

Regression and Causation

Week 7

Yunkyu Sohn

School of Political Science and Economics

Waseda University

2019 Spring Statistics I

Summary of Statistics I

- ▶ Causality
 - ▶ Counterfactuals
 - ▶ Causal Effects
 - ▶ Fundamental Problem of Causal Inference
- ▶ Experimental Research
 - ▶ Role of Randomization
 - ▶ Sample Average Treatment Effect
- ▶ Internal and External Validities

Summary of Statistics I

- ▶ Observational studies
 - ▶ Confounding bias
 - ▶ Cross-section design
 - ▶ Before-and-after design
 - ▶ Difference-in-differences design

Summary of Statistics I

- ▶ Survey sampling
 - ▶ in order to get a representative sample of the target population
 - ▶ many sources of bias
 - ▶ Frame bias
 - ▶ non-response bias
 - ▶ response bias (e.g. social desirability bias)
 - ▶ Statistical remedy: e.g. Item count technique
- ▶ Univariate distribution
 - ▶ can be summarized using central/spread tendency indices
 - ▶ or graphs (e.g. histogram, boxplot)

Summary of Statistics I

- ▶ Correlation
 - ▶ z-score: standardized measure of deviation from the mean
 - ▶ Correlation coefficient: correlational measure of z-scores
- ▶ Linear regression
 - ▶ Least squares estimation => Coefficient estimation
 - ▶ Residuals: residual plot / used to calculate RMSE and R^2
 - ▶ Coefficient of determination: Explained variation by the model

Contents (4.3.1 - 4.3.3)

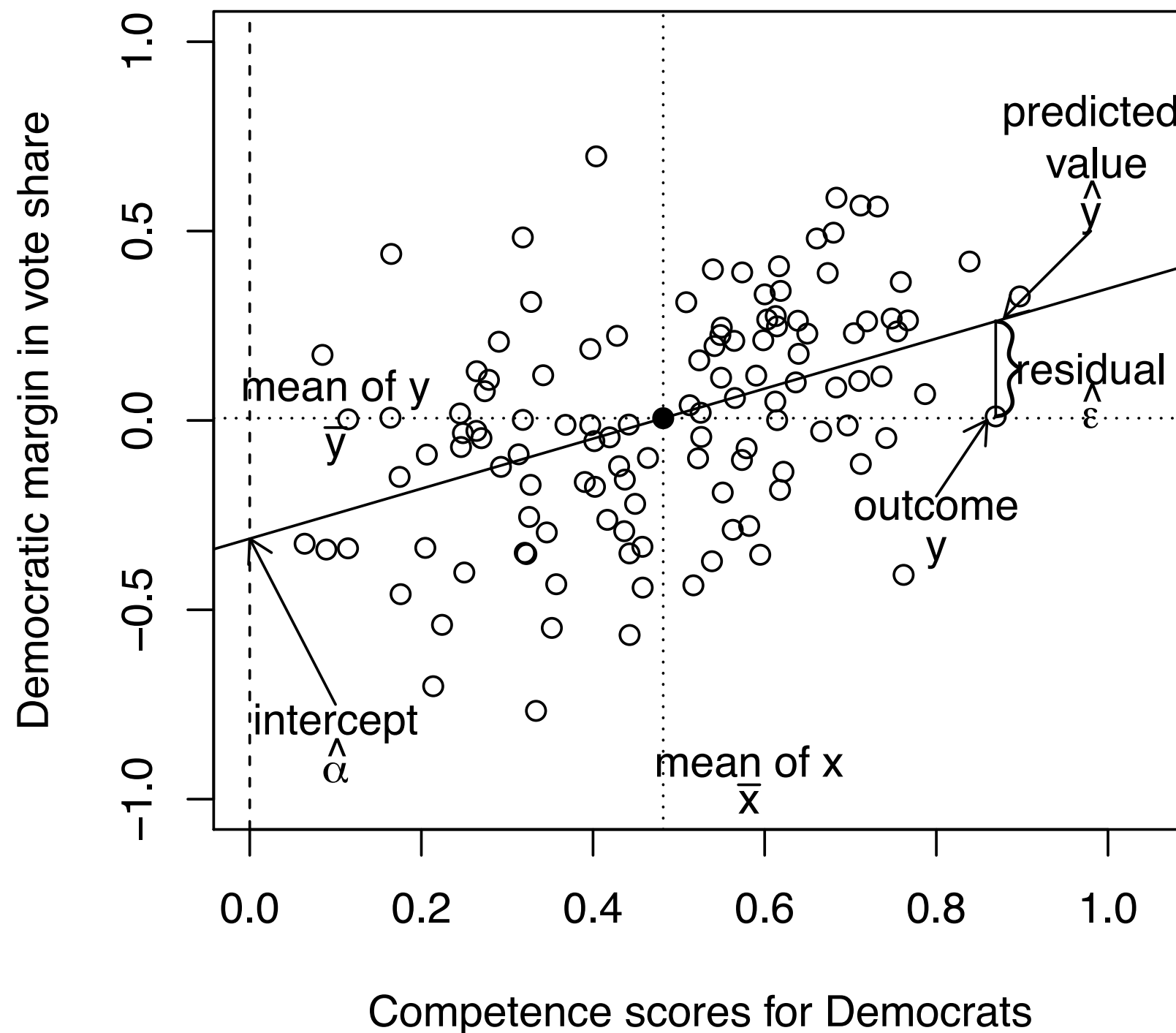
- ▶ Review: Linear regression as a **prediction model**
 - ▶ Coefficient estimates' connection to correlation coefficient
 - ▶ Measure of fit: Coefficient of determination
- ▶ **Using linear regression for causal inference**
 - ▶ How can you associate coefficients with treatment effects?
 - ▶ Example: Randomized natural experiment in India
 - ▶ Regression with multiple predictors
 - ▶ Multiple regression
 - ▶ Heterogeneous treatment effects
 - ▶ Linear regression with interaction terms

