# More Multiple Linear Regression

Prof. Wells

Math 209, 2/27/23

# Outline

In this lecture, we will...

- Discuss framework for multiple linear regression and compare to simple linear regression
- Use the moderndive packages to create multiple regression models.
- Investigate the geometry of multilinear regression models

# Section 1

# Multiple Linear Regression

• Often, several explanatory variables could be used to predict values of a single response variable.

- Often, several explanatory variables could be used to predict values of a single response variable.
  - Response: Penguin bill length
  - Potential Explanatory: body mass, species, bill depth, age

- Often, several explanatory variables could be used to predict values of a single response variable.
  - Response: Penguin bill length
  - Potential Explanatory: body mass, species, bill depth, age
  - Response: Home prices
  - **Potential Explanatory**: square feet, # bedrooms, # bathrooms, neighborhood

- Often, several explanatory variables could be used to predict values of a single response variable.
  - Response: Penguin bill length
  - Potential Explanatory: body mass, species, bill depth, age
  - **Response**: Home prices
  - **Potential Explanatory**: square feet, # bedrooms, # bathrooms, neighborhood
  - Response: State graduation rate
  - Potential Explanatory: poverty rate, per capita tax revenue, region, teen pregnancy rate

- Often, several explanatory variables could be used to predict values of a single response variable.
  - Response: Penguin bill length
  - Potential Explanatory: body mass, species, bill depth, age
  - Response: Home prices
  - **Potential Explanatory**: square feet, # bedrooms, # bathrooms, neighborhood
  - Response: State graduation rate
  - Potential Explanatory: poverty rate, per capita tax revenue, region, teen pregnancy rate
- In each case, we could create simple linear regression models for each explanatory variable.

- Often, several explanatory variables could be used to predict values of a single response variable.
  - Response: Penguin bill length
  - Potential Explanatory: body mass, species, bill depth, age
  - Response: Home prices
  - **Potential Explanatory**: square feet, # bedrooms, # bathrooms, neighborhood
  - Response: State graduation rate
  - Potential Explanatory: poverty rate, per capita tax revenue, region, teen pregnancy rate
- In each case, we could create simple linear regression models for each explanatory variable.
  - But the results may be misleading:

- Often, several explanatory variables could be used to predict values of a single response variable.
  - Response: Penguin bill length
  - Potential Explanatory: body mass, species, bill depth, age
  - Response: Home prices
  - Potential Explanatory: square feet, # bedrooms, # bathrooms, neighborhood
  - Response: State graduation rate
  - Potential Explanatory: poverty rate, per capita tax revenue, region, teen pregnancy rate
- In each case, we could create simple linear regression models for each explanatory variable.
  - But the results may be misleading:
  - Some individual models may be stronger than others.

- Often, several explanatory variables could be used to predict values of a single response variable.
  - Response: Penguin bill length
  - Potential Explanatory: body mass, species, bill depth, age
  - Response: Home prices
  - Potential Explanatory: square feet, # bedrooms, # bathrooms, neighborhood
  - Response: State graduation rate
  - Potential Explanatory: poverty rate, per capita tax revenue, region, teen pregnancy rate
- In each case, we could create simple linear regression models for each explanatory variable.
  - But the results may be misleading:
  - Some individual models may be stronger than others.
  - Results may be correlated, so we can't easily quantify uncertainty

- Often, several explanatory variables could be used to predict values of a single response variable.
  - Response: Penguin bill length
  - Potential Explanatory: body mass, species, bill depth, age
  - Response: Home prices
  - Potential Explanatory: square feet, # bedrooms, # bathrooms, neighborhood
  - Response: State graduation rate
  - Potential Explanatory: poverty rate, per capita tax revenue, region, teen pregnancy rate
- In each case, we could create simple linear regression models for each explanatory variable.
  - But the results may be misleading:
  - Some individual models may be stronger than others.
  - Results may be correlated, so we can't easily quantify uncertainty
- Could we get better predictive power by including all explanatory variables in the *same* model?

