

Course Name :Basic Statistics using GUI-R (RKWard)
Module : Terms of Statistics
Week 4 Lecture : 1

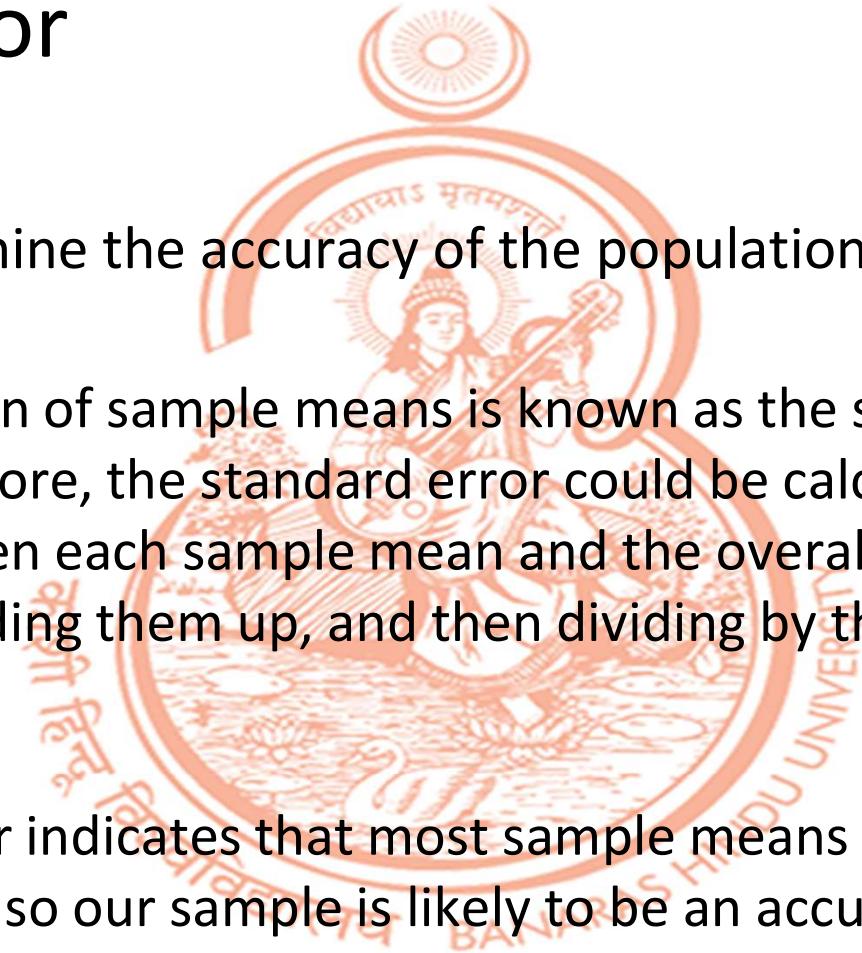
Harsh Pradhan, Assistant Professor,
Institute of Management Studies, BHU
https://bhu.ac.in/Site/FacultyProfile/1_5?FA000562

Standard Error

So how do we determine the accuracy of the population mean?

The standard deviation of sample means is known as the standard error of the mean (SE). Therefore, the standard error could be calculated by taking the difference between each sample mean and the overall mean, squaring these differences, adding them up, and then dividing by the number of samples.

A small standard error indicates that most sample means are similar to the population mean and so our sample is likely to be an accurate reflection of the population



Central Limit Theorem

In many practical problems, the exact knowledge of the parameter may not be necessary. It is quite adequate if an interval along with a probability statement is specified such that the probability that the random interval will cover the unknown parameter is a specified number

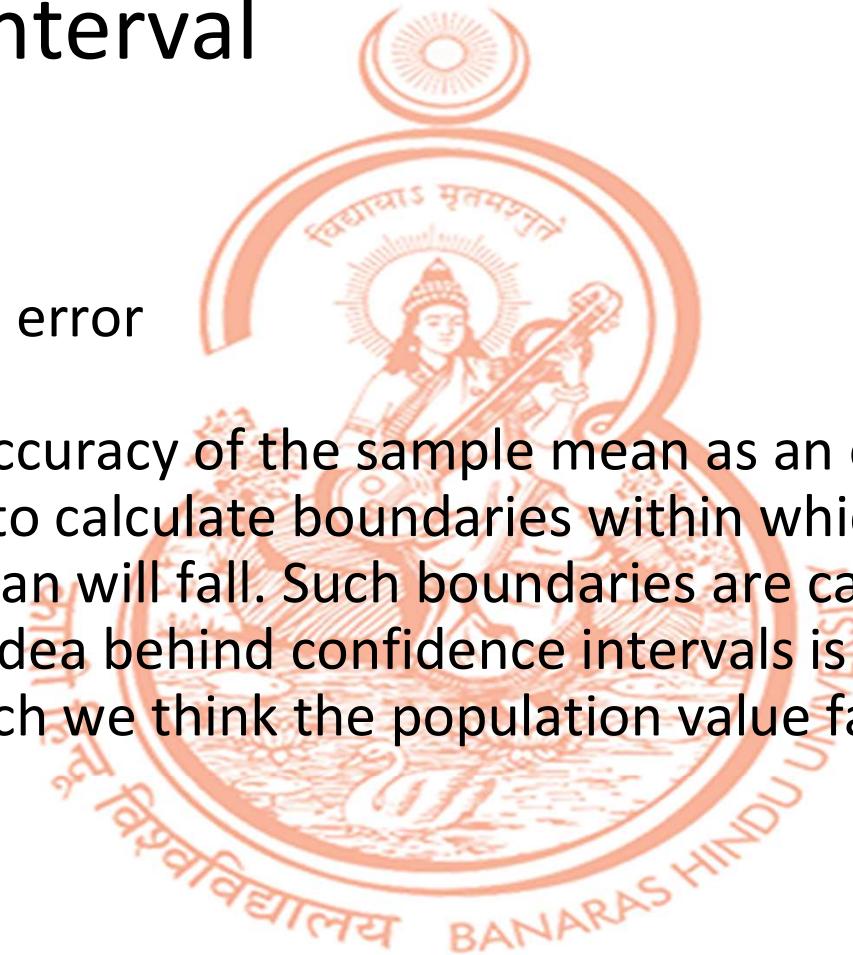
Central limit theorem when the sample is large ($N>30$) then sampling distribution has a normal distribution with a mean equal to the population mean, and a standard deviation of:

$$\sigma = s/\sqrt{N}$$

s standard deviation of population

Confidence Interval

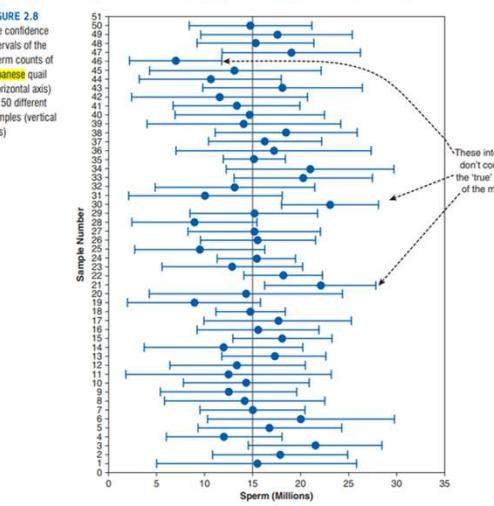
Apart from standard error



We can assess the accuracy of the sample mean as an estimate of the mean in the population is to calculate boundaries within which we believe the true value of the mean will fall. Such boundaries are called confidence intervals. The basic idea behind confidence intervals is to construct a range of values within which we think the population value falls.

Confidence Interval

So, when you see a 95% confidence interval for a mean, think of it like this: if we'd collected 100 samples, calculated the mean and then calculated a confidence interval for that mean (Figure) then for 95 of these samples, the confidence intervals we constructed would contain the true value of the mean in the population



Confidence Interval for μ

- The confidence interval can be expressed as

$\left(\bar{x} - Z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}, \bar{x} + Z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \right)$ is a $100(1-\alpha)\%$ CI for μ

or

$\left(\bar{x} - Z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x} + Z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \right)$ with $100(1-\alpha)\%$ confidence.

Confidence Interval

To calculate the confidence interval, we need to know the limits within which 95% of means will fall.

$$\text{lower boundary of confidence interval} = \bar{X}_{avg} - (1.96 \times SE)$$

$$\text{upper boundary of confidence interval} = \bar{X}_{avg} + (1.96 \times SE)$$

If the mean represents the true mean well, then the confidence interval of that mean should be small. We know that 95% of confidence intervals contain the true mean, so we can assume this confidence interval contains the true mean; therefore, if the interval is small, the sample mean must be very close to the true mean. Conversely, if the confidence interval is very wide then the sample mean could be very different from the true mean, indicating that it is a bad representation of the population.

A confidence interval for the mean is a range of scores constructed such that the population mean will fall within this range in 95% of samples.

Statistic

A numerical quantity calculated from a sample of data, summarizing characteristics such as distribution, central tendency, variability, or relationships.

Examples:

Mean

Median

Standard Deviation

Correlation Coefficient

Frequency

Proportion

Parameter

An unknown characteristic of a population that is being estimated using sample data.

Examples:

Population Mean

Population Median

Population Standard Deviation

Population Correlation Coefficient

Population Frequency

Population Proportion

Estimate

An approximation or prediction of the population parameter based on sample data.

