

GLM Concept

- Fundamental questions to consider when trying to understand the relationship between two variables: Are the two variables related?, What is the direction of the relationship between two variables?, How strong is the relationship between two variables?
- Statistical models:

$$\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3$$

Y: dependent (outcome) variable  
 $\hat{Y}$ : predicted score of Y  
Y-  $\hat{Y}$ : error

- Measuring the error of the model:
  - The basic measure of a model’s error is the **Sum of Squared Errors (SSE)** where you square each individual’s error then add up all the squared errors
- Model comparison:** make incremental changes then see if the more complex model predicts Y sufficiently better, and if it does then we conclude that adding that increased complexity (adding the new predictor X) benefited the model
- General Linear modeling/general multivariate regression model:** a compact way of writing several multiple linear regression models in one linear model
  - \*Used when the data is normally distributed\*
  - Examples: t-test, regression, ANOVA

Data = Model + Error

$$Y = \beta_0 + \varepsilon$$

Descriptive Statistics

		Notation	Definition	Formula	Excel Formula/Steps
Central Tendency	Mean	$\bar{Y}$	The average of the numbers	$\frac{\sum x}{n} = \frac{Sum}{sample\ size}$	=AVERAGE(B2:B6)
	Median	-	The middle score in a set of data arranged numerically	-	=MEDIAN(B2:B6)
	Mode	-	The most frequent score in a data set	-	=MODE(B2:B6)
Dispersion	Range	-	The difference between the max & min score	-	=MIN(B2:B6) =MAX(B2:B6)
	Standard Deviation	S	A measure of how dispersed the data is in relation to the mean	$(x_i - \bar{x})^2$ = the difference between the individual value and the mean  $S = \sqrt{Variance}$	1. Calculate the <b>Mean</b> 2. Calculate the <b>Deviation</b> (=Data cell – Mean cell) 3. Calculate the <b>Squared Deviation</b> (=Deviation cell^2) 4. Calculate the <b>Sum of Squares (SS)</b> - sq. dev. (=SUM(F2:F6)) 5. Input the <b>Degrees of Freedom</b> (=n-1) 6. Calculate the <b>Variance</b> (=SS/df) 7. Calculate the <b>Standard Deviation</b> (=SQRT(F10)) 8. Calculate the Standard Error (=Std. Dev./(SQRT(n)))
	Variance	$\sigma^2 / S^2$	A measure of how far each number in a dataset is from the mean	$\frac{\sum(x_i-\bar{x})^2}{n-1} = \frac{square\ of\ deviation}{degree\ of\ freedom}$  $\frac{SSE}{df}$	

Sum of Squared Errors (SSE or SSEA): basic measure of a model’s error where you square each individual’s error and add them all up; Quantifies error as an area; allows us to calculate variance & standard deviation

Conceptual Formula	Computational Formula (preferred)
$SSE = \sum(Y_i - \hat{Y}_i)^2$	$SSE = \sum Y^2 - \frac{(\sum Y)^2}{N}$
1. Calculate each ind. <b>Deviation</b> (=Data cell -Mean cell) 2. Calculate each ind. <b>Squared Deviation</b> (=Deviation cell ^2) 3. <b>Sum</b> the Squared Deviations (=SUM(F2:F6))	1. Calculate each individual <b>Y^2</b> (=Data cell^2) 2. <b>Sum</b> all the Y^2 values 3. <b>Sum</b> all the Y values 4. Add values from Steps 1-3 into formula

t-Tests

Hypothesis Testing	
Definition	A method for testing a hypothesis about a parameter or population, using data measured in a sample
Purpose	To test claims or ideas about a group or population
Goal	To determine the likelihood that a population parameter, such as the mean, is likely to be true

