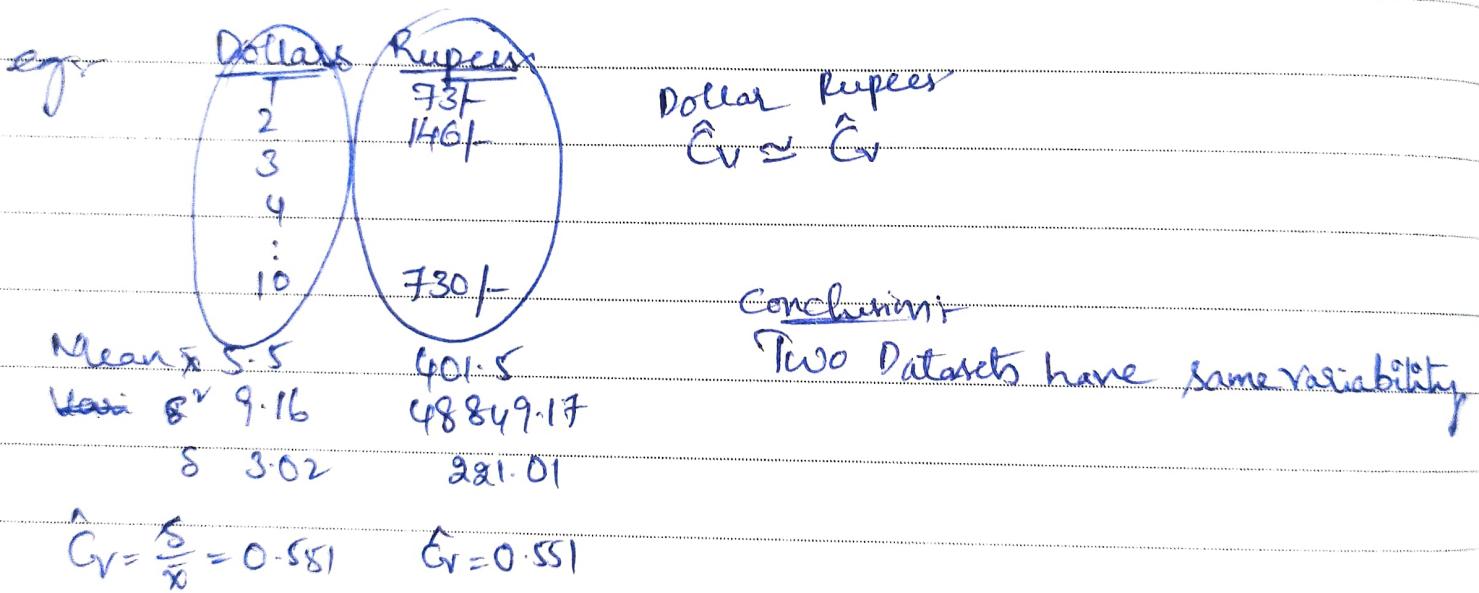
$$C.V = \frac{\text{Std. Deviation}}{\text{Mean}} = \frac{\sigma}{\bar{x}} \quad (\text{or}) \quad \frac{s}{\bar{x}}$$

population      Sample

$$C_v = \frac{\sigma}{\bar{x}} \quad \hat{C}_v = \frac{s}{\bar{x}}$$

or Std. Deviation is the most common measure of variability for a Single Dataset.

$C_v$  = Coef. of Variation in comparing Two or More Datasets



eg:- Income of People in U.S. V/S. Denmark Sample.

Currency is different. \$

We can compare  $\hat{C}_v$ .

Kr

## Measures of Relationship b/w Variables

1. Covariance

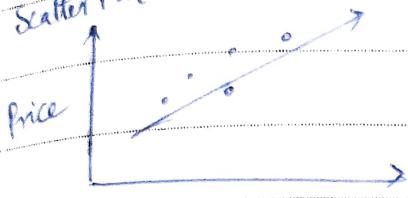
NOTES

2. Linear Correlation Coef.

\* Covariance :-

Land Size (A)	Price (Y)
650	472,000
780	998,000
1200	1,200,000
720	800,000
975	895,000

Scatter Plot



$$\text{Cov. Sample} = 33,491,250$$

Sample Covariance

$$S_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

Pop Covariance N

$$\sigma_{xy} = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{N}$$

During  $\Sigma$  (summation)

+ - also summed  
cup sign also.