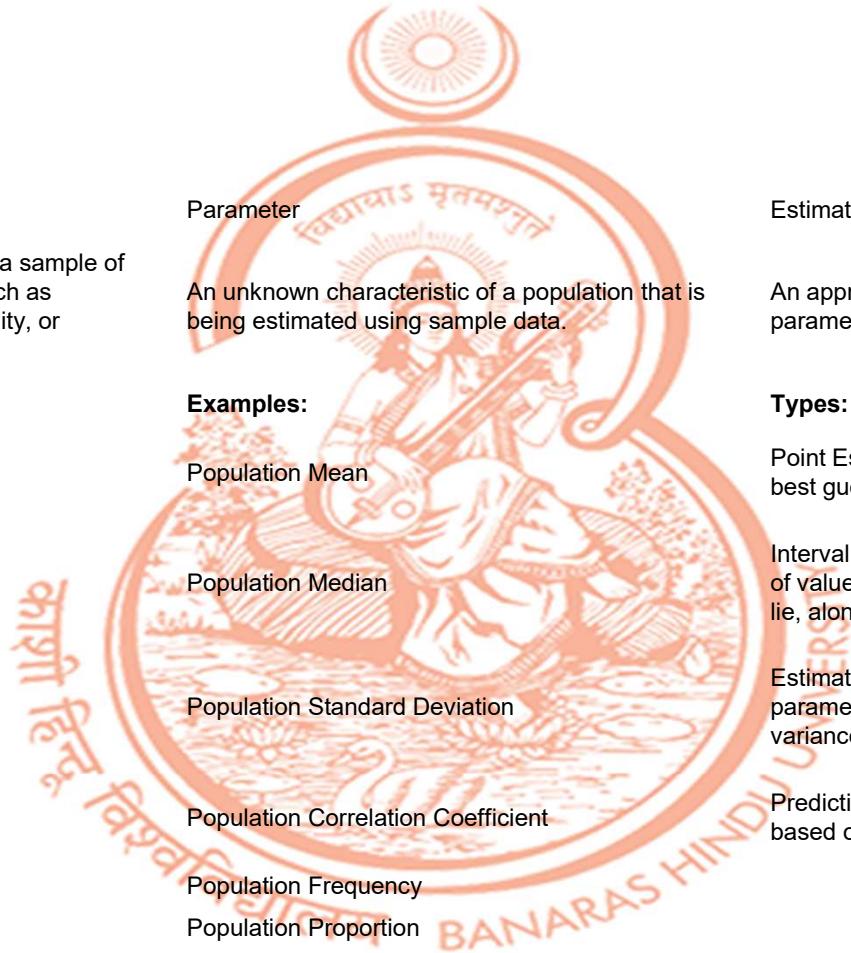
Types:

Point Estimate: A single value that serves as the best guess for the parameter being estimated.

Interval Estimate: An estimate that provides a range of values within which the parameter is believed to lie, along with a level of confidence.

Estimation of Parameters: Estimating various parameters of a population, such as the mean, variance, proportion, etc., using sample data.

Prediction: Estimating future values or outcomes based on historical data or trends.



- Total Variance = Variance Explained + Unexplained variance

Explained Variance
Systematic Variance
Model Variance
Predictive Variance
Treatment Variance
Regulated Variance
Deterministic Variance

Unexplained Variance
Residual Variance
Random Variance
Error Variance
Noise Variance
Unsystematic Variance
Intrinsic Variance

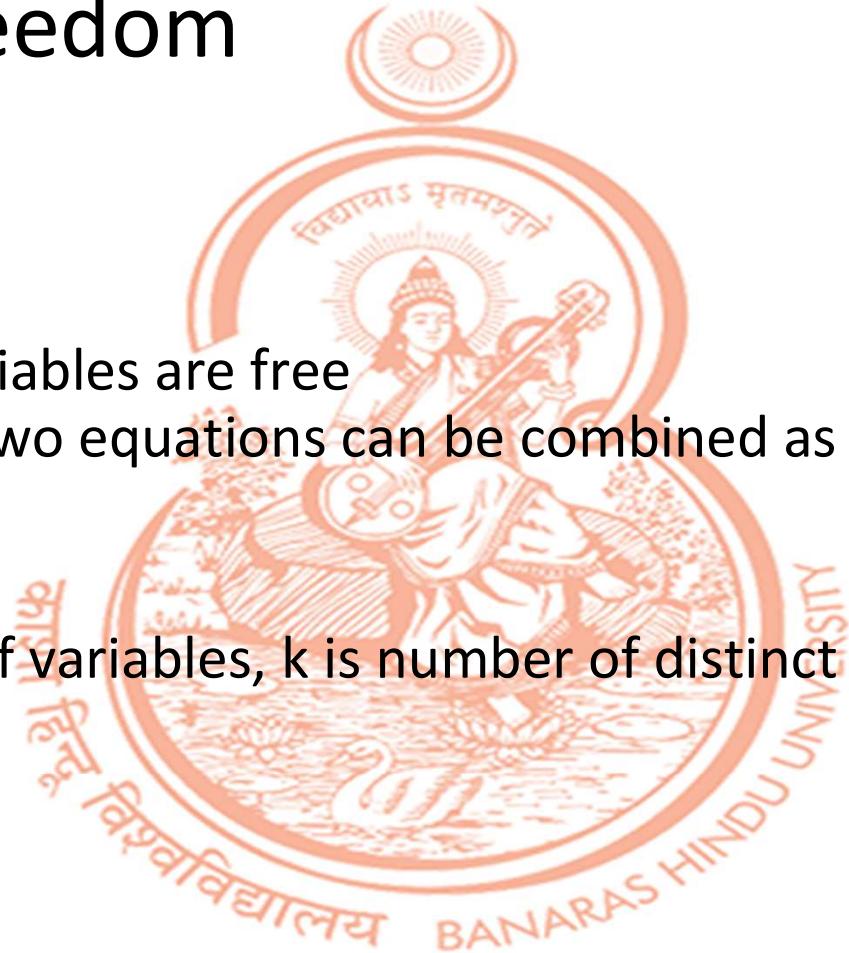


Degree of Freedom

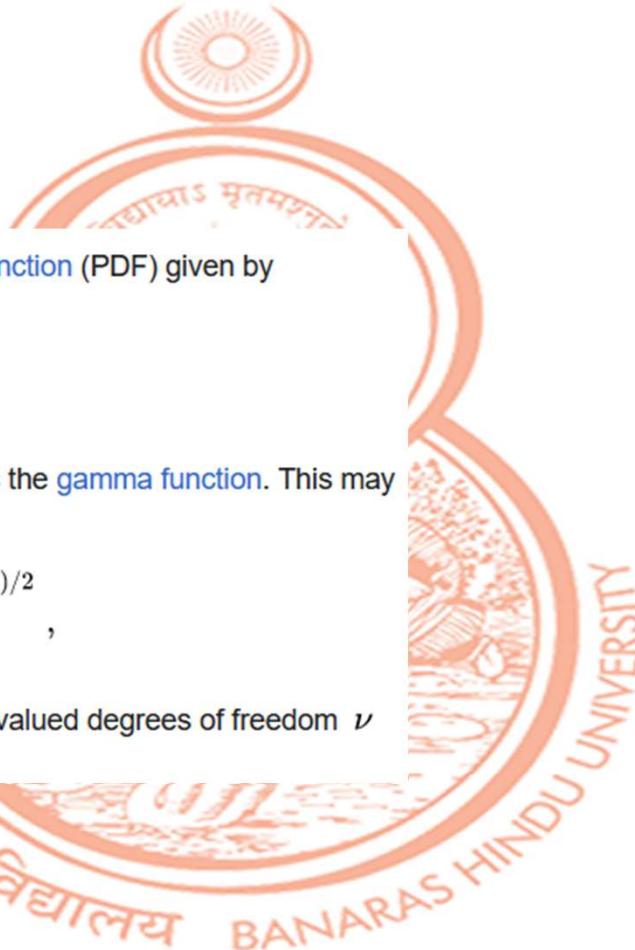
$z=3-x-y$, only two variables are free

$x+y+z=3$, $x+y-z=4 \Rightarrow$ two equations can be combined as $x+y=3.5$. Hence only one free variable.

$df = n - k$ (n - number of variables, k is number of distinct equations).



Student T test



Student's t distribution has the probability density function (PDF) given by

$$f(t) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\pi\nu} \Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{t^2}{\nu}\right)^{-(\nu+1)/2},$$

where ν is the number of *degrees of freedom* and Γ is the *gamma function*. This may also be written as

$$f(t) = \frac{1}{\sqrt{\nu} B\left(\frac{1}{2}, \frac{\nu}{2}\right)} \left(1 + \frac{t^2}{\nu}\right)^{-(\nu+1)/2},$$

where B is the *Beta function*. In particular for integer valued degrees of freedom ν we have:

Student T test



- When \bar{X} is the mean of a random sample of n from a normal distribution with mean μ , the rv

$$T = \frac{\bar{X} - \mu}{S / \sqrt{n}}$$

has a probability distribution called a t distribution with $n-1$ degrees of freedom.



Confidence Interval for T test

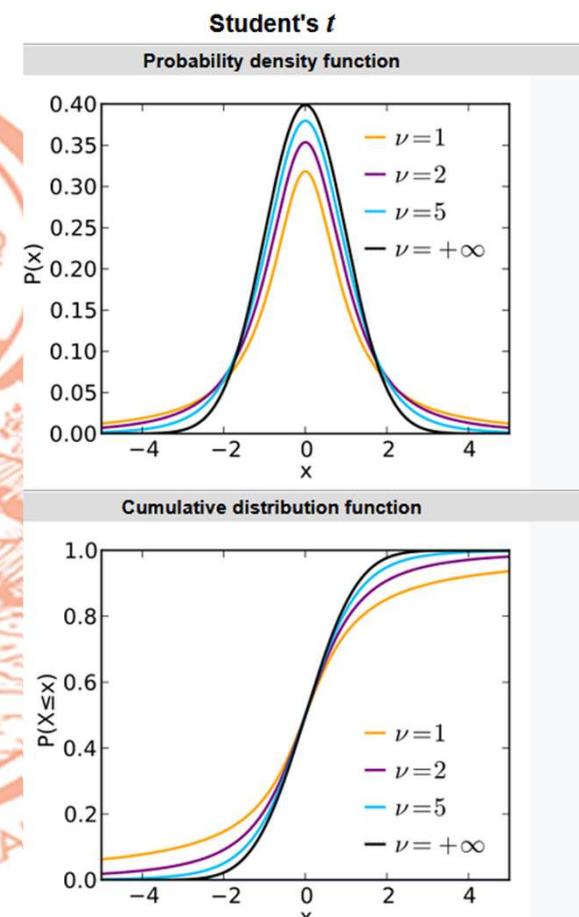
- Let \bar{x} and s be the sample mean and sample standard deviation computed from the results of a random sample from a normal population with mean μ . Then a $100(1-\alpha)\%$ confidence interval for μ is

$$\left(\bar{x} - t_{\alpha/2, n-1} \frac{s}{\sqrt{n}}, \bar{x} + t_{\alpha/2, n-1} \frac{s}{\sqrt{n}} \right)$$