Goal: Visualize quantitative response variable and 2 quantitative explanatory variables.

Goal: Visualize quantitative response variable and 2 quantitative explanatory variables.

• Option 1: 2D scatterplot with explanatory variables on x and y axes, color for response:

Goal: Visualize quantitative response variable and 2 quantitative explanatory variables.

• Option 1: 2D scatterplot with explanatory variables on x and y axes, color for response:



2D scatterplot with color for response

9

6 3 0

Goal: Visualize quantitative response variable and 2 quantitative explanatory variables.

Goal: Visualize quantitative response variable and 2 quantitative explanatory variables.

• Option 2: 3D scatterplot with explanatory variables on x and y axes, response on z axis:

Goal: Visualize quantitative response variable and 2 quantitative explanatory variables.

• Option 2: 3D scatterplot with explanatory variables on x and y axes, response on z axis:



• An interactive 3D plot is available on schedule page of course website.

• In a simple linear regression model (SLR), we express the response variable Y as a linear function of one explanatory variable X:

$$\hat{Y} = \beta_0 + \beta_1 \cdot X$$

• In a **simple linear regression model** (SLR), we express the response variable *Y* as a linear function of one explanatory variable *X*:

$$\hat{Y} = \beta_0 + \beta_1 \cdot X$$

 In a multiple linear regression model (MLR), we express the response variable Y as a linear combination of p explanatory variables X<sub>1</sub>, X<sub>2</sub>,..., X<sub>p</sub>:

$$\hat{Y} = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \dots + \beta_p \cdot X_p$$

• In a **simple linear regression model** (SLR), we express the response variable *Y* as a linear function of one explanatory variable *X*:

$$\hat{Y} = \beta_0 + \beta_1 \cdot X$$

 In a multiple linear regression model (MLR), we express the response variable Y as a linear combination of p explanatory variables X<sub>1</sub>, X<sub>2</sub>,..., X<sub>p</sub>:

$$\hat{Y} = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \dots + \beta_p \cdot X_p$$

• In the MLR model, explanatory variables can either be quantitative or binary categorical

• In a **simple linear regression model** (SLR), we express the response variable *Y* as a linear function of one explanatory variable *X*:

$$\hat{Y} = \beta_0 + \beta_1 \cdot X$$

• In a multiple linear regression model (MLR), we express the response variable Y as a linear combination of p explanatory variables  $X_1, X_2, \ldots, X_p$ :

$$\hat{Y} = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \dots + \beta_p \cdot X_p$$

- In the MLR model, explanatory variables can either be quantitative or binary categorical
  - If we want to use categoricals with more than 2 levels, we need to first create indicators for each level.

• In a simple linear regression model (SLR), we express the response variable Y as a linear function of one explanatory variable X:

$$\hat{Y} = \beta_0 + \beta_1 \cdot X$$

• In a multiple linear regression model (MLR), we express the response variable Y as a linear combination of p explanatory variables  $X_1, X_2, \ldots, X_p$ :

$$\hat{Y} = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \dots + \beta_p \cdot X_p$$

- In the MLR model, explanatory variables can either be quantitative or binary categorical
  - If we want to use categoricals with more than 2 levels, we need to first create indicators for each level.
- We do lose a nice 2D graphical representation (although higher dimensional graphics are possible), but statistical software allows us to estimate coefficients of the model.

• To perform simple linear regression, we found a formula for the model that minimized the sum of squared residuals:

Minimize 
$$\sum_{i=1}^{n} e_i^2$$
 where  $e = y - \hat{y} = y - (\beta_0 + \beta_1 x)$ 

• To perform simple linear regression, we found a formula for the model that minimized the sum of squared residuals:

Minimize 
$$\sum_{i=1}^{n} e_i^2$$
 where  $e = y - \hat{y} = y - (\beta_0 + \beta_1 x)$ 

• To create an MLR model, we do the exact same thing!