Covariance

- > 0, the two variables moving together
- < 0, " " " opposite
- = 0, " " " are independent.

Covariance  
30000

50  
0.23

How Can We Know, Covariance More or just more?

So: We go for Correlation Coef.

\* Correlation Coef:-

-1 to 1

$$\frac{\text{Cov}(x,y)}{\text{std.dev}(x) * \text{std.dev}(y)}$$

Cor. Coef = 0.87

Highly correlated  
positive indicate

if size of land ↑ → Price ↑

1 → Perfectly Correlated ↑↑ ↗ ↘

0 → Independent

-1 → Perfectly Correlated ↑ ↓

Decrease of Umbrella in Seaton Row

\* Causality:-

Size causes Price, but not Vice Versa

Correlation doesn't imply causation.

# Descriptive Statistics

## NOTES

### Product Table

- I.D → Categorical
- Type of Property → Qualitative (Nominal)
- Area → Numerical
- Price → Numerical
- Status → Categorical
- Cust. ID → Numerical
- Name → Categorical
- Age → Numerical (Ratio)
- Gender → Categorical
- Country → Quantitative (Numerical) (Ratio)
- State → Categorical

### Descriptive Statistics      Real Example

### Customer Table

Set →  
Pending, Empty - No one has bought.

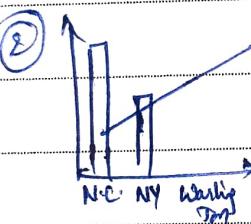
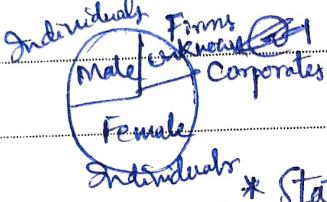
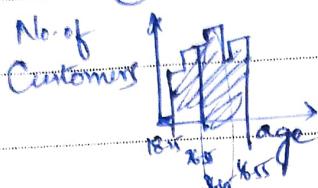
Age  
Mean 46.15  
Median 45  
Mode 48  
Skew 0.24  
→ Right Skew

Sample Dataset

Sample formula's

Variance 164  
St.dev 12.84

~~Descriptive~~



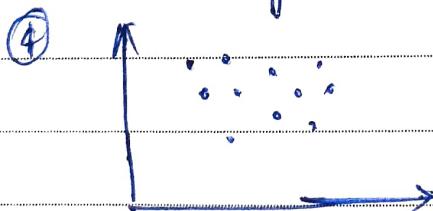
\* Age → Histogram → (Numerical)

@ What age more customers are buying  
Hence?

\* Gender → Pie Chart. (Categorical)

\* State → Pareto diagram

@ Age Vs Price (Scatter Plot)



Covariance -176.381.87

Correlation Coef. -0.17

Inference: Price Not related to age.

1. Females are buying more. (Gender)

2. Most sales in N.C. state (State)

3. Age 45 are buying more with Avg mean & Std dev of Right Skewness → Younger people buy more than older people

4. No relation b/w Age & Price they are buying.

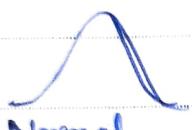
## Inferential Statistics

- Infer → To give some proper conclusions.  
Probability Theory & Distributions based on sample data

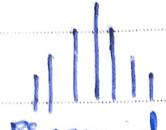
### NOTES

- Distribution
- Point Estimates
- Confidence Intervals.

Distribution :- Usually Means Probability distribution.



Normal



Binomial



Uniform Distributions

A distribution is a function that shows the possible values for a variable

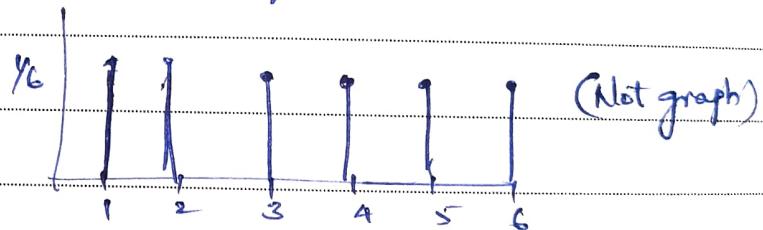
and how often they occur.

Numerical - Discrete Distribution :-

Die & prob. of getting of 1

Outcome	Probability
1	$\frac{1}{6}$ = 0.17
2	$\frac{1}{6}$
3	$\frac{1}{6}$
4	$\frac{1}{6}$
5	$\frac{1}{6}$
6	$\frac{1}{6}$
7	0

⇒ Discrete uniform Distribution.