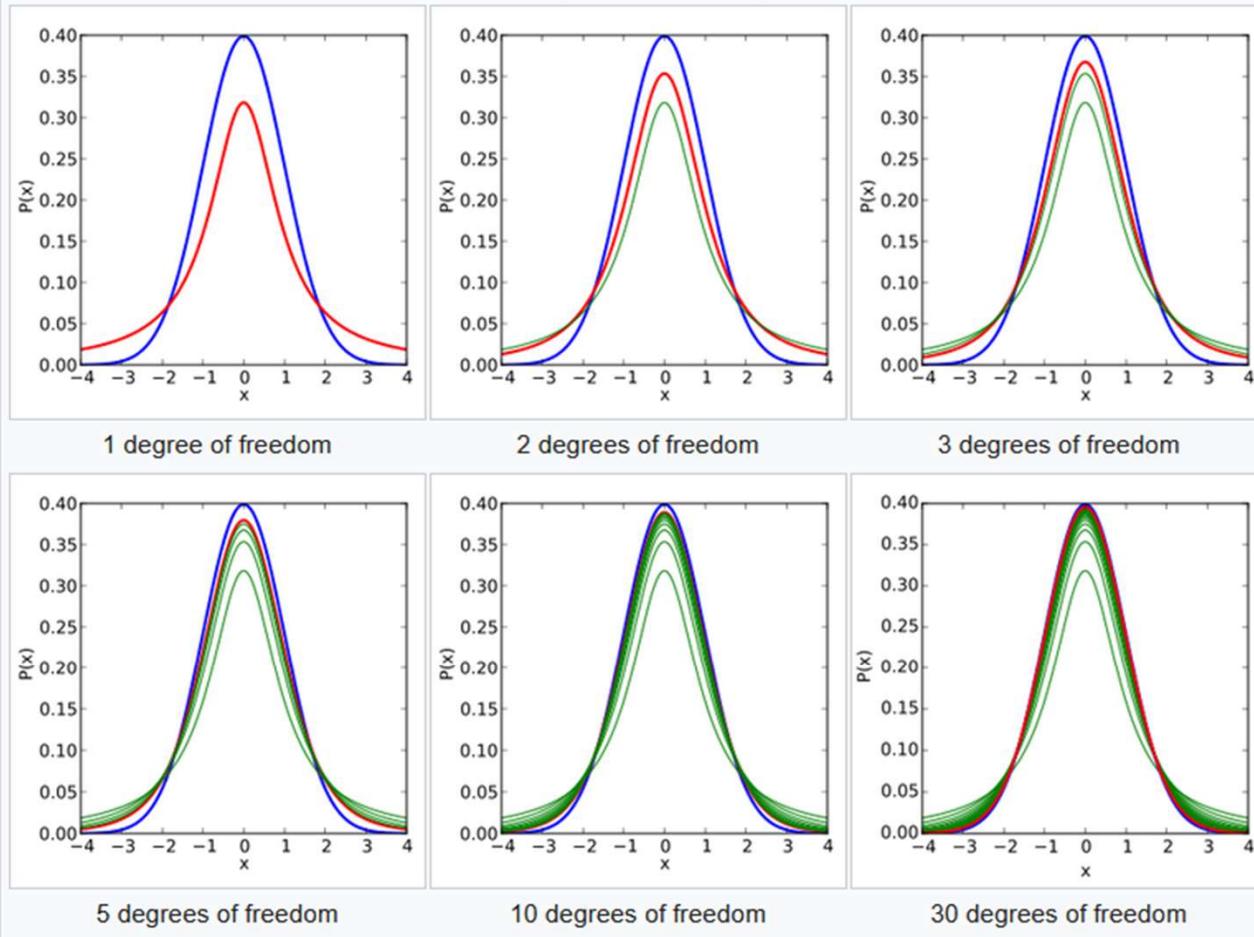
Link with Normal Distribution

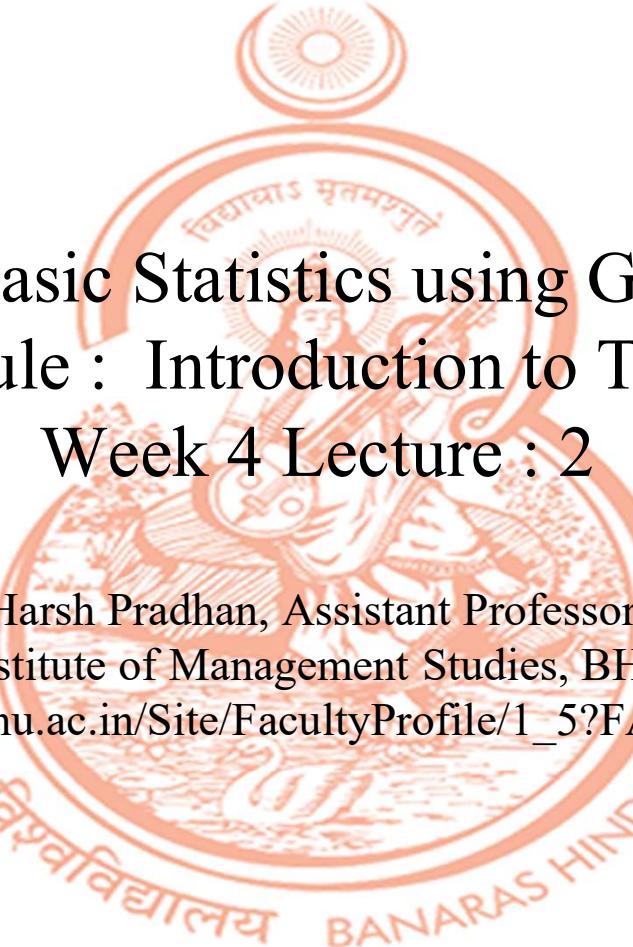
T as degree of freedom increases
student t approached normal



Density of the t distribution (red) for 1, 2, 3, 5, 10, and 30 degrees of freedom compared to the standard normal distribution (blue).

Previous plots shown in green.



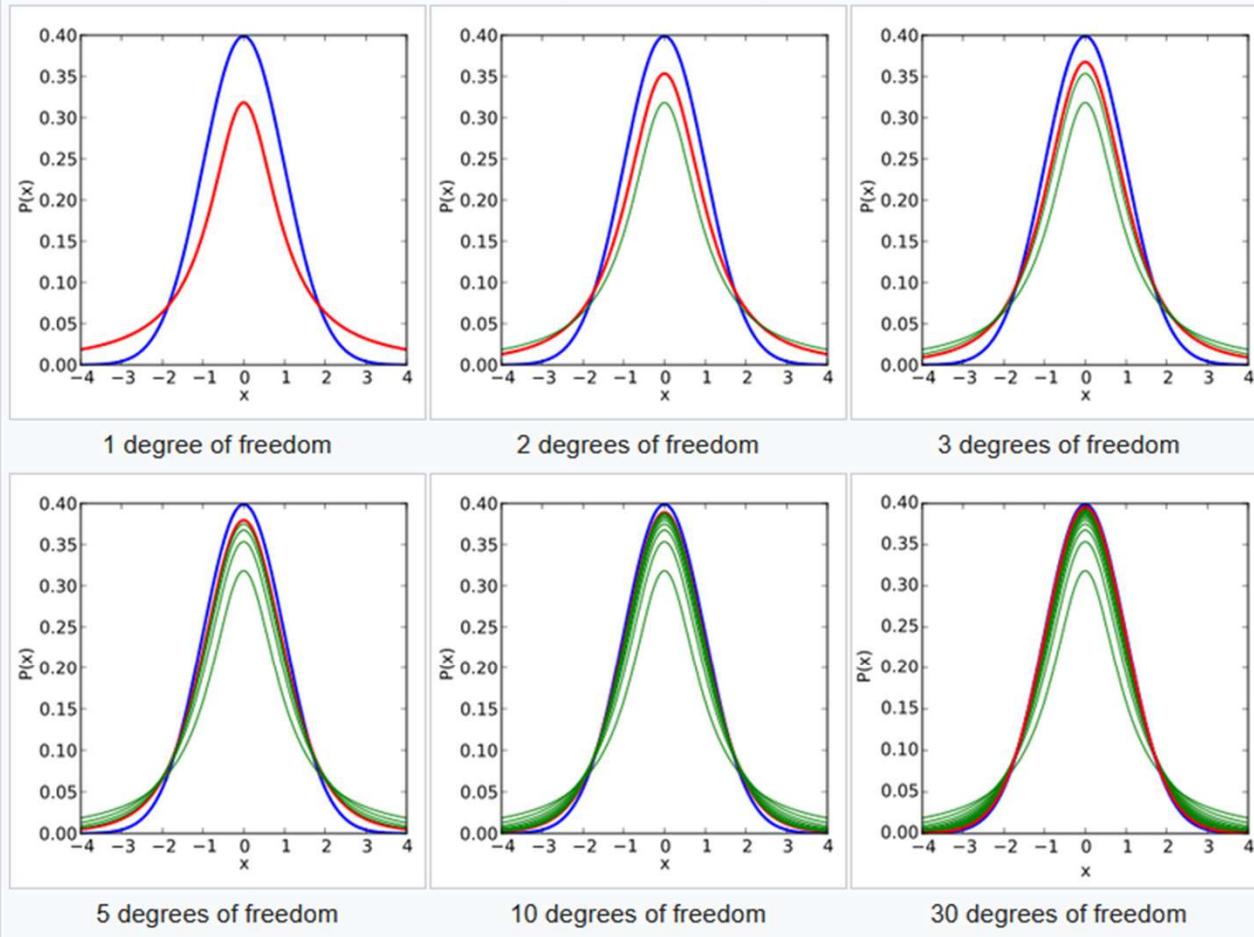


Course Name :Basic Statistics using GUI-R (RKWard)
Module : Introduction to T-Test
Week 4 Lecture : 2

Harsh Pradhan, Assistant Professor,
Institute of Management Studies, BHU
https://bhu.ac.in/Site/FacultyProfile/1_5?FA000562

Density of the t distribution (red) for 1, 2, 3, 5, 10, and 30 degrees of freedom compared to the standard normal distribution (blue).