Steps for Hypothesis Testing	
Step 1: State the null & alternative	<ul style="list-style-type: none"><li>Null Hypothesis (H<sub>0</sub>) - statement about the population parameter that is <u>assumed to be true</u><ul style="list-style-type: none"><li>Example: Children in the US watch an average of 3 hrs. of TV a week (H<sub>0</sub>: μ = 3)</li></ul></li><li>Alternative Hypothesis (H<sub>1</sub>) – statement that <u>contradicts the null</u> hypothesis<ul style="list-style-type: none"><li>Example: Children in the US watch more or less than 3 hrs. of TV a week (H<sub>1</sub>: μ ≠ 3)</li></ul></li><li>Possible Hypotheses to Test<ul style="list-style-type: none"><li>Nondirectional or Two-tailed:<ul style="list-style-type: none"><li>H<sub>0</sub>: μ = 0</li><li>H<sub>1</sub>: μ ≠ 0</li></ul></li><li>Directional or One-tailed:<ul style="list-style-type: none"><li>H<sub>0</sub>: μ = 0</li><li>H<sub>1</sub>: μ &gt; 0 OR H<sub>1</sub>: μ &lt; 0</li></ul></li></ul></li></ul>
Step 2: Set criteria for the decision/level of significance	<ul style="list-style-type: none"><li>Criterion of judgement upon which a decision is made regarding the value stated in the null hypothesis</li><li>Typically set at 5% (.05)</li></ul>
Step 3: Compute the test statistic (t-score)	<div><math display="block">t_{obt} = \frac{M - \mu}{s_M}</math><ul style="list-style-type: none"><li>M: population mean</li><li>μ: mean of means (the number you are suspecting that your mean is different from in your hypotheses)</li><li>s<sub>M</sub>: standard error (standard deviation of the mean)</li></ul></div> <ul style="list-style-type: none"><li>Calculating t<sub>obt</sub> in Excel<ul style="list-style-type: none"><li>Sum the Y values</li><li>Calculate each individual Y<sup>2</sup> (=Data cell^2)</li><li>Sum all the Y<sup>2</sup> values</li><li>Calculate the Mean</li><li>Calculate the SSE</li><li>Calculate the Variance (=SSE cell/ df (N-1))</li><li>Calculate the Std. Dev. (=SQRT(variance cell))</li><li>Calculate the Std. Error (= std. dev. Cell/SQRT(N))</li><li>Calculate t<sub>obt</sub> =(mean – hypothesis value)/std. error)</li></ul></li><li>Calculate the 95% CI (M ± (t<sub>crit</sub> &amp; s<sub>M</sub> )<ul style="list-style-type: none"><li>M: sample mean</li><li>T<sub>crit</sub>: use t Table<ul style="list-style-type: none"><li>Need to know: df (N-1), alpha level/level of significance, &amp; 1-tailed or 2-tailed</li></ul></li><li>s<sub>M</sub>: standard error (std. dev./SQRT(n))</li><li>In Excel<ul style="list-style-type: none"><li>Calculate 95% CI LB (=mean-(t<sub>crit</sub>*std. error))</li><li>Calculate 95% CI UB (=mean+( t<sub>crit</sub>*std. error))</li></ul></li></ul></li></ul>
Step 4: Make a decision	<ul style="list-style-type: none"><li>Compare the t<sub>obt</sub> &amp; t<sub>crit</sub> values<ul style="list-style-type: none"><li>If t<sub>obt</sub> &gt; t<sub>crit</sub> = reject the null (the sample mean is associated with a low probability of occurrence; <b>results are significant</b>)</li><li>If t<sub>obt</sub> &lt; t<sub>crit</sub> = retain/fail to reject the null (the sample mean is associated with a high probability of occurrence; <b>results aren't significant</b>)</li></ul></li><li>Reporting/Interpretation<ul style="list-style-type: none"><li>Example: According to a one sample t-test, the average soda consumption of the sample is/is not statistically different from 2 cans per day (t(df) = x.xx)</li></ul></li></ul>

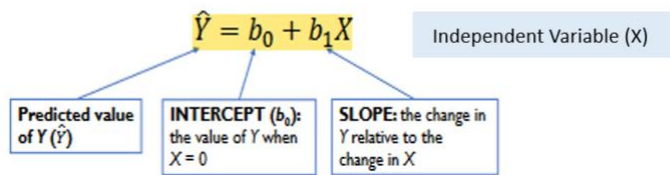
Bivariate Regression

- Simple model (mean-only model) isn't very informative since it has a lot of error – can add a predictor variable (X) to improve prediction & reduce error (provides an equation of a line that best predicts Y from X)
- Bivariate regression has two types of variables: Predictor (X; known) and Outcome/Criterion (Y; to-be-predicted)
- A population parameter (β), estimate of population parameter using sample data (b)
- Role of Error

