Theory-Based Methods 0000 Multiple Linear Regression 000000000

Inference for Multiple Linear Regression

Nate Wells

Math 141, 4/20/22

Nate Wells

Conditions for Inference

Theory-Based Methods 0000 Multiple Linear Regression

Outline

In this lecture, we will...

- Review conditions for inference in simple linear models
- Create confidence intervals for parameters of linear models
- Discuss theory-based methods for regression
- Review framework for multilinear regression
- Discuss inference procedures for MLR models

Theory-Based Methods 0000 Multiple Linear Regression 000000000

Section 1

Confidence Intervals

Nate Wells

Confidence Intervals ○●○○○○	Conditions for Inference	Theory-Based Methods 0000	Multiple Linear Regression
Reed Thesis			

• Earlier this year, Math 141 students collected data on several hundred senior theses from thesis tower.

Confidence Intervals ○●○○○○	Conditions for Inference	Theory-Based Methods 0000	Multiple Linear Regression

Reed Thesis

- Earlier this year, Math 141 students collected data on several hundred senior theses from thesis tower.
 - Page Count and Year Published for several of these theses are shown below:



Reed Theses

Theory-Based Methods 0000 Multiple Linear Regression

Confidence Intervals for Linear Models

• A hypothesis test allows us to assess the strength of evidence of a claim, while a confidence interval allows us to assess the magnitude of an effect.

Theory-Based Methods 0000 Multiple Linear Regression

Confidence Intervals for Linear Models

- A hypothesis test allows us to assess the strength of evidence of a claim, while a confidence interval allows us to assess the magnitude of an effect.
- Suppose page count can be perfectly predicted by year (with no deviations or errors). What slope would we expect to find in the regression model?

Theory-Based Methods 0000 Multiple Linear Regression

Confidence Intervals for Linear Models

- A hypothesis test allows us to assess the strength of evidence of a claim, while a confidence interval allows us to assess the magnitude of an effect.
- Suppose page count can be perfectly predicted by year (with no deviations or errors). What slope would we expect to find in the regression model?
 - It's hard to say without knowing the variability in the year and in the page count data.
 - Remember that slope tells us the average increase in the response variable per unit increase in the explanatory variable

Theory-Based Methods 0000 Multiple Linear Regression

Confidence Intervals for Linear Models

- A hypothesis test allows us to assess the strength of evidence of a claim, while a confidence interval allows us to assess the magnitude of an effect.
- Suppose page count can be perfectly predicted by year (with no deviations or errors). What slope would we expect to find in the regression model?
 - It's hard to say without knowing the variability in the year and in the page count data.
 - Remember that slope tells us the average increase in the response variable per unit increase in the explanatory variable
- If we want to estimate the strength of the linear relationship between the two variables, we should instead create a confidence interval for the correlation *R*.

Confidence	Intervals
000000	

Theory-Based Methods 0000 Multiple Linear Regression

Bootstrapping for confidence intervals

- To approximate variablity in the correlation statistic *R*, we create a bootstrap sample by resampling the paired data and then calculation correlation
 - This corresponds to sampling with replacement from the columns of the original sample

Confidence	Intervals
000000	

Theory-Based Methods 0000 Multiple Linear Regression 000000000

Bootstrapping for confidence intervals

- To approximate variablity in the correlation statistic *R*, we create a bootstrap sample by resampling the paired data and then calculation correlation
 - This corresponds to sampling with replacement from the columns of the original sample

```
theses samp %>%
  specify(n_pages~year) %>%
  generate(1, type = "bootstrap")
    A tibble: 6 x 3
               replicate [1]
##
   # Groups:
##
     replicate n_pages year
##
         <int>
                  <dbl> <dbl>
## 1
                     51
                        1991
             1
## 2
                     78
                        1987
             1
## 3
             1
                    103 2010
             1
                     81
                         2008
## 4
## 5
             1
                     36
                        1964
## 6
             1
                     37
                         1973
```

Conditions for Inferend

Theory-Based Methods 0000 Multiple Linear Regression

Bootstrapping for confidence intervals

- To approximate variablity in the correlation statistic *R*, we create a bootstrap sample by resampling the paired data and then calculation correlation
 - This corresponds to sampling with replacement from the columns of the original sample

```
theses samp %>%
  specify(n_pages~year) %>%
  generate(1, type = "bootstrap")
     A tibble: 6 x 3
                replicate [1]
##
     Groups:
##
     replicate n_pages
                          year
          <int>
                   <dbl> <dbl>
##
## 1
                      51
                           1991
## 2
                      78
                          1987
## 3
                     103
                          2010
                      81
                           2008
## 4
              1
## 5
                      36
                          1964
               1
## 6
                      37
                           1973
##
   #
     A tibble:
                1 \times 2
##
     replicate
                    cor
                 <dbl>
##
          <int>
## 1
              1 - 0.382
```



- Dashed red line indicates regression line for original sample
- Darker points correspond to observations included in bootstrap more than once

Theory-Based Methods 0000 Multiple Linear Regression

Bootstrap Distribution for correlation

Now we generate 1000 replicates, and compute the correlation for each

Theory-Based Methods 0000 Multiple Linear Regression

Bootstrap Distribution for correlation

Now we generate 1000 replicates, and compute the correlation for each

```
theses_samp %>%
specify(n_pages~year) %>%
generate(1000, type = "bootstrap") %>%
calculate(stat = "correlation")
```

Theory-Based Methods 0000 Multiple Linear Regression

Bootstrap Distribution for correlation

Now we generate 1000 replicates, and compute the correlation for each

```
theses_samp %>%
specify(n_pages~year) %>%
generate(1000, type = "bootstrap") %>%
calculate(stat = "correlation")
```

```
## Response: n_pages (numeric)
## Explanatory: year (numeric)
## # A tibble: 6 x 2
##
     replicate
                  stat
         <int> <dbl>
##
## 1
             1 - 0.294
## 2
             2 - 0.242
## 3
             3 -0.235
## 4
          4 -0.0830
## 5
           5 -0.268
## 6
             6 - 0.407
```

Confidence Intervals		Theory-Based Methods	
00000	000000	0000	00000000

The Bootstrap Distribution for R

correlation_ci <- boot_slope %>% get_ci(level = .95, type = "percentile")
correlation_ci

A tibble: 1 x 2
lower_ci upper_ci
<dbl> <dbl>
1 -0.542 -0.0829

boot_slope %>% visualize()+shade_ci(endpoints =correlation_ci)



Confidence Intervals		Theory-Based Methods	
000000	00000	0000	00000000

The Bootstrap Distribution for R

correlation_ci <- boot_slope %>% get_ci(level = .95, type = "percentile")
correlation_ci

A tibble: 1 x 2
lower_ci upper_ci
<dbl> <dbl>
1 -0.542 -0.0829

boot_slope %>% visualize()+shade_ci(endpoints =correlation_ci)



Simulation–Based Bootstrap Distribution

- The original sample had correlation R = -0.3
 - It is possible the true relationship between page count and year has between very weak (-0.08) and moderate (-0.54) negative correlation.