Previous plots shown in green.



Test Statistic = explained / error

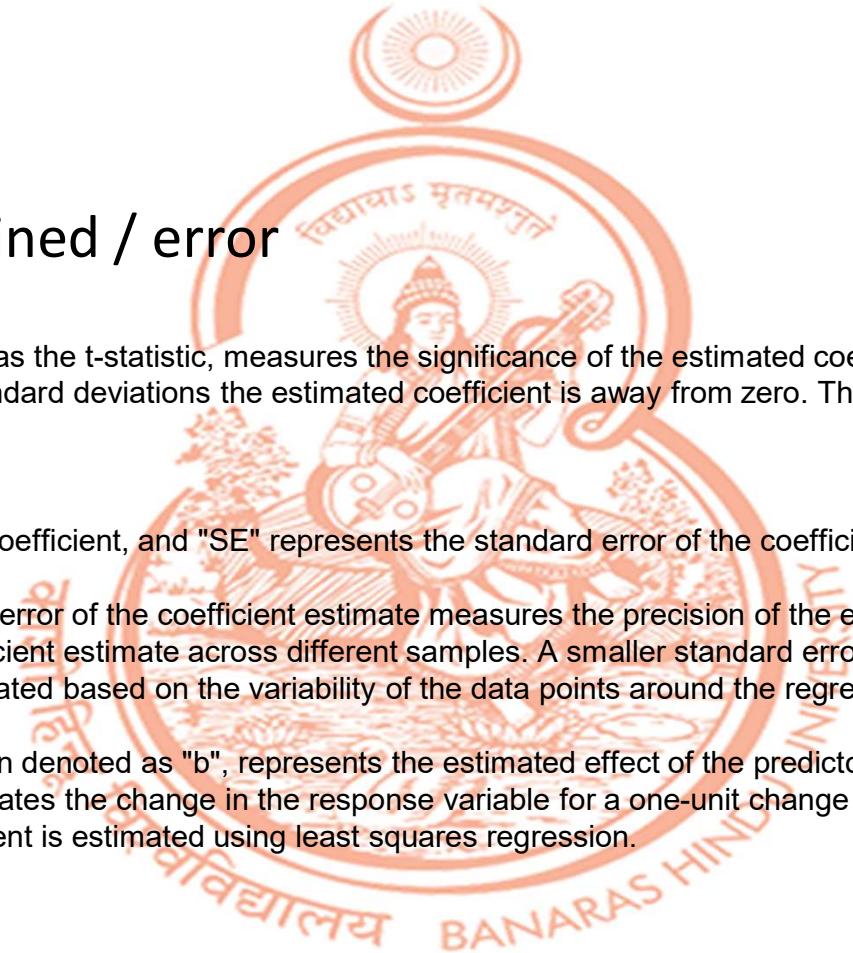
- **t-value (t):** The t-value, also known as the t-statistic, measures the significance of the estimated coefficient (b) relative to its standard error (SE). It indicates how many standard deviations the estimated coefficient is away from zero. The formula for calculating the t-value is:

$$\bullet t=b/SE$$

• Here, " b " represents the estimated coefficient, and "SE" represents the standard error of the coefficient estimate.

- **Standard Error (SE):** The standard error of the coefficient estimate measures the precision of the estimated coefficient (b). It represents the variability of the coefficient estimate across different samples. A smaller standard error indicates a more precise estimate. The standard error is calculated based on the variability of the data points around the regression line.

- **Coefficient (b):** The coefficient, often denoted as " b ", represents the estimated effect of the predictor variable on the response variable in the linear regression model. It indicates the change in the response variable for a one-unit change in the predictor variable, holding other variables constant. The coefficient is estimated using least squares regression.



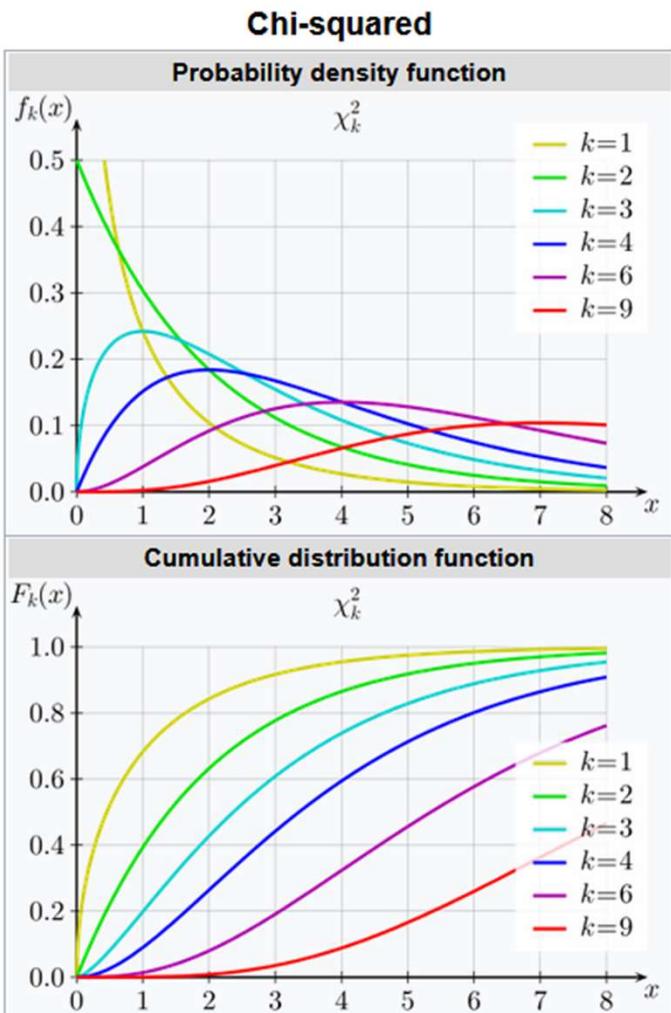
Example

- whether the mean weight of a certain breed of dogs is significantly different from the hypothesized population mean of 25 kg
22, 24, 27, 26, 28, 23, 25, 29, 21, 26, 24, 27
- Step 1: Calculate the sample mean (\bar{x}) $\bar{x} = (22 + 24 + 27 + 26 + 28 + 23 + 25 + 29 + 21 + 26 + 24 + 27) / 12 = 25.17$ kg
- Step 2: Calculate the sample standard deviation (s)
 - Subtract the mean from each data point to get the deviations: -3.17, -1.17, 1.83, 0.83, 2.83, -2.17, -0.17, 3.83, -4.17, 0.83, -1.17, 1.83
 - Square the deviations: 10.05, 1.37, 3.35, 0.69, 8.01, 4.71, 0.03, 14.67, 17.39, 0.69, 1.37, 3.35
 - Sum of squared deviations = $10.05 + 1.37 + 3.35 + 0.69 + 8.01 + 4.71 + 0.03 + 14.67 + 17.39 + 0.69 + 1.37 + 3.35 = 65.68$
 - Variance = Sum of squared deviations / $(n - 1) = 65.68 / (12 - 1) = 5.97$
 - Sample standard deviation (s) = $\sqrt{5.97} = 2.44$ kg
- Step 3: Calculate the standard error of the mean (SE) $SE = s / \sqrt{n} = 2.44 / \sqrt{12} = 0.705$ kg
- Step 4: Calculate the t-statistic Hypothesized population mean (μ_0) = 25 kg $t = (\bar{x} - \mu_0) / SE = (25.17 - 25) / 0.705 = 0.24$
- Step 5: Determine the degrees of freedom Degrees of freedom = $n - 1 = 12 - 1 = 11$

T-test

- Step 6: Determine the critical value from the t-distribution Let's assume a significance level (α) of 0.05 (two-tailed test). The critical value for the t-distribution with 11 degrees of freedom and $\alpha = 0.05$ is ± 2.201 .
- Step 7: Compare the t-statistic to the critical value The calculated t-statistic (0.24) is less than the critical value (2.201) in absolute value.
- Step 8: Make a decision Since the absolute value of the t-statistic (0.24) is less than the critical value (2.201), we fail to reject the null hypothesis. We do not have enough evidence to conclude that the mean weight of this breed of dogs is significantly different from the hypothesized population mean of 25 kg.
- In this example, we calculated the sample mean, sample standard deviation, standard error of the mean, and t-statistic. We then compared the t-statistic to the critical value from the t-distribution to determine whether the difference between the sample mean and the hypothesized population mean was statistically significant or not.

Chi Square



If Z_1, \dots, Z_k are independent, standard normal random variables, then the sum of their squares,

$$Q = \sum_{i=1}^k Z_i^2,$$

is distributed according to the chi-squared distribution with k degrees of freedom. This is usually denoted as

$$Q \sim \chi^2(k) \text{ or } Q \sim \chi_k^2.$$

The chi-squared distribution has one parameter: a positive integer k that specifies the number of degrees of freedom (the number of random variables being summed, Z_i s).

Chi Square Table

F Distribution



A random variate of the F -distribution with parameters d_1 and d_2 arises as the ratio of two appropriately scaled chi-squared variates:^[8]

$$X = \frac{U_1/d_1}{U_2/d_2}$$

where

- U_1 and U_2 have chi-squared distributions with d_1 and d_2 degrees of freedom respectively, and
- U_1 and U_2 are independent.

In instances where the F -distribution is used, for example in the analysis of variance, independence of U_1 and U_2 might be demonstrated by applying Cochran's theorem.