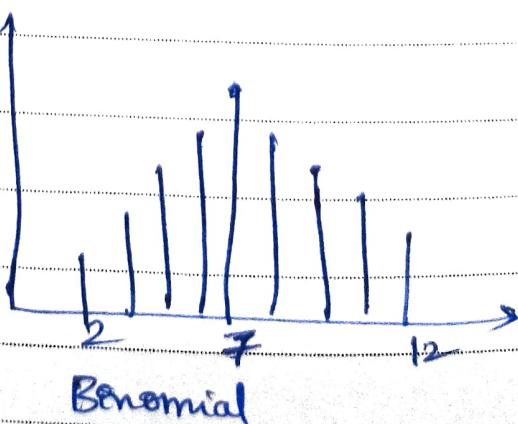


Die & Rolling Two Dice.

- |       |       |     |       |
|-------|-------|-----|-------|
| (1,1) | (1,2) | ... | (1,6) |
| (2,1) |       | ... |       |
| ⋮     | ⋮     |     |       |
| (6,1) | ...   | ... | (6,6) |

Outcome sum of Prob.

2	$\frac{1}{36}$
3	$\frac{2}{36}$
4	$\frac{3}{36}$
⋮	⋮
12	$\frac{1}{36}$



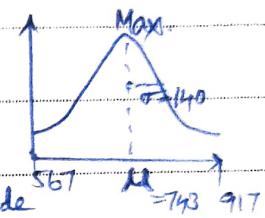
Continuous DistributionNormal Distribution:

- \* They approximate a wide variety of Random Variables
- \* Distribution of sample means with large enough sample sizes could be approximated to Normal
- \* All compatible Decisions based on Normal Distribution is good.

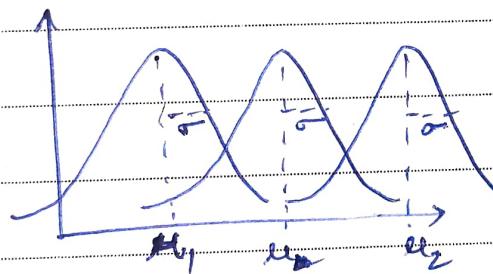
$$N \sim (\mu, \sigma^2)$$

Distribution Mean Variance

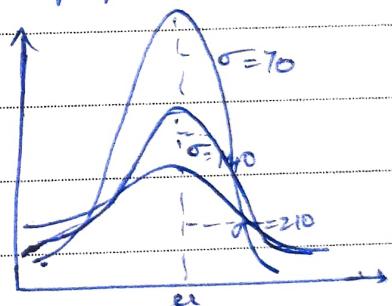
Bell Curve  
Normal Distribution  
Gaussian Distribution  
Symmetrical  $\Rightarrow$  Skewness = 0  
 $\Rightarrow$  Mean = Median = Mode



## Controlling for Std. Deviation:



## Controlling for Mean - By changing Std. dev.



More data in Middle  
less in thinner tails

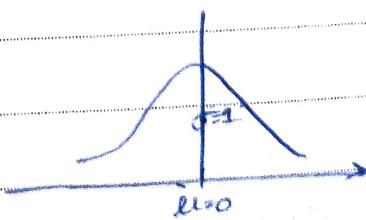
Higher Std. deviation  $\Rightarrow$  Make graph flat  $\Rightarrow$  fatter tails