Multiple Linear Regression 000000000

Section 2

Conditions for Inference

Nate Wells

Theory-Based Methods 0000 Multiple Linear Regression

Conditions for Inference: LINE!

Theory-Based Methods 0000 Multiple Linear Regression

Conditions for Inference: LINE!

- The relationship between explanatory and response variables must be approximately linear. (Linear)
 - Check using scatterplot/residual plot

Multiple Linear Regression

Conditions for Inference: LINE!

- The relationship between explanatory and response variables must be approximately linear. (Linear)
 - Check using scatterplot/residual plot
- **2** The observations should be independent of one another. (Independence)
 - Check using scatterplot/residual plot, as well as sample design
- O The distribution of residuals should be bell-shaped, unimodal, symmetric, and centered at 0. (Normal)
 - Check using histogram of residuals

Multiple Linear Regression

Conditions for Inference: LINE!

- The relationship between explanatory and response variables must be approximately linear. (Linear)
 - Check using scatterplot/residual plot
- **2** The observations should be independent of one another. (Independence)
 - Check using scatterplot/residual plot, as well as sample design
- O The distribution of residuals should be bell-shaped, unimodal, symmetric, and centered at 0. (Normal)
 - Check using histogram of residuals
- The variability of residuals should be roughly constant across entire data set. (Equal Variability)
 - Check using residual plot.

Conditions for Inference

Theory-Based Methods 0000 Multiple Linear Regression

Checking Conditions: Linear



Data is not tightly clustered around line of best fit

Conditions for Inference

Theory-Based Methods 0000 Multiple Linear Regression

Checking Conditions: Linear



Data is not tightly clustered around line of best fit

• But this doesn't mean data is not linear. Just that residuals have high variance

Conditions for Inference

Theory-Based Methods 0000 Multiple Linear Regression

Checking Conditions: Linear



Data is not tightly clustered around line of best fit

• But this doesn't mean data is not linear. Just that residuals have high variance

```
## # A tibble: 1 x 1
## cor
## <dbl>
## 4 0.005
```

1 -0.295

Theory-Based Methods 0000 Multiple Linear Regression

Checking Conditions: Independence

• When students were tasked with sampling theses, they were asked to consider whether their method represented an SRS. Here are some methods used:

Theory-Based Methods 0000 Multiple Linear Regression

Checking Conditions: Independence

- When students were tasked with sampling theses, they were asked to consider whether their method represented an SRS. Here are some methods used:
- Sort theses in the online library catalog by year published and title. Generate 10 random numbers between 1 and 16159, and use these to select theses from catalog.
- **2** Use the library database with no order specified. Randomly generate a letter of the alphabet and pick the first thesis in the list whose title included the letter.
- e Generate 3 random letters of the alphabet, and choose 10 theses whose author's last name begins with the given letter.
- O Divide the thesis tower into 6 sections of approx. equal size. Randomly choose 1 section using 6-sided die. Then randomly choose a shelf in this section, followed by a row, and then a thesis on the row (using appropriately sized dice)

Theory-Based Methods 0000 Multiple Linear Regression 000000000



Theory-Based Methods 0000 Multiple Linear Regression

Checking Conditions: Normal



• The distribution does appears somewhat right-skewed, with a notable outliers on the right.

Theory-Based Methods 0000 Multiple Linear Regression



- The distribution does appears somewhat right-skewed, with a notable outliers on the right.
- This provides some evidence residuals are not Normally disributed.

Theory-Based Methods 0000 Multiple Linear Regression



- The distribution does appears somewhat right-skewed, with a notable outliers on the right.
- This provides some evidence residuals are not Normally disributed.
- Do we discard conclusions entirely?

Theory-Based Methods 0000 Multiple Linear Regression



- The distribution does appears somewhat right-skewed, with a notable outliers on the right.
- This provides some evidence residuals are not Normally disributed.
- Do we discard conclusions entirely?
 - No. But this does warrant further research.

Theory-Based Methods 0000 Multiple Linear Regression

Checking Conditions: Equal Variability



Residual Plot

Residuals appear to have constant variability between 1975 and 2020

Theory-Based Methods 0000 Multiple Linear Regression

Checking Conditions: Equal Variability



Residual Plot

Residuals appear to have constant variability between 1975 and 2020

• However, theses prior to 1975 appear to have more spread (and almost all outliers come from this region of sparser data)

Section 3

Theory-Based Methods

Confidence Intervals	Conditions for Inference	Theory-Based Methods	Multiple Linear Regression
000000	000000	○●○○	

Inference for Slope

• Can we make inference about the slope β_1 of a linear model without using simulation?
Confidence Intervals	Conditions for Inference	Theory-Based Methods	Multiple Linear Regression
000000	000000	○●○○	

- Can we make inference about the slope β_1 of a linear model without using simulation?
 - We need to know the mean, standard error, and shape of the sampling distribution for \hat{eta}_1

Confidence Intervals	Conditions for Inference	Theory-Based Methods	Multiple Linear Regression
000000	000000	○●○○	

- Can we make inference about the slope β_1 of a linear model without using simulation?
 - We need to know the *mean*, *standard error*, and *shape* of the sampling distribution for \hat{eta}_1
- If LINE conditions are satisfied, then $\hat{\beta}_1$ is Normally distributed with mean β_1 .