• To perform simple linear regression, we found a formula for the model that minimized the sum of squared residuals:

Minimize 
$$\sum_{i=1}^{n} e_i^2$$
 where  $e = y - \hat{y} = y - (\beta_0 + \beta_1 x)$ 

- To create an MLR model, we do the exact same thing!
  - That is, we find the model involving sums of the variables that minimize the squared sum of residuals:

Minimize 
$$\sum_{i=1}^{n} e_i^2 \quad \text{where } e = y - \hat{y} = y - (\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p)$$

• To perform simple linear regression, we found a formula for the model that minimized the sum of squared residuals:

Minimize 
$$\sum_{i=1}^{n} e_i^2$$
 where  $e = y - \hat{y} = y - (\beta_0 + \beta_1 x)$ 

- To create an MLR model, we do the exact same thing!
  - That is, we find the model involving sums of the variables that minimize the squared sum of residuals:

Minimize 
$$\sum_{i=1}^{n} e_i^2 \quad \text{where } e = y - \hat{y} = y - (\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p)$$

• The only difference is that instead of the equation describing a line, the equation describes a "plane" in higher dimensional space.

• To perform simple linear regression, we found a formula for the model that minimized the sum of squared residuals:

Minimize 
$$\sum_{i=1}^{n} e_i^2$$
 where  $e = y - \hat{y} = y - (\beta_0 + \beta_1 x)$ 

- To create an MLR model, we do the exact same thing!
  - That is, we find the model involving sums of the variables that minimize the squared sum of residuals:

Minimize  $\sum_{i=1}^{n} e_i^2 \quad \text{where } e = y - \hat{y} = y - (\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p)$ 

- The only difference is that instead of the equation describing a line, the equation describes a "plane" in higher dimensional space.
- There is a formula for the coefficients of the multilinear model. But we will use lm in R, rather than the formula.

```
mlr_mod <- lm(y ~ x1 + x2 + ... + xp, data = my_data)
get_regression_table(mlr_mod)</pre>
```

# Visualizing Regression Plane

• The regression plane in 3D space minimizes the sum of squared residuals:

# Visualizing Regression Plane

• The regression plane in 3D space minimizes the sum of squared residuals:



- An interactive 3D plot is available on schedule page of course website.
- Regression Equation:  $\hat{y} = -0.8 + 0.67x_1 + 0.83x_2$

$$\hat{Y} = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \dots + \beta_p \cdot X_p$$

• Consider a multilinear model with equation

$$\hat{Y} = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \dots + \beta_p \cdot X_p$$

• The intercept  $\beta_0$  of the MLR is the predicted value of the response when *all* explanatory values take the value 0

$$\hat{Y} = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \dots + \beta_p \cdot X_p$$

- The **intercept**  $\beta_0$  of the MLR is the predicted value of the response when *all* explanatory values take the value 0
  - Whether it is reasonable to make this prediction depends on whether it is plausible for all explanatory variables to be 0.

$$\hat{Y} = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \dots + \beta_p \cdot X_p$$

- The **intercept**  $\beta_0$  of the MLR is the predicted value of the response when *all* explanatory values take the value 0
  - Whether it is reasonable to make this prediction depends on whether it is plausible for all explanatory variables to be 0.
- A slope β<sub>i</sub> is the average change in the response Y per 1 unit change in X<sub>i</sub>, while holding *all* other variables in the model constant.

$$\hat{Y} = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \dots + \beta_p \cdot X_p$$

- The **intercept**  $\beta_0$  of the MLR is the predicted value of the response when *all* explanatory values take the value 0
  - Whether it is reasonable to make this prediction depends on whether it is plausible for all explanatory variables to be 0.
- A slope β<sub>i</sub> is the average change in the response Y per 1 unit change in X<sub>i</sub>, while holding *all* other variables in the model constant.
  - Positive values of β<sub>i</sub> indicate that increases in the corresponding explanatory variable X<sub>i</sub> are associated with increases in the response, while other variables are held constant.