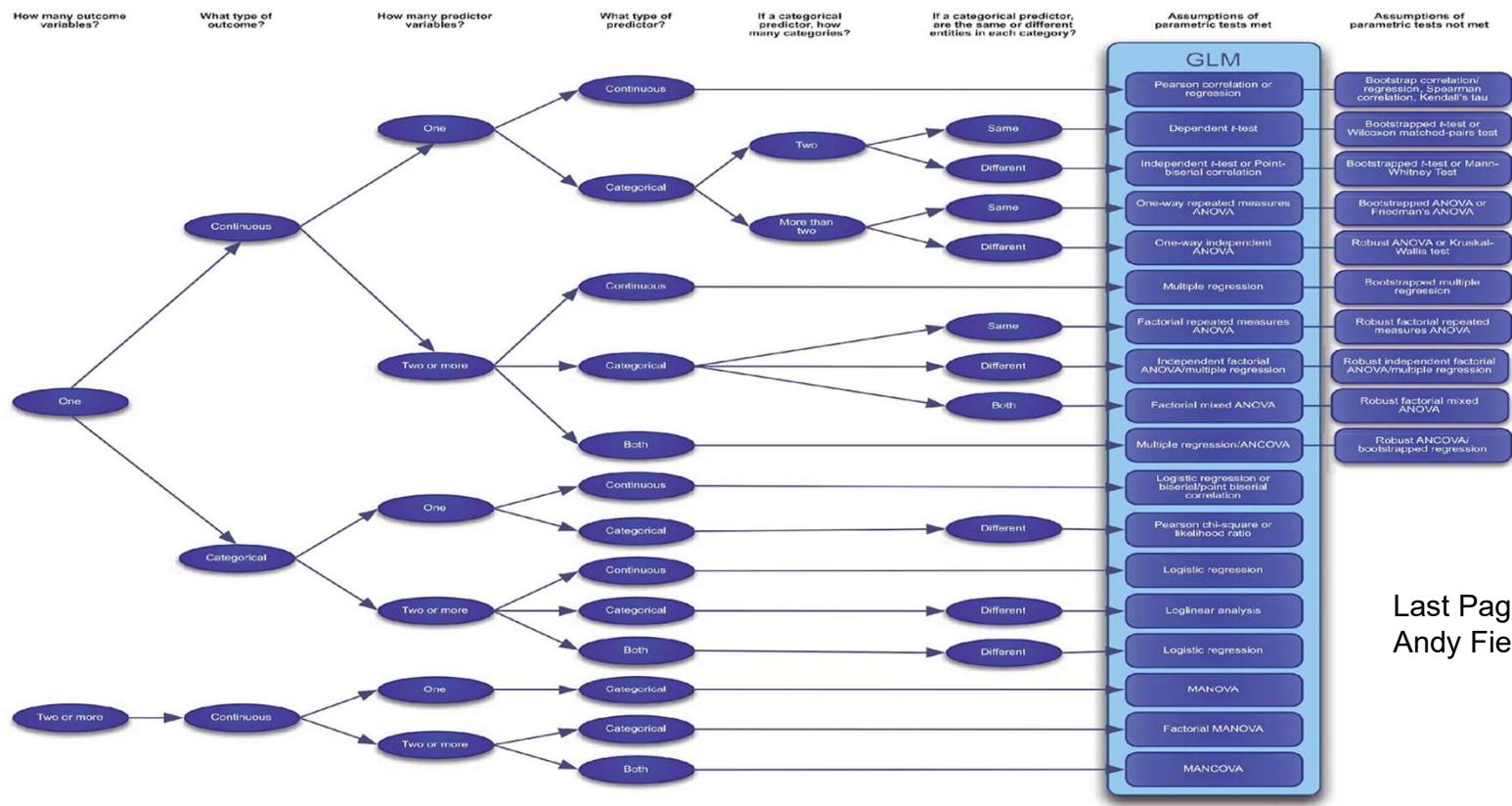
Equivalently, the random variable of the F -distribution may also be written

$$X = \frac{s_1^2}{\sigma_1^2} \div \frac{s_2^2}{\sigma_2^2},$$

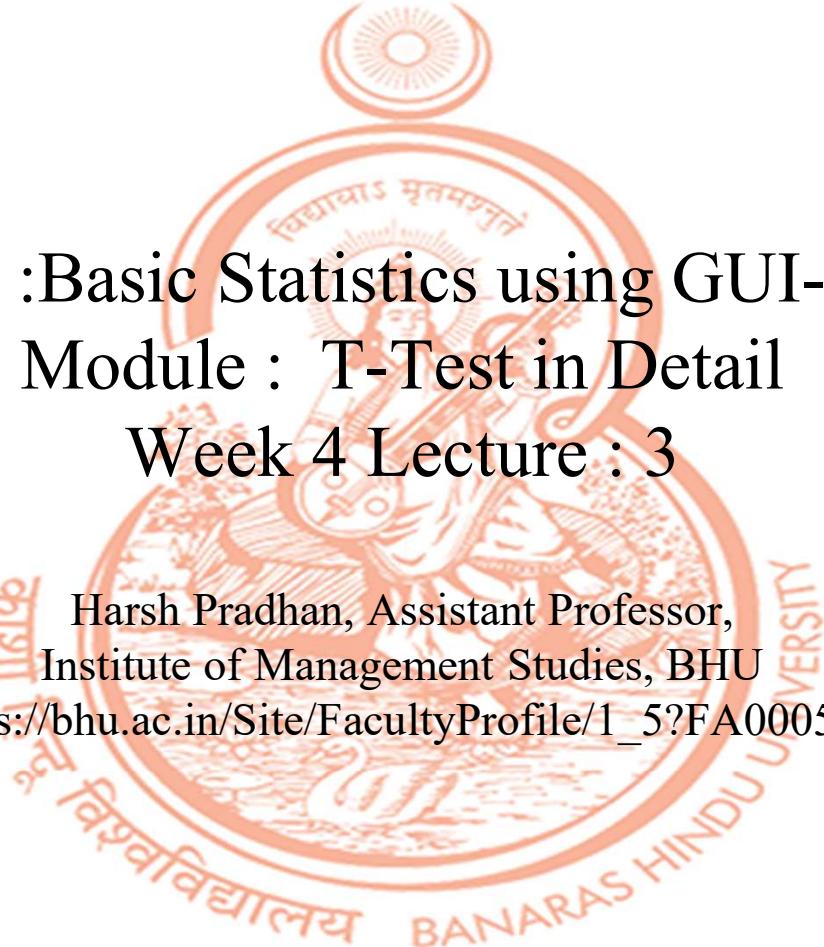
where $s_1^2 = \frac{S_1^2}{d_1}$ and $s_2^2 = \frac{S_2^2}{d_2}$, S_1^2 is the sum of squares of d_1 random variables from normal distribution $N(0, \sigma_1^2)$ and S_2^2 is the sum of squares of d_2 random variables from normal distribution $N(0, \sigma_2^2)$. [\[discuss\]](#) [\[citation needed\]](#)

F Value Table

Univariate	Bivariate	Multivariate
It only summarize single variable at a time.	It only summarize two variables	It only summarize more than 2 variables.
It does not deal with causes and relationships.	It does deal with causes and relationships and analysis is done.	It does not deal with causes and relationships and analysis is done.
It does not contain any dependent variable.	It does contain only one dependent variable.	It is similar to bivariate but it contains more than 2 variables.
The main purpose is to describe.	The main purpose is to explain.	The main purpose is to study the relationship among them.
The example of a univariate can be height.	The example of bivariate can be temperature and ice sales in summer vacation.	Example, Suppose an advertiser wants to compare the popularity of four advertisements on a website. Then their click rates could be measured for both men and women and relationships between variable can be examined

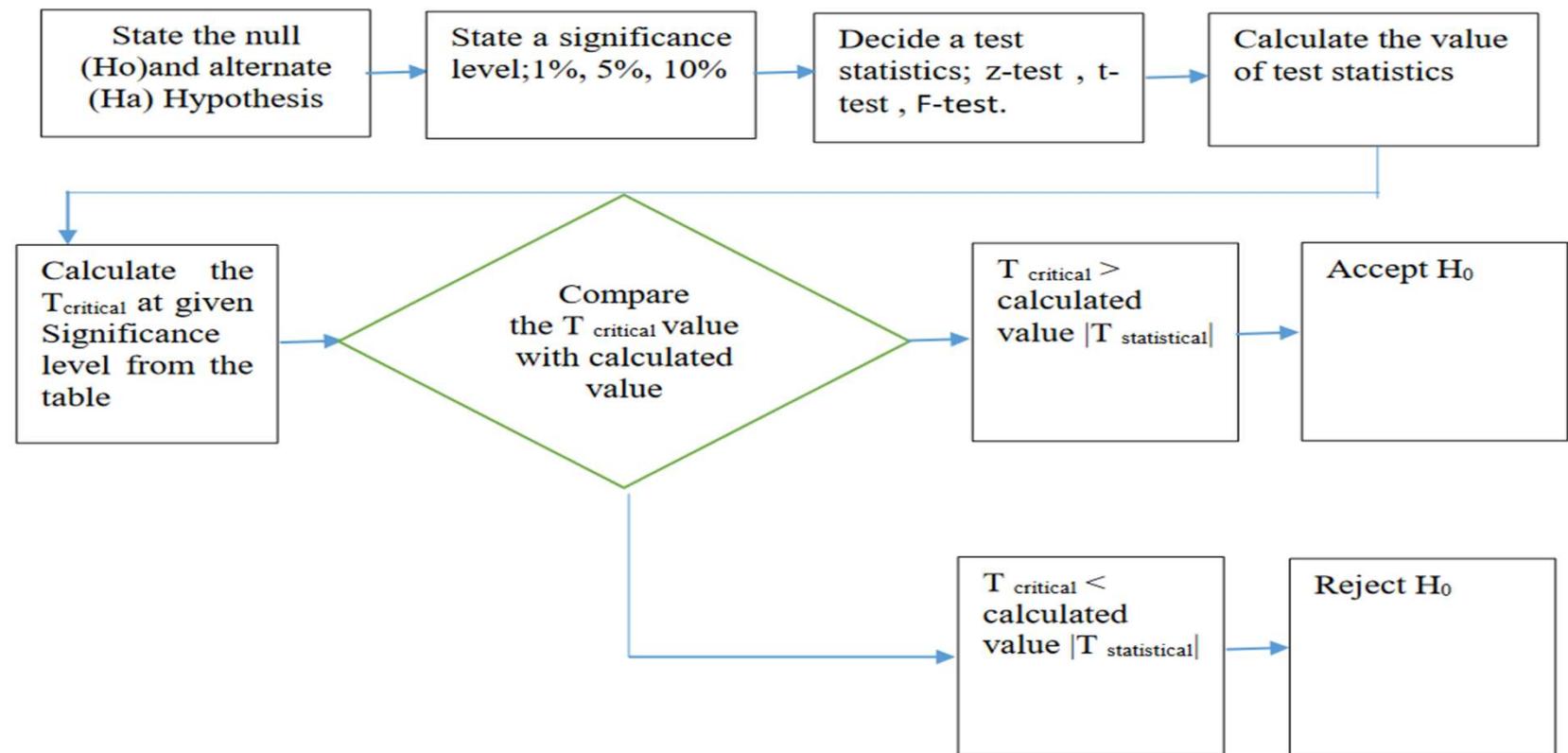


Last Page of
Andy Field



Course Name :Basic Statistics using GUI-R (RKWard)
Module : T-Test in Detail
Week 4 Lecture : 3