Confidence Intervals	Conditions for Inference	Theory-Based Methods	Multiple Linear Regression
000000	000000	○●○○	

- Can we make inference about the slope β_1 of a linear model without using simulation?
 - We need to know the mean, standard error, and shape of the sampling distribution for \hat{eta}_1
- If LINE conditions are satisfied, then $\hat{\beta}_1$ is Normally distributed with mean β_1 .
 - And the standard error is given by:

Confidence Intervals 000000	Conditions for Inference	Theory-Based Methods ○●○○	Multiple Linear Regression

- Can we make inference about the slope β_1 of a linear model without using simulation?
 - We need to know the mean, standard error, and shape of the sampling distribution for \hat{eta}_1
- If LINE conditions are satisfied, then $\hat{\beta}_1$ is Normally distributed with mean β_1 .
 - And the standard error is given by:

$$SE(\hat{\beta}_1) = \sqrt{\frac{1}{n-2} \frac{\sum_{i=1}^{n} (y_i - (\beta_0 + \beta_1 x_i))^2}{\sum_{i=1}^{n} (x_i - \bar{x})^2}} \quad (\text{DON'T MEMORIZE!})$$

Confidence Intervals 000000	Conditions for Inference	Theory-Based Methods O●OO	Multiple Linear Regression

- Can we make inference about the slope β_1 of a linear model without using simulation?
 - We need to know the mean, standard error, and shape of the sampling distribution for \hat{eta}_1
- If LINE conditions are satisfied, then $\hat{\beta}_1$ is Normally distributed with mean β_1 .
 - And the standard error is given by:

$$SE(\hat{\beta}_1) = \sqrt{\frac{1}{n-2} \frac{\sum_{i=1}^{n} (y_i - (\beta_0 + \beta_1 x_i))^2}{\sum_{i=1}^{n} (x_i - \bar{x})^2}} \quad (\text{DON'T MEMORIZE!})$$

• In practice, we estimate β_0, β_1 in the formula using $\hat{\beta}_0, \hat{\beta}_1$.

Confidence Intervals	Conditions for Inference	Theory-Based Methods	Multiple Linear Regression
000000	000000	○●○○	

- Can we make inference about the slope β_1 of a linear model without using simulation?
 - We need to know the mean, standard error, and shape of the sampling distribution for \hat{eta}_1
- If LINE conditions are satisfied, then $\hat{\beta}_1$ is Normally distributed with mean β_1 .
 - And the standard error is given by:

$$SE(\hat{\beta}_1) = \sqrt{\frac{1}{n-2} \frac{\sum_{i=1}^{n} (y_i - (\beta_0 + \beta_1 x_i))^2}{\sum_{i=1}^{n} (x_i - \bar{x})^2}} \quad (\text{DON'T MEMORIZE!})$$

• In practice, we estimate β_0, β_1 in the formula using $\hat{\beta}_0, \hat{\beta}_1$.

• We perfom a hypothesis test of H_0 : $\beta_1 = 0$ using the test statistic

$$t = \frac{\text{sample stat} - \text{null value}}{SE} = \frac{\hat{\beta}_1 - 0}{SE}$$

Confidence Intervals 000000	Conditions for Inference	Theory-Based Methods ○●○○	Multiple Linear Regression

- Can we make inference about the slope β_1 of a linear model without using simulation?
 - We need to know the mean, standard error, and shape of the sampling distribution for \hat{eta}_1
- If LINE conditions are satisfied, then $\hat{\beta}_1$ is Normally distributed with mean β_1 .
 - And the standard error is given by:

$$SE(\hat{\beta}_1) = \sqrt{\frac{1}{n-2} \frac{\sum_{i=1}^{n} (y_i - (\beta_0 + \beta_1 x_i))^2}{\sum_{i=1}^{n} (x_i - \bar{x})^2}} \quad (\text{DON'T MEMORIZE!})$$

• In practice, we estimate β_0, β_1 in the formula using $\hat{\beta}_0, \hat{\beta}_1$.

• We perfom a hypothesis test of H_0 : $\beta_1 = 0$ using the test statistic

$$t = \frac{\text{sample stat} - \text{null value}}{SE} = \frac{\hat{\beta}_1 - 0}{SE}$$

• And we create a confidence interval for β_1 using

sample stat
$$\pm t^* \cdot SE = \hat{\beta}_1 \pm t^* \cdot SE$$

Confidence Intervals	Conditions for Inference	Theory-Based Methods	Multiple Linear Regression
000000	000000	○●○○	

- Can we make inference about the slope β_1 of a linear model without using simulation?
 - We need to know the mean, standard error, and shape of the sampling distribution for \hat{eta}_1
- If LINE conditions are satisfied, then $\hat{\beta}_1$ is Normally distributed with mean β_1 .
 - And the standard error is given by:

$$SE(\hat{\beta}_1) = \sqrt{\frac{1}{n-2} \frac{\sum_{i=1}^{n} (y_i - (\beta_0 + \beta_1 x_i))^2}{\sum_{i=1}^{n} (x_i - \bar{x})^2}} \quad (\text{DON'T MEMORIZE!})$$

• In practice, we estimate β_0, β_1 in the formula using $\hat{\beta}_0, \hat{\beta}_1$.

• We perfom a hypothesis test of H_0 : $\beta_1 = 0$ using the test statistic

$$t = \frac{\text{sample stat} - \text{null value}}{SE} = \frac{\hat{\beta}_1 - 0}{SE}$$

• And we create a confidence interval for β_1 using

sample stat
$$\pm t^* \cdot SE = \hat{\beta}_1 \pm t^* \cdot SE$$

• In both cases, the reference distribution is the *t*-distribution with n - 2 degrees of freedom.

Theory-Based Methods

Multiple Linear Regression

Calculating test statistics and confidence intervals

• Can we get test statisics and confidence intervals for β_1 without tedious calculation?

Theory-Based Methods

Multiple Linear Regression

Calculating test statistics and confidence intervals

- Can we get test statisics and confidence intervals for β_1 without tedious calculation?
 - Yes! Using the 1m function in R.

Theory-Based Methods

Multiple Linear Regression

Calculating test statistics and confidence intervals

• Can we get test statisics and confidence intervals for β_1 without tedious calculation?

```
    Yes! Using the lm function in R.
    thesis_mod <- lm(n_pages ~ year, data = theses_samp)
get_regression_table(theses_mod)
```

A tibble: 2 x 7 estimate std_error statistic p_value lower_ci upper_ci ## term <chr>> <dbl> <dbl> <dbl> <dbl> <dbl> ## <db1> ## 1 intercept 1292. 394. 3.28 0.001 509. 2074. ## 2 year -0.605 0.198 -3.06 0.003 -0.998 -0.212

Theory-Based Methods

Multiple Linear Regression

Calculating test statistics and confidence intervals

• Can we get test statisics and confidence intervals for β_1 without tedious calculation?

```
• Yes! Using the lm function in R.
thesis_mod <- lm(n_pages ~ year, data = theses_samp)
get_regression_table(theses_mod)
```

A tibble: 2 x 7 estimate std error statistic p value lower ci upper ci ## term <dbl> <dbl> <dbl> <dbl> <dbl> ## <chr> <db1> ## 1 intercept 1292. 394. 3.28 0.001 509. 2074. ## 2 year -0.605 0.198 -3.06 0.003 -0.998 -0.212

• The theory-based standard error is std_error, the test statistic is statistic, and the corresponding p-value in the t-distribution with n-2 df is p_value.