Contents (4.3.1 - 4.3.3)

- ▶ Review: Linear regression as a prediction model
 - ▶ Coefficient estimates' connection to correlation coefficient
 - ▶ Measure of fit: Coefficient of determination
- ▶ Using linear regression for causal inference
 - ▶ How can you associate coefficients with treatment effects?
 - ▶ Example: Randomized natural experiment in India
 - ▶ Regression with multiple predictors
 - ▶ Multiple regression
 - ▶ Heterogeneous treatment effects
 - ▶ Linear regression with interaction terms

Linear Regression Model

$$Y = \underbrace{\alpha}_{\text{intercept}} + \underbrace{\beta}_{\text{slope}} X + \underbrace{\epsilon}_{\text{error term}}$$



Linear Regression Model

- ▶ Model:
$$Y = \underbrace{\alpha}_{\text{intercept}} + \underbrace{\beta}_{\text{slope}} X + \underbrace{\epsilon}_{\text{error term}}$$
- ▶ Y : dependent/outcome/response variable
- ▶ X : independent/explanatory variable, predictor
- ▶ α, β : coefficients (parameters of the model)
- ▶ ϵ : unobserved error/disturbance term (mean zero)
- ▶ Interpretation
 - ▶ $\alpha + \beta X$: mean of Y given the value of X
 - ▶ α : the value of Y when X is zero
 - ▶ β : increase in Y associated with one unit increase in X

Linear Regression Model

- ▶ Linear regression draws a single line w. the highest level of explanation

$$Y = \underbrace{\alpha}_{\text{intercept}} + \underbrace{\beta}_{\text{slope}} X + \underbrace{\epsilon}_{\text{error term}}$$

- ▶ Estimation (hat notation for estimated)
 - ▶ $\hat{\alpha}, \hat{\beta}$: estimated coefficients
 - ▶ $\hat{Y} = \hat{\alpha} + \hat{\beta}X$: predicted/fitted value
 - ▶ $\hat{\epsilon} = Y - \hat{Y}$: residuals

Linear Regression Model

- Estimated coefficients (will be covered next week)

$$\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X}$$

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

$$\hat{\beta} = \text{correlation of } X \text{ and } Y \times \frac{\text{standard deviation of } Y}{\text{standard deviation of } X}$$