$$\hat{Y} = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \dots + \beta_p \cdot X_p$$

- The **intercept**  $\beta_0$  of the MLR is the predicted value of the response when *all* explanatory values take the value 0
  - Whether it is reasonable to make this prediction depends on whether it is plausible for all explanatory variables to be 0.
- A slope  $\beta_i$  is the average change in the response Y per 1 unit change in  $X_i$ , while holding *all* other variables in the model constant.
  - Positive values of β<sub>i</sub> indicate that increases in the corresponding explanatory variable X<sub>i</sub> are associated with increases in the response, while other variables are held constant.
  - The multilinear model allows us to isolate the effect of one variable on the response
# Section 2

# Application of Multiple Linear Regression

### House Prices

• What factors determine the sale price of a house?

### House Prices

- What factors determine the sale price of a house?
  - We'll consider a subset of 1000 homes from the house\_price dataset in the moderndive package, which contains sale prices for homes in King County, WA between May 2014 and May 2015.

### House Prices

- What factors determine the sale price of a house?
  - We'll consider a subset of 1000 homes from the house\_price dataset in the moderndive package, which contains sale prices for homes in King County, WA between May 2014 and May 2015.

```
## Rows: 1,000
```

```
## Columns: 17
```

##	\$ price	<dbl></dbl>	241, 262, 765, 430, 215, 675, 885, 907, 395, 650, 300, 6~
##	\$ bedrooms	<dbl></dbl>	3, 4, 4, 2, 3, 2, 4, 3, 3, 3, 2, 3, 4, 4, 2, 3, 5, 3, 3,~
##	\$ bathrooms	<dbl></dbl>	1.8, 2.0, 1.0, 2.2, 2.0, 1.8, 2.5, 1.5, 1.5, 2.8, 1.5, 2~
##	\$ sqft_living	<dbl></dbl>	1350, 1540, 2520, 1040, 1280, 2140, 2830, 1340, 1120, 16~
##	\$ sqft_lot	<dbl></dbl>	7588, 5110, 5500, 1516, 6994, 5000, 5000, 6000, 7000, 13~
##	\$ floors	<dbl></dbl>	1.0, 1.0, 1.5, 2.0, 1.0, 1.0, 2.0, 1.5, 1.0, 3.0, 1.0, 2~
##	\$ waterfront	<lgl></lgl>	FALSE, FALSE, FALSE, FALSE, FALSE, FALSE, FALSE, ~
##	\$ view	<dbl></dbl>	0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 2, 0, 0, 0, 0, 0, 0, 0, ~
##	\$ condition	<dbl></dbl>	3, 3, 5, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3,
##	\$ grade	<dbl></dbl>	7, 7, 8, 8, 7, 7, 9, 9, 7, 9, 6, 9, 7, 8, 8, 7, 9, 6, 7,~
##	\$ sqft_above	<dbl></dbl>	1350, 1540, 1820, 1040, 1280, 1000, 2830, 1340, 1120, 13~
##	\$ sqft_basement	<dbl></dbl>	0, 0, 700, 0, 0, 1140, 0, 0, 0, 320, 480, 0, 890, 0, 0, ~
##	\$ yr_built	<dbl></dbl>	1993, 1957, 1912, 2008, 1991, 1930, 1995, 1927, 1955, 20~
##	\$ <pre>yr_renovated</pre>	<dbl></dbl>	0, 0, 0, 0, 0, 1991, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
##	\$ zipcode	<dbl></dbl>	98010, 98118, 98144, 98122, 98038, 98112, 98105, 98105, ~
##	\$ lat	<dbl></dbl>	47, 48, 48, 48, 47, 48, 48, 48, 48, 48, 48, 48, 48, 48, 47, ~
##	\$ long	<dbl></dbl>	-122, -122, -122, -122, -122, -122, -122, -122, -122, -1~

