

Relationship between two numerical variables	
Correlation	<p>Definition: Linear association between two variables, denoted as R</p> <p>Properties</p> <ol style="list-style-type: none"> <li>1. Magnitude of R measures strength of linear association</li> <li>2. Sign of R indicate slope</li> <li>3. R is between -1 and 1</li> <li>4. R is unitless (not affect by changes in center or scale such as unit conversion)</li> <li>5. Correlation of X with Y = Y with X</li> <li>6. R is sensitive to outliers</li> </ol>
Residuals	<p>Definintion: The sum of the difference between the data points and the linear regression.</p> <p style="text-align: center;">► difference between the observed and predicted y</p> <div style="border: 1px solid black; padding: 10px; width: fit-content; margin-left: auto; margin-right: auto;"> <p style="text-align: center;"><b>residual:</b> <math>e_i = y_i - \hat{y}_i</math></p> </div> <p>Interpret:</p> <ul style="list-style-type: none"> <li>- Positive residual: observed value higher than predicted</li> <li>- Negative residual: observed value is lower than predicted</li> </ul>
Least Squares Line (line of best fit)	<p>To minimize residuals. Calculated by finding the slope, x-intercept, y-intercept.</p> <p>Features:</p> <ul style="list-style-type: none"> <li>• Slope: For each unit increase in x, the y is expected to be higher/lower on average by the slope</li> <li>• X-intercept - serves to adjust the height of the line.</li> </ul> <p>Slope:</p> <div style="border: 1px solid black; padding: 10px; width: fit-content; margin-left: auto; margin-right: auto;"> <math display="block">\text{slope: } b_1 = \frac{s_y}{s_x} R \quad \begin{array}{l} s_x : \text{SD of } x \\ s_y : \text{SD of } y \\ R = \text{cor}(x, y) \end{array}</math> </div> <p>Notation:</p>

	parameter	point estimate
intercept	$\beta_0$	$b_0$
slope	$\beta_1$	$b_1$

### Least Squares Line Example:

The standard deviation of % living poverty is 3.1%, and the standard deviation of % HS graduates is 3.73%. Given that the correlation between these variables is -0.75, what is the slope of the regression line for predicting % living poverty from % HS graduates?

Step 1: Identify what's given

$$s_y = 3.1\%$$

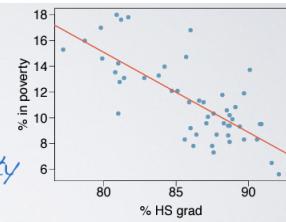
$$s_x = 3.73\%$$

$$R = -0.75$$

Step 2: Find the Slope

$$b_1 = \frac{s_y}{s_x} R = \frac{3.1}{3.73} \times -0.75 \approx -0.62$$

For each % point increase in % HS graduate rate, we would expect the % living in poverty to be lower on average by 0.62% points.



Step 3: Estimate the regression using points given

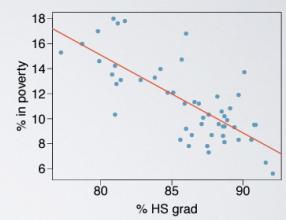
Given that the average % living in poverty is 11.35%, and the average % HS graduates is 86.01%, what is the intercept of the regression line for predicting % living poverty from % HS graduates?

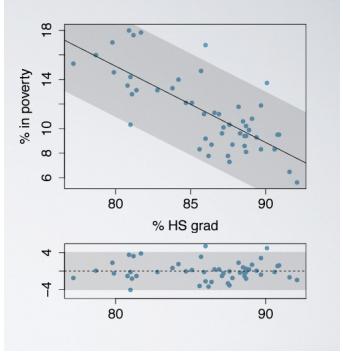
$$\bar{y} = 11.35\%$$

$$\bar{x} = 86.01\%$$

$$b_0 = \bar{y} - b_1 \bar{x} = 11.35 - (-0.62) 86.01 = 64.68$$

States with no HS graduates are expected on average to have 64.68% of their residents living below the poverty line.



Purpose of Least Squares Line	<p><b>Prediction</b></p> <ul style="list-style-type: none"> <li>- Predict the value of the response variable for a given value of the explanatory variable</li> </ul> <p><b>Extrapolation</b></p> <ul style="list-style-type: none"> <li>- Applying a model estimate to values outside of the realm of the original data</li> <li>- Sometimes the intercept might be an extrapolation</li> </ul>
Conditions for Linear Regression	<p><a href="https://gallery.shinyapps.io/sl_r_diag/">https://gallery.shinyapps.io/sl_r_diag/</a></p> <p><b>Linearity</b></p> <ul style="list-style-type: none"> <li>- Linear relationship</li> <li>- Methods for fitting a model to non-linear relationships exist</li> <li>- Check using a scatter plot of the data or a residual plot.</li> </ul> <p><b>Nearly normal residuals</b></p> <ul style="list-style-type: none"> <li>- Normally distributed centered at 0</li> <li>- May not be satisfied if there are unusual observations (outliers)</li> <li>- Check using a histogram or normal probability plot</li> </ul> <p><b>Constant Variability</b></p> <ul style="list-style-type: none"> <li>- Variability of points around regression line should be roughly constant</li> <li>- Implies that the variability of residuals around 0 line should be roughly constant as well</li> <li>- Aka homoscedasticity</li> <li>- Check using a residuals plot</li> </ul> 
R squared	<p><b>Definition:</b></p> <ul style="list-style-type: none"> <li>- Strength of the fit of a linear model calculated as the square of the correlation coefficient</li> <li>- Tells us what percent of variability in the response variable is explained by the model</li> <li>- The rest of the variability is explained by variables not included in the model</li> <li>- Always between 0 and 1</li> </ul> <p><b>Example: <math>R^2 = 0.5625</math></b></p> <ul style="list-style-type: none"> <li>- Means 56.25% of the variability in the % of residents living in poverty among the states is explained by the model</li> </ul>