Model	Source	Label	Error
$Y_i = \beta_0 + \varepsilon_i$	Model C	Total error	$SST (SSE_C)$
$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$	Model A	Residual error	$SSE_A$
	Model C- Model A	Explained error	$SSR$

\*SST is almost always greater than  $SSE_A$

$SST = SSE_A + SSR$

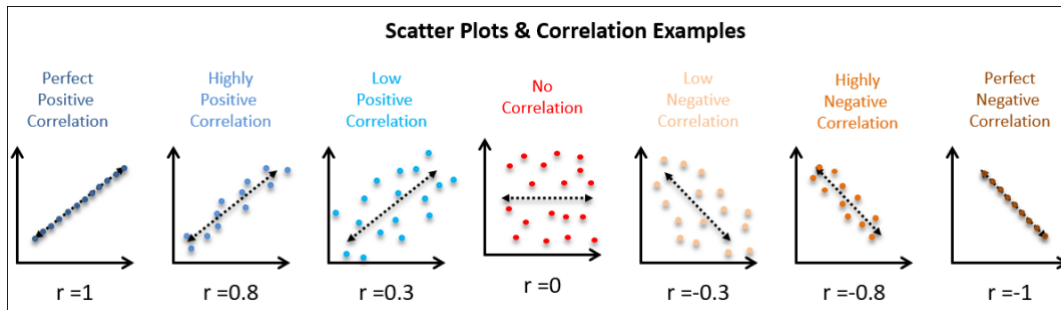


10 Steps to Model Comparison	
Step 1: State the Augmented & Compact Models	<ul style="list-style-type: none"> <li>Need to identify the DV &amp; IV</li> <li>Compact Model: <math>DV = b_0 + \varepsilon</math></li> <li>Augmented Model: <math>DV = b_0 + b_1 IV + \varepsilon</math></li> </ul>
Step 2: Identify the Null Hypothesis	<ul style="list-style-type: none"> <li>Testing to see if <math>b_1</math> is statistically different from 0               <ul style="list-style-type: none"> <li><math>H_0: b_1 = 0</math></li> <li><math>H_1: b_1 \neq 0</math></li> </ul> </li> </ul>
Step 3: Count the Number of Parameters Estimated by Each Model	<ul style="list-style-type: none"> <li>Compact Model (PC = 1)               <ul style="list-style-type: none"> <li>Parameter: <math>b_0</math></li> </ul> </li> <li>Augmented Model (PA = 2)               <ul style="list-style-type: none"> <li>parameters: <math>b_0</math> and <math>b_1</math></li> </ul> </li> </ul>
Step 4: Calculate the Regression Equation	<ul style="list-style-type: none"> <li>Bivariate Regression Coefficient: <math>Y = b_0 + b_1 X + \varepsilon</math></li> <li><math>b_1</math> <ul style="list-style-type: none"> <li>Conceptual Formula: <math>\frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2} = \frac{\text{sum of (x deviation * y deviation)}}{\text{sum of squared x deviation}}</math></li> <li>Computational Formula: <math display="block">\frac{\sum XY - \frac{(\sum X)(\sum Y)}{n}}{\sum X^2 - \frac{(\sum x)^2}{n}}</math></li> </ul> </li> <li><math>b_0</math>: <math>b_0 = \bar{Y} - b_1 \bar{X}</math></li> <li>Calculating in Excel (NOTE: order data as Y then X)               <ul style="list-style-type: none"> <li><b>Conceptual Formula (preferred)</b> <ul style="list-style-type: none"> <li>Calculate the <b>Mean</b> for X &amp; Y</li> <li>Calculate <b>Y_DEV</b></li> <li>Calculate <b>X_DEV</b></li> <li>Calculate <b>X_DEV*Y_DEV</b></li> <li>Calculate the <b>Sum</b> for X_DEV*Y_DEV</li> <li>Calculate <b>X_DEV^2</b></li> <li>Calculate the <b>Sum</b> for X_DEV^2</li> <li>Calculate <b>b1</b> (=Sum for X_DEV*Y_DEV/Sum for X_DEV^2)</li> <li>Calculate <b>b0</b> (=Y Mean-(b1*X Mean))</li> </ul> </li> <li>Computational Formula                   <ul style="list-style-type: none"> <li>Calculate <b>Sum</b> for X &amp; Y</li> <li>Calculate <b>X*Y</b></li> <li>Calculate the <b>Sum</b> for X*Y</li> <li>Calculate <b>X^2</b></li> <li>Calculate the <b>Sum</b> for X^2</li> <li>Calculate <b>b1</b> (=Sum of X*Y – (Sum of X &amp; Sum of Y/# of Cases))/(Sum of X^2 – (Sum of X^2/ # of Cases))</li> <li>Calculate <b>b0</b> (=Y Mean-(b1*X Mean))</li> </ul> </li> </ul> </li> <li>Assumptions that need to be checked – see <i>SPSS Output Slides</i></li> </ul>
Step 5: Compute $SSE_T$	<ul style="list-style-type: none"> <li>Also SST; represents Total Error</li> <li>Comes from the compact model: <math>SSE_T = \sum(Y_i - \bar{Y})^2</math></li> <li>Calculating in Excel (NOTES: order data as Y then X; will need to move over b1 and b0 values or calculate them)               <ul style="list-style-type: none"> <li>Calculate <b>Mean</b> of Y</li> <li>Calculate <b>Y_DEV</b></li> <li>Calculate <b>Y_DEV^2</b></li> <li>Calculate <b>SST</b> (=SUM for Y_DEV^2 data column)</li> </ul> </li> </ul>
Step 6: Compute $SSE_A$	<ul style="list-style-type: none"> <li>Residual Error; Error for the Augmented Model</li> <li><math>SSE_A = \sum(Y_i - \hat{Y}_i)^2</math> (<math>\hat{Y}_i</math> is the score predicted by the regression equation)</li> <li>Calculating in Excel               <ul style="list-style-type: none"> <li>Calculate <b>Y_HAT or Y Predicted</b> (=b0+(b1*X value))</li> <li>Calculate <b>DEV</b> (=Y Value – Y Predicted Value)</li> <li>Calculate <b>DEV^2</b></li> <li>Calculate <b>SSE_A</b> (=SUM for DEV^2 data column)</li> </ul> </li> </ul>
Step 7: Compute SSR	<ul style="list-style-type: none"> <li>Explained Error; Reduced Sum of Squares</li> <li>Calculating in Excel               <ul style="list-style-type: none"> <li>Calculate <b>SSR</b> (=SSE_T-SSE_A)</li> </ul> </li> </ul>