Theory-Based Methods

Multiple Linear Regression

Calculating test statistics and confidence intervals

• Can we get test statistics and confidence intervals for β_1 without tedious calculation?

```
• Yes! Using the lm function in R.
thesis_mod <- lm(n_pages ~ year, data = theses_samp)
get_regression_table(theses_mod)
```

A tibble: 2 x 7 estimate std error statistic p value lower ci upper ci ## term <dbl> <dbl> <dbl> <dbl> <dbl> ## <chr>> <db1> ## 1 intercept 1292. 394. 3.28 0.001 509. 2074. ## 2 year -0.605 0.198 -3.06 0.003 -0.998 -0.212

- The theory-based standard error is std_error, the test statistic is statistic, and the corresponding p-value in the t-distribution with n-2 df is p_value.
- The upper and lower bounds for the 95% confidence interval are lower_ci and upper_ci

Theory-Based Methods

Multiple Linear Regression

Calculating test statistics and confidence intervals

• Can we get test statistics and confidence intervals for β_1 without tedious calculation?

```
    Yes! Using the lm function in R.
    thesis_mod <- lm(n_pages ~ year, data = theses_samp)</li>
    get_regression_table(theses_mod)
```

A tibble: 2 x 7 estimate std error statistic p value lower ci upper ci ## term <dbl> <dbl> <dbl> <dbl> <dbl> ## <chr> <db1> ## 1 intercept 1292. 394. 3.28 0.001 509. 2074. ## 2 year -0.605 0.198 -3.06 0.003 -0.998 -0.212

- The theory-based standard error is std_error, the test statistic is statistic, and the corresponding p-value in the t-distribution with n-2 df is p_value.
- The upper and lower bounds for the 95% confidence interval are lower_ci and upper_ci
- The table also gives similar information for the intercept and hypothesis test $H_0: \beta_0 = 0$ (but this is less useful in practice)

Confidence Intervals	Conditions for Inference	Theory-Based Methods	Multiple Linear Regression
000000	000000	○○○●	

• Suppose we are interested in investigating the correlation ρ between two variables

	Theory-Based Methods	Multiple Linear Regression
	0000	

- Suppose we are interested in investigating the correlation ρ between two variables
- The standard error for the sample correlation R when $\rho = 0$ is

$$SE(R) = \sqrt{\frac{1-R^2}{n-2}}$$

	Theory-Based Methods	Multiple Linear Regression
	0000	

- Suppose we are interested in investigating the correlation ρ between two variables
- The standard error for the sample correlation R when $\rho = 0$ is

$$SE(R) = \sqrt{\frac{1-R^2}{n-2}}$$

• To test the hypothesis $H_0: \rho = 0$ against $H_a: \rho \neq 0$, use the test statistic

$$t = \frac{\text{sample stat} - \text{null value}}{SE} = \frac{R - 0}{\sqrt{\frac{1 - R^2}{n - 2}}}$$

	Theory-Based Methods	Multiple Linear Regression
	0000	

- Suppose we are interested in investigating the correlation ho between two variables
- The standard error for the sample correlation R when $\rho = 0$ is

$$SE(R) = \sqrt{\frac{1-R^2}{n-2}}$$

• To test the hypothesis $H_0: \rho = 0$ against $H_a: \rho \neq 0$, use the test statistic

$$t = \frac{\text{sample stat} - \text{null value}}{SE} = \frac{R - 0}{\sqrt{\frac{1 - R^2}{n - 2}}}$$

where *t* follows the *t*-distribution with n - 2 degrees of freedom.

• There is a formula for confidence intervals, but it is considerably more complicated.

	Theory-Based Methods	Multiple Linear Regression
	0000	

- Suppose we are interested in investigating the correlation ho between two variables
- The standard error for the sample correlation R when $\rho = 0$ is

$$SE(R) = \sqrt{\frac{1-R^2}{n-2}}$$

• To test the hypothesis $H_0: \rho = 0$ against $H_a: \rho \neq 0$, use the test statistic

$$t = \frac{\text{sample stat} - \text{null value}}{SE} = \frac{R - 0}{\sqrt{\frac{1 - R^2}{n - 2}}}$$

- There is a formula for confidence intervals, but it is considerably more complicated.
 - This is because the sampling distribution for R is highly skewed unless R is close to 0

	Theory-Based Methods	Multiple Linear Regression
	0000	

- Suppose we are interested in investigating the correlation ho between two variables
- The standard error for the sample correlation R when $\rho = 0$ is

$$SE(R) = \sqrt{\frac{1-R^2}{n-2}}$$

• To test the hypothesis $H_0: \rho = 0$ against $H_a: \rho \neq 0$, use the test statistic

$$t = \frac{\text{sample stat} - \text{null value}}{SE} = \frac{R - 0}{\sqrt{\frac{1 - R^2}{n - 2}}}$$

- There is a formula for confidence intervals, but it is considerably more complicated.
 - This is because the sampling distribution for R is highly skewed unless R is close to 0
 - Therefore, we can't use the Normal approximation for *R* unless either the sample size is very large, or *R* is close to 0.

	Theory-Based Methods	Multiple Linear Regression
	0000	

- Suppose we are interested in investigating the correlation ρ between two variables
- The standard error for the sample correlation R when $\rho = 0$ is

$$SE(R) = \sqrt{\frac{1-R^2}{n-2}}$$

• To test the hypothesis $H_0: \rho = 0$ against $H_a: \rho \neq 0$, use the test statistic

$$t = \frac{\text{sample stat} - \text{null value}}{SE} = \frac{R - 0}{\sqrt{\frac{1 - R^2}{n - 2}}}$$

- There is a formula for confidence intervals, but it is considerably more complicated.
 - This is because the sampling distribution for R is highly skewed unless R is close to 0
 - Therefore, we can't use the Normal approximation for *R* unless either the sample size is very large, or *R* is close to 0.
 - This is one situation where the simulation-based method clearly outperforms the theory-based method

Section 4

Multiple Linear Regression

		Theory-Based Methods	Multiple Linear Regression
000000	000000	0000	00000000

 In a multiple linear regression model (MLR), we express the response variable Y as a linear combination of k explanatory variables X₁, X₂,..., X_k:

$$\hat{Y} = eta_0 + eta_1 \cdot X_1 + eta_2 \cdot X_2 + \dots + eta_k \cdot X_k$$