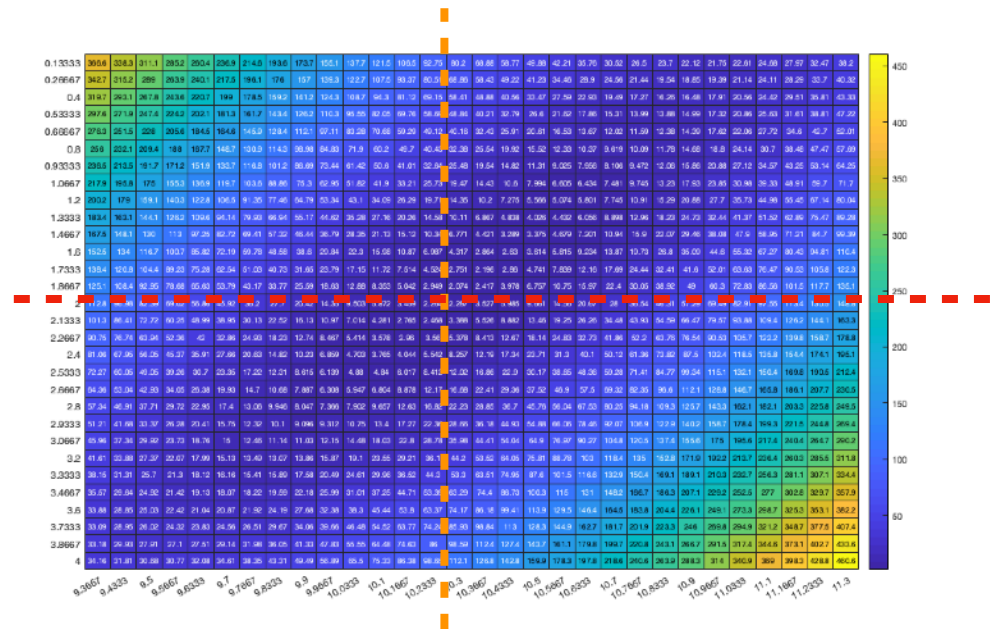
- Linear regression draws a single line w. the highest level of explanation
 - = Estimating the coefficients $\hat{\alpha}, \hat{\beta}$ (via least squares)
 - How?: Minimize the sum of squared residuals (SSR)

$$\text{SSR} = \sum_{i=1}^n \hat{\epsilon}_i^2 = \sum_{i=1}^n (Y_i - \hat{\alpha} - \hat{\beta}X_i)^2$$

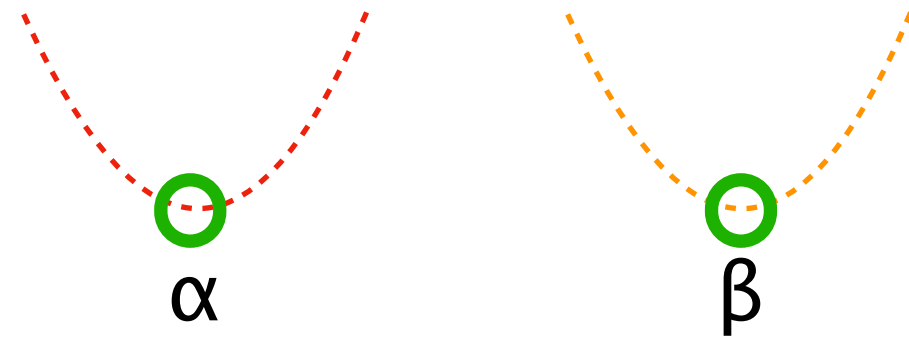
LEAST Squares Estimation

- Minimize SSR

$$\text{SSR} = \sum_{i=1}^n \hat{\epsilon}_i^2 = \sum_{i=1}^n (Y_i - \hat{\alpha} - \hat{\beta}X_i)^2$$

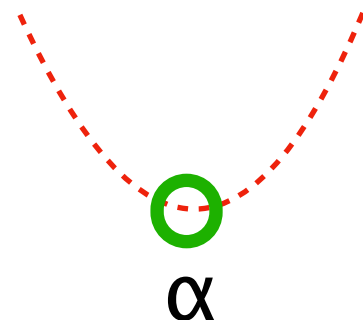


- Linear Regression: U-shape guaranteed

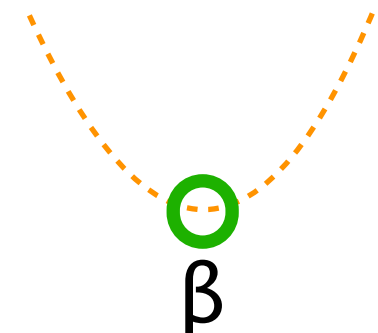


- Our objective: Find corresponding values of α and β at stationary points
- Derivative: Sensitivity of change in SSR w.r.t. change in α or β
- At stationary points, partial derivatives become 0

$$\frac{\partial \text{SSR}}{\partial \alpha} = 0$$



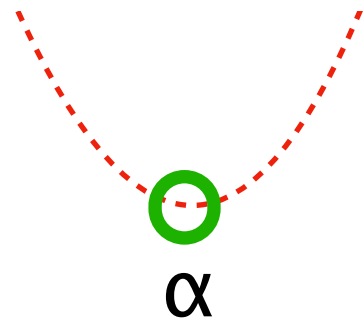
$$\frac{\partial \text{SSR}}{\partial \beta} = 0$$