Harsh Pradhan, Assistant Professor,
Institute of Management Studies, BHU
https://bhu.ac.in/Site/FacultyProfile/1_5?FA000562



Example

- whether the mean weight of a certain breed of dogs is significantly different from the hypothesized population mean of 25 kg
22, 24, 27, 26, 28, 23, 25, 29, 21, 26, 24, 27
- Step 1: Calculate the sample mean (\bar{x}) $\bar{x} = (22 + 24 + 27 + 26 + 28 + 23 + 25 + 29 + 21 + 26 + 24 + 27) / 12 = 25.17$ kg
- Step 2: Calculate the sample standard deviation (s)
 - Subtract the mean from each data point to get the deviations: -3.17, -1.17, 1.83, 0.83, 2.83, -2.17, -0.17, 3.83, -4.17, 0.83, -1.17, 1.83
 - Square the deviations: 10.05, 1.37, 3.35, 0.69, 8.01, 4.71, 0.03, 14.67, 17.39, 0.69, 1.37, 3.35
 - Sum of squared deviations = $10.05 + 1.37 + 3.35 + 0.69 + 8.01 + 4.71 + 0.03 + 14.67 + 17.39 + 0.69 + 1.37 + 3.35 = 65.68$
 - Variance = Sum of squared deviations / $(n - 1) = 65.68 / (12 - 1) = 5.97$
 - Sample standard deviation (s) = $\sqrt{5.97} = 2.44$ kg
- Step 3: Calculate the standard error of the mean (SE) $SE = s / \sqrt{n} = 2.44 / \sqrt{12} = 0.705$ kg
- Step 4: Calculate the t-statistic Hypothesized population mean (μ_0) = 25 kg $t = (\bar{x} - \mu_0) / SE = (25.17 - 25) / 0.705 = 0.24$
- Step 5: Determine the degrees of freedom Degrees of freedom = $n - 1 = 12 - 1 = 11$

T-test

- Step 6: Determine the critical value from the t-distribution Let's assume a significance level (α) of 0.05 (two-tailed test). The critical value for the t-distribution with 11 degrees of freedom and $\alpha = 0.05$ is ± 2.201 .
- Step 7: Compare the t-statistic to the critical value The calculated t-statistic (0.24) is less than the critical value (2.201) in absolute value.
- Step 8: Make a decision Since the absolute value of the t-statistic (0.24) is less than the critical value (2.201), we fail to reject the null hypothesis. We do not have enough evidence to conclude that the mean weight of this breed of dogs is significantly different from the hypothesized population mean of 25 kg.
- In this example, we calculated the sample mean, sample standard deviation, standard error of the mean, and t-statistic. We then compared the t-statistic to the critical value from the t-distribution to determine whether the difference between the sample mean and the hypothesized population mean was statistically significant or not.

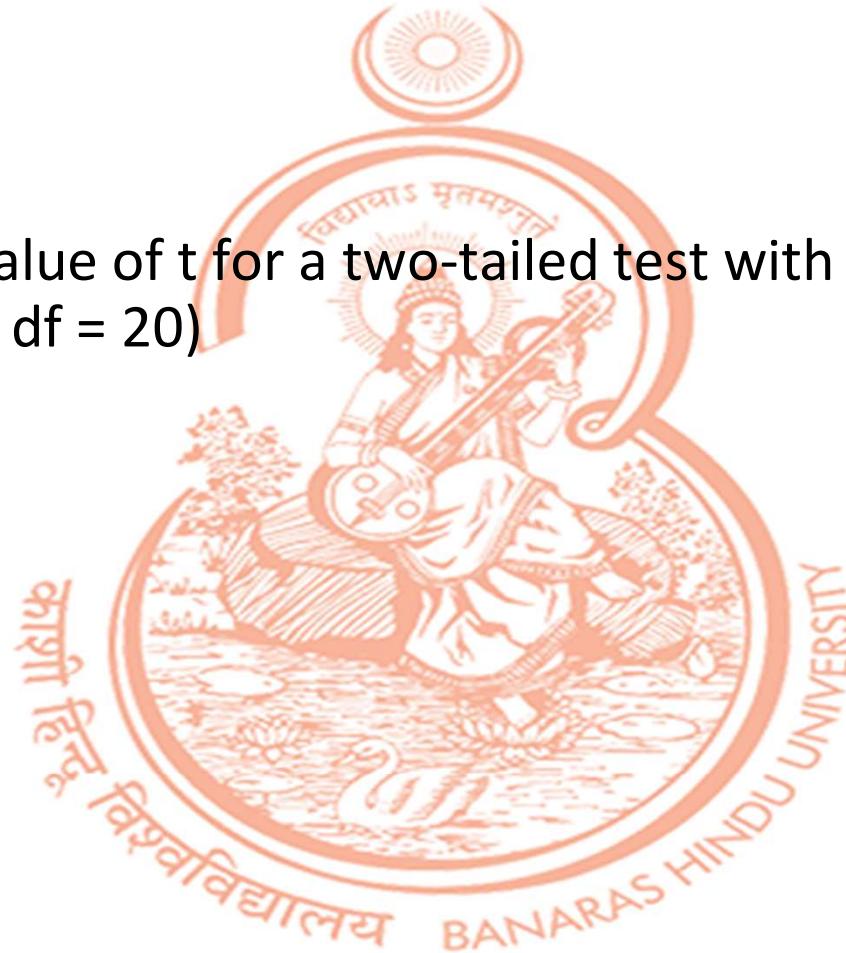
Confidence Interval for T test

- Let \bar{x} and s be the sample mean and sample standard deviation computed from the results of a random sample from a normal population with mean μ . Then a $100(1-\alpha)\%$ confidence interval for μ is

$$\left(\bar{x} - t_{\alpha/2, n-1} \frac{s}{\sqrt{n}}, \bar{x} + t_{\alpha/2, n-1} \frac{s}{\sqrt{n}} \right)$$



- # Find critical value of t for a two-tailed test with alpha = 0.05 and df = 20 , $qt(0.025, df = 20)$





	One-sample t-test	Two-sample t-test	Paired t-test
Synonyms	Student's t-test	Independent groups t-test	Paired groups t-test
		Independent samples t-test	Dependent samples t-test
		Equal variances t-test	
		Pooled t-test	
		Unequal variances t-test	
Number of variables	One	Two	Two
Type of variable	Continuous measurement	Continuous measurement	Continuous measurement
		Categorical or Nominal to define groups	Categorical or Nominal to define pairing within group
Purpose of test	Decide if the population mean is equal to a specific value or not	Decide if the population means for two different groups are equal or not	Decide if the difference between paired measurements for a population is zero or not
Example: test if...	Mean heart rate of a group of people is equal to 65 or not	Mean heart rates for two groups of people are the same or not	Mean difference in heart rate for a group of people before and after exercise is zero or not
Degrees of freedom	Number of observations in sample minus 1, or: $n-1$	Sum of observations in each sample minus 2, or: $n_1 + n_2 - 2$	Number of paired observations in sample minus 1, or: $n-1$

Two Sample Test

Aspect	Independent Samples t-test	Paired t-test
Nature of Data	Data consists of two separate, unrelated groups	Data consists of paired observations
Dependence/Independence	Observations in each group are independent of each other	Observations within pairs are dependent/correlated
Assumption	Population distributions are normally distributed, and variances are equal	Differences between paired observations are normally distributed
Hypotheses	H ₀ : No difference between population means H _a : Significant difference between population means	H ₀ : No difference between related groups H _a : Significant difference between related groups
Formula	Based on means, standard deviations, and sample sizes	Based on mean difference and standard error of the mean difference
Example	Comparing two groups receiving different treatments	Comparing before-and-after measurements

Assumptions for Parametric Tests

Dependent variable is a scale variable interval or ratio

- If the dependent variable is ordinal or nominal, it is a non-parametric test

Participants are randomly selected

- If there is no randomization, it is a non-parametric test

The shape of the population of interest is approximately normal

- If the shape is not normal, it is a non-parametric test

Effect

- Having a statistically significant difference between means does not mean that the difference is large!
- To measure the size of the difference, we calculate effect sizes: the difference between the two means divided by the standard deviation.