		Theory-Based Methods	Multiple Linear Regression
000000	000000	0000	00000000

 In a multiple linear regression model (MLR), we express the response variable Y as a linear combination of k explanatory variables X₁, X₂,..., X_k:

$$\hat{Y} = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \dots + \beta_k \cdot X_k$$

• We use the following R code to fit and summarize a linear model: mod<-lm(Y ~ X1 + X2 + X3, data = my_data) get_regression_table(mod)

##	#	A tibble:	4 x 7					
##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	intercept	3.26	7.94	0.41	0.686	-13.3	19.8
##	2	X1	-1.24	0.313	-3.95	0.001	-1.89	-0.584
##	3	X2	2.68	1.94	1.38	0.182	-1.36	6.72
##	4	XЗ	3.20	0.397	8.06	0	2.37	4.02

		Theory-Based Methods	Multiple Linear Regression
000000	000000	0000	00000000

 In a multiple linear regression model (MLR), we express the response variable Y as a linear combination of k explanatory variables X₁, X₂,..., X_k:

$$\hat{Y} = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \dots + \beta_k \cdot X_k$$

• We use the following R code to fit and summarize a linear model: mod<-lm(Y ~ X1 + X2 + X3, data = my_data) get_regression_table(mod)

##	#	A tibble:	4 x 7					
##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	intercept	3.26	7.94	0.41	0.686	-13.3	19.8
##	2	X1	-1.24	0.313	-3.95	0.001	-1.89	-0.584
##	3	X2	2.68	1.94	1.38	0.182	-1.36	6.72
##	4	XЗ	3.20	0.397	8.06	0	2.37	4.02

Which gives us our linear regression formula:

 $\hat{Y} = 3.26 - 1.24 \cdot X_1 + 2.68 \cdot X_2 + 3.2 \cdot X_3$

		Theory-Based Methods	Multiple Linear Regression
000000	000000	0000	00000000

 In a multiple linear regression model (MLR), we express the response variable Y as a linear combination of k explanatory variables X₁, X₂,..., X_k:

$$\hat{Y} = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \cdots + \beta_k \cdot X_k$$

• We use the following R code to fit and summarize a linear model: mod<-lm(Y ~ X1 + X2 + X3, data = my_data) get_regression_table(mod)

##	#	A tibble:	4 x 7					
##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	intercept	3.26	7.94	0.41	0.686	-13.3	19.8
##	2	X1	-1.24	0.313	-3.95	0.001	-1.89	-0.584
##	3	X2	2.68	1.94	1.38	0.182	-1.36	6.72
##	4	XЗ	3.20	0.397	8.06	0	2.37	4.02

Which gives us our linear regression formula:

 $\hat{Y} = 3.26 - 1.24 \cdot X_1 + 2.68 \cdot X_2 + 3.2 \cdot X_3$

• The slope on each variable indicates the changed in the predicted value of Y per unit change in that variable, with all other variables held constant

Confidence Intervals 000000	Conditions for Inference	Theory-Based Methods 0000	Multiple Linear Regression

• How does page count vary across year of publication and division?

Confidence Intervals	Conditions for Inference	Theory-Based Methods	Multiple Linear Regression
000000	000000	0000	

• How does page count vary across year of publication and division?



Confidence Intervals	Conditions for Inference	Theory-Based Methods	Multiple Linear Regression
000000	000000	0000	

• How does page count vary across year of publication and division?

Theory-Based Methods 0000 Multiple Linear Regression

Reed Theses

 How does page count vary across year of publication and division? theses_mlr <- lm(n_pages ~ year + division, data = theses) get_regression_table(theses_mlr)

```
## # A tibble: 7 \times 7
                  estimate std_error statistic p_value lower_ci upper_ci
##
     term
##
     <chr>
                     <dbl>
                                <dbl>
                                          <dbl>
                                                  <dbl>
                                                           <dbl>
                                                                     <dbl>
## 1 intercept
                   1149.
                             126.
                                           9.14
                                                  0
                                                         902.
                                                                 1396.
## 2 vear
                     -0.54
                               0.063
                                          -8.58
                                                          -0.664
                                                                    -0.416
                                                  0
## 3 divisionHSS
                     34.5
                               5.61
                                           6.16
                                                  0
                                                          23.5
                                                                   45.5
## 4 divisionID
                     18.7
                              7.38
                                           2.53
                                                  0.011
                                                         4.22
                                                                   33.2
## 5 divisionLL
                     12.6
                               5.84
                                           2.16
                                                  0.031
                                                           1.16
                                                                    24.1
## 6 divisionMNS
                    -11.2
                               5.50
                                          -2.03
                                                  0.042 - 22.0
                                                                   -0.389
## 7 divisionPRPL
                      8.28
                               5.92
                                           1.40
                                                  0.162
                                                          -3.34
                                                                   19.9
```

Confidence Intervals 000000	Conditions for Inference	Theory-Based Methods 0000	Multiple Linear Regression

 How does page count vary across year of publication and division? theses_mlr <- lm(n_pages ~ year + division, data = theses) get_regression_table(theses_mlr)

```
## # A tibble: 7 \times 7
               estimate std_error statistic p_value lower_ci upper_ci
##
    term
##
    <chr>
                  <dbl>
                           <dbl>
                                   <dbl>
                                          <dbl>
                                                  <dbl>
                                                          <dbl>
## 1 intercept
                1149.
                         126.
                                    9.14
                                          0
                                                902.
                                                       1396.
## 2 vear
                -0.54
                         0.063
                                  -8.58 0
                                                -0.664
                                                         -0.416
## 3 divisionHSS
                  34.5
                          5.61
                                    6.16 0
                                                 23.5
                                                        45.5
                      7.38
## 4 divisionTD
                18.7
                                  2.53 0.011 4.22
                                                        33.2
## 5 divisionLL
               12.6 5.84
                               2.16 0.031 1.16
                                                         24.1
## 6 divisionMNS
                -11.2 5.50 -2.03 0.042 -22.0
                                                        -0.389
## 7 divisionPRPL
                   8.28
                          5.92
                                  1.40
                                          0.162 -3.34
                                                        19.9
```

 $\mathrm{Pages} = 1149 - 0.54 \cdot \mathrm{Year} + 34.5 \cdot \mathrm{HSS} + 18.7 \cdot \mathrm{ID} + 12.6 \cdot \mathrm{LL} - 11.2 \cdot \mathrm{MNS} + 8.3 \cdot \mathrm{PRPL}$

Confidence Intervals 000000	Conditions for Inference	Theory-Based Methods 0000	Multiple Linear Regression