Step 8: Compute PRE/R <sup>2</sup>	<ul style="list-style-type: none"><li>Proportion of Reduced Error Variance; Percentage of Reduced Sum of Squares of Error</li><li>Referred to as the “variance explained”</li><li>Calculating in Excel<ul style="list-style-type: none"><li>Calculate <b>R<sup>2</sup></b> (=SSR/SST)</li></ul></li></ul>																								
Step 9: Complete Summary Table	<table><tr><th>Source</th><th>SS</th><th>df</th><th>MS</th><th>F</th><th>R<sup>2</sup></th></tr><tr><td>Model Comparison</td><td>SSR</td><td>PA - PC</td><td><math>MS_{model} = \frac{SSR}{PA - PC}</math></td><td><math>F = \frac{MS_{model}}{MS_{residual}}</math></td><td><math>R^2 = \frac{SSR}{SSE_T}</math></td></tr><tr><td>Residual</td><td>SSE<sub>A</sub></td><td>N - PA</td><td><math>MS_{residual} = \frac{SSE_A}{N - PA}</math></td><td></td><td></td></tr><tr><td>Total</td><td>SSE<sub>T</sub></td><td>N - PC</td><td></td><td></td><td></td></tr></table>	Source	SS	df	MS	F	R <sup>2</sup>	Model Comparison	SSR	PA - PC	$MS_{model} = \frac{SSR}{PA - PC}$	$F = \frac{MS_{model}}{MS_{residual}}$	$R^2 = \frac{SSR}{SSE_T}$	Residual	SSE <sub>A</sub>	N - PA	$MS_{residual} = \frac{SSE_A}{N - PA}$			Total	SSE <sub>T</sub>	N - PC			
Source	SS	df	MS	F	R <sup>2</sup>																				
Model Comparison	SSR	PA - PC	$MS_{model} = \frac{SSR}{PA - PC}$	$F = \frac{MS_{model}}{MS_{residual}}$	$R^2 = \frac{SSR}{SSE_T}$																				
Residual	SSE <sub>A</sub>	N - PA	$MS_{residual} = \frac{SSE_A}{N - PA}$																						
Total	SSE <sub>T</sub>	N - PC																							
Step 10: Make Decision Regarding H <sub>0</sub>	<ul style="list-style-type: none"><li>Look at the <i>F</i> statistic (calculated) to make a decision<ul style="list-style-type: none"><li><i>F</i> critical value can be obtained from <i>F</i> table</li></ul></li><li>When you have <i>F</i><sub>obt</sub> and <i>F</i><sub>crit</sub><ul style="list-style-type: none"><li>If <i>F</i><sub>obt</sub> &lt; <i>F</i><sub>crit</sub>, we do not reject H<sub>0</sub></li><li>If <i>F</i><sub>obt</sub> &gt; <i>F</i><sub>crit</sub>, we reject H<sub>0</sub></li></ul></li><li>If you have SPSS output, you are given a <i>p</i> value – compare to alpha level of .05<ul style="list-style-type: none"><li>If <i>p</i> &gt; .05, we do not reject H<sub>0</sub></li><li>If <i>p</i> &lt; .05, we reject H<sub>0</sub></li></ul></li></ul>																								