## Regression with Categorical variables

Set 0 and 1 for each categorical variable, if it's the one you want, set the rest to 0 and the desired variable to 1.

Next, we use a new region variable (`region4`) with four levels: northeast, midwest, west, south. Write the linear regression model based on the regression output below.

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	9.50	0.87	10.94	0.00
region4:midwest	0.03	1.15	0.02	0.98
region4:west	1.79	1.13	1.59	0.12
region4:south	4.16	1.07	3.87	0.00

$$\hat{\% \text{ in poverty}} = 9.50 + 0.03 \text{ reg4:mw} + 1.79 \text{ reg4:w} + 4.16 \text{ reg4:s}$$

Calculate the predicted poverty rate for western states.

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	9.50	0.87	10.94	0.00
region4:midwest	0.03	1.15	0.02	0.98
region4:west	1.79	1.13	1.59	0.12
region4:south	4.16	1.07	3.87	0.00

$$\begin{aligned}\hat{\% \text{ in poverty}} &= 9.50 + 0.03 \text{ reg4:mw} + 1.79 \text{ reg4:w} + 4.16 \text{ reg4:s} \\ &= 9.50 + 0 + 1.79 + 0 \\ &= 11.29\end{aligned}$$

## Outliers and Regression

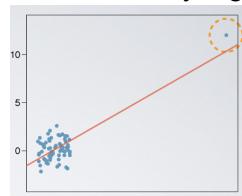
### Outliers in Regression

Two Types:

1. Leverage Points: outliers that fall horizontally away from the center of the cloud but doesn't influence the slope of the regression line



2. Influential Points: influence the slope of the regression line
  - a. Usually high leverage points



<p>Hypothesis Testing for Linear Regression</p>	<p><b>Hypothesis</b></p> <ul style="list-style-type: none"> <li>- <math>H_0 = \text{slope is } 0</math> (no relationship)</li> <li>- <math>H_A = \text{slope } \neq 0</math>, there is a relationship</li> </ul> <p><b>T-statistic and df for Slope Calculation</b></p> <div style="border: 1px solid black; padding: 10px; width: fit-content; margin-left: auto; margin-right: auto;"> <math display="block">\text{t-statistic for the slope: } T = \frac{b_1 - 0}{SE_{b_1}} \quad df = n - 2</math> </div> <ul style="list-style-type: none"> <li>-</li> <li>- <math>Df = n-2</math> because we are estimating for two parameters, <math>b_1</math> and <math>b_2</math></li> </ul> <p><b>Confidence Interval for Slope</b></p> <div style="border: 1px solid black; padding: 10px; width: fit-content; margin-left: auto; margin-right: auto;"> <math display="block">b_1 \pm t_{df}^* SE_{b_1}</math> </div> <ul style="list-style-type: none"> <li>-</li> <li>- SE usually given</li> </ul> <p>Regression output gives <math>b_1</math>, <math>SE_{b_1}</math>, two-tailed p-value for the t-test for the slope where the null value is 0</p>															
<p><b>Hypothesis Test for Linear Regression Example:</b></p>	<p>Calculate the 95% confidence interval for the slope of the relationship between biological and foster twins' IQs?</p> <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th></th> <th>Estimate</th> <th>Std. Error</th> <th>t value</th> <th>Pr(&gt; t )</th> </tr> </thead> <tbody> <tr> <td>(Intercept)</td> <td>9.2076</td> <td>9.2999</td> <td>0.99</td> <td>0.3316</td> </tr> <tr> <td>bioIQ</td> <td>0.9014</td> <td>0.0963</td> <td>9.36</td> <td>0.0000</td> </tr> </tbody> </table> <p>(Also given that 27 subjects were studied)</p> <p>Step 1: Find df</p> <div style="border: 1px solid black; padding: 5px; width: fit-content; margin-left: auto; margin-right: auto;"> <math display="block">df = 27 - 2 = 25</math> </div> <p>Step 2: Find <math>t^*</math></p> <ul style="list-style-type: none"> <li>- By hand...using t-table locate 0.05 for two tail and degree of freedom of 25</li> <li>- By R...</li> </ul> <div style="border: 1px solid black; padding: 10px; width: fit-content; margin-left: auto; margin-right: auto;"> <pre>R &gt; qt(0.025, df = 25) [1] -2.059539</pre> </div> <p>Step 3: Find the Confidence Interval</p> <div style="border: 1px solid black; padding: 5px; width: fit-content; margin-left: auto; margin-right: auto;"> <math display="block">0.9014 \pm 2.059539 \times 0.0963 = (0.7, 1.1)</math> </div> <p>Step 4: Interpret</p> <div style="border: 1px solid black; padding: 10px; width: fit-content; margin-left: auto; margin-right: auto;"> <p>We are 95% confident that for each additional point on the biological twins' IQs, the foster twins' IQs are expected on average to be higher by 0.7 to 1.1 points.</p> </div>		Estimate	Std. Error	t value	Pr(> t )	(Intercept)	9.2076	9.2999	0.99	0.3316	bioIQ	0.9014	0.0963	9.36	0.0000
	Estimate	Std. Error	t value	Pr(> t )												
(Intercept)	9.2076	9.2999	0.99	0.3316												
bioIQ	0.9014	0.0963	9.36	0.0000												