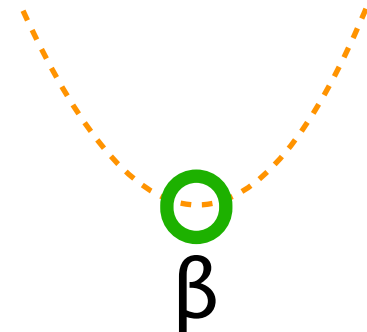
LEAST Squares Estimation

$$\text{SSR} = \sum_{i=1}^n \hat{\epsilon}_i^2 = \sum_{i=1}^n (Y_i - \hat{\alpha} - \hat{\beta}X_i)^2$$

$$\frac{\partial \text{SSR}}{\partial \alpha} = 0$$

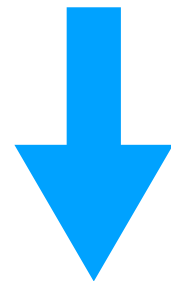


$$\frac{\partial \text{SSR}}{\partial \beta} = 0$$



$$-2 \sum_{i=1}^n (Y_i - \alpha - \beta X_i) = 0$$

$$-2 \sum_{i=1}^n (Y_i - \alpha - \beta X_i)X_i = 0$$



Solving simultaneous equations

$$\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X}$$

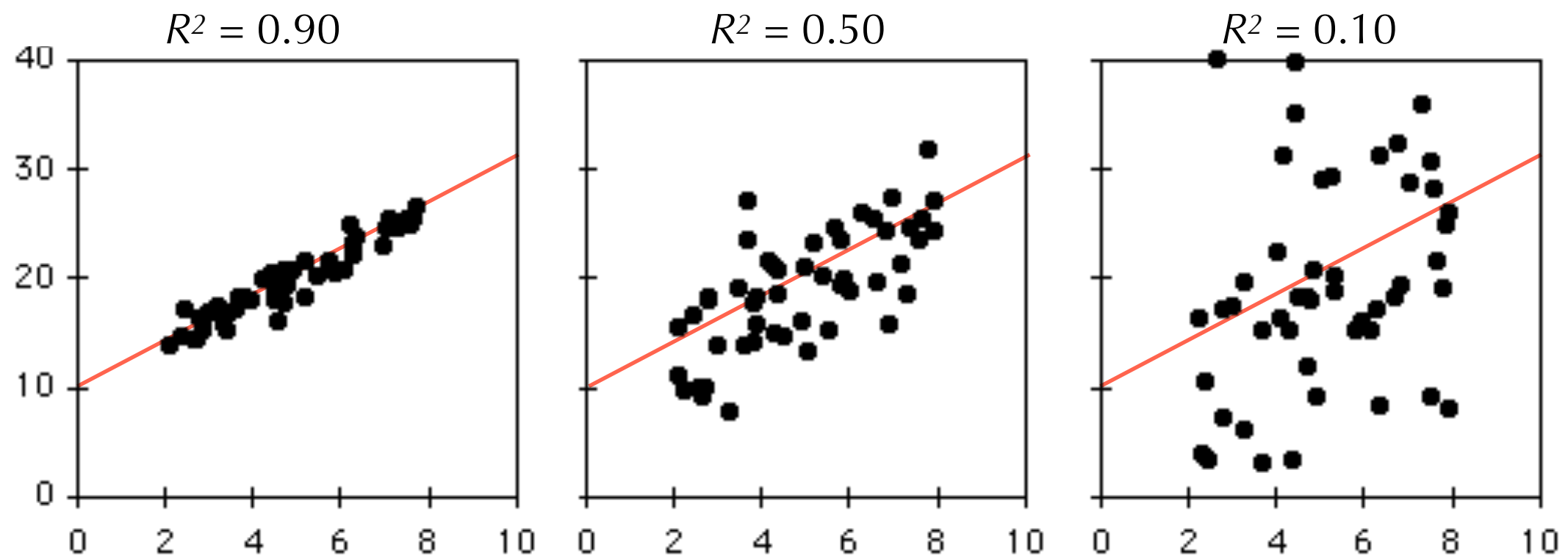
$$\hat{\beta} = \frac{\sum_{i=1}^n (Y_i - \bar{Y})(X_i - \bar{X})}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

Model Fit

residual variation of Y
left unexplained by X

$$R^2 = 1 - \frac{\text{SSR}}{\text{Total sum of squares (TSS)}} = 1 - \frac{\sum_{i=1}^n \hat{\epsilon}_i^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$$

- Square of correlation between Y and \hat{Y} : magnitude of association



- Look at the residuals!