 How does page count vary across year of publication and division? theses_mlr <- lm(n_pages ~ year + division, data = theses) get_regression_table(theses_mlr)

```
## # A tibble: 7 \times 7
               estimate std_error statistic p_value lower_ci upper_ci
##
    term
##
    <chr>
                 <dbl>
                          <dbl>
                                  <dbl>
                                         <dbl>
                                                 <dbl>
                                                        <dbl>
## 1 intercept 1149.
                        126.
                                   9.14
                                         0
                                               902.
                                                      1396.
## 2 vear
                -0.54
                        0.063
                               -8.58 0
                                              -0.664
                                                       -0.416
## 3 divisionHSS
               34.5
                          5.61
                                  6.16 0
                                                23.5
                                                       45.5
## 4 divisionTD
               18.7 7.38
                                 2.53 0.011 4.22 33.2
## 5 divisionLL
              12.6 5.84
                              2.16 0.031 1.16
                                                       24.1
## 6 divisionMNS
               -11.2 5.50 -2.03 0.042 -22.0 -0.389
## 7 divisionPRPL
                  8.28
                          5.92
                                1.40 0.162 -3.34
                                                       19.9
```

 $Pages = 1149 - 0.54 \cdot Year + 34.5 \cdot HSS + 18.7 \cdot ID + 12.6 \cdot LL - 11.2 \cdot MNS + 8.3 \cdot PRPL$

• Which division is used as the baseline?

Confidence Intervals 000000	Conditions for Inference	Theory-Based Methods 0000	Multiple Linear Regression

 How does page count vary across year of publication and division? theses_mlr <- lm(n_pages ~ year + division, data = theses) get_regression_table(theses_mlr)

```
## # A tibble: 7 \times 7
               estimate std_error statistic p_value lower_ci upper_ci
##
    term
##
    <chr>
                 <dbl>
                          <dbl>
                                  <dbl>
                                         <dbl>
                                                <dbl>
                                                        <dbl>
## 1 intercept 1149.
                        126.
                                   9.14
                                         0
                                               902.
                                                     1396.
## 2 vear
               -0.54
                        0.063
                               -8.58 0
                                              -0.664
                                                      -0.416
## 3 divisionHSS
               34.5
                         5.61
                                 6.16 0
                                               23.5
                                                      45.5
## 4 divisionTD
               18.7 7.38
                                2.53 0.011 4.22 33.2
## 5 divisionLL
             12.6 5.84 2.16 0.031 1.16
                                                      24.1
## 6 divisionMNS
               -11.2 5.50 -2.03 0.042 -22.0 -0.389
## 7 divisionPRPL
                  8.28
                         5.92
                                1.40 0.162 -3.34
                                                      19.9
```

 $Pages = 1149 - 0.54 \cdot Year + 34.5 \cdot HSS + 18.7 \cdot ID + 12.6 \cdot LL - 11.2 \cdot MNS + 8.3 \cdot PRPL$

- Which division is used as the baseline?
 - Arts (because it's first alphabetically)

Confidence Intervals 000000	Conditions for Inference	Theory-Based Methods 0000	Multiple Linear Regression

Hypothesis Testing

• The regression table provides p-values for each slope in the model.

Confidence Intervals	Conditions for Inference	Theory-Based Methods	Multiple Linear Regression
000000	000000	0000	

Hypothesis Testing

- The regression table provides p-values for each slope in the model.
 - But what hypotheses are being tested?

Confidence Intervals 000000	Conditions for Inference	Theory-Based Methods 0000	Multiple Linear Regression

Hypothesis Testing

- The regression table provides p-values for each slope in the model.
 - But what hypotheses are being tested?
- In a **MLR model**, we are still interested in determining whether a slope β_i is 0.
| Confidence Intervals | Conditions for Inference | Theory-Based Methods | Multiple Linear Regression |
|----------------------|--------------------------|----------------------|----------------------------|
| 000000 | 000000 | 0000 | |
| | | | |

- The regression table provides p-values for each slope in the model.
 - But what hypotheses are being tested?
- In a **MLR model**, we are still interested in determining whether a slope β_i is 0.
 - But we want to investigate this slope in light of the other variables in the model.

Confidence Intervals	Conditions for Inference	Theory-Based Methods	Multiple Linear Regression
000000	000000	0000	

- The regression table provides p-values for each slope in the model.
 - But what hypotheses are being tested?
- In a **MLR model**, we are still interested in determining whether a slope β_i is 0.
 - But we want to investigate this slope in light of the other variables in the model.
- Each row corresponds to a hypothesis test of the form

 $H_0: \beta_i = 0$, given that other variables are included in the model

Confidence Intervals	Conditions for Inference	Theory-Based Methods	Multiple Linear Regression
000000	000000	0000	

- The regression table provides p-values for each slope in the model.
 - But what hypotheses are being tested?
- In a **MLR model**, we are still interested in determining whether a slope β_i is 0.
 - But we want to investigate this slope in light of the other variables in the model.
- Each row corresponds to a hypothesis test of the form

 $H_0: \beta_i = 0$, given that other variables are included in the model

• I.e. The year row corresponds to the test of

 $H_0: \beta_{year} = 0$, given each division is included in the model

Confidence Intervals 000000	Conditions for Inference	Theory-Based Methods 0000	Multiple Linear Regression

- The regression table provides p-values for each slope in the model.
 - But what hypotheses are being tested?
- In a **MLR model**, we are still interested in determining whether a slope β_i is 0.
 - But we want to investigate this slope in light of the other variables in the model.
- Each row corresponds to a hypothesis test of the form

 $H_0: \beta_i = 0$, given that other variables are included in the model

• I.e. The year row corresponds to the test of

 $H_0: \beta_{year} = 0$, given each division is included in the model

• Reminder: The p-value is the probability of obtaining a statistic as extreme as the observed statistic, **if the null hypothesis were true**.

Confidence Intervals 000000	Conditions for Inference	Theory-Based Methods 0000	Multiple Linear Regression

- The regression table provides p-values for each slope in the model.
 - But what hypotheses are being tested?
- In a **MLR model**, we are still interested in determining whether a slope β_i is 0.
 - But we want to investigate this slope in light of the other variables in the model.
- Each row corresponds to a hypothesis test of the form

 $H_0: \beta_i = 0$, given that other variables are included in the model

• I.e. The year row corresponds to the test of

 $H_0: \beta_{year} = 0$, given each division is included in the model

- Reminder: The p-value is the probability of obtaining a statistic as extreme as the observed statistic, **if the null hypothesis were true**.
- The standard error, statistic, and p-values are all calculated using theory-based methods.