## Variability Partitioning (ANOVA)

Partitioning variability is the process of breaking down the total variability in a dataset into different components that explain the observed data.

Use ANOVA in context of linear regression

### Sum of Squares (Calculated)

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
anova	1	5231.13	5231.13	87.56	0.0000
output	25	1493.53	59.74		
	Total	26	6724.66		

#### sum of squares

$$\text{total variability in } y: SS_{Tot} = \sum(y - \bar{y})^2 = 6724.66$$

$$\text{unexplained variability in } y \text{ (residuals): } SS_{Res} = \sum(y - \hat{y})^2 = \sum e_i^2 = 1493.53$$

$$\text{explained variability in } y: SS_{Reg} = 6724.66 - 1493.53 = 5231.13$$

### Degrees of Freedom

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
anova	1	5231.13	5231.13	87.56	0.0000
output	25	1493.53	59.74		
	Total	26	6724.66		

#### degrees of freedom

$$\text{total degrees of freedom: } df_{Tot} = 27 - 1 = 26$$

$$\text{regression degrees of freedom: } df_{Reg} = 1 \text{ only 1 predictor}$$

$$\text{residual degrees of freedom: } df_{Res} = 26 - 1 = 25$$

### Mean Squares and F-Statistic

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
anova	1	5231.13	5231.13	87.56	0.0000
output	25	1493.53	59.74		
	Total	26	6724.66		

$$\text{mean squares} \quad \text{MS regression: } MS_{Reg} = \frac{SS_{Reg}}{df_{Reg}} = \frac{5231.13}{1} = 5231.13$$

$$\text{MS residual: } MS_{Res} = \frac{SS_{Res}}{df_{Res}} = \frac{1493.53}{25} = 59.74$$

$$\text{F statistic} \quad \text{ratio of explained to unexplained variability} \quad F_{(1,25)} = \frac{MS_{Reg}}{MS_{Res}} = \frac{5231.13}{1493.53} = 87.56$$

### P-Value

anova

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
bioIQ	1	5231.13	5231.13	87.56	0.0000
Residuals	25	1493.53	59.74		
Total	26	6724.66			

$$H_0 : \beta_1 = 0$$

small p-value → reject  $H_0$

$$H_A : \beta_1 \neq 0$$

The data provide convincing evidence that the slope is significantly different than 0, i.e. the explanatory variable is a significant predictor of the response variable.

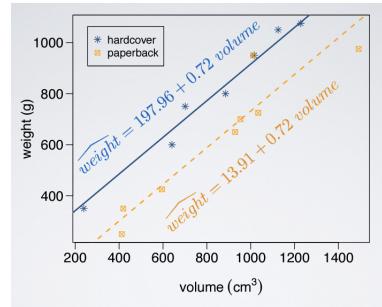
$R^2$  also derived from Proportion of Sum of Squares  
(remember, the other way is the square of correlation coefficient)

$$(2) \quad R^2 = \frac{\text{explained variability}}{\text{total variability}} = \frac{SS_{Reg}}{SS_{Tot}} = \frac{5231.13}{6724.66} \approx 0.78$$

### Multiple Regressions

Multiple Predictors

We are comparing hardcover and paperback books



Plug into R to find the linear model output ( lm() ). ANOVA is referencing hardcover as the base and then layers on the paperback to do the paperback calculation.

The '+cover' adds on the second item to be considered

```
# fit model
> book_mlr = lm(weight ~ volume + cover, data = allbacks)
> summary(book_mlr)
```

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	197.96	59.19	3.34	0.01
volume	0.72	0.06	11.67	0.00
cover:pb	-184.05	40.49	-4.55	0.00

$$\widehat{weight} = 197.96 + 0.72 \text{ volume} - 184.05 \text{ cover : pb}$$

This means that to look at ONLY hardcover, must plug in 0 for the cover:pb

- For hardcover books: plug in 0 for cover:

$$\begin{aligned}\hat{\text{weight}} &= 197.96 + 0.72 \text{ volume} - 184.05 \times 0 \\ &= 197.96 + 0.72 \text{ volume}\end{aligned}$$

To look at ONLY paperback, must plug in 1 for the cover:pb to layer on the difference

- For paperback books: plug in 1 for cover:

$$\begin{aligned}\hat{\text{weight}} &= 197.96 + 0.72 \text{ volume} - 184.05 \times 1 \\ &= 13.91 + 0.72 \text{ volume}\end{aligned}$$

We can predict the weight of a paperback book that is 600cm^3 in volume.

$$197.96 + 0.72 \times 600 - 184.05 \times 1 = 445.91 \text{ grams}$$

### Interpret the output:

#### Slope

- Slope of Volume: for each 1cm^3 increase in volume, the model predicts the books to be heavier on average by 0.72grams
- Slope of Cover:pb: paperback books weigh 184.05g lower than hardcover books

#### Intercept

- Hardcover books with n volume are expected to weigh 198g  
-> makes no sense, just used to adjust the height of the line

Note that the model assumes the same slope for both hardcover and paperback. If that isn't reasonable, we would include an interaction variable in the model.