• Consider price as function of square footage, and above ground square footage

 Consider price as function of square footage, and above ground square footage Price vs Square Footage



 Consider price as function of square footage, and above ground square footage Price vs Square Footage



 $Price = 158.66 + 0.16 \cdot sqft$  R = 0.56



 $Price = 158.66 + 0.16 \cdot sqft$  R = 0.56



 $Price = 158.66 + 0.16 \cdot sqft$  R = 0.56  $Price = 236.53 + 0.13 \cdot abv$  R = 0.45

• Both models have some explanatory power for price.

### The Regression Plane

• How do total square footage and above ground square footage together explain price?



Price vs Total / Above Ground Square Footage

### The Regression Plane

• How do total square footage and above ground square footage together explain price?



Price vs Total / Above Ground Square Footage

• What does the upper diagonal line correspond to?

## The Regression Plane

• How do total square footage and above ground square footage together explain price?



Price vs Total / Above Ground Square Footage

- What does the upper diagonal line correspond to?
- Which type of houses tend to have the highest price?

• Let's find the MLR model

house\_sqft\_abv\_mod <-lm(price ~ sqft\_living + sqft\_above, data = house)</pre>

Let's find the MLR model

house\_sqft\_abv\_mod <-lm(price ~ sqft\_living + sqft\_above, data = house)</pre>

And investigate the regression table get\_regression\_table(house\_sqft\_abv\_mod)

## # A tibble: 3 x 7 ## estimate std error statistic p value lower ci upper ci term ## <chr>> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> ## 1 intercept 161. 14.8 10.9 132. 190. 0 0.145 ## 2 sqft living 0.1720.014 12.6 0 0.199 ## 3 saft above -0.0170.014 -1.170.243 -0.045 0.011

• Let's find the MLR model

house\_sqft\_abv\_mod <-lm(price ~ sqft\_living + sqft\_above, data = house)</pre>

```
And investigate the regression table
get_regression_table(house_sqft_abv_mod)
```

##	#	A tibble: 3	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	161.	14.8	10.9	0	132.	190.
##	2	sqft_living	0.172	0.014	12.6	0	0.145	0.199
##	З	sqft_above	-0.017	0.014	-1.17	0.243	-0.045	0.011

• Which gives us the regression equation:

 $Price = 160.924 + 0.172 \cdot sqft - 0.017 \cdot abv$ 

• Let's find the MLR model

house\_sqft\_abv\_mod <-lm(price ~ sqft\_living + sqft\_above, data = house)</pre>

```
And investigate the regression table
get_regression_table(house_sqft_abv_mod)
```

##	#	A tibble: 3	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	161.	14.8	10.9	0	132.	190.
##	2	sqft_living	0.172	0.014	12.6	0	0.145	0.199
##	З	sqft_above	-0.017	0.014	-1.17	0.243	-0.045	0.011

• Which gives us the regression equation:

 $Price = 160.924 + 0.172 \cdot sqft - 0.017 \cdot abv$ 

• Increasing total footage 1 ft, while keeping above ground fixed, *increases* Price by an average of \$0.1724.

• Let's find the MLR model

house\_sqft\_abv\_mod <-lm(price ~ sqft\_living + sqft\_above, data = house)</pre>

```
And investigate the regression table
get_regression_table(house_sqft_abv_mod)
```

##	#	A tibble: 3	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	161.	14.8	10.9	0	132.	190.
##	2	sqft_living	0.172	0.014	12.6	0	0.145	0.199
##	З	sqft_above	-0.017	0.014	-1.17	0.243	-0.045	0.011

• Which gives us the regression equation:

 $Price = 160.924 + 0.172 \cdot sqft - 0.017 \cdot abv$ 

- Increasing total footage 1 ft, while keeping above ground fixed, *increases* Price by an average of \$0.1724.
- Increasing above ground footage 1 ft, while keeping total footage fixed, *decreases* Price by an average of \$0.017.