It is affected by

- the distance between means of two distributions (larger distance means larger effect size)
- the standard deviations of the two distributions (a smaller standard deviation means a larger effect size)

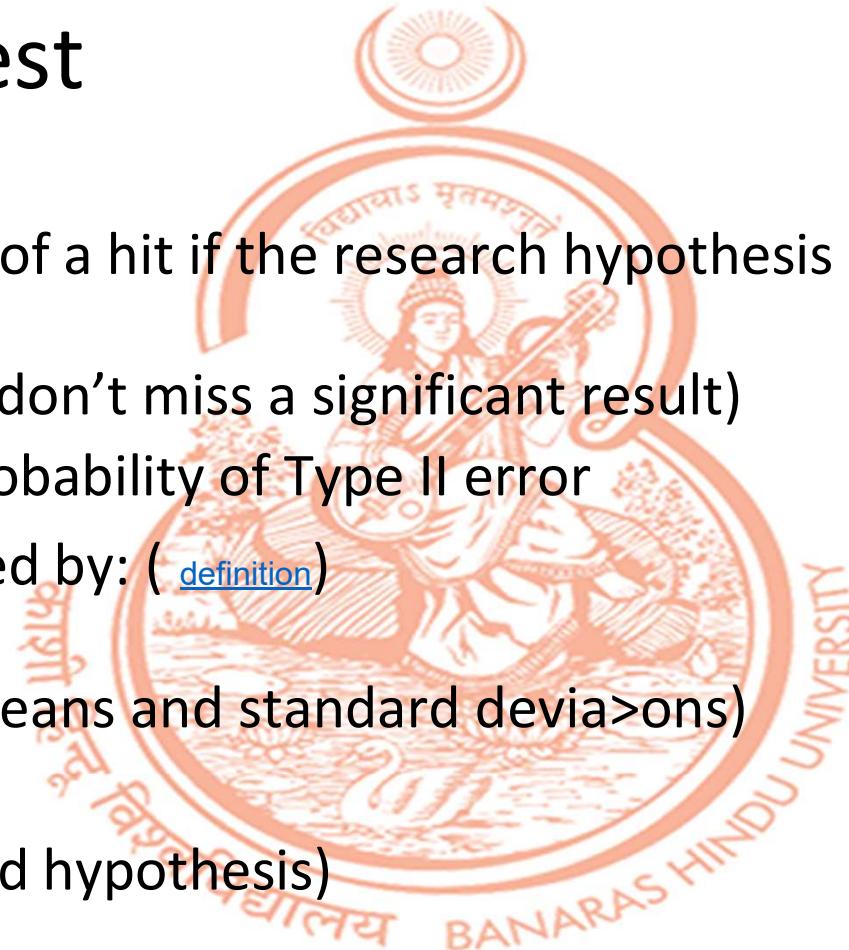
Unaffected by sample size

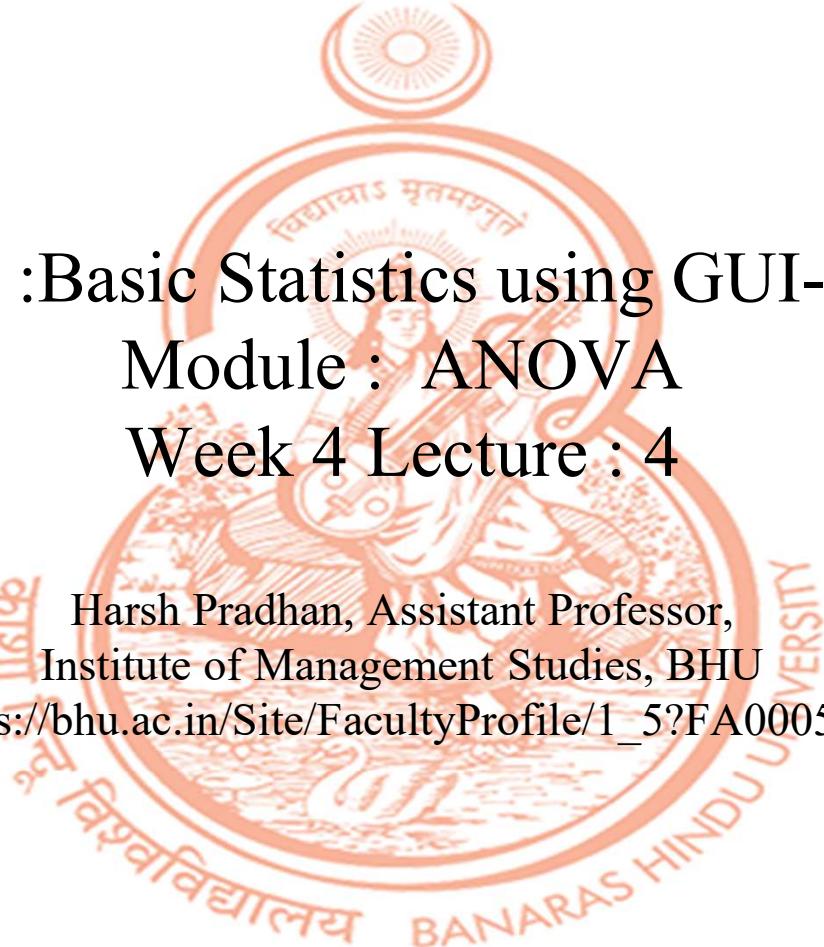
- If we use standardised effect sizes, we can compare studies using different measurement scales (Meta-analyses)

Power of Test

- The probability of a hit if the research hypothesis is true (i.e. that we don't miss a significant result)
Power = $100 - \text{probability of Type II error}$
- Power is affected by:
 - ([definition](#)) sample size
 - effect size (i.e. means and standard deviations)
 - (alpha)
 - (one or two-tailed hypothesis)

[p-hacking 2:38](#)



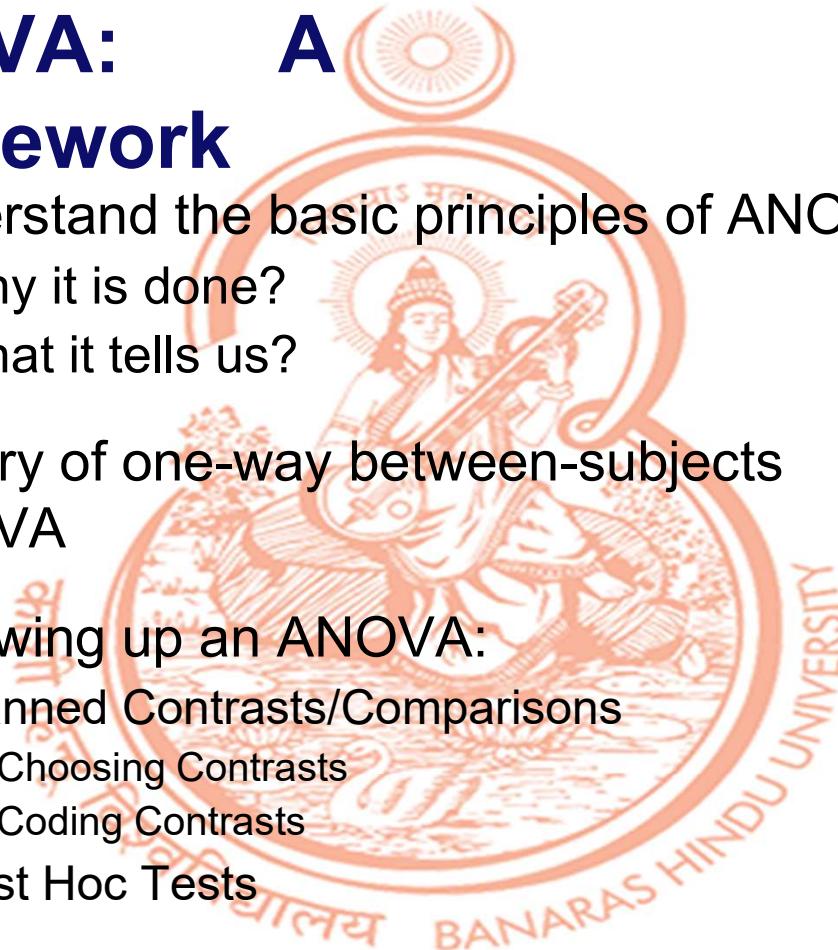


Course Name :Basic Statistics using GUI-R (RKWard)
Module : ANOVA
Week 4 Lecture : 4

Harsh Pradhan, Assistant Professor,
Institute of Management Studies, BHU
https://bhu.ac.in/Site/FacultyProfile/1_5?FA000562

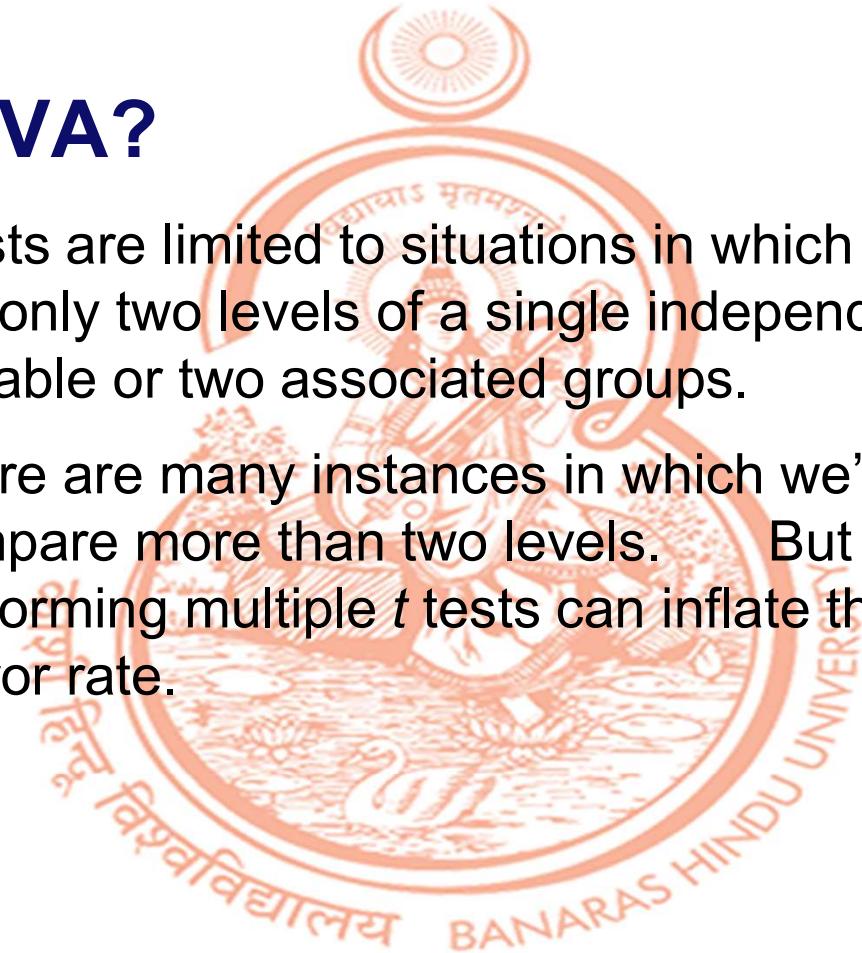
ANOVA: A Framework

- Understand the basic principles of ANOVA
 - Why it is done?
 - What it tells us?
- Theory of one-way between-subjects ANOVA
- Following up an ANOVA:
 - Planned Contrasts/Comparisons
 - Choosing Contrasts
 - Coding Contrasts
 - Post Hoc Tests
- Writing up Results



Why ANOVA?