### Adjusted R-squared

Definition: Applies a penalty to R-squared for the number of predictors included in the model. The magnitude of the penalty depends on how k compares to our sample size n. The larger the sample size, the more predictors the model can handle and the less penalty is going to be for additional predictors being added to the model.

$$\text{adjusted } R^2: R_{adj}^2 = 1 - \left( \frac{SSE}{SST} \times \frac{n-1}{n-k-1} \right) \quad k : \text{number of predictors}$$

Example:

## predicting poverty from % female householder + % white

R																														
> pov_mlr = lm(poverty ~ female_house + white, data = states)																														
> summary(pov_mlr)																														
<table border="1"> <thead> <tr> <th></th> <th>Estimate</th> <th>Std. Error</th> <th>t value</th> <th>Pr(&gt; t )</th> </tr> </thead> <tbody> <tr> <td>(Intercept)</td> <td>-2.58</td> <td>5.78</td> <td>-0.45</td> <td>0.66</td> </tr> <tr> <td>female_house</td> <td>0.89</td> <td>0.24</td> <td>3.67</td> <td>0.00</td> </tr> <tr> <td>white</td> <td>0.04</td> <td>0.04</td> <td>1.08</td> <td>0.29</td> </tr> </tbody> </table>		Estimate	Std. Error	t value	Pr(> t )	(Intercept)	-2.58	5.78	-0.45	0.66	female_house	0.89	0.24	3.67	0.00	white	0.04	0.04	1.08	0.29										
	Estimate	Std. Error	t value	Pr(> t )																										
(Intercept)	-2.58	5.78	-0.45	0.66																										
female_house	0.89	0.24	3.67	0.00																										
white	0.04	0.04	1.08	0.29																										
R																														
> anova(pov_mlr)																														
<table border="1"> <thead> <tr> <th></th> <th>Df</th> <th>Sum Sq</th> <th>Mean Sq</th> <th>F value</th> <th>Pr(&gt;F)</th> </tr> </thead> <tbody> <tr> <td>female_house</td> <td>1</td> <td>132.57</td> <td>132.57</td> <td>18.74</td> <td>0.00</td> </tr> <tr> <td>white</td> <td>1</td> <td>8.21</td> <td>8.21</td> <td>1.16</td> <td>0.29</td> </tr> <tr> <td>Residuals</td> <td>48</td> <td>339.47</td> <td>7.07</td> <td></td> <td></td> </tr> <tr> <td>Total</td> <td>50</td> <td>480.25</td> <td></td> <td></td> <td></td> </tr> </tbody> </table>		Df	Sum Sq	Mean Sq	F value	Pr(>F)	female_house	1	132.57	132.57	18.74	0.00	white	1	8.21	8.21	1.16	0.29	Residuals	48	339.47	7.07			Total	50	480.25			
	Df	Sum Sq	Mean Sq	F value	Pr(>F)																									
female_house	1	132.57	132.57	18.74	0.00																									
white	1	8.21	8.21	1.16	0.29																									
Residuals	48	339.47	7.07																											
Total	50	480.25																												
$R^2 = \frac{132.57 + 8.21}{480.25} = 0.29$																														

We've partitioned the predicted poverty level into two predictors: female householder and white, decreasing the residual (unexplained). However, we need to adjust the  $R^2$  for the additional predictor.

Calculate adjusted  $R^2$  for the multiple linear regression model predicting % living in poverty from % female householders and % white.  
Remember  $n = 51$  (50 states + DC).

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
female_house	1	132.57	132.57	18.74	0.00
white	1	8.21	8.21	1.16	0.29
Residuals	48	339.47	7.07		
Total	50	480.25			

$$\begin{aligned}
 R_{adj}^2 &= 1 - \left( \frac{SSE}{SST} \times \frac{n-1}{n-k-1} \right) \\
 &= 1 - \left( \frac{339.47}{480.25} \times \frac{51-1}{51-2-1} \right) = 0.26
 \end{aligned}$$

As a result, adjusted  $R^2$  is less than the original  $R^2$ .

Why is this important?

- When any variable is added to the model,  $R^2$  increases but that's not necessarily true.
- If the added variable doesn't really provide any new information or is completely unrelated, the adjusted  $R^2$  does not increase
- Therefore we choose adjusted  $R^2$  for the penalty it includes. We choose models with the highest adjusted  $R^2$ .

Collinearity and Parsimony

Collinearity

- Two predictors are correlated with each other, which complicates a model estimation because predictors should be independent
- The addition of collinear variables can result in biased estimates of the regression parameters.

