# Lecture 6: Correlation and Linear Regression Basics

## ECON 480 - Econometrics - Fall 2018

Ryan Safner

September 17, 2018

# Covariance and Correlation

- We looked at single variables for descriptive statistics

- We looked at single variables for descriptive statistics
- Most uses of statistics in economics and business investigate relationships *between* variables

- We looked at single variables for descriptive statistics
- Most uses of statistics in economics and business investigate relationships *between* variables
    - # of police & crime rates

- We looked at single variables for descriptive statistics
- Most uses of statistics in economics and business investigate relationships *between* variables
    - # of police & crime rates
    - healthcare spending & life expectancy

- We looked at single variables for descriptive statistics
- Most uses of statistics in economics and business investigate relationships *between* variables
    - # of police & crime rates
    - healthcare spending & life expectancy
    - government spending & GDP growth

HOOD
COLLEGE

- We looked at single variables for descriptive statistics
- Most uses of statistics in economics and business investigate relationships *between* variables

    - # of police & crime rates
    - healthcare spending & life expectancy
    - government spending & GDP growth
    - carbon dioxide emissions & temperatures

- We looked at single variables for descriptive statistics
- Most uses of statistics in economics and business investigate relationships *between* variables

    - # of police & crime rates
    - healthcare spending & life expectancy
    - government spending & GDP growth
    - carbon dioxide emissions & temperatures

- We will begin with bivariate data for relationships between *X* and *Y*

HOOD
COLLEGE

- We looked at single variables for descriptive statistics
- Most uses of statistics in economics and business investigate relationships *between* variables

  - # of police & crime rates
  - healthcare spending & life expectancy
  - government spending & GDP growth
  - carbon dioxide emissions & temperatures

- We will begin with bivariate data for relationships between $X$ and $Y$

  - Immediate aim is to explore associations between variables, quantified with **correlation** and **linear regression**

HOOD COLLEGE

- We looked at single variables for descriptive statistics
- Most uses of statistics in economics and business investigate relationships *between* variables

  - # of police & crime rates
  - healthcare spending & life expectancy
  - government spending & GDP growth
  - carbon dioxide emissions & temperatures

- We will begin with bivariate data for relationships between $X$ and $Y$

  - Immediate aim is to explore associations between variables, quantified with **correlation** and **linear regression**
  - Later we want to develop more sophisticated tools to argue for **causation**

```
econfreedom<-read.csv("~/Dropbox/Teaching/Hood College/ECON 480 - Econometrics/D
head(econfreedom)
```

```
##   ISO.Code       Country Economic.Freedom.Summary.Index GDP.Per.Capita
## 1      AGO        Angola                           5.08       4153.146
## 2      ALB       Albania                           7.40       4543.088
## 3      ARE Unit. Arab Em.                          7.98      39313.274
## 4      ARG     Argentina                           4.81      10501.660
## 5      ARM       Armenia                           7.71       3796.517
## 6      AUS     Australia                           7.93      54688.446
```

- **Rows** are individual observations

```
econfreedom<-read.csv("~/Dropbox/Teaching/Hood College/ECON 480 - Econometrics/D
head(econfreedom)
```

```
##   ISO.Code        Country Economic.Freedom.Summary.Index GDP.Per.Capita
## 1     AGO          Angola                           5.08       4153.146
## 2     ALB         Albania                           7.40       4543.088
## 3     ARE Unit. Arab Em.                            7.98      39313.274
## 4     ARG       Argentina                           4.81      10501.660
## 5     ARM         Armenia                           7.71       3796.517
## 6     AUS       Australia                           7.93      54688.446
```

- **Rows** are individual observations
- **Columns** are variables on all individuals

```
econfreedom<-read.csv("~/Dropbox/Teaching/Hood College/ECON 480 - Econometrics/D
head(econfreedom)
```

```
##   ISO.Code        Country Economic.Freedom.Summary.Index GDP.Per.Capita
## 1     AGO         Angola                            5.08       4153.146
## 2     ALB        Albania                            7.40       4543.088
## 3     ARE Unit. Arab Em.                            7.98      39313.274
## 4     ARG      Argentina                            4.81      10501.660
## 5     ARM        Armenia                            7.71       3796.517
## 6     AUS      Australia                            7.93      54688.446
```

- **Rows** are individual observations
- **Columns** are variables on all individuals
- Let $X$ be Economic Freedom and $Y$ be GDP per capita

4

```r
str(econfreedom)
```

```
## 'data.frame':    152 obs. of  4 variables:
##  $ ISO.Code                      : Factor w/ 152 levels "AGO","ALB","ARE",..:
##  $ Country                       : Factor w/ 152 levels "Albania","Algeria",
##  $ Economic.Freedom.Summary.Index: num  5.08 7.4 7.98 4.81 7.71 7.93 7.56 6.5
##  $ GDP.Per.Capita                : num  4153 4543 39313 10502 3797 ...
```

```r
summary(econfreedom)
```

```
##     ISO.Code        Country      Economic.Freedom.Summary.Index
##  AGO    :  1    Albania  :  1    Min.   :4.800
##  ALB    :  1    Algeria  :  1    1st Qu.:6.430
##  ARE    :  1    Angola   :  1    Median :7.050
##  ARG    :  1    Argentina:  1    Mean   :6.909
##  ARM    :  1    Armenia  :  1    3rd Qu.:7.428
##  AUS    :  1    Australia:  1    Max.   :9.030
##  (Other):146    (Other)  :146
##  GDP.Per.Capita
##  Min.   :   206.7
##  1st Qu.:  1588.3
##  Median :  5719.3
```

```
# syntax for plotting is similar to hist() and boxplot()
# just tell R "plot(df$x,df$y)"
plot(econfreedom$Economic.Freedom.Summary.Index, econfreedom$GDP.Per.Capita)
```