Numerical Exercise

Observation	1st sibling age (X)	2nd sibling age (Y)
1	20	18
2	18	16
3	20	16
4	20	18

$$\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X}$$

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

$$\text{RMSE} = \sqrt{\frac{1}{n}\text{SSR}} = \sqrt{\frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2}$$

$$\text{TSS} = \sum_{i=1}^n (Y_i - \bar{Y})^2 \quad R^2 = 1 - \frac{\text{SSR}}{\text{TSS}}$$

- Predict Y using X: $\hat{\alpha}, \hat{\beta}$
- SSR
- RMSE
- TTS
- R^2

Numerical Exercise

Observation	1st sibling age (X)	2nd sibling age (Y)
1	20	18
2	18	16
3	20	16
4	20	18

$$\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X}$$

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

$$\text{RMSE} = \sqrt{\frac{1}{n}\text{SSR}} = \sqrt{\frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2}$$

$$\text{TSS} = \sum_{i=1}^n (Y_i - \bar{Y})^2 \quad R^2 = 1 - \frac{\text{SSR}}{\text{TSS}}$$

Challenge:

Can we say that the coefficient implies a causal relationship?

Contents (4.3.1 - 4.3.3)

- ▶ Review: Linear regression as a prediction model
 - ▶ Coefficient estimates' connection to correlation coefficient
 - ▶ Measure of fit: Coefficient of determination
- ▶ Using linear regression for causal inference
 - ▶ How can you associate coefficients with treatment effects?
 - ▶ Example: Randomized natural experiment in India
 - ▶ Regression with multiple predictors
 - ▶ Multiple regression
 - ▶ Heterogeneous treatment effects
 - ▶ Linear regression with interaction terms

Example: Women as Policy Makers

- ▶ Research question:
 - ▶ Do women promote different policies compared to men?
- ▶ Usual approach: observational studies
 - ▶ Problems: confounding (pre-treatment) variables
- ▶ Ideal design: randomized experiment
 - ▶ Special case: natural (randomized) experiment
- ▶ Data: Randomized policy experiment in India
 - ▶ one third of village council heads **randomly reserved** for women
 - ▶ assigned at the level of Gram Panchayat (GP) since mid-1990s
 - ▶ each GP has multiple villages

Example: Women as Policy Makers

- ▶ Hypothesis:
 - ▶ Female politicians represent the interests of female voters
 - ▶ Female voters complain about drinking water while male voters complain about irrigation

Name	Description
GP	An identifier for the Gram Panchayat (GP)
village	identifier for each village
reserved	binary variable indicating whether the GP was reserved for women leaders or not
female	binary variable indicating whether the GP had a female leader or not
irrigation	variable measuring the number of new or re-paired irrigation facilities in the village since the reserve policy started
water	variable measuring the number of new or re-paired drinking-water facilities in the village since the reserve policy started

Example: Women as Policy Makers

- ▶ Each GP has the same number of villages.
- ▶ Compliance check:
 - ▶ Does the reservation policy increase number of female politicians?

```
mean(women$female[women$reserved == 1])  
## [1] 1  
  
mean(women$female[women$reserved == 0])  
## [1] 0.0748
```

- ▶ Policy outcomes: difference in the means estimator

```
## drinking-water facilities  
mean(women$water[women$reserved == 1]) -  
  mean(women$water[women$reserved == 0])  
  
## [1] 9.25
```

Example: Women as Policy Makers

- Now let us instead use regression for the difference in means estimator

$$\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X}$$

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

- When X is binary: $X_i=0$ (control) 1 (treated); n_0 : #control; n_1 : #treated

$$\hat{\alpha} = \underbrace{\frac{1}{n_0} \sum_{i=1}^n (1 - X_i) Y_i}_{\text{mean outcome among the control}},$$

$$\hat{\beta} = \underbrace{\frac{1}{n_1} \sum_{i=1}^n X_i Y_i}_{\text{mean outcome among the treated}} - \underbrace{\frac{1}{n_0} \sum_{i=1}^n (1 - X_i) Y_i}_{\text{mean outcome among the control}}$$