Parsimony

- |  |   |
|--|---|
|  | <ul style="list-style-type: none"> <li>- Want to avoid adding predictors associated with each other because it would bring nothing new to the table so we prefer the simplest best model: <b><i>parsimonious model</i></b> <ul style="list-style-type: none"> <li>- Occam's razor: among competing hypotheses, the one with the fewest assumption should be selected</li> </ul> </li> </ul> |
|--|---|

## Inference for Multiple Regression and Model Selection

### Inference for Multiple Linear Regression (MLR)

Using R: We have a set of data and we want to know if different factors (explanatory variables) for the mom's impact their kid's IQ score.

Using R, we input the full model

```
# full model
> cog_full = lm(kid_score ~ mom_hs + mom_iq + mom_work + mom_age, data = cognitive)
> summary(cog_full)
```

This yields the output below:

```
Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 19.59241   9.21906   2.125   0.0341 *  
mom_hs:yes  5.09482   2.31450   2.201   0.0282 *  
mom_iq      0.56147   0.06064   9.259 <2e-16 *** 
mom_work:yes 2.53718   2.35067   1.079   0.2810    
mom_age     0.21802   0.33074   0.659   0.5101
```

```
Residual standard error: 18.14 on 429 degrees of freedom
Multiple R-squared:  0.2171, Adjusted R-squared:  0.2098 
F-statistic: 29.74 on 4 and 429 DF,  p-value: < 2.2e-16
```

Step 1: Take a F-test to determine the model as a whole to know if it is worthwhile to hypothesis test for slopes

- Set the Hypothesis for F-test

$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$   
 $H_A: \text{At least one } \beta_j \text{ is different than 0}$

F-statistic: 29.74 on 4 and 429 DF, p-value: < 2.2e-16

- $H_0 = \text{all the slopes are 0, there is no relation}$
- $H_A = \text{atleast one slope is not 0}$
- Analysis:
  - Since the p-value is samaller, we reject the null hypothesis and know that at least one slope is not 0

Step 2: Perform hypothesis testing (T-Test) on each of the factors to determine whether or not they are a significant predictor of kids IQ. Example below is with mom who went to high school.

- Set the Hypothesis for T-Test

$H_0: \beta_1 = 0$ , when all other variables are included in the model  
 $H_A: \beta_1 \neq 0$ , when all other variables are included in the model

- $H_0$  = slope is 0, no relation
- $H_A$  = slope is not 0, there is a relationship
- Look at the p-value to see if we can reject the null hypothesis or not
  - 
  - Mom\_hs, mom\_iq are significant
  - mom\_work , mom\_age are not significant

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	19.59241	9.21906	2.125	0.0341
mom_hs:yes	5.09482	2.31450	2.201	0.0282
mom_iq	0.56147	0.06064	9.259	<2e-16
mom_work:yes	2.53718	2.35067	1.079	0.2810
mom_age	0.21802	0.33074	0.659	0.5101

## Inference for Multiple Linear Regression (MLR) Done by Hand

Verify the T score and the p-value for the slope of mom\_hs.

### Step 1: T-score and df

- Formula

$$\text{t-statistic for the slope: } T = \frac{b_1 - 0}{SE_{b_1}} \quad df = n - k - 1$$

$$\begin{aligned} T &= \frac{5.095 - 0}{2.315} \\ &= 2.201 \\ df &= n - k - 1 \\ &= 434 - 4 - 1 \\ &= 429 \end{aligned}$$

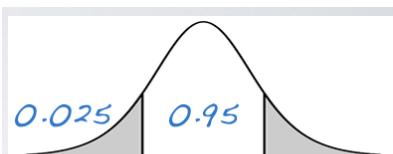
- Output

```
R
> pt(2.201, df = 429, lower.tail = FALSE) * 2
[1] 0.0282
```

Calculate the 95% confidence interval for the slope of mom\_work.