Confidence Intervals 000000	Conditions for Inference	Theory-Based Methods 0000	Multiple Linear Regression

- The regression table provides p-values for each slope in the model.
 - But what hypotheses are being tested?
- In a **MLR model**, we are still interested in determining whether a slope β_i is 0.
 - But we want to investigate this slope in light of the other variables in the model.
- Each row corresponds to a hypothesis test of the form

 $H_0: \beta_i = 0$, given that other variables are included in the model

• I.e. The year row corresponds to the test of

 $H_0: \beta_{year} = 0$, given each division is included in the model

- Reminder: The p-value is the probability of obtaining a statistic as extreme as the observed statistic, **if the null hypothesis were true**.
- The standard error, statistic, and p-values are all calculated using theory-based methods.
 - But the formula is very complicated, requiring matrices and linear algebra (If interested, take Math 392)

Theory-Based Methods 0000 Multiple Linear Regression

Analysis

• Consider the regression table...

##	#	A tibble: /	X /					
##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	intercept	1149.	126.	9.14	0	902.	1396.
##	2	year	-0.54	0.063	-8.58	0	-0.664	-0.416
##	3	divisionHSS	34.5	5.61	6.16	0	23.5	45.5
##	4	divisionID	18.7	7.38	2.53	0.011	4.22	33.2
##	5	divisionLL	12.6	5.84	2.16	0.031	1.16	24.1
##	6	divisionMNS	-11.2	5.50	-2.03	0.042	-22.0	-0.389
##	7	divisionPRPI	8.28	5.92	1.40	0.162	-3.34	19.9

Theory-Based Methods 0000 Multiple Linear Regression

Analysis

• Consider the regression table...

##	Ŧ	A tibble: /	X /					
##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	intercept	1149.	126.	9.14	0	902.	1396.
##	2	year	-0.54	0.063	-8.58	0	-0.664	-0.416
##	3	divisionHSS	34.5	5.61	6.16	0	23.5	45.5
##	4	divisionID	18.7	7.38	2.53	0.011	4.22	33.2
##	5	divisionLL	12.6	5.84	2.16	0.031	1.16	24.1
##	6	divisionMNS	-11.2	5.50	-2.03	0.042	-22.0	-0.389
##	7	divisionPRPI	8.28	5.92	1.40	0.162	-3.34	19.9

• For which hypothesis tests would you reject *H*₀?

Theory-Based Methods 0000 Multiple Linear Regression

Analysis

• Consider the regression table...

##	#	A tibble: /	X /					
##		term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
##		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	intercept	1149.	126.	9.14	0	902.	1396.
##	2	year	-0.54	0.063	-8.58	0	-0.664	-0.416
##	3	divisionHSS	34.5	5.61	6.16	0	23.5	45.5
##	4	divisionID	18.7	7.38	2.53	0.011	4.22	33.2
##	5	divisionLL	12.6	5.84	2.16	0.031	1.16	24.1
##	6	divisionMNS	-11.2	5.50	-2.03	0.042	-22.0	-0.389
##	7	divisionPRPI	8.28	5.92	1.40	0.162	-3.34	19.9

- For which hypothesis tests would you reject H₀?
 - What does this mean in context?

Theory-Based Methods 0000 Multiple Linear Regression

Analysis

- Consider the regression table...
- ## # A tibble: 7 x 7 ## estimate std_error statistic p_value lower_ci upper_ci term ## <chr>> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> 902. ## 1 intercept 1149. 126. 9.14 0 1396. ## 2 year -0.540.063 -8.58 -0.664-0.4160 ## 3 divisionHSS 34.5 5.61 6.16 23.5 45.5 0 ## 4 divisionID 18.7 7.38 2.53 0.011 4.22 33.2 ## 5 divisionLL 12.6 5.84 2.16 0.031 1.16 24.1 ## 6 divisionMNS -11.25.50 -2.030.042 -22.0 -0.389 ## 7 divisionPRPL 8.28 5.92 1.40 0.162 -3.34 19.9
 - For which hypothesis tests would you reject *H*₀?
 - What does this mean in context?
 - For which would you fail to reject H₀?

Theory-Based Methods 0000 Multiple Linear Regression

Analysis

- Consider the regression table...
- ## # A tibble: 7 x 7 estimate std error statistic p value lower ci upper ci ## term <chr>> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> ## 902. ## 1 intercept 1149. 126. 9.14 0 1396. ## 2 year -0.540.063 -8.58 -0.664-0.4160 ## 3 divisionHSS 34.5 5.61 6.16 23.5 45.5 0 ## 4 divisionID 18.7 7.38 2.53 0.011 4.22 33.2 ## 5 divisionLL 12.6 5.84 2.16 0.031 1.16 24.1 ## 6 divisionMNS -11.2 5.50 -2.030.042 -22.0 -0.389 ## 7 divisionPRPL 8.28 5.92 1.40 0.162 -3.34 19.9
 - For which hypothesis tests would you reject *H*₀?
 - What does this mean in context?
 - For which would you fail to reject H₀?
 - What does this mean in context?

Theory-Based Methods 0000 Multiple Linear Regression

Analysis

- Consider the regression table...
- ## # A tibble: 7 x 7 estimate std error statistic p value lower ci upper ci ## term <chr>> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> ## 902. ## 1 intercept 1149. 126. 9.14 0 1396. -0.540.063 -8.58 -0.664-0.416## 2 year 0 ## 3 divisionHSS 34.5 5.61 6.16 23.5 45.5 0 ## 4 divisionID 18.7 7.38 2.53 0.011 4.22 33.2 5.84 ## 5 divisionLL 12.6 2.16 0.031 1.16 24.1 ## 6 divisionMNS -11.25.50 -2.03 0.042 -22.0 -0.389## 7 divisionPRPL 8.28 5.92 1.40 0.162 -3.34 19.9
 - For which hypothesis tests would you reject *H*₀?
 - What does this mean in context?
 - For which would you fail to reject H₀?
 - What does this mean in context?
 - What is the relationship between the estimate, the std_error and lower_ci, upper_ci?

Theory-Based Methods 0000 Multiple Linear Regression

Model Assumptions for MLR

• In order to responsibly use MLR to make inference, we need...