- The best way to visualize an association between two variables is with a scatterplot

7

```
# syntax for plotting is similar to hist() and boxplot()
# just tell R "plot(df$x,df$y)"
plot(econfreedom$Economic.Freedom.Summary.Index, econfreedom$GDP.Per.Capita)
```



- Each point is a pair of variable values $(X_i, Y_i)$ for observation $i$

```r
library("ggplot2")
ggplot(econfreedom, aes(x=Economic.Freedom.Summary.Index,y=GDP.Per.Capita))+
  geom_point(color="blue")+theme_bw()+
  xlab("Economic Freedom Index (2014)")+ylab("GDP per Capita (2014 USD)")
```

- Look for **association** between independent and dependent variables

- Look for **association** between independent and dependent variables

  1. *Direction*: is the trend positive or negative?

- Look for **association** between independent and dependent variables
    1. *Direction*: is the trend positive or negative?
    2. *Form*: is the trend linear, quadratic, something else, or no pattern?

- Look for **association** between independent and dependent variables

    1. *Direction*: is the trend positive or negative?
    2. *Form*: is the trend linear, quadratic, something else, or no pattern?
    3. *Strength*: is the association strong or weak?

HOOD
COLLEGE

- Look for **association** between independent and dependent variables

  1. *Direction*: is the trend positive or negative?
  2. *Form*: is the trend linear, quadratic, something else, or no pattern?
  3. *Strength*: is the association strong or weak?
  4. *Outliers*: do any observations break the trends above?

HOOD
COLLEGE

- For any two variables, we can measure their sample covariance, $cov(X, Y)$ or $s_{X,Y}$ to quantify how they vary *together*[1]

$$s_{X,Y} = E\big[(X - \bar{X})(Y - \bar{Y})\big]$$

---

[1]Henceforth we limit to samples, for convenience. Population covariance is denoted $\sigma_{X,Y}$

- For any two variables, we can measure their sample covariance, $cov(X, Y)$ or $s_{X,Y}$ to quantify how they vary *together*[1]

$$s_{X,Y} = E\big[(X - \bar{X})(Y - \bar{Y})\big]$$

- Intuition: if $X$ is above its mean, would we expect $Y$:

---

[1]Henceforth we limit to samples, for convenience. Population covariance is denoted $\sigma_{X,Y}$

- For any two variables, we can measure their sample covariance, $cov(X, Y)$ or $s_{X,Y}$ to quantify how they vary *together*[1]

$$s_{X,Y} = E\big[(X - \bar{X})(Y - \bar{Y})\big]$$

- Intuition: if $X$ is above its mean, would we expect $Y$:

  - to be *above* its mean also ($X$ and $Y$ covary *positively*)

---

[1] Henceforth we limit to samples, for convenience. Population covariance is denoted $\sigma_{X,Y}$

- For any two variables, we can measure their sample covariance, $cov(X, Y)$ or $s_{X,Y}$ to quantify how they vary *together*[1]

$$s_{X,Y} = E\big[(X - \bar{X})(Y - \bar{Y})\big]$$

- Intuition: if $X$ is above its mean, would we expect $Y$:

  - to be *above* its mean also ($X$ and $Y$ covary *positively*)
  - to be *below* its mean ($X$ and $Y$ covary *negatively*)

HOOD
C O L L E G E

---

[1]Henceforth we limit to samples, for convenience. Population covariance is denoted $\sigma_{X,Y}$

- For any two variables, we can measure their sample covariance, $cov(X, Y)$ or $s_{X,Y}$ to quantify how they vary *together*[1]

$$s_{X,Y} = E\big[(X - \bar{X})(Y - \bar{Y})\big]$$

- Intuition: if $X$ is above its mean, would we expect $Y$:
  - to be *above* its mean also ($X$ and $Y$ covary *positively*)
  - to be *below* its mean ($X$ and $Y$ covary *negatively*)

- Covariance is a common measure, but the units are meaningless, thus we rarely need to use it so **don't worry about learning the formula**

HOOD
COLLEGE

---

[1]Henceforth we limit to samples, for convenience. Population covariance is denoted $\sigma_{X,Y}$

- More convenient to standardize covariance into a more intuitive concept: correlation ($\rho$ or $r$), normalized to be between -1 and 1

$$r_{X,Y} = \frac{s_{X,Y}}{s_X s_Y} = \frac{cov(X, Y)}{sd(X)sd(Y)}$$

- More convenient to standardize covariance into a more intuitive concept: correlation ($\rho$ or $r$), normalized to be between -1 and 1

$$r_{X,Y} = \frac{s_{X,Y}}{s_X s_Y} = \frac{cov(X, Y)}{sd(X)sd(Y)}$$

- Alternatively, sample correlation can be found by standardizing (finding the $Z$-score) $X$ and $Y$ and multiplying, for each $(X, Y)$ pair, and then averaging (over $n - 1$, due to sampling df, again):

$$r = \frac{1}{n - 1} \sum_{i=1}^{n} \left( \frac{X_i - \bar{X}}{s_X} \right) \left( \frac{Y_i - \bar{Y}}{s_Y} \right)$$