### Step 1: Find critical point $t^*$ and df

- Using R/t-table



```
R
> qt(0.025, df = 429)
[1] -1.97
```

### Step 2: Find the confidence interval

	<table border="1"> <thead> <tr> <th></th><th>Estimate</th><th>Std. Error</th><th>t value</th><th>Pr(&gt; t )</th></tr> </thead> <tbody> <tr> <td>(Intercept)</td><td>19.59241</td><td>9.21906</td><td>2.125</td><td>0.0341</td></tr> <tr> <td>mom_hs:yes</td><td>5.09482</td><td>2.31450</td><td>2.201</td><td>0.0282</td></tr> <tr> <td>mom_iq</td><td>0.56147</td><td>0.06064</td><td>9.259</td><td>&lt;2e-16</td></tr> <tr> <td>mom_work:yes</td><td>2.53718</td><td>2.35067</td><td>1.079</td><td>0.2810</td></tr> <tr> <td>mom_age</td><td>0.21802</td><td>0.33074</td><td>0.659</td><td>0.5101</td></tr> </tbody> </table> <p>- Residual standard error: 18.14 on 429 degrees of freedom</p> <p style="text-align: center;"><math>2.54 \pm 1.97 \times 2.35 \approx (-2.09, 7.17)</math></p> <p><b>Step 3: Interpret</b></p> <ul style="list-style-type: none"> <li>We are 95% confident that, all else being equal, the model predicts that children whose moms worked during the first 3 years of their lives score 2.09 points lower to 7.17 points higher than those whose mom did not work.</li> </ul>		Estimate	Std. Error	t value	Pr(> t )	(Intercept)	19.59241	9.21906	2.125	0.0341	mom_hs:yes	5.09482	2.31450	2.201	0.0282	mom_iq	0.56147	0.06064	9.259	<2e-16	mom_work:yes	2.53718	2.35067	1.079	0.2810	mom_age	0.21802	0.33074	0.659	0.5101																																																																																																												
	Estimate	Std. Error	t value	Pr(> t )																																																																																																																																							
(Intercept)	19.59241	9.21906	2.125	0.0341																																																																																																																																							
mom_hs:yes	5.09482	2.31450	2.201	0.0282																																																																																																																																							
mom_iq	0.56147	0.06064	9.259	<2e-16																																																																																																																																							
mom_work:yes	2.53718	2.35067	1.079	0.2810																																																																																																																																							
mom_age	0.21802	0.33074	0.659	0.5101																																																																																																																																							
Model Selection - Stepwise Model Selection	<p><b>Backwards Elimination - adjusted R<sup>2</sup></b></p> <ul style="list-style-type: none"> <li>Start with the full model</li> <li>Drop one variable at a time and record adjusted r<sup>2</sup> of each smaller model</li> <li>Pick the model with the highest increase in adjusted R<sup>2</sup></li> <li>Repeat until none of the models yield an increase in adjusted R<sup>2</sup></li> </ul> <table border="1"> <thead> <tr> <th>step</th> <th>variables included</th> <th>removed</th> <th>adjusted R</th> </tr> </thead> <tbody> <tr> <td>FULL</td> <td>kid_score ~ mom_hs + mom_iq + mom_work + mom_age</td> <td></td> <td><b>0.2098</b></td> </tr> <tr> <td>STEP 1</td> <td>kid_score ~ mom_iq + mom_work + mom_age</td> <td>[-mom_hs]</td> <td>0.2027</td> </tr> <tr> <td></td> <td>kid_score ~ mom_hs + mom_work + mom_age</td> <td>[-mom_iq]</td> <td>0.0541</td> </tr> <tr> <td></td> <td>kid_score ~ mom_hs + mom_iq + mom_age</td> <td>[-mom_work]</td> <td>0.2095</td> </tr> <tr> <td></td> <td><b>kid_score ~ mom_hs + mom_iq + mom_work</b></td> <td>[-mom_age]</td> <td><b>0.2109</b></td> </tr> <tr> <td>STEP 2</td> <td>kid_score ~ mom_iq + mom_work</td> <td>[-mom_hs]</td> <td>0.2024</td> </tr> <tr> <td></td> <td>kid_score ~ mom_hs + mom_work</td> <td>[-mom_iq]</td> <td>0.0546</td> </tr> <tr> <td></td> <td>kid_score ~ mom_hs + mom_iq</td> <td>[-mom_work]</td> <td>0.2105</td> </tr> </tbody> </table> <p><b>Backwards Elimination - p-value</b></p> <ul style="list-style-type: none"> <li>Start with the full model</li> <li>Drop the variable with the highest p-value and refit a smaller model</li> <li>Repeat until all variables left in the model are significant</li> </ul> <table border="1"> <thead> <tr> <th colspan="6">backwards elimination - p-value</th> </tr> <tr> <th></th> <th>Estimate</th> <th>Std. Error</th> <th>t value</th> <th>Pr(&gt; t )</th> <th></th> </tr> </thead> <tbody> <tr> <td>FULL</td> <td>19.5924</td> <td>9.2191</td> <td>2.13</td> <td>0.0341</td> <td></td> </tr> <tr> <td></td> <td>(Intercept)</td> <td>19.5924</td> <td>9.2191</td> <td>2.13</td> <td>0.0341</td> </tr> <tr> <td></td> <td>mom_hs:yes</td> <td>5.0948</td> <td>2.3145</td> <td>2.20</td> <td>0.0282</td> </tr> <tr> <td></td> <td>mom_iq</td> <td>0.5615</td> <td>0.0606</td> <td>9.26</td> <td>0.0000</td> </tr> <tr> <td></td> <td>mom_work:yes</td> <td>2.5372</td> <td>2.3507</td> <td>1.08</td> <td>0.2810</td> </tr> <tr> <td></td> <td>mom_age</td> <td>0.2180</td> <td>0.3307</td> <td>0.66</td> <td>0.5101</td> </tr> <tr> <td>STEP 1</td> <td>24.1794</td> <td>6.0432</td> <td>4.00</td> <td>0.0001</td> <td></td> </tr> <tr> <td></td> <td>(Intercept)</td> <td>24.1794</td> <td>6.0432</td> <td>4.00</td> <td>0.0001</td> </tr> <tr> <td></td> <td>mom_hs:yes</td> <td>5.3823</td> <td>2.2716</td> <td>2.37</td> <td>0.0183</td> </tr> <tr> <td></td> <td>mom_iq</td> <td>0.5628</td> <td>0.0606</td> <td>9.29</td> <td>0.0000</td> </tr> <tr> <td></td> <td>mom_work:yes</td> <td>2.5664</td> <td>2.3487</td> <td>1.09</td> <td>0.2751</td> </tr> <tr> <td>STEP 2</td> <td>25.7315</td> <td>5.8752</td> <td>4.38</td> <td>0.0000</td> <td></td> </tr> <tr> <td></td> <td>(Intercept)</td> <td>25.7315</td> <td>5.8752</td> <td>4.38</td> <td>0.0000</td> </tr> <tr> <td></td> <td>mom_hs:yes</td> <td>5.9501</td> <td>2.2118</td> <td>2.69</td> <td>0.0074</td> </tr> <tr> <td></td> <td>mom_iq</td> <td>0.5639</td> <td>0.0606</td> <td>9.31</td> <td>0.0000</td> </tr> </tbody> </table> <p><b>Forward