Example: Women as Policy Makers

- Now let us instead use regression for the difference in means estimator

$$\hat{\alpha} = \underbrace{\frac{1}{n_0} \sum_{i=1}^n (1 - X_i) Y_i}_{\text{mean outcome among the control}},$$

$$\hat{\beta} = \underbrace{\frac{1}{n_1} \sum_{i=1}^n X_i Y_i}_{\text{mean outcome among the treated}} - \underbrace{\frac{1}{n_0} \sum_{i=1}^n (1 - X_i) Y_i}_{\text{mean outcome among the control}}$$

- Compare with

$$\mathbf{SATE} = \frac{1}{n} \sum_{i=1}^n \{Y_i(1) - Y_i(0)\}$$
$$D = \frac{1}{|\{T_i = 1\}|} \sum_{i \in \{T_i=1\}} Y_i - \frac{1}{|\{T_i = 0\}|} \sum_{i \in \{T_i=0\}} Y_i$$

Example: Women as Policy Makers

$$\widehat{Y(X)} = \hat{\alpha} + \hat{\beta}X$$

► Why?

► Potential outcomes

$$\widehat{Y(0)} = \hat{\alpha} \qquad \widehat{Y(1)} = \hat{\alpha} + \hat{\beta}$$

► Difference in the means estimator

$$\widehat{Y(1)} - \widehat{Y(0)} = (\hat{\alpha} + \hat{\beta}) - \hat{\alpha} = \hat{\beta}$$

```
lm(water ~ reserved, data = women)
```

```
##
```

```
## Call:
```

```
## lm(formula = water ~ reserved, data = women)
```

```
##
```

```
## Coefficients:
```

```
## (Intercept)      reserved
```

```
##          14.74          9.25
```

```
## drinking-water facilities
```

```
mean(women$water[women$reserved == 1]) -  
mean(women$water[women$reserved == 0])
```