Wait...

Wait...

• The SLR for Price and Above Ground Square Footage was

 $\hat{\mathrm{Price}} = 236.53 + 0.13 \cdot \mathrm{abv}$ 

Wait...

• The SLR for Price and Above Ground Square Footage was

 $\hat{\mathrm{Price}} = 236.53 + 0.13 \cdot \mathrm{abv}$ 

• That is, increasing above ground square footage by 1 ft INCREASED price by \$0.13.

Wait...

$$Price = 236.53 + 0.13 \cdot abv$$

- That is, increasing above ground square footage by 1 ft INCREASED price by \$0.13.
- But the MLR is

$$Price = 160.924 + 0.172 \cdot sqft - 0.017 \cdot abv$$

Wait...

• The SLR for Price and Above Ground Square Footage was

$$\hat{\rm Price} = 236.53 + 0.13 \cdot {\rm abv}$$

- That is, increasing above ground square footage by 1 ft INCREASED price by \$0.13.
- But the MLR is

$$Price = 160.924 + 0.172 \cdot sqft - 0.017 \cdot abv$$

 Not only has MLR given us a new rate of change, but it's completely switched the direction!

Wait...

$$\hat{\rm Price} = 236.53 + 0.13 \cdot {\rm abv}$$

- That is, increasing above ground square footage by 1 ft INCREASED price by \$0.13.
- But the MLR is

$$Price = 160.924 + 0.172 \cdot sqft - 0.017 \cdot abv$$

- Not only has MLR given us a new rate of change, but it's completely switched the direction!
- How is this possible?

Wait...

$$\hat{\rm Price} = 236.53 + 0.13 \cdot {\rm abv}$$

- That is, increasing above ground square footage by 1 ft INCREASED price by \$0.13.
- But the MLR is

$$Price = 160.924 + 0.172 \cdot sqft - 0.017 \cdot abv$$

- Not only has MLR given us a new rate of change, but it's completely switched the direction!
- How is this possible?
  - Basements are expensive in Seattle. Why?

Wait...

$$\hat{\rm Price} = 236.53 + 0.13 \cdot {\rm abv}$$

- That is, increasing above ground square footage by 1 ft INCREASED price by \$0.13.
- But the MLR is

$$Price = 160.924 + 0.172 \cdot sqft - 0.017 \cdot abv$$

- Not only has MLR given us a new rate of change, but it's completely switched the direction!
- How is this possible?
  - Basements are expensive in Seattle. Why?
  - Seattle is hilly, with firm clay soil, making it more difficult to excavate
  - Could basements be associated with other desirable housing attributes?

• Let's consider the relationship between above ground and total square footage

• Let's consider the relationship between above ground and total square footage



• Let's consider the relationship between above ground and total square footage



Total vs Above Ground Square Footage

• In a vacuum, as total square footage increases, so too does above ground square footage

• Let's consider the relationship between above ground and total square footage



- In a vacuum, as total square footage increases, so too does above ground square footage
  - So in the SLR model, when we look at change in price due to increase in above ground square footage, we are implicitly also increasing total square footage too.

• Let's consider the relationship between above ground and total square footage



- In a vacuum, as total square footage increases, so too does above ground square footage
  - So in the SLR model, when we look at change in price due to increase in above ground square footage, we are implicitly also increasing total square footage too.
  - We could say total square footage is a confounding variable in the SLR model.