Selection - adjusted R<sup>2</sup></b></p>	step	variables included	removed	adjusted R	FULL	kid_score ~ mom_hs + mom_iq + mom_work + mom_age		<b>0.2098</b>	STEP 1	kid_score ~ mom_iq + mom_work + mom_age	[-mom_hs]	0.2027		kid_score ~ mom_hs + mom_work + mom_age	[-mom_iq]	0.0541		kid_score ~ mom_hs + mom_iq + mom_age	[-mom_work]	0.2095		<b>kid_score ~ mom_hs + mom_iq + mom_work</b>	[-mom_age]	<b>0.2109</b>	STEP 2	kid_score ~ mom_iq + mom_work	[-mom_hs]	0.2024		kid_score ~ mom_hs + mom_work	[-mom_iq]	0.0546		kid_score ~ mom_hs + mom_iq	[-mom_work]	0.2105	backwards elimination - p-value							Estimate	Std. Error	t value	Pr(> t )		FULL	19.5924	9.2191	2.13	0.0341			(Intercept)	19.5924	9.2191	2.13	0.0341		mom_hs:yes	5.0948	2.3145	2.20	0.0282		mom_iq	0.5615	0.0606	9.26	0.0000		mom_work:yes	2.5372	2.3507	1.08	0.2810		mom_age	0.2180	0.3307	0.66	0.5101	STEP 1	24.1794	6.0432	4.00	0.0001			(Intercept)	24.1794	6.0432	4.00	0.0001		mom_hs:yes	5.3823	2.2716	2.37	0.0183		mom_iq	0.5628	0.0606	9.29	0.0000		mom_work:yes	2.5664	2.3487	1.09	0.2751	STEP 2	25.7315	5.8752	4.38	0.0000			(Intercept)	25.7315	5.8752	4.38	0.0000		mom_hs:yes	5.9501	2.2118	2.69	0.0074		mom_iq	0.5639	0.0606	9.31	0.0000
step	variables included	removed	adjusted R																																																																																																																																								
FULL	kid_score ~ mom_hs + mom_iq + mom_work + mom_age		<b>0.2098</b>																																																																																																																																								
STEP 1	kid_score ~ mom_iq + mom_work + mom_age	[-mom_hs]	0.2027																																																																																																																																								
	kid_score ~ mom_hs + mom_work + mom_age	[-mom_iq]	0.0541																																																																																																																																								
	kid_score ~ mom_hs + mom_iq + mom_age	[-mom_work]	0.2095																																																																																																																																								
	<b>kid_score ~ mom_hs + mom_iq + mom_work</b>	[-mom_age]	<b>0.2109</b>																																																																																																																																								
STEP 2	kid_score ~ mom_iq + mom_work	[-mom_hs]	0.2024																																																																																																																																								
	kid_score ~ mom_hs + mom_work	[-mom_iq]	0.0546																																																																																																																																								
	kid_score ~ mom_hs + mom_iq	[-mom_work]	0.2105																																																																																																																																								
backwards elimination - p-value																																																																																																																																											
	Estimate	Std. Error	t value	Pr(> t )																																																																																																																																							
FULL	19.5924	9.2191	2.13	0.0341																																																																																																																																							
	(Intercept)	19.5924	9.2191	2.13	0.0341																																																																																																																																						
	mom_hs:yes	5.0948	2.3145	2.20	0.0282																																																																																																																																						
	mom_iq	0.5615	0.0606	9.26	0.0000																																																																																																																																						
	mom_work:yes	2.5372	2.3507	1.08	0.2810																																																																																																																																						
	mom_age	0.2180	0.3307	0.66	0.5101																																																																																																																																						
STEP 1	24.1794	6.0432	4.00	0.0001																																																																																																																																							
	(Intercept)	24.1794	6.0432	4.00	0.0001																																																																																																																																						
	mom_hs:yes	5.3823	2.2716	2.37	0.0183																																																																																																																																						
	mom_iq	0.5628	0.0606	9.29	0.0000																																																																																																																																						
	mom_work:yes	2.5664	2.3487	1.09	0.2751																																																																																																																																						
STEP 2	25.7315	5.8752	4.38	0.0000																																																																																																																																							
	(Intercept)	25.7315	5.8752	4.38	0.0000																																																																																																																																						
	mom_hs:yes	5.9501	2.2118	2.69	0.0074																																																																																																																																						
	mom_iq	0.5639	0.0606	9.31	0.0000																																																																																																																																						