## Correlation

- Statistical procedure used to describe the strength and direction of the linear relationship between 2 variables (degree to which two variables are associated)
  - Values range from -1, to 1, with both of those values insinuating a perfect correlation (rare). Values of *r* = 0 doesn't mean the variables aren't related, it just means they aren't LINEARLY related.
  - Correlation doesn't equal causation
  - Key Descriptors



<b>Perfectly Correlated</b>
• Either -1 or 1
<b>Strongly Correlated</b>
• $r >  .75 $
<b>Moderately Correlated</b>
• $ .5  < r <  .75 $
<b>Weakly Correlated</b>
• $ .25  < r <  .5 $
<b>Uncorrelated</b>
• $0 < r <  .25 $
<b>Sign</b>
• Positive or Negative

- Significance

Two-tailed Testing	One-tailed Testing
H <sub>0</sub> : <i>r</i> = 0; there is no association between the two variables H <sub>1</sub> : <i>r</i> ≠ 0; there is an association between the two variables	H <sub>0</sub> : <i>r</i> = 0; there is no association between the two variables H <sub>1</sub> : <i>r</i> < 0 or <i>r</i> > 0; there is a positive/negative association between the two variables

- Pearson's *r* (Most common, 2 continuous variables)

Calculating in Excel (Note: Will have data for X & Y)	Interpretation
<ul style="list-style-type: none"> <li>Calculate the <b>Mean</b> for X &amp; Y</li> <li>Calculate <b>X_DEV</b></li> <li>Calculate <b>X_DEV^2</b></li> <li>Calculate the <b>Sum</b> for X_DEV^2</li> <li>Calculate <b>Y_DEV</b></li> <li>Calculate the <b>Y_DEV^2</b></li> <li>Calculate the <b>Sum</b> for Y_DEV^2</li> <li>Calculate <b>X_DEV*Y_DEV</b></li> <li>Calculate the <b>Sum</b> for X_DEV*Y_DEV</li> <li>Calculate <b>Pearson's r</b> (=Sum for X_DEV*Y_DEV/SQRT(X_DEV^2*Y_DEV^2))</li> <li>Calculate <b>R^2</b> (=Pearson's r^2)</li> </ul>	<ul style="list-style-type: none"> <li>Obtain Critical Value <ul style="list-style-type: none"> <li>Need to know: df (<i>n</i>-2), level of significance, &amp; one- or two-tailed</li> </ul> </li> <li>Compare Obtained to Critical Values <ul style="list-style-type: none"> <li>Reject the Null (statistically significant) <ul style="list-style-type: none"> <li>Obtained &gt; Critical</li> </ul> </li> <li>Retain the Null (not significant) <ul style="list-style-type: none"> <li>Obtained &lt; Critical</li> </ul> </li> </ul> </li> <li>Write Up <ul style="list-style-type: none"> <li>Example: Using the Pearson correlation coefficient, there is a statistically significant relationship between X and Y, <i>r</i> = -.897, <i>p</i> &lt; .05.</li> </ul> </li> </ul>