Standard Normal Distribution

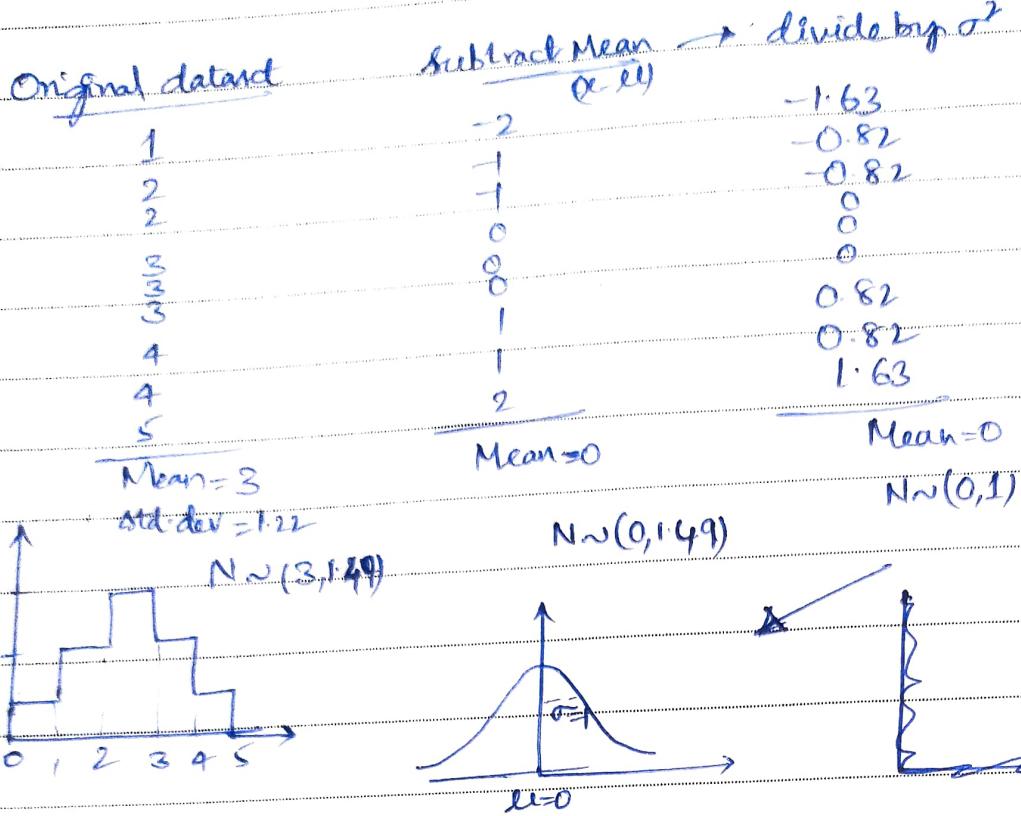
$$N \sim (\mu, \sigma^2) \sim (0, 1)$$

\*\*\*

$$Z = \frac{x-\mu}{\sigma} \quad Z \sim N(0, 1)$$



## NOTES



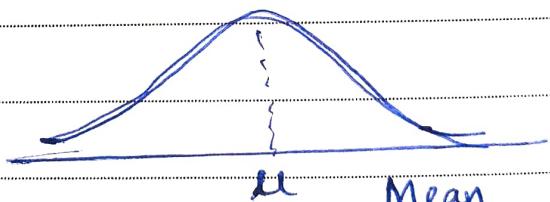
### \* Central Limit Theorem:-

Sample 1      Mean 2600/-  
 Sample 2      3201/-  
 Sample 3      2844/-

~~We see different values~~  
 We get for different samples.

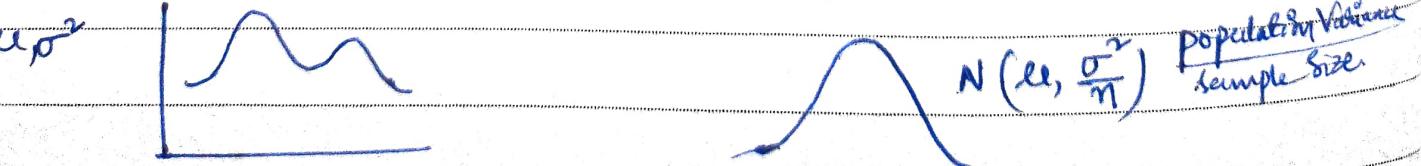
### Sampling Distribution of the Mean

We take some more samples and calculate Means



- We are not calculating Population Mean.
- We Calculate Means of different samples in population. (~~Average of Sample Means~~)
- All together gives good result precisely.

### Original Distribution



### Sampling Distribution

$$N(\mu, \frac{\sigma^2}{n}) \quad \begin{matrix} \text{Population Variance} \\ \text{sample size} \end{matrix}$$

(The bigger the sample  
we get less variance)      n > 30

Central Limit Theorem allows us to perform tests, solve problems and make inferences using Normal Distribution, even when the population is not Normally Distributed.

### Standard Error:

$$S.E. = \frac{\sigma}{\sqrt{n}} = \frac{\sigma}{\sqrt{m}}$$

Stand. Deviation

As sample size  $\uparrow$ , Error  $\downarrow$ .