• Let's consider the relationship between above ground and total square footage



- In a vacuum, as total square footage increases, so too does above ground square footage
  - So in the SLR model, when we look at change in price due to increase in above ground square footage, we are implicitly also increasing total square footage too.
  - We could say total square footage is a confounding variable in the SLR model.
  - The MLR model allows us to *control* for this confounding variable

• Let's convert above ground square footage to a categorical variable (by grouping into 7 levels with roughly the same number of houses each)

• Let's convert above ground square footage to a categorical variable (by grouping into 7 levels with roughly the same number of houses each)



• Let's convert above ground square footage to a categorical variable (by grouping into 7 levels with roughly the same number of houses each)



• While price has a positive overall relationship with above ground square footage, within each band of total square footage, price has a weakly negative relationship

• Let's convert above ground square footage to a categorical variable (by grouping into 7 levels with roughly the same number of houses each)



- While price has a positive overall relationship with above ground square footage, within each band of total square footage, price has a weakly negative relationship
  - This is an example of **Simpson's Paradox**: a trend present in the aggregate data can reverse itself when data is considered by group.

## Assessing Strength of Multilinear Models

• For SLR, we used the correlation coefficient R to assess model strength.
- For SLR, we used the correlation coefficient R to assess model strength.
- We also saw that R<sup>2</sup> had a natural interpretation: the percentage of variability in the response due to linear relationship with explanatory variable.

- For SLR, we used the correlation coefficient R to assess model strength.
- We also saw that R<sup>2</sup> had a natural interpretation: the percentage of variability in the response due to linear relationship with explanatory variable.
- For MLR, we cannot define the correlation coefficient, because we have multiple explanatory variables.

- For SLR, we used the correlation coefficient R to assess model strength.
- We also saw that R<sup>2</sup> had a natural interpretation: the percentage of variability in the response due to linear relationship with explanatory variable.
- For MLR, we cannot define the correlation coefficient, because we have multiple explanatory variables.
- However, we can still define  $R^2$ !

 $R^2 = \frac{\text{variability in response explained by model}}{\text{variability in response}} = \frac{s_y^2 - s_{\text{res}}^2}{s_y^2}$ 

- For SLR, we used the correlation coefficient R to assess model strength.
- We also saw that R<sup>2</sup> had a natural interpretation: the percentage of variability in the response due to linear relationship with explanatory variable.
- For MLR, we cannot define the correlation coefficient, because we have multiple explanatory variables.
- However, we can still define R<sup>2</sup>!

$$R^{2} = \frac{\text{variability in response explained by model}}{\text{variability in response}} = \frac{s_{y}^{2} - s_{\text{res}}^{2}}{s_{y}^{2}}$$

• Usually, we use software to compute  $R^2$  for multivariate models

- For SLR, we used the correlation coefficient R to assess model strength.
- We also saw that R<sup>2</sup> had a natural interpretation: the percentage of variability in the response due to linear relationship with explanatory variable.
- For MLR, we cannot define the correlation coefficient, because we have multiple explanatory variables.
- However, we can still define R<sup>2</sup>!

$$R^{2} = \frac{\text{variability in response explained by model}}{\text{variability in response}} = \frac{s_{y}^{2} - s_{\text{res}}^{2}}{s_{y}^{2}}$$

• Usually, we use software to compute  $R^2$  for multivariate models house\_sqft\_abv\_mod <- lm(price ~ sqft\_living + sqft\_above, data = house) get\_regression\_summaries(house\_sqft\_abv\_mod)

```
## # A tibble: 1 x 9
    r_squared adj_r_squared mse rmse sigma statistic p_value
##
                                                                 df nobs
        <dbl>
                    <dbl> <dbl> <dbl> <dbl> <dbl>
                                                 <dbl>
                                                         <dbl> <dbl> <dbl>
##
        0.309
                     0.308 24397. 156. 156.
                                                  223.
                                                                  2
                                                                    1000
## 1
                                                             0
```

• Can we build a multivariate model that explains a higher proportion of the variability in price?

Can we build a multivariate model that explains a higher proportion of the variability in price?

Can we build a multivariate model that explains a higher proportion of the variability in price?