- Point Biserial (1 continuous & 1 categorical variable) – see *SPSS Output Slides*
- Spearman's *p* (rho) (have ordinal variables)
  - Relationship must be monotonic (as 1 variable increases, the other either increases or decreases)
  - Can't handle normal data, not sensitive to outliers
- Kendall's tau-b

- Alternative to Spearman's when sample size is small – therefore, assumptions are the same

## One-Way ANOVA

- **Analysis of Variance (ANOVA):** omnibus test used to determine whether at least 1 of 3+ group means are different from one another
  - DV (Y) is continuous
  - IV (X) is categorical (the grouping variable)
- Equivalent to a t-test but can have more levels (groups) tested at once
- If the means of k different groups are all the same, then the variance between them is 0
- Hypothesis
  - $H_0: \mu_1 = \mu_2 \dots$  [ $\mu$  represents the population mean for each group]
  - $H_1$ : At least one mean is different
- F statistic: Ratio of Variances (=Between/Within)
  - Should be a positive value.
  - $F = 1$  means no mean difference because the between and within group variabilities would be the same. The higher your F the greater chance of it being significant (would mean there is a much higher variability between groups than you see within groups)
  - Calculating in Excel

Mean	$SS_{BG}$	$SS_{WG}$	SST	Complete Summary Table
<ul style="list-style-type: none"> <li>• Calculate the <b>Mean</b> for each group</li> <li>• Calculate the <b>Grand Mean</b> (includes all subjects regardless of group)</li> </ul>	<ul style="list-style-type: none"> <li>• Calculate the <b>DEV.</b></li> <li>• Calculate the <b>SQ. DEV.</b></li> <li>• Calculate the <b>Sum</b> for SQ. DEV.</li> <li>• Calculate <b>SSR (Between)</b> <math>(= (G1 \cdot n * (G1 \text{ mean} - \text{grand mean})^2) + (G2 \cdot n * (G2 \text{ mean} - \text{grand mean})^2)</math></li> </ul>	<ul style="list-style-type: none"> <li>• Calculate <b>SSE (Within)</b> <math>(= G1 \text{ Sum for SQ. DEV.} + G2 \text{ Sum for SQ. DEV.})</math></li> </ul>	<ul style="list-style-type: none"> <li>• Calculate <b>SST</b> <math>(= SSR + SSE)</math></li> </ul>	<ul style="list-style-type: none"> <li>• df Column               <ul style="list-style-type: none"> <li>○ BG: # of groups - 1</li> <li>○ WG: Total # of subjects - # of groups</li> <li>○ Total: Total # of subjects - 1</li> </ul> </li> <li>• MS column               <ul style="list-style-type: none"> <li>○ BG: <math>SS_{BG} - BG \text{ df}</math></li> <li>○ WG: <math>SS_{WG} - WG \text{ df}</math></li> </ul> </li> <li>• F Statistic               <ul style="list-style-type: none"> <li>○ <math>MS_{BG} / MS_{WG}</math></li> </ul> </li> <li>• eta-squared (<math>\eta^2</math>)               <ul style="list-style-type: none"> <li>○ <math>SS_{BG} / SST</math></li> <li>○ Small Effect: <math>.01 &lt; \eta^2 &lt; .06</math></li> <li>○ Medium Effect: <math>.06 &lt; \eta^2 &lt; .14</math></li> <li>○ Large Effect: <math>.14 &lt; \eta^2</math></li> </ul> </li> </ul>

Source	SS	df	MS	F
Between Groups	$SS_{BG}$	$k - 1$	$MS_{BG} = \frac{SS_{BG}}{k - 1}$	$F = \frac{MS_{BG}}{MS_{WG}}$
Within Groups	$SS_{WG}$	$N - k$	$MS_{WG} = \frac{SS_{WG}}{N - k}$	
Total	SST	$N - 1$		

**\*\*Where N is total number of subjects and k is number of groups.**

- Interpretation
  - Obtain F crit value (from F table)
    - df between – horizontal
    - df within - vertical
  - If F value is GREATER than F crit then we can say there is a significant difference between the means of the groups (reject the null)

## Miscellaneous

- If a test result is **NOT SIGNIFICANT**, then the alpha level should be written in the following way in a reporting/write up:  $p > .05$ 
  - If result is significant:  $p < .05$