• Kan Pradeep

### Estimators and Estimates

- \* Point Estimate: Exact point of Estimate • Middle of Confidence Interval
- \* Confidence Estimate:  $\boxed{\quad}$  Range: Good for Inferences

### \* Point Estimators and Estimates

#### Estimator of Parameter

$\bar{x}$  of  $\mu$

$s^2$  of  $\sigma^2$

$\overline{\bar{x}}$   
 $\overline{\text{average}}$

#### \* Unbiased Estimator

$\bar{x}$  of  $\mu$

#### \* Biased Estimator

$\bar{x} + 1$  ft of  $\bar{x}$

#### \* Efficiency:

Most Efficient Estimator is Unbiased Estimator with Smallest Variance

### \* Estimators: A type of statistic

e.g. Avg. pizza cost \* 100% Confidence  $\Rightarrow$  We have to consider entire population  
But not possible

160/- 200/- 250/-

\*  $\alpha$  Confidence level  
5%  $\leftarrow$  95%

## NOTES



(Point Estimate  $\pm$  Reliability Factor \* Standard Error, P.E. = R.E \* SE)

$$[\bar{x} - R.F. \cdot \frac{\sigma}{\sqrt{n}}, \bar{x} + R.F. \cdot \frac{\sigma}{\sqrt{n}}]$$

90% Confidence Interval 99%

$$\alpha = 0.1$$

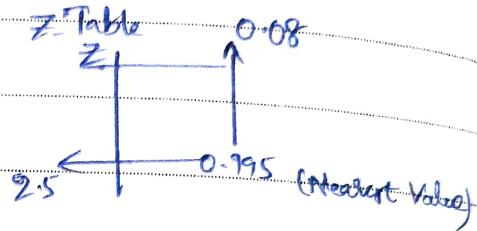
$$\alpha = 0.01$$

$$\alpha_1 = 0.05$$

$$1 - \alpha_1 = 0.95$$

$$1 - \alpha_1 = 0.995$$

$$\frac{\alpha}{2} = 0.005$$

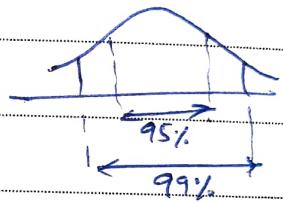


$$[\bar{x} - Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}, \bar{x} + Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}] \quad R.F. = 2.5 + 0.08 = 2.58$$

[Sample Mean  $\pm$  R.F. \* Deviation,  $\bar{x}$ ,  $\bar{x} +$  ]  
From Z-table.

Confidence 95%, Narrower

[Low Value, High Value]



e.g. Average of Salaries

Min.

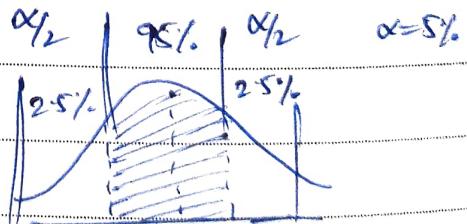
Max.

100% confidence [Min, Max]

Confidence Intervals:-

95% of Confidence

We are 95% confident that the true population means falls within this interval.



Majority of observations are around Mean.

## NOTES

Std-Normal Distribution Mean = 0

Limit  $[-z, +z]$



Confidence Interval

$1 - \alpha$  % is smaller

$\Rightarrow (1-\alpha)$  smaller

$\Rightarrow$  Confidence Interval is Narrower.

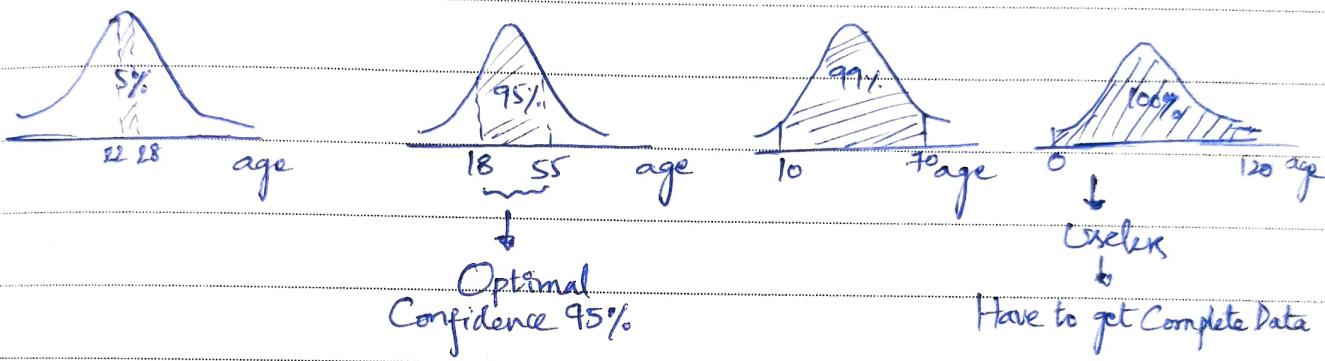
95% is Larger

Larger.