- Start with a single predictor regression of the response vs explanatory variable
- Pick the model with the highest adjusted R<sup>2</sup>
- Add the remaining variables one at a time to the existing model and pick the one with highest adjusted R<sup>2</sup>
- Repeat until addition of variables does not result in higher adjusted R<sup>2</sup>

step	variables included	adjusted R
STEP 1	kid_score ~ mom_hs	0.0539
	kid_score ~ mom_work	0.0097
	kid_score ~ mom_age	0.0062
	kid_score ~ mom_iq	0.1991
STEP 2	kid_score ~ mom_iq + mom_work	0.2024
	kid_score ~ mom_iq + mom_age	0.1999
	kid_score ~ mom_iq + mom_hs	0.2105
STEP 3	kid_score ~ mom_iq + mom_hs + mom_age	0.2095
	kid_score ~ mom_iq + mom_hs + mom_work	0.2109
STEP 4	kid_score ~ mom_hs + mom_iq + mom_work + mom_age	0.2098

### Forward Selection - p-value

- Start with a single predictor regression of the response vs explanatory variable
- Pick the variable with lowest p-value
- Add the remaining variable one at a time to the existing model and pick the variable with lowest p-value

### Model Selection - Backward Elimination Example

The following model uses data from the American Community Survey to predict income from hours worked per week, race, and gender. Which variable (if any) should be dropped from the model first when doing backwards elimination using the p-value approach?

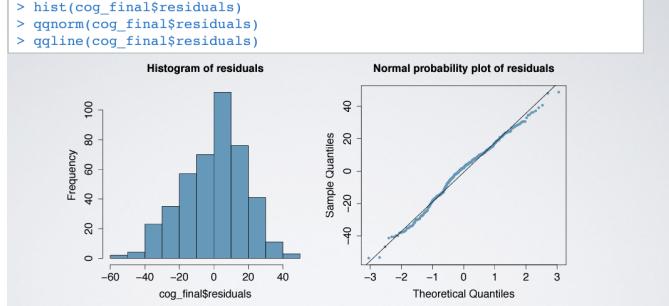
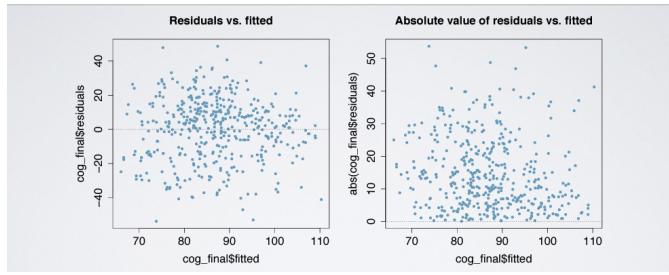
	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	2782.5726	6676.5534	0.42	0.6770
hrs_work	1247.2128	146.2013	8.53	0.0000 ✓
race:black	-9565.3090	6393.2168	-1.50	0.1350
race:asian	35816.6156	8690.3484	4.12	0.0000
race:other	-11112.8617	7213.3220	-1.54	0.1238
gender:female	-16430.0916	3803.4700	-4.32	0.0000 ✓

*don't drop any variables*

**NOTE:** if you have categorical variables with various levels and you are doing backwards elimination. Do not drop variable Race because the level asian's p-value is significant

### Model Selection - R<sup>2</sup> vs P-value

Adjusted R<sup>2</sup> is more reliable prediction but p-value has more significant predictors

	<pre>R &gt; cog_final = lm(kid_score ~ mom_hs + mom_iq + mom_work, data = cognitive) &gt; summary(cog_final)  Coefficients:             Estimate Std. Error t value Pr(&gt; t )     (Intercept) 24.17944   6.04319   4.001 7.42e-05 *** mom_hsyes    5.38225   2.27156   2.369   0.0183 *   mom_iq        0.56278   0.06057   9.291 &lt; 2e-16 *** mom_workyes   2.56640   2.34871   1.093   0.2751     Residual standard error: 18.13 on 430 degrees of freedom Multiple R-squared:  0.2163, Adjusted R-squared:  0.2109  F-statistic: 39.57 on 3 and 430 DF,  p-value: &lt; 2.2e-16</pre> <p>We selected this model using the adjusted R^2 method, which tells us that including the mom_workyes gives us higher predictive power even though it may not be statistically significant than the other variables (looking at p-values)</p>
Conditions required for Multiple Linear Regressions to be Valid	<ol style="list-style-type: none"> <li>Linear relationship between numerical explanatory and response variable           <pre>&gt; hist(cog_final\$residuals) &gt; qqnorm(cog_final\$residuals) &gt; qqline(cog_final\$residuals)</pre>  </li> <li>Constant variability of residuals           <pre>&gt; plot(cog_final\$residuals ~ cog_final\$fitted) &gt; plot(abs(cog_final\$residuals) ~ cog_final\$fitted)</pre>  </li> <li>Independent residuals (independent observations)</li> </ol>