Multiple Linear Regression

Model Assumptions for MLR

- In order to responsibly use MLR to make inference, we need...
- The relationship between explanatory and response variables must be approximately multilinear linear. (Linear)
- **2** The observations should be independent of one another. (Independence)
- O The distribution of residuals should be bell-shaped, unimodal, symmetric, and centered at 0. (Normal)
- The variability of residuals should be roughly constant across entire data set. (Equal Variability)

Multiple Linear Regression

Model Assumptions for MLR

- In order to responsibly use MLR to make inference, we need...
- The relationship between explanatory and response variables must be approximately multilinear linear. (Linear)
- **2** The observations should be independent of one another. (Independence)
- O The distribution of residuals should be bell-shaped, unimodal, symmetric, and centered at 0. (Normal)
- The variability of residuals should be roughly constant across entire data set. (Equal Variability)
- How do we check some of these conditions? Why can't we create a scatterplot of residuals as we did for SLR?

Multiple Linear Regression

Model Assumptions for MLR

- In order to responsibly use MLR to make inference, we need...
- The relationship between explanatory and response variables must be approximately multilinear linear. (Linear)
- **2** The observations should be independent of one another. (Independence)
- O The distribution of residuals should be bell-shaped, unimodal, symmetric, and centered at 0. (Normal)
- The variability of residuals should be roughly constant across entire data set. (Equal Variability)
- How do we check some of these conditions? Why can't we create a scatterplot of residuals as we did for SLR?
 - Instead, we will use a scatterplot of residuals vs predicted values

Confidence	
000000	

Theory-Based Methods 0000 Multiple Linear Regression

Residuals vs Fitted Values

```
mod_res <- get_regression_points(theses_mlr)</pre>
mod res %>%
  select(n_pages, n_pages_hat, residual)
## # A tibble: 772 x 3
##
      n_pages n_pages_hat residual
##
        <dbl>
                    <dbl>
                             <dbl>
                     72.5
                           11.5
##
    1
           84
##
   2
          139
                     88.2
                           50.8
##
   3
           50
                     82.2
                           -32.2
##
    4
           79
                     53.6
                            25.4
##
   5
           36
                     47.1
                            -11.1
##
   6
           80
                     54.1
                            25.9
   7
          134
                     67.6
                           66.4
##
##
   8
           58
                     75.2
                           -17.2
           69
                     75.7
                           -6.74
##
   9
##
  10
           74
                     85.4
                             -11.4
         with 762 more rows
##
   #
```

	Theory-Based Methods	Multiple Linear Regression
		000000000

Residuals vs Fitted Values

```
mod_res <- get_regression_points(theses_mlr
mod_res %>%
```

```
select(n_pages, n_pages_hat, residual)
```

A tibble: 772 x 3

##		n_pages	n_pages_hat	residual
##		<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	84	72.5	11.5
##	2	139	88.2	50.8
##	3	50	82.2	-32.2
##	4	79	53.6	25.4
##	5	36	47.1	-11.1
##	6	80	54.1	25.9
##	7	134	67.6	66.4
##	8	58	75.2	-17.2
##	9	69	75.7	-6.74
##	10	74	85.4	-11.4
##	# .	with	762 more row	IS

```
mod_res %>% ggplot(aes(x = n_pages_hat, y = residual)
geom_point()+
geom_smooth(method = "lm", se = F)+
labs(title = "Residual vs Predicted")
```



Confidence Intervals	Conditions for Inference	Theory-Based Methods	Multiple Linear Regression
000000	000000	0000	0000000●0

Residuals vs Fitted Values

```
mod res \% ggplot(aes(x = n pages hat, y = residual)
mod res <- get regression points(theses mlr
                                                geom point()+
mod res %>%
                                                geom smooth(method = "lm", se = F)+
  select(n_pages, n_pages_hat, residual)
                                                labs(title = "Residual vs Predicted")
                                                         Residual vs Predicted
     A tibble: 772 x 3
##
      n_pages n_pages_hat residual
##
                                                                  •
        <db1>
                     <dbl>
                               <dbl>
##
                                                                                       ••
                      72.5
                             11.5
##
           84
    1
##
    2
          139
                      88.2
                             50.8
                            -32.2
##
    3
           50
                     82.2
                             25.4
##
    4
           79
                     53.6
                                                      levidual
##
    5
           36
                     47.1
                            -11.1
##
    6
           80
                     54.1
                            25.9
    7
          134
                     67.6
                            66.4
##
           58
                     75.2
                            -17.2
##
    8
                      75.7
                             -6.74
##
    9
           69
           74
                      85.4
                              -11.4
   10
##
      ... with 762 more rows
                                                                        n pages hat
```

When analyzing residual vs. predicted plots, look for...

Confidence Intervals	Conditions for Inference	Theory-Based Methods	Multiple Linear Regression
000000	000000	0000	0000000●0

Residuals vs Fitted Values

```
mod res \% ggplot(aes(x = n pages hat, y = residual)
mod res <- get regression points(theses mlr
                                                geom point()+
mod res %>%
                                                geom smooth(method = "lm", se = F)+
  select(n_pages, n_pages_hat, residual)
                                                labs(title = "Residual vs Predicted")
                                                        Residual vs Predicted
     A tibble: 772 x 3
##
      n_pages n_pages_hat residual
##
                                                                 •
        <db1>
                     <dbl>
                              <dbl>
##
                      72.5
                             11.5
##
           84
    1
          139
                     88.2
                             50.8
##
    2
##
    3
           50
                     82.2
                            -32.2
##
    4
           79
                     53.6
                            25.4
                                                     esidual
##
    5
           36
                     47.1
                            -11.1
##
           80
                     54.1
                            25.9
    6
          134
                     67.6
                            66.4
##
    7
           58
                     75.2 -17.2
##
    8
                      75.7
                            -6.74
           69
##
    9
           74
                      85.4
                             -11.4
   10
         with 762 more rows
                                                                        n pages hat
```

- When analyzing residual vs. predicted plots, look for...
 - Non-linear patterns
 - Increasing variability across range of predicted values
 - Outliers with atypical predicted value or large residual

••

Confidence Intervals	Conditions for Inference	Theory-Based Methods	Multiple Linear Regression
000000	000000	0000	0000000●

Distribution of Residuals

• We can still look at the histogram and QQ Plot of residuals, as we did for SLR:

Theory-Based Methods 0000 Multiple Linear Regression

Distribution of Residuals

• We can still look at the histogram and QQ Plot of residuals, as we did for SLR:





```
ggplot(mod_res, aes(sample = residual))+
geom_point(stat = "qq")+
labs(title = "QQ Plot of Residuals")
```



Confidence Intervals	Conditions for Inference	Theory-Based Methods	Multiple Linear Regression
000000		0000	00000000●

Distribution of Residuals

• We can still look at the histogram and QQ Plot of residuals, as we did for SLR:



• We see some evidence that residuals are not Normally distributed (right-skew)