Broader.



b) Age of 100 persons taking Statistics Course.



Have to get Complete Data

T-tablet:-

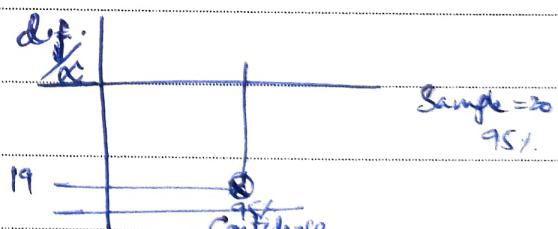
After 30° df (degree freedom) t-statistics  $\approx$  Z-statistics.

\* Student T-distribution :-

Inference through small samples.

Unknown population Variance

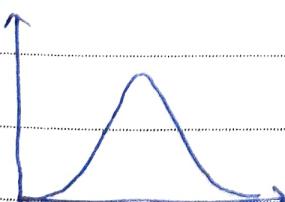
Huge Real Life Application.



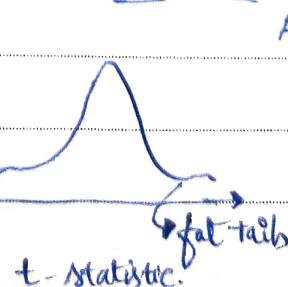
$$t_{n-1, \alpha} = \frac{\bar{x} - \mu}{\left( \frac{s}{\sqrt{n}} \right)}$$

n - sample size  
 $n-1$  freedom degree.  
 $\alpha$  significance level

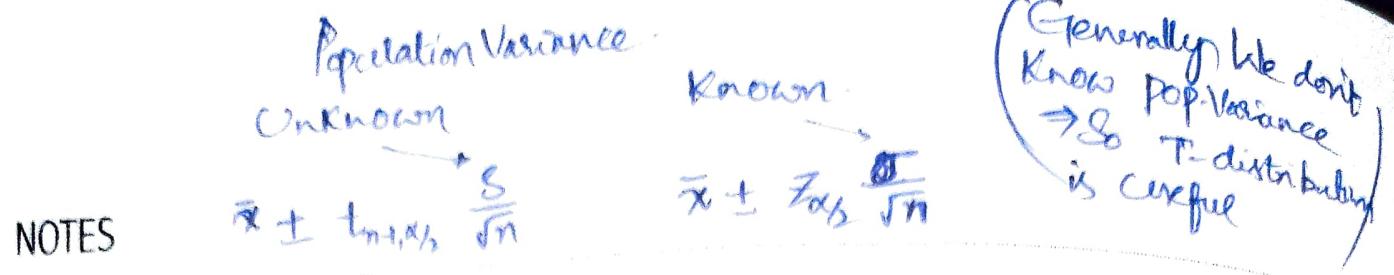
Normal distribution



Student T distribution



App. to Normal.



e.g.: Data Scientist Sample

S No Salaries (Dataset)

- 1 →
- 2 →
- 3 →
- 4 →
- 5 →
- 6 →

[Student T-table Range to analyze]

~~Conf.~~ (

\* C.I. <sub>1 - \alpha</sub> unknown  $\hat{z} \equiv (LL, UL) \equiv$  <sup>Width</sup> More

\* C.I. <sub>95%</sub> Known  $\hat{z} \equiv (LL, UL) \equiv$  less

If population variance is unknown, our boundaries are more little less accurate:



MARGIN OF ERROR:

M.E.

$$\bar{x} \pm \left( Z_{\alpha/2} * \frac{\sigma}{\sqrt{n}} \right)$$

M.E.

$$\bar{x} \pm \left[ t_{n-1, \alpha/2} * \frac{s}{\sqrt{n}} \right]$$

Conf. Int. =  $\bar{x} \pm M.E.$

\* M.E.  $\downarrow^{\text{est}}$   $\Rightarrow$  Conf. Int.  $\downarrow^{\text{est}}$   $\Rightarrow$  Std. Dev.  $\downarrow^{\text{est}}$

less error  
more

\* M.E.  $\downarrow^{\text{est}}$   $\Rightarrow n \uparrow^{\text{est}}$

Higher no. of samples,  
We get near to True Mean.