• Can we build a multivariate model that explains a higher proportion of the variability in price?

get\_regression\_table(price\_big\_mod)

```
## # A tibble: 9 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	3661.	404.	9.06	0	2868.	4453.
2	bedrooms	-19.7	6.82	-2.88	0.004	-33.0	-6.29
3	bathrooms	30.0	11.2	2.67	0.008	7.91	52.0
4	sqft_living	0.158	0.016	9.93	0	0.127	0.189
5	sqft_above	0.039	0.014	2.74	0.006	0.011	0.066
6	sqft_lot	-0.014	0.002	-8.57	0	-0.017	-0.011
7	view	50.4	8.61	5.85	0	33.5	67.3
8	condition	11.8	7.48	1.58	0.114	-2.84	26.5
9	yr_built	-1.78	0.205	-8.67	0	-2.18	-1.38
	1 2 3 4 5 6 7 8 9	<pre>term <chr> 1 intercept 2 bedrooms 3 bathrooms 4 sqft_living 5 sqft_above 6 sqft_lot 7 view 8 condition 9 yr_built</chr></pre>	term         estimate <chr> <dbl>           1 intercept         3661.           2 bedrooms         -19.7           3 bathrooms         30.0           4 sqft_living         0.158           5 sqft_above         0.039           6 sqft_lot         -0.014           7 view         50.4           8 condition         11.8           9 yr_built         -1.78</dbl></chr>	term         estimate std_error <ch>&gt; <dbl>           1 intercept         3661.         404.           2 bedrooms         -19.7         6.82           3 bathrooms         30.0         11.2           4 sqft_living         0.158         0.016           5 sqft_above         0.039         0.014           6 sqft_lot         -0.014         0.002           7 view         50.4         8.61           8 condition         11.8         7.48           9 yr_built         -1.78         0.205</dbl></ch>	term         estimate std_error statistic <chr> <dbl><dbl><dbl><dbl><dbl></dbl>            1 intercept         3661.         404.         9.06           2 bedrooms         -19.7         6.82         -2.88           3 bathrooms         30.0         11.2         2.67           4 sqft_living         0.158         0.016         9.93           5 sqft_above         0.039         0.014         2.74           6 sqft_lot         -0.014         0.002         -8.57           7 view         50.4         8.61         5.85           8 condition         11.8         7.48         1.58           9 yr_built         -1.78         0.205         -8.67</dbl></dbl></dbl></dbl></chr>	term         estimate std_error statistic         p_value <ch>&gt; <dbl> <dbl> <dbl>           1 intercept         3661.         404.         9.06         0           2 bedrooms         -19.7         6.82         -2.88         0.004           3 bathrooms         30.0         11.2         2.67         0.008           4 sqft_living         0.158         0.016         9.93         0           5 sqft_above         0.039         0.014         2.74         0.006           6 sqft_lot         -0.014         0.002         -8.57         0           7 view         50.4         8.61         5.85         0           8 condition         11.8         7.48         1.58         0.114           9 yr_built         -1.78         0.205         -8.67         0</dbl></dbl></dbl></ch>	term         estimate std_error statistic         p_value         lower_ci <chr> <chr> <dbl></dbl> <dbl></dbl> <dbl></dbl> <dbl></dbl>           1 intercept         3661.         404.         9.06         0         2868.           2 bedrooms         -19.7         6.82         -2.88         0.004         -33.0           3 bathrooms         30.0         11.2         2.67         0.008         7.91           4 sqft_living         0.158         0.016         9.93         0         0.127           5 sqft_above         0.039         0.014         2.74         0.006         0.011           6 sqft_lot         -0.014         0.002         -8.57         0         -0.017           7 view         50.4         8.61         5.85         0         33.5           8 condition         11.8         7.48         1.58         0.114         -2.84           9 yr_built         -1.78         0.205         -8.67         0         -2.18</chr></chr>

get\_regression\_summaries(price\_big\_mod)

##	#	A tibble:	1 x 9							
##		r_squared	adj_r_squared	mse	rmse	sigma	statistic	p_value	df	nobs
##		<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	0.434	0.429	19987.	141.	142.	94.9	0	8	1000