- t tests are limited to situations in which there are only two levels of a single independent variable or two associated groups.
- There are many instances in which we'd like to compare more than two levels. But ... performing multiple t tests can inflate the Type I error rate.



One-Way ANOVA

- The one-way analysis of variance is used to test the null hypothesis that three or more population means are equal
 - more precisely: test the null hypothesis that the means of the groups are not significantly different from the grand mean of all participants

One-Way ANOVA

- The *response variable* is the variable you're comparing, i.e., dependent variable
- The *factor variable* is the categorical variable being used to define the groups, i.e., independent variable
 - Usually called k samples (groups)
- The *one-way* is because there is one independent variable

Assumptions



- Scale dependent variable
- Normal distribution
- Equal variances

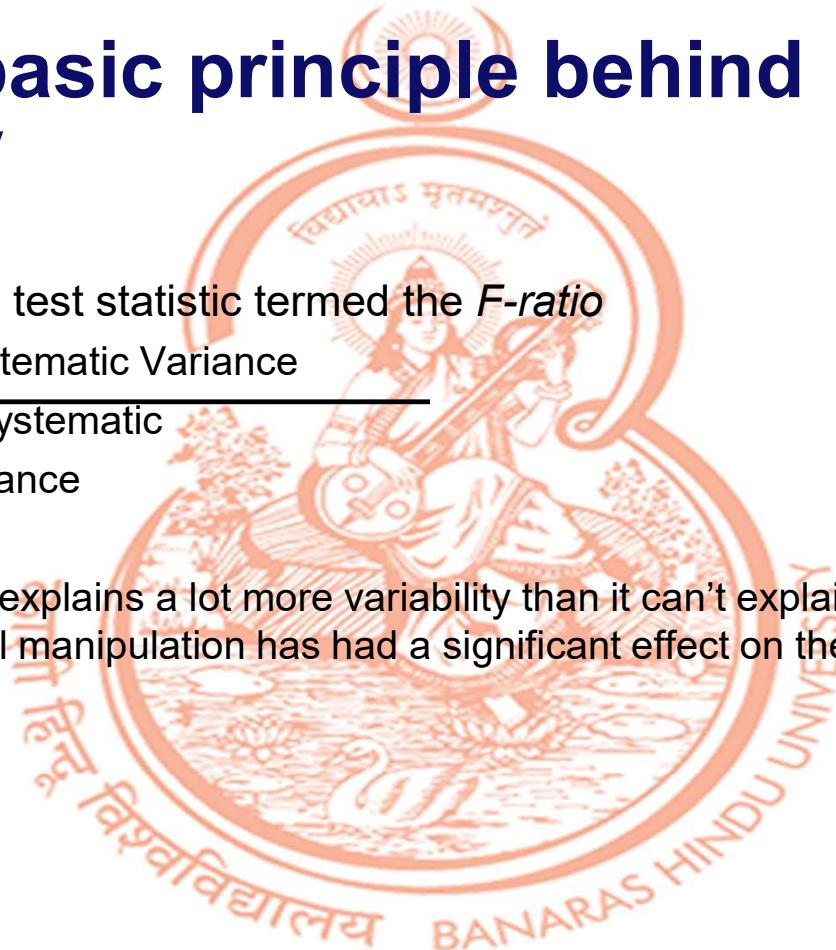
The basic principle behind ANOV

A

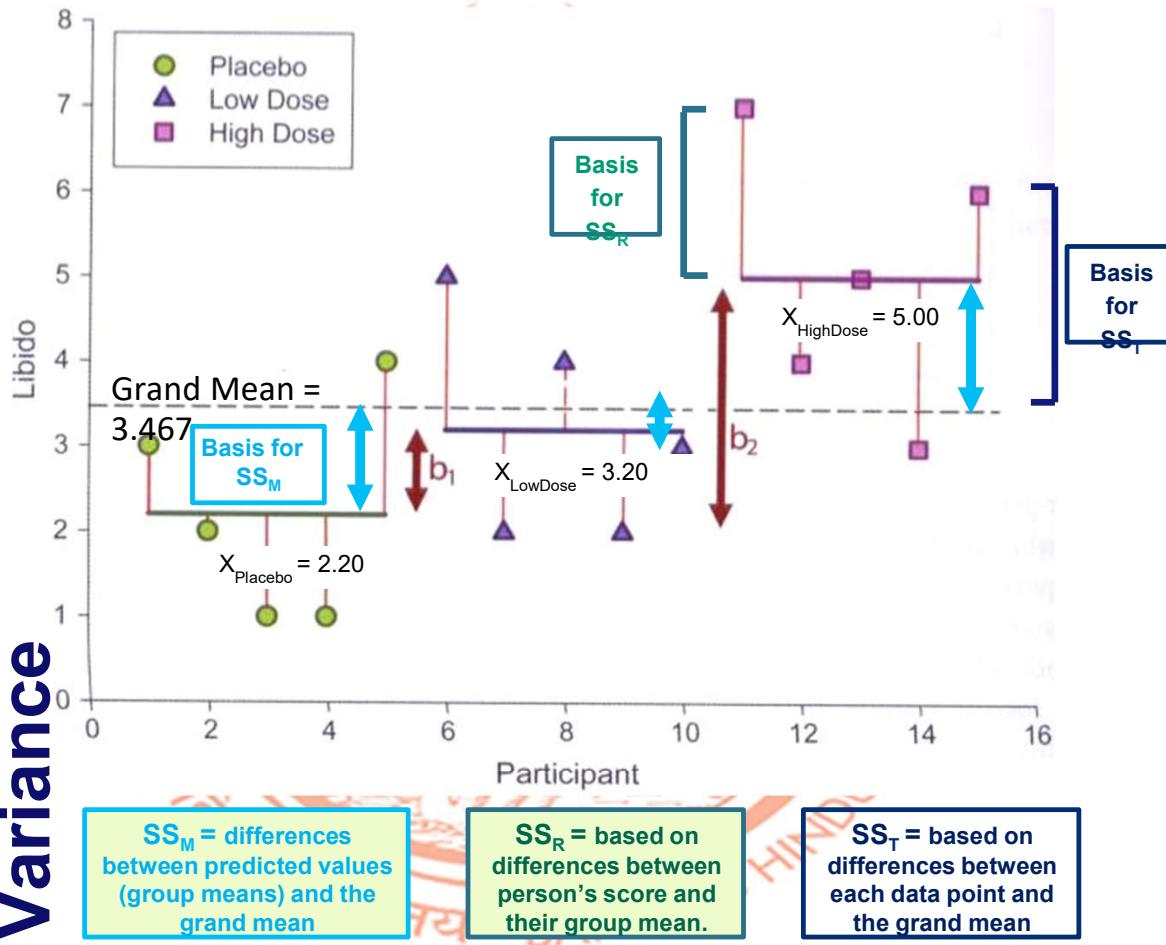
Produces a test statistic termed the *F-ratio*

$$\frac{\text{Systematic Variance}}{\text{Unsystematic Variance}}$$

If the model explains a lot more variability than it can't explain, then the experimental manipulation has had a significant effect on the outcome (DV).



Partitioning the Variance



Partitioning the variance



- SS_M = differences between predicted values (group means) and the grand mean (i.e., variation across groups) –
if there is a significant difference between the groups, this should be a large number relative to SS_R
- SS_R = based on differences between person's score and their group mean (i.e., variation with each group)
- SS_T = based on differences between each data point and the grand mean (the total variation in the entire data set)

Mean Squares (MS_M and MS_R)

SS_M = amount of variation explained by the model (exp. manipulation).

SS_R = amount of variation due to extraneous factors.

These are “summed” scores and will therefore be influenced by the number of scores. To eliminate this bias we calculate the average sum of squares (mean squares) by dividing by the appropriate degrees of freedom.

Calculating Degrees of Freedom

(for one-way independent groups ANOVA)

$$df_{\text{total}} = N - 1 \text{ (number of all scores minus 1)}$$

$$df_{M / \text{between}} = k - 1 \text{ (number of groups minus 1)}$$

$$df_{R / \text{within}} = N - k \text{ (number of all scores minus number of groups)}$$

The F-ratio

- We compare the amount of variability explained by the Model (MS_M), to the error in the model [individual differences] (MS_R)
 - This ratio is called the *F*-ratio
- If the model explains a lot more variability than it can't explain, then the experimental manipulation has had a significant effect on the outcome (DV).

$$F \equiv \frac{MS_M}{MS_R}$$

One-Way ANOVA

- Here is the basic one-way ANOVA table

Source	SS	df	MS	F	p
Between					
Within					
Total					

An example: Fairness in different types of societies

Fairness score: proportion of money shared in a game

	Hunter-gatherer	Farming	Natural resources	Industrial
P1	28	32	47	40
P2	36	33	43	47
P3	38	40	52	45
P4	31			
Mean	33.25	35.0	47.33	44.0
<i>N</i>	4	3	3	3

Grand Mean = 39.385 (The sum of all scores divided by the total N)

Total SS

The sum of the squared deviation of each score from the grand mean



$$SS_T = \sum (x_i - M_{grand})^2$$



Error (within-group) sum of squares (SS_R)

- Calculate the squared deviation of each score from its group mean
- Add the results



$$SS_R = \sum (x_i - M_i)^2$$



Model (between-group) sum of squares (SS_M)

1. Calculate the difference between the mean of each group and the grand mean.
The grand mean is the mean of all scores
2. Square each of these differences
3. Multiply each result by the number of participants within that group – this is a correction (or “weighting”): a smaller sample will have less “weight” in the equation, a larger sample will have more “weight”.
4. Add the values for each group together.

$$SS_M = \sum n_i (M_i - M_{grand})^2$$

One-Way ANOVA

- After filling in the sum of squares, we have ...

Source	SS	df	MS	F	p
Between	461.64				
Within	167.42				
Total	629.08				

Degrees of freedom

- The between group df is one less than the number of groups, $k - 1$
 - We have four groups, so $df_M = 3$
- The within group df is the sum of the individual df's of each group, which equals $N - k$
 - The sample sizes are 4, 3, 3, and 3
 - $df_R = 3 + 2 + 2 + 2 = 13 - 4 = 9$
- The total df is one less than the sample size, $N - 1$
 - $df(\text{Total}) = 13 - 1 = 12 = (n+n+n) - 1$