## \* Dependent

\* Before & After situation  
Cause and Effect

## Independent

- \* Population Variance Known
- \* Population Variance Unknown, but assumed to be equal.
- \* Population Variance is Unknown, but assumed to be different

NOTES

## ① Calculating Confidential Intervals for 2 Means with Dependent Samples

e.g. Weight loss - Blood samples of a person are Dependent

\* Person is same here.

eg:- Sugar Test:

Before Breakfast 90

After Breakfast 160 → Why increasing?

Reason - Cause - Had breakfast.

What happened - Effect - Increased Sugar level in the body.

Here we do testing on same person.

Testing Approaches:-

1. Confidence Intervals for Dependent Samples
2. Statistical Methods like Regressions.

Patient	Drug taken		Difference
	Before	After	
1	10	15	-0.3
2	15	20	-5
3	20	24	-4
:			
n	10	15	0.9
			Mean = $\frac{\sum \text{diff}}{n}$
			Std. Dev = $\sqrt{\frac{\sum (\text{diff} - \bar{\text{diff}})^2}{n-1}}$

Confidence Interval

$$\bar{d} \pm t_{n-1, \alpha/2} * \frac{s_d}{\sqrt{n}}$$

for eg Result: (0.01, 0.65)

\* True Mean will fall in this interval.

\* Whole interval is  $t^{\text{re}}$   $\Rightarrow$

\* Level of drug is  $t^{\text{re}}$  (bcz re interval)

$\Rightarrow$  Drug is Effective.

## ② Confidence Interval for the difference of 2 Means, Independent Samples, Variance Known:

eg:-	Engg' x		M.B.A' y		Diff x-y
	Size	Sample Mean	Size	Sample Mean	
	100	58 pts	70	65 pts	-7
	Pop. Std. Dev	10 pts	5 pts		

Considerations:

1. Populations are Normally Distributed.
2. Population Variances are Known
3. Sample sizes are different.

## NOTES

- \* The grade getting in Engg. does not effect the grade of a M.B.A. Student.
- \* the samples are truly independent.

### Considerations

- ↳ Both samples are big.  $\Rightarrow Z$  Statistics
- ↳ Population Variance Known
- ↳ Populations are assumed to follow Normal Dist.

Variance of difference:

$$\sigma_{\text{diff}}^2 = \frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}$$

$$= \frac{10^2}{100} + \frac{5^2}{70} = 1.36$$

Mean diff = -7.

C.I:  $(\bar{x} - \bar{y}) \pm Z_{\alpha/2} \sqrt{\frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}}$

different point  
estimator

95% confidence  
(-9.28, -4.72)

- Inferences:-
1. True Mean falls in this interval. We are 95% confident.
  2. Whole interval is -ve  $\rightarrow$  Engineers consistently gets lower grades.
  3. We calculate diff b/w Management - Engineers  $\Rightarrow$  (-4.72, -9.28)

## NOTES

③ Confidence Interval for difference of 2 Means, Independent Samples, Variance Unknown but assumed to be Equal.

New York Apples

1. 3.8/-
2. 3.4/-
3. 4/-
4. 3.2/-
5. 3.62/-

L.A. Apples

1. 3.1/-
2. 8/-
3. 6.5/-
4. 3.2/-

	NY	L.A.
Mean	3.9	3.25
Std. Dev.	0.18	0.27
Sample Size	10	8

We assume variance is same for both

Pop Variance Unknown  
Small Samples

→ T-distr

$$\text{Pool Variance } S_p^2 = \frac{(n_x - 1) s_x^2 + (n_y - 1) s_y^2}{(n_x - 1) + (n_y - 1)}$$

$$\text{Pooled Variance} = \frac{(10-1)(0.18)^2 + (8-1)(0.27)^2}{(10-1) + (8-1)} = 0.05$$

$$\text{Pooled Std. Dev} = \sqrt{0.05} = 0.22$$

$$\text{C.I. } (\bar{x} - \bar{y}) \pm t_{n_x+n_y-2, \alpha/2} \sqrt{\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}}$$

↳ pooled std. dev.

$$\text{C.I. } (0.47, 0.92)$$

Inferences: Apple Cost  $>$  Apple Cost in L.A. because diff C.I. value is +ve  
N.Y.