$$= \frac{1}{n - 1} \sum_{i=1}^{n} Z_X Z_Y$$

Example

$$(1,1), (2,2), (3,4), (4,9)$$

```
corr.example<-data.frame(x=c(1,2,3,4),
                         y=c(1,2,4,9))
ggplot(corr.example,aes(x=x,y=y))+geom_point()
```

```r
mean(corr.example$x) #find mean of x
```

```
## [1] 2.5
```

```r
mean(corr.example$y) #find mean of y
```

```
## [1] 4
```

```r
sd(corr.example$x) #find sd of x
```

```
## [1] 1.290994
```

```r
sd(corr.example$y) #find sd of y
```

```
## [1] 3.559026
```

```
#take z score of x,y for each pair and multiply them
corr.example$z.product<-(((corr.example$x-2.5)/1.291)*
                          ((corr.example$y-4)/3.559))

corr.example


##   x y z.product
## 1 1 1 0.9793959
## 2 2 2 0.2176435
## 3 3 4 0.0000000
## 4 4 9 1.6323265
```

```r
(sum(corr.example$z.product)/3) #average z products over n-1
```

```
## [1] 0.943122
```

```r
cor(corr.example$x, corr.example$y) #compare our answer to cor() command
```

```
## [1] 0.9431191
```

```r
cov(corr.example$x, corr.example$y) #just for kicks - covariance
```

```
## [1] 4.333333
```

- Correlation is standardized to $-1 \leq r \leq 1$

- Correlation is standardized to $-1 \leq r \leq 1$
  - Negative values $\implies$ negative association

- Correlation is standardized to $-1 \leq r \leq 1$
  - Negative values $\implies$ negative association
  - Positive values $\implies$ positive association

- Correlation is standardized to $-1 \leq r \leq 1$
  - Negative values $\implies$ negative association
  - Positive values $\implies$ positive association
  - Correlation of 0 $\implies$ no association

- Correlation is standardized to $-1 \leq r \leq 1$
    - Negative values $\implies$ negative association
    - Positive values $\implies$ positive association
    - Correlation of 0 $\implies$ no association
    - As $|r| \to 1 \implies$ the stronger the association



| *Perfect Positive Correlation* | *High Positive Correlation* | *Low Positive Correlation* | *No Correlation* | *Low Negative Correlation* | *High Negative Correlation* | *Perfect Negative Correlation* |
|---|---|---|---|---|---|---|
| 1 | 0.9 | 0.5 | 0 | -0.5 | -0.9 | -1 |

- Correlation is standardized to $-1 \leq r \leq 1$
  - Negative values $\implies$ negative association
  - Positive values $\implies$ positive association
  - Correlation of 0 $\implies$ no association
  - As $|r| \to 1 \implies$ the stronger the association
  - Correlation of $|r| = 1 \implies$ a perfect linear relationship



| *Perfect Positive Correlation* | *High Positive Correlation* | *Low Positive Correlation* | *No Correlation* | *Low Negative Correlation* | *High Negative Correlation* | *Perfect Negative Correlation* |
|---|---|---|---|---|---|---|
| 1 | 0.9 | 0.5 | 0 | -0.5 | -0.9 | -1 |

Guess The Correlation Game

- Reminder: Correlation does not imply causation!

- Reminder: Correlation does not imply causation!
- See the **Handout** for more on Covariance and Correlation

## Example



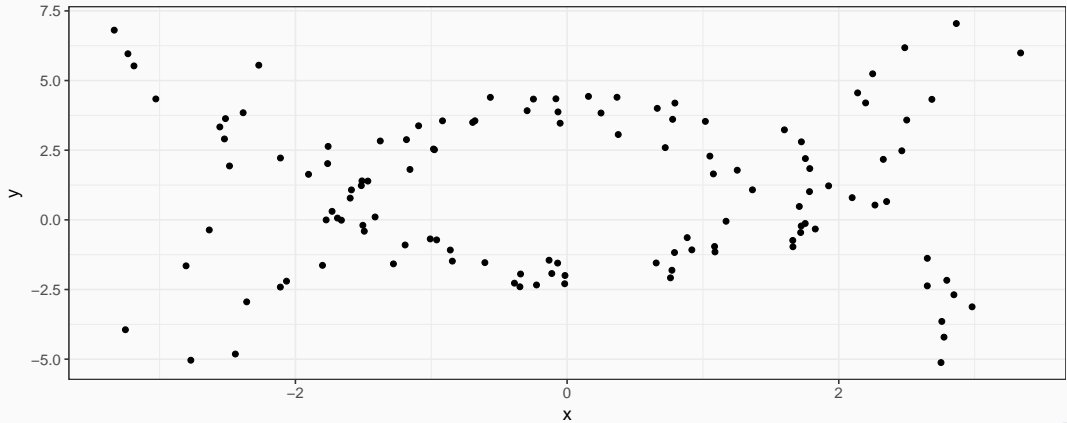- The correlation between Life Expectancy and Doctors Per Person is 0.705.

### Example



- The correlation between Life Expectancy and Doctors Per Person is 0.705.
- So should we send more doctors to developing countries to increase their life expectancy?

### Example



- The correlation between Life Expectancy and Doctors Per Person is 0.705.
- So should we send more doctors to developing countries to increase their life expectancy?
- Properly interpreting relationships requires both statistical *and* economic intuition!

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##       4.0     6.5     9.0     9.0    11.5    14.0
```

```
##    dataset              x              y
## Length:1846     Min.   :15.56   Min.   : 0.01512
## Class :character 1st Qu.:41.07   1st Qu.:22.56107
## Mode  :character Median :52.59   Median :47.59445
##                  Mean   :54.27   Mean   :47.83510
##                  3rd Qu.:67.28   3rd Qu.:71.81078
##                  Max.   :98.29   Max.   :99.69468
```

See the Datasaurus

# Population Linear Regression Model

- If an association appears linear, we can estimate the equation of a line that would "fit" the data

- If an association appears linear, we can estimate the equation of a line that would "fit" the data

$$Y = a + bX$$

- Recall a linear equation describing a line contains: - $a$: vertical intercept - $b$: slope - Note we will use different symbols for $a$ and $b$, in line with standard econometric notation

HOOD
COLLEGE

- How do we choose the equation that best fits the data? Process is called linear regression

- Linear regression lets us estimate the slope of the population regression line between *X* and *Y*

- Linear regression lets us estimate the slope of the population regression line between *X* and *Y*
- We can make **inferences** about the population slope coefficient

- Linear regression lets us estimate the slope of the population regression line between *X* and *Y*
- We can make **inferences** about the population slope coefficient
    - eventually, a causal interpretation

- Linear regression lets us estimate the slope of the population regression line between *X* and *Y*
- We can make **inferences** about the population slope coefficient
    - eventually, a causal interpretation
    - slope $= \frac{\Delta y}{\Delta x}$: for a 1-unit change in *X*, how many units will this *cause Y* to change?

- Statistically, we want to use the population regression model for:

· Statistically, we want to use the population regression model for:

1. Estimation of the marginal effect of $X$ on $Y$ (slope of population regression line)

- Statistically, we want to use the population regression model for:

1. Estimation of the marginal effect of $X$ on $Y$ (slope of population regression line)
2. Hypothesis Testing of the value of the marginal effect (slope)

- Statistically, we want to use the population regression model for:

1. Estimation of the marginal effect of $X$ on $Y$ (slope of population regression line)
2. Hypothesis Testing of the value of the marginal effect (slope)
3. Confidence Interval construction of a range for the true effect (slope)

**Example**

What is the relationship between class size and educational performance?

- Policy question: What is the effect of reducing class sizes by 1 student per class on test scores? 10 students?

```r
library("foreign") #for importing .dta files
CASchool<-read.dta("~/Dropbox/Teaching/Hood College/ECON 480 - Econometrics/Data

ca.scatter<-ggplot(CASchool, aes(str,testscr))+
  geom_point(color="blue",fill="blue")+
  xlab("Student to Teacher Ratio")+
  ylab("Test Score")+theme_bw()
```

- If we *change* ($\Delta$) the class size by an amount, what would we expect the *change* in test scores to be?

$$\beta_{ClassSize} = \frac{\text{change in test score}}{\text{change in class size}} = \frac{\Delta\text{test score}}{\Delta\text{class size}}$$

· If we *change* ($\Delta$) the class size by an amount, what would we expect the *change* in test scores to be?

$$\beta_{ClassSize} = \frac{\text{change in test score}}{\text{change in class size}} = \frac{\Delta \text{test score}}{\Delta \text{class size}}$$

· If we knew $\beta_{ClassSize}$, we could say that changing class size by 1 student will change test scores by $\beta_{ClassSize}$

- Rearranging:

$$\Delta\text{test score} = \beta_{ClassSize} \times \Delta\text{class size}$$

- Rearranging:

$$\Delta\text{test score} = \beta_{ClassSize} \times \Delta\text{class size}$$

- Suppose $\beta_{ClassSize} = -0.6$. If we shrank class size by 2 students, our model predicts:

- Rearranging:

$$\Delta\text{test score} = \beta_{ClassSize} \times \Delta\text{class size}$$

- Suppose $\beta_{ClassSize} = -0.6$. If we shrank class size by 2 students, our model predicts:

$$\Delta\text{test score} = -0.6$$
$$\Delta\text{test score} = \times -2 = 1.2$$

$$\text{test score} = \beta_0 + \beta_{ClassSize} \times \text{class size}$$

- The line relating class size and test scores has the above equation

$$\text{test score} = \beta_0 + \beta_{ClassSize} \times \text{class size}$$

- The line relating class size and test scores has the above equation
  - $\beta_0$ is the vertical-intercept, test score where class size is 0

$$\text{test score} = \beta_0 + \beta_{ClassSize} \times \text{class size}$$

- The line relating class size and test scores has the above equation

  - $\beta_0$ is the vertical-intercept, test score where class size is 0
  - $\beta_{ClassSize}$ is the slope of the regression line

$$\text{test score} = \beta_0 + \beta_{ClassSize} \times \text{class size}$$

- The line relating class size and test scores has the above equation

  - $\beta_0$ is the vertical-intercept, test score where class size is 0
  - $\beta_{ClassSize}$ is the slope of the regression line

- This relationship only holds **on average** for all districts in the population, individual districts are also affected by other factors

- To get an equation that holds for *each* district, we need to include other factors

$$\text{test score} = \beta_0 + \beta_{ClassSize} \times \text{class size} + \text{other factors}$$

- To get an equation that holds for *each* district, we need to include other factors

$$\text{test score} = \beta_0 + \beta_{ClassSize} \times \text{class size} + \text{other factors}$$

- For now, we will ignore these until the next lesson

- To get an equation that holds for *each* district, we need to include other factors

$$\text{test score} = \beta_0 + \beta_{ClassSize} \times \text{class size} + \text{other factors}$$

- For now, we will ignore these until the next lesson
- Thus, $\beta_0 + \beta_{ClassSize} \times$ class size gives the **average effect** of class sizes on scores

- To get an equation that holds for *each* district, we need to include other factors

$$\text{test score} = \beta_0 + \beta_{ClassSize} \times \text{class size} + \text{other factors}$$

- For now, we will ignore these until the next lesson
- Thus, $\beta_0 + \beta_{ClassSize} \times$ class size gives the **average effect** of class sizes on scores
- Later, we will want to estimate the **marginal effect** (**causal effect**) of each factor on an individual district's test score, holding all other factors constant

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

- $y$ is the dependent variable of interest
  - AKA "response variable," "regressand," "Left-hand side (LHS) variable"
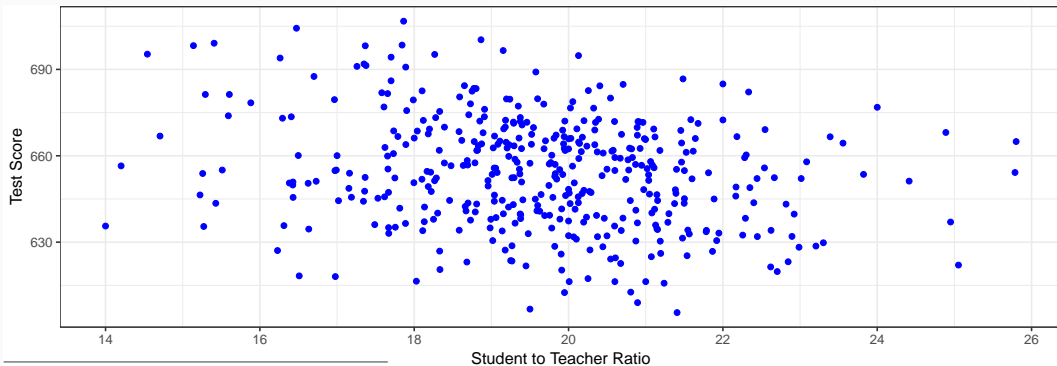
$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

- $y$ is the dependent variable of interest
  - AKA "response variable," "regressand," "Left-hand side (LHS) variable"
- $x_1$ and $x_2$ are independent variables
  - AKA "explanatory variables," "regressors," "Right-hand side (RHS) variables," "covariates," "control variables"

HOOD
COLLEGE

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

- $y$ is the dependent variable of interest
  - AKA "response variable," "regressand," "Left-hand side (LHS) variable"
- $x_1$ and $x_2$ are independent variables
  - AKA "explanatory variables," "regressors," "Right-hand side (RHS) variables," "covariates," "control variables"
- We have observed values of $y$, $x_1$, and $x_2$ & "regress $y$ on $x_1$ and $x_2$"

HOOD
COLLEGE

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

- $y$ is the dependent variable of interest
  - AKA "response variable," "regressand," "Left-hand side (LHS) variable"
- $x_1$ and $x_2$ are independent variables
  - AKA "explanatory variables," "regressors," "Right-hand side (RHS) variables," "covariates," "control variables"
- We have observed values of $y$, $x_1$, and $x_2$ & "regress $y$ on $x_1$ and $x_2$"
- $\beta_0$, $\beta_1$, and $\beta_2$ are unknown parameters to *estimate*

HOOD
COLLEGE

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

- $y$ is the dependent variable of interest
  - AKA "response variable," "regressand," "Left-hand side (LHS) variable"
- $x_1$ and $x_2$ are independent variables
  - AKA "explanatory variables," "regressors," "Right-hand side (RHS) variables," "covariates," "control variables"
- We have observed values of $y$, $x_1$, and $x_2$ & "regress $y$ on $x_1$ and $x_2$"
- $\beta_0$, $\beta_1$, and $\beta_2$ are unknown parameters to *estimate*
- $\epsilon$ is the error term
  - It is **stochastic** (random)
  - We can never measure the error term

HOOD
COLLEGE

· How do we draw a line through the scatterplot? We do not know the true $\beta_{ClassSize}$

- How do we draw a line through the scatterplot? We do not know the true $\beta_{ClassSize}$
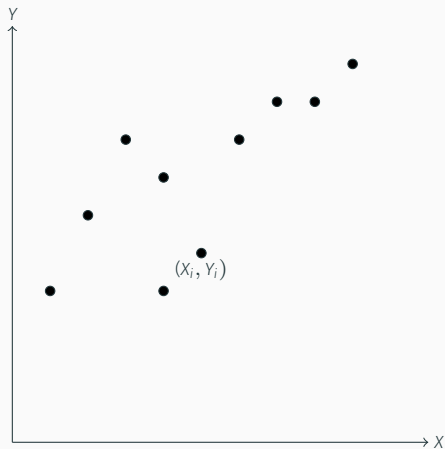- We do have data from a *sample* of class sizes and test scores[2]



[2] Data is student-teacher-ratio and average test scores on Stanford 9 Achievement Test for 5th grade students for 420 K-6 and K-8 school districts in California in 1999, (Stock and Watson, 2015: p. 141)
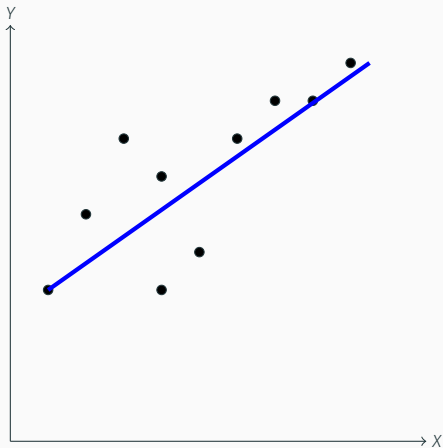
- How do we draw a line through the scatterplot? We do not know the true $\beta_{ClassSize}$
- We do have data from a *sample* of class sizes and test scores[2]
- So the real question is, **how can we estimate $\beta_0$ and $\beta_1$?**



[2] Data is student-teacher-ratio and average test scores on Stanford 9 Achievement Test for 5th grade students for 420 K-6 and K-8 school districts in California in 1999, (Stock and Watson, 2015: p. 141)

# OLS Estimators and Sample Regression Model

· Suppose we have a scatter plot of points $(X_i, Y_i)$

- Suppose we have a scatter plot of points $(X_i, Y_i)$
- We can draw a "line of best fit" through our scatterplot

- Suppose we have a scatter plot of points $(X_i, Y_i)$

- We can draw a "line of best fit" through our scatterplot

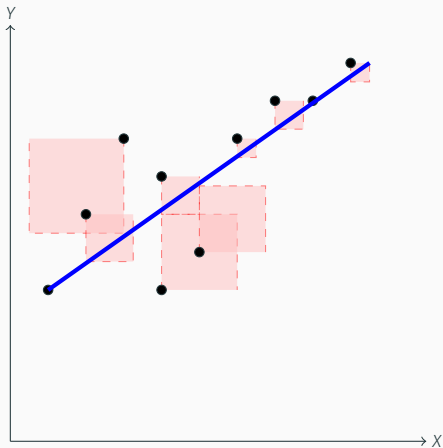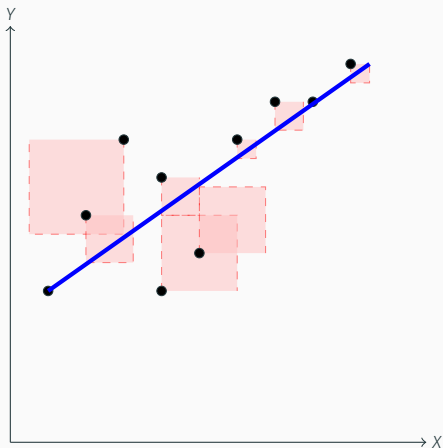- The residual ($\epsilon_i$) of each data point is the difference between **actual** and **predicted value** of $Y$ given $X$

$$\epsilon_i = Y_i - \hat{Y}_i$$

- Suppose we have a scatter plot of points $(X_i, Y_i)$
- We can draw a "line of best fit" through our scatterplot
- The residual ($\epsilon_i$) of each data point is the difference between **actual** and **predicted value** of $Y$ given $X$

$$\epsilon_i = Y_i - \hat{Y}_i$$

- If we were to **square** each residual and add them all up, this is Sum of Squared Errors (SSE)

$$SSE = \sum_{i=1}^{n} \epsilon_i^2$$

- Suppose we have a scatter plot of points $(X_i, Y_i)$

- We can draw a "line of best fit" through our scatterplot

- The residual ($\epsilon_i$) of each data point is the difference between **actual** and **predicted value** of $Y$ given $X$

$$\epsilon_i = Y_i - \hat{Y}_i$$

- If we were to **square** each residual and add them all up, this is Sum of Squared Errors (SSE)

$$SSE = \sum_{i=1}^{n} \epsilon_i^2$$

- The line of best fit **minimizes SSE**

- I coded an example (using an application of R called `shiny`) to demonstrate how OLS tries to solve the problem by picking optimal line parameters

- The ordinary least squares (OLS) estimators of the unknown population parameters $\beta_0$ and $\beta_1$, solve the calculus problem:

- The ordinary least squares (OLS) estimators of the unknown population parameters $\beta_0$ and $\beta_1$, solve the calculus problem:

$$\min_{\beta_0, \beta_1} \sum_{i=1}^{n} [Y_i - (\underbrace{\beta_0 + \beta_1 X_i}_{\hat{Y}_i})]^2$$

- OLS estimators minimize the average squared distance between the actual values ($Y_i$) and the predicted values ($\hat{Y}_i$) along the estimated regression line

HOOD COLLEGE

- The OLS regression line or sample regression line is the linear function constructed using the OLS estimators:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

- The OLS regression line or sample regression line is the linear function constructed using the OLS estimators:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

- $\hat{\beta}_0$ and $\hat{\beta}_1$ ("beta 0 hat" & "beta 1 hat") are the OLS estimators of population parameters $\beta_0$ and $\beta_1$ using sample data

- The OLS regression line or sample regression line is the linear function constructed using the OLS estimators:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

- $\hat{\beta}_0$ and $\hat{\beta}_1$ ("beta 0 hat" & "beta 1 hat") are the OLS estimators of population parameters $\beta_0$ and $\beta_1$ using sample data
- The **predicted value** of Y given X, based on the regression, is $E(Y_i|X_i) = \hat{Y}_i$

- The OLS regression line or sample regression line is the linear function constructed using the OLS estimators:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

- $\hat{\beta}_0$ and $\hat{\beta}_1$ ("beta 0 hat" & "beta 1 hat") are the OLS estimators of population parameters $\beta_0$ and $\beta_1$ using sample data
- The **predicted value** of Y given X, based on the regression, is $E(Y_i|X_i) = \hat{Y}_i$
- The **residual** or **prediction error** for the $i^{th}$ observation is the difference between observed $Y_i$ and its predicted value, $\hat{\epsilon}_i = Y_i - \hat{Y}_i$

HOOD
COLLEGE

- The solution to the SSE minimization problem yields:[3]

---

[3]See **Handout** on Blackboard for proofs.

- The solution to the SSE minimization problem yields:[3]
- For $\hat{\beta}_0$:

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

---

[3]See **Handout** on Blackboard for proofs.

- The solution to the SSE minimization problem yields:[3]
- For $\hat{\beta}_0$:

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$
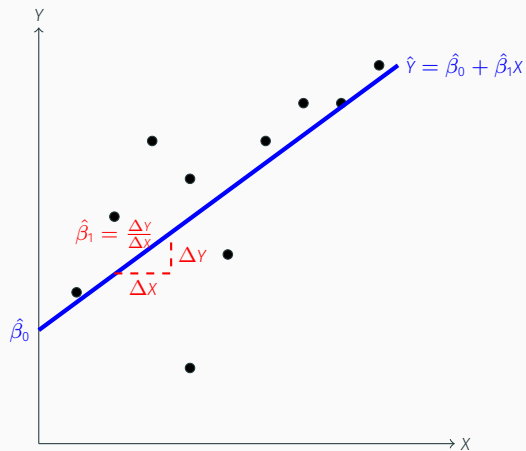
- For $\hat{\beta}_1$:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{n}(X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^{n}(X_i - \bar{X})^2}$$

[3]See **Handout** on Blackboard for proofs.

- The solution to the SSE minimization problem yields:[3]
- For $\hat{\beta}_0$:

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

- For $\hat{\beta}_1$:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{n}(X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^{n}(X_i - \bar{X})^2} = \frac{s_{XY}}{s_X^2}$$

[3] See **Handout** on Blackboard for proofs.

- The solution to the SSE minimization problem yields:[3]
- For $\hat{\beta}_0$:

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

- For $\hat{\beta}_1$:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{n}(X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^{n}(X_i - \bar{X})^2} = \frac{s_{XY}}{s_X^2} = \frac{cov(X, Y)}{var(X)}$$
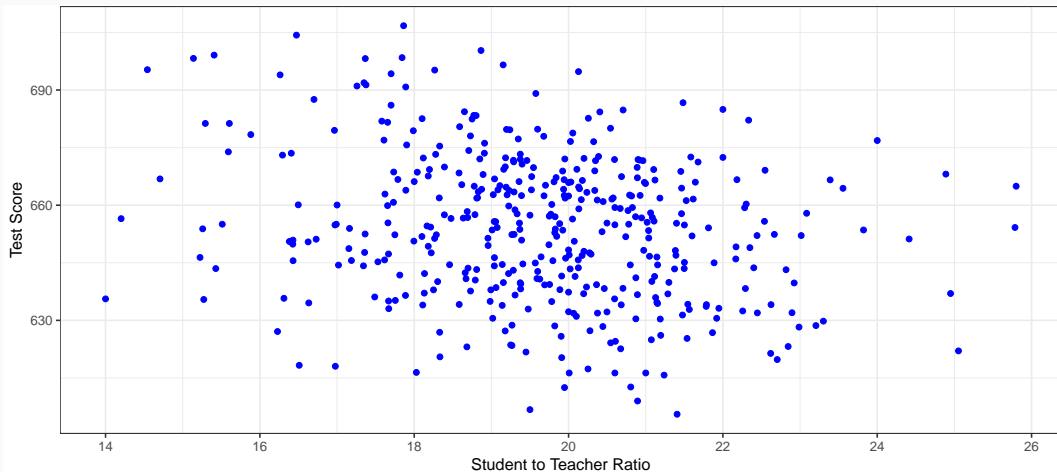
---

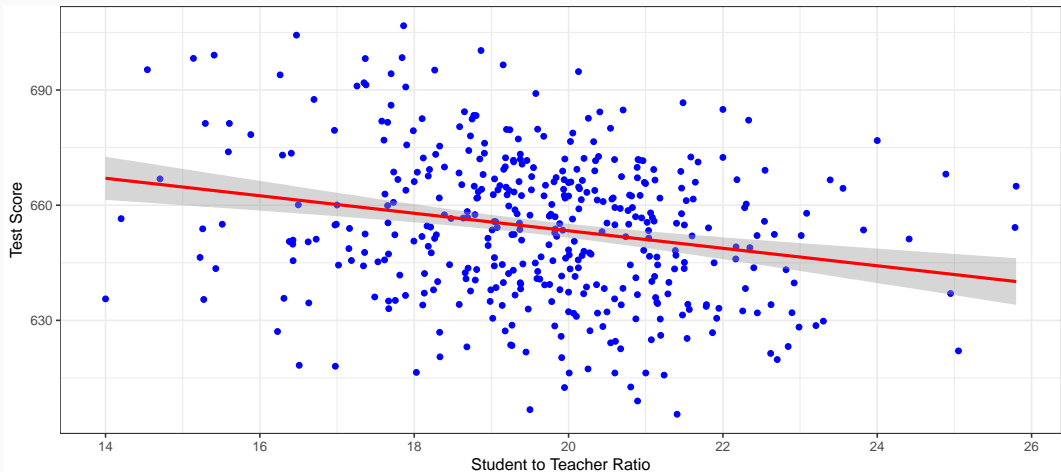[3]See **Handout** on Blackboard for proofs.

HOOD
COLLEGE

- There is some true (unknown) population relationship:

$$\text{Test Score} = \beta_0 + \beta_1 \times STR$$

44

- There is some true (unknown) population relationship:

$$\text{Test Score} = \beta_0 + \beta_1 \times STR$$

44

- Using OLS, we find:

$$\widehat{\text{Test Score}} = 689.9 - 2.28 \times STR$$

$$\widehat{\text{Test Score}} = 689.9 - 2.28 \times STR$$

· Estimated slope: $\hat{\beta}_1 = \frac{\Delta \text{test score}}{\Delta \text{STR}} = -2.28$

$$\widehat{\text{Test Score}} = 689.9 - 2.28 \times STR$$

- Estimated slope: $\hat{\beta}_1 = \frac{\Delta \text{test score}}{\Delta \text{STR}} = -2.28$
- Estimated intercept: $\hat{\beta}_0 = 689.9$

$$\widehat{\text{Test Score}} = 689.9 - 2.28 \times \textit{STR}$$

- Estimated slope: $\hat{\beta}_1 = \frac{\Delta \text{test score}}{\Delta \text{STR}} = -2.28$
- Estimated intercept: $\hat{\beta}_0 = 689.9$
    - Not always economically meaningful

HOOD
COLLEGE

$$\widehat{\text{Test Score}} = 689.9 - 2.28 \times STR$$

- Estimated slope: $\hat{\beta}_1 = \frac{\Delta \text{test score}}{\Delta \text{STR}} = -2.28$
- Estimated intercept: $\hat{\beta}_0 = 689.9$

    - Not always economically meaningful
    - Literally: "districts with 0 students have a predicted test score of 689.9"

$$\widehat{\text{Test Score}} = 689.9 - 2.28 \times STR$$

- We can now make simple predictions with our model:

$$\widehat{\text{Test Score}} = 689.9 - 2.28 \times STR$$

- We can now make simple predictions with our model:

  - For a district with 20 students per teacher, the predicted test score is:

$$689.9 - 2.28(20) = 644.3$$

$$\widehat{\text{Test Score}} = 689.9 - 2.28 \times STR$$

- We can now make simple predictions with our model:

    - For a district with 20 students per teacher, the predicted test score is:

$$689.9 - 2.28(20) = 644.3$$

- Is this big or small? How **economically** meaningful is 644?

· Syntax for running a regression in R is simple:

```
# name an object e.g. "regression.name", "lm" stands for "linear model"
regression.name<-lm(y~x, data=data.frame.name)

# get simple (beta) coefficients by calling the object
regression.name

# get more detailed information with summary()
summary(regression.name)
```

HOOD
COLLEGE

```r
school.regression<-lm(testscr~str, data=CASchool)
school.regression
```

```
##
## Call:
## lm(formula = testscr ~ str, data = CASchool)
##
## Coefficients:
## (Intercept)         str
##      698.93       -2.28
```
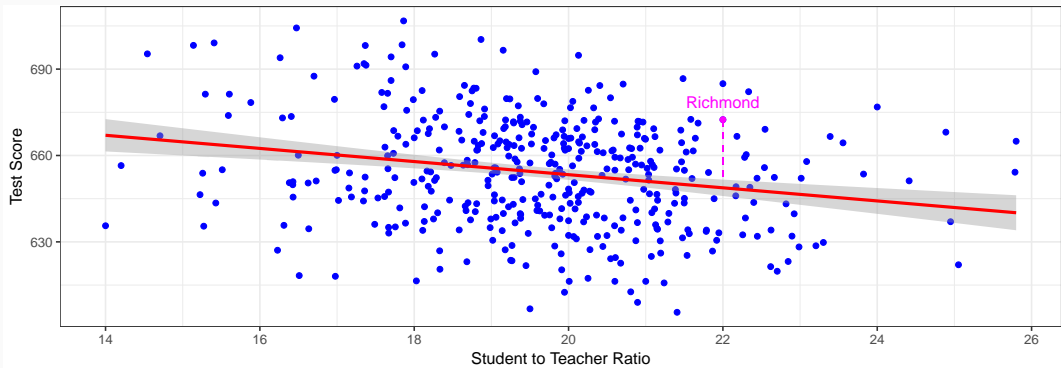
```r
summary(school.regression)
```
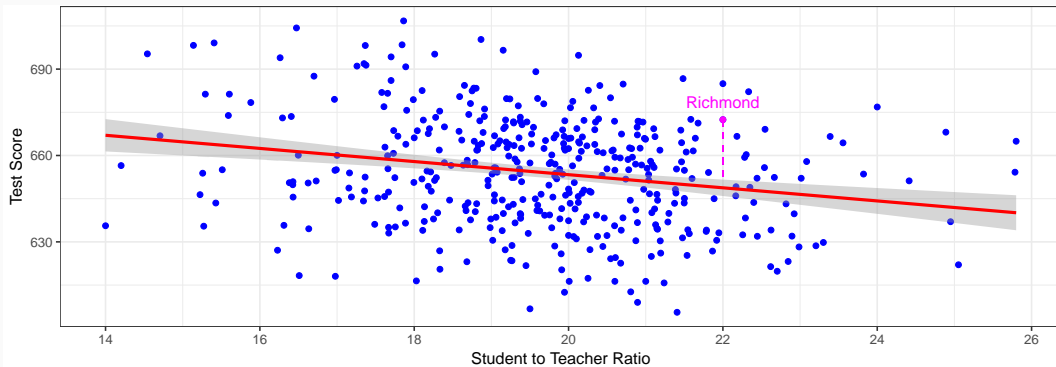
```
##
## Call:
## lm(formula = testscr ~ str, data = CASchool)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -47.727 -14.251   0.483  12.822  48.540
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 698.9330     9.4675  73.825  < 2e-16 ***
## str          -2.2798     0.4798  -4.751 2.78e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 18.58 on 418 degrees of freedom
## Multiple R-squared:  0.05124,    Adjusted R-squared:  0.04897
## F-statistic: 22.58 on 1 and 418 DF,  p-value: 2.783e-06
```

HOOD
COLLEGE

- One district in our sample is Richmond, CA with STR=22, Test Score=672

- One district in our sample is Richmond, CA with STR=22, Test Score=672
- Predicted value: $\widehat{\text{Test Score}}_{Richmond} = 698 - 2.28(22) \approx 647$

- One district in our sample is Richmond, CA with STR=22, Test Score=672
- Predicted value: $\widehat{\text{Test Score}}_{Richmond} = 698 - 2.28(22) \approx 647$
- Residual: $\widehat{\epsilon_{Richmond}} = 672 - 647 = 25$