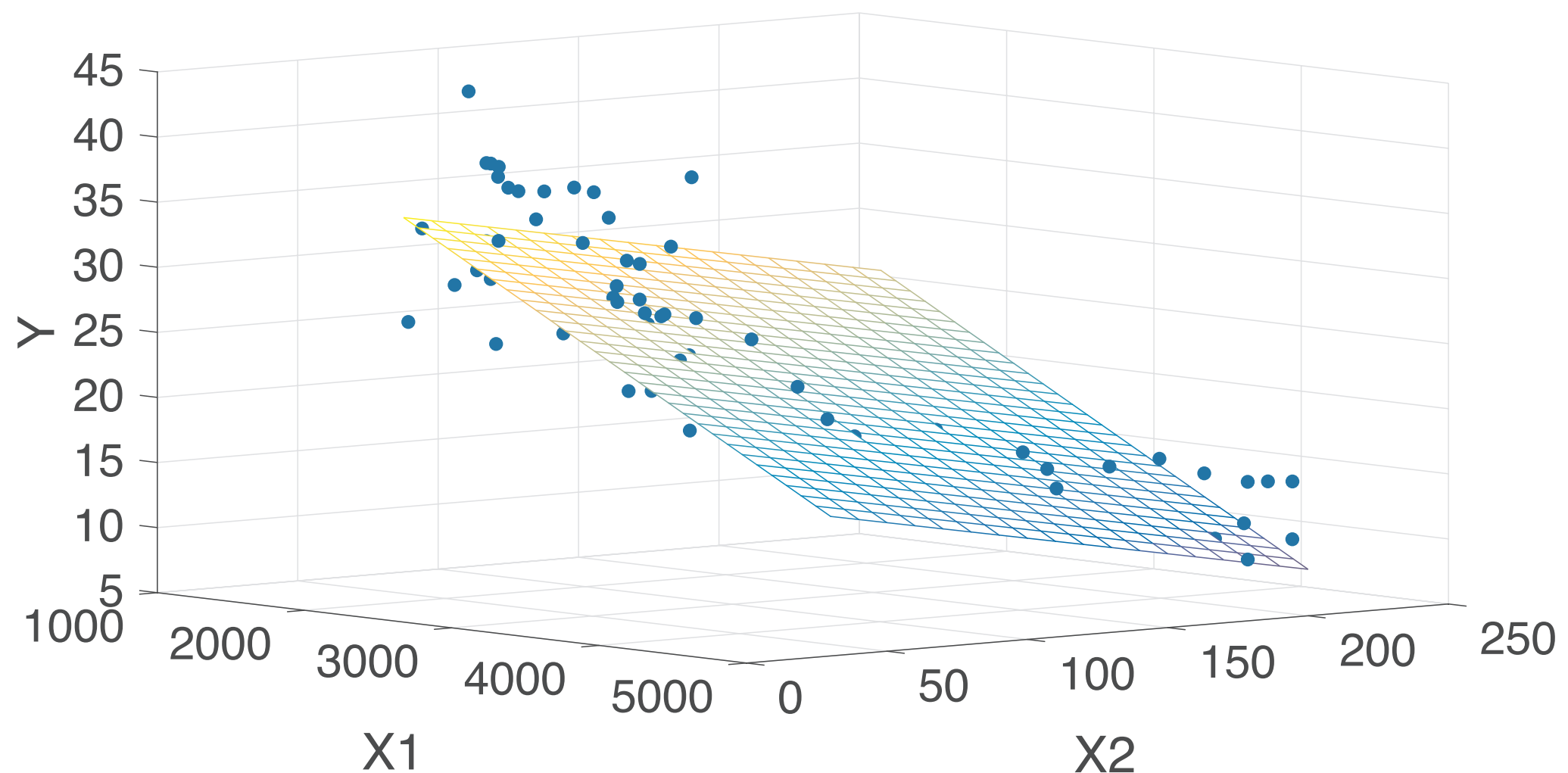
```
## [1] 9.25
```

Linear Regression With Multiple Predictors

- Generalizing linear regression for multiple predictors

$$Y = \alpha + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p + \epsilon$$

- Find a hyperplane instead of a line (e.g. for 2 predictors: simple plane)

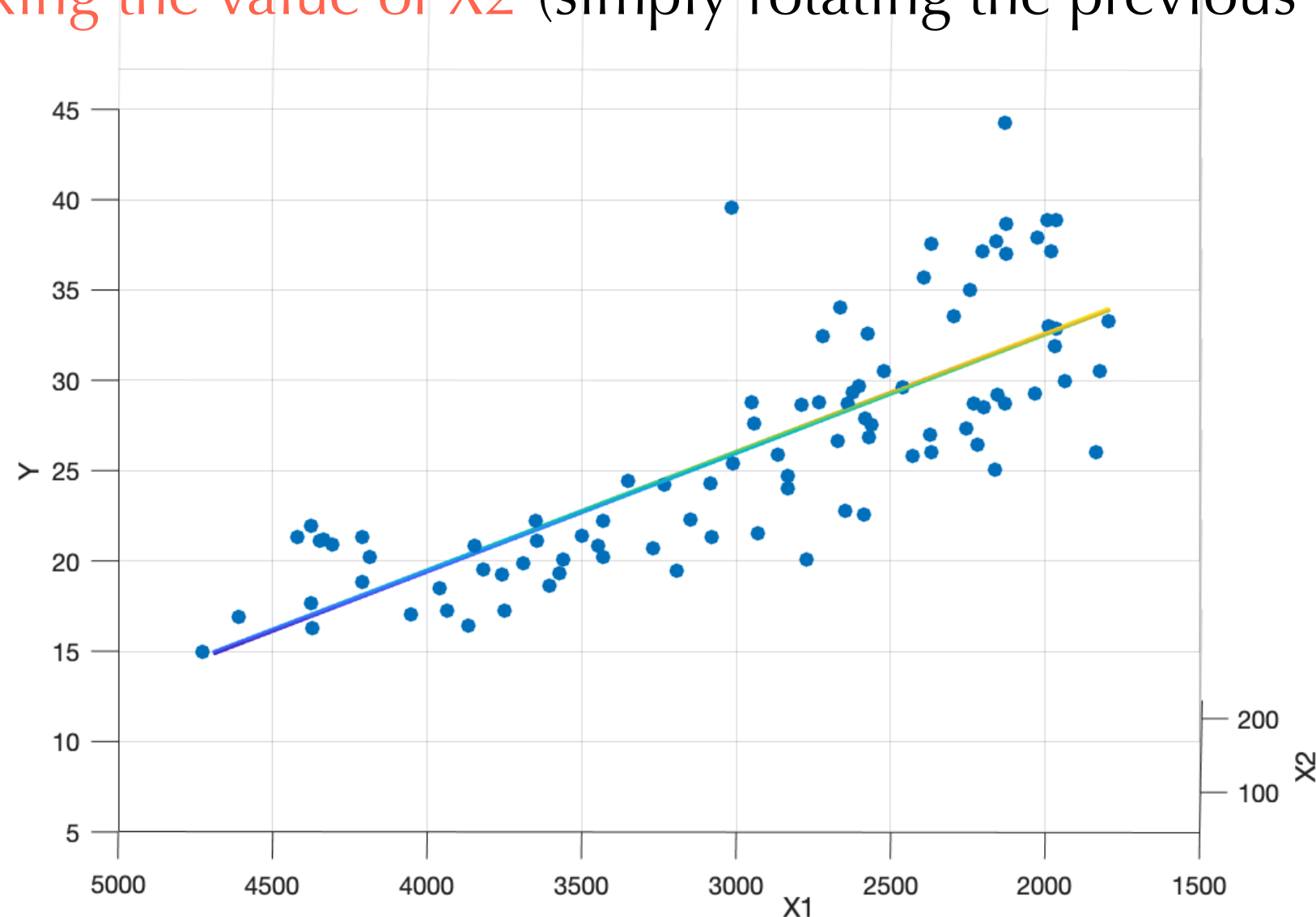


Linear Regression With Multiple Predictors

- Generalizing linear regression

$$Y = \alpha + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p + \epsilon$$

- Find a hyperplane instead of a line (e.g. for 2 predictors: simple plane)
- after fixing the value of X_2 (simply rotating the previous plot)



Linear Regression With Multiple Predictors

- Generalizing linear regression

$$Y = \alpha + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p + \epsilon$$

- We need to minimize SSR as same as the case $p=1$

$$\text{SSR} = \sum_{i=1}^n \hat{\epsilon}_i^2 = \sum_{i=1}^n (Y_i - \hat{\alpha} - \hat{\beta}_1 X_{i1} - \hat{\beta}_2 X_{i2} - \cdots - \hat{\beta}_p X_{ip})^2$$

Social Pressure Turnout Experiment (Book Chapter 2.4.2)

- ▶ Get-out-the-vote (GOTV) messaged postcard with social pressure
 - ▶ naming and shaming strategy

Dear Registered Voter:

WHAT IF YOUR NEIGHBORS KNEW WHETHER YOU VOTED?

Why do so many people fail to vote? We've been talking about the problem for years, but it only seems to get worse. This year, we're taking a new approach. We're sending this mailing to you and your neighbors to publicize who does and does not vote.

The chart shows the names of some of your neighbors, showing which have voted in the past. After the August 8 election, we intend to mail an updated chart. You and your neighbors will all know who voted and who did not.

DO YOUR CIVIC DUTY – VOTE!

MAPLE DR	Aug 04	Nov 04	Aug 06
9995 JOSEPH JAMES SMITH	Voted	Voted	_____
995 JENNIFER KAY SMITH		Voted	_____
9997 RICHARD B JACKSON		Voted	_____
9999 KATHY MARIE JACKSON		Voted	_____

Social Pressure Turnout Experiment (Book Chapter 2.4.2)

- ▶ **Research Question:** Does peer pressure facilitate turnout?
- ▶ **Unit:** voter in a primary election in the state of Michigan
- ▶ **Treatment variable** T : GOTV + social pressure messaged postcard sent
- ▶ **Treatment group** (treated units): voters receiving certain postcards
- ▶ **Control group** (untreated units): voters receiving none
- ▶ **Outcome variable** (response variable) Y : turnout (publicly available)
- ▶ **Potential outcomes:** *turnout*(received) and *turnout*(not-received)
- ▶ **Causal effect:** Difference in the sample means estimator
- ▶ Further details: 2 additional treatment conditions
 - ▶ GOTV w/o social pressure but civic duty messaged postcard
 - ▶ Hawthorne effect group
 - ▶ YOU ARE BEING STUDIED! w/o social pressure

Social Pressure Turnout Experiment (Book Chapter 2.4.2)

<i>Variable</i>	<i>Description</i>
hhsiz	household size of the voter
messages	GOTV messages the voter received (Civic Duty, Control, Neighbors, Hawthorne)
sex	sex of the voter (female or male)
yearofbirth	year of birth of the voter
primary2004	whether the voter voted in the 2004 primary election (1=voted, 0=abstained)
primary2006	whether the voter turned out in the 2006 primary election (1=voted, 0=abstained)

► Model

$$Y = \alpha + \beta_1 \text{Control} + \beta_2 \text{Hawthorne} + \beta_3 \text{Neighbors} + \epsilon$