④ C.I. for diff. of 2 Means, Indep. Samples, Vari. Unknown but assumed to be different.

Apples | Oranges

$$\text{C.I. } (\bar{x} - \bar{y}) \pm t_{v, \alpha/2} \sqrt{\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}}$$

Here Main Problem is Calculating Degree of Freedom

$$v = \frac{\left(\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}\right)^2}{\left(\frac{s_x^2}{n_x}\right)^2/(n_x-1) + \left(\frac{s_y^2}{n_y}\right)^2/(n_y-1)}$$

# Inferential Statistics

## NOTES

Al Bandy Shoe Shop in US: Operating in UK, Canada, Germany.

1. Some shoes are no one buying
2. Some shoes are more buying and failed to supply as per demand

Here Comes Confidential Intervals.

Shop Last 20 yrs Data, We are about to Consider <sup>Last</sup> 3 yrs ~~Last~~ Data

Size Conversion Table

Invoice No	Date	Country	Product	Shop	Gender	Shoe Size (US)	Shoe size (UK)	Shoe size (Canada)	Shoe size (Germany)	Unit Price	Discount %
					M						

## Inferences Regarding Dataset

### Shoe Size Conversion Table

Men	Women
UK PS (Euros)	US PS (Inches)

### 1. Sample Data

2. Two Main groups Men group & Women  
(Almost Independent) - ~~different sizes~~  
~~different models~~

Segment the data by 1. Shoe Size:

2. Country
3. Gender

### Men Group

US men size	Canada	US	UK	Germany	Total
6	15	5	6	30	105
Total					

### Women Group

US Women Group size	Canada	US	UK	Germany	Total
Total					

NOTES

Sum or Average of No.s  $\Rightarrow$  We assume Normality  $\Rightarrow$  Central Limit Theorem

Game Plan:

95% of Confidence Interval

1. Last 12 months of sales.

2. Only for Men shoes

3. Only for U.S.A.

(last 2 yrs, last 3 yrs)

(size 6, size 6.5 ...)

(All Countries Canada, U.K. ...)

U.S. 2016

U.S. Size	1	2	3	...	12 Dec.	Mean	Std. Error	M.E.	95% CI	No. of Pairs
6	4	1	3	...	0	Avg. $\frac{4+1+3}{12} = 2.92$		12	1.8 to 4.04	4
6.5								t-dist	0.46 to 2.88	3
7									2.90	
7.5									2.50 to 3.50	
16									25.00 to 35.67	36
Total									0.00 to 0.00	0

$$n = 12$$

$$t_{11, 0.025} = 2.2$$

Shoe size of "6", we require 4 pairs May.

Shoe size of 7.5", we require 36 pairs

Shoe size of 16", @ Nothing to Manufacture.

### Another Application of Confidential Intervals

Two shops are selling the same no. of Shoes

95% Confidentiality of 2 Shops Sales in Germany

U.S. Size	1	2	3	...	12 Dec.	GCR1
4	1	2	3	...	0	
5						
6						
11.5						

U.S. Size	1	2	3	...	12 Dec.	GCR2
4						
5						
6						
11.5						

logically, same year, same people, won't buy same brand shoes in different shops  
 So, two samples are independent

## NOTES

Same Country, same People.

We assume population variance of Sample 1  $\equiv$  Pop Variance of Sample 2.

Here The 2 samples are Independent, Population Variance Unknown, but assumed  
 to be equal.

	Mean		Sample Variance		Pooled Variance	ME	95% CI
	Gen1	Gen2	Gen1	Gen2			
4	-	-	-	-	-	-	0.00 0.00
4.5	-	-	-	-	-	-	0.23 0.09
:	-	-	-	-	-	-	-
11.5 <del>11.5</del>	-	-	-	-	-	-	-

95% confidence

$$t_{12+12-2, 0.025} = 2.07$$

We can't conclude  
 One shop has more no. of Pairshoes than Other on all sizes.  
 Reason Below

All confidence Intervals Start in Negatives and finish in +ve.

$$[L-L, U+U]$$

↓  
 '0' is included

~~Gen1~~

$$(-0.23, 0.90)$$

$$(1.11, 0.27)$$

$$(-0.72, 3.55)$$

Mostly Negative

Mostly +ve

In terms of Sales, Both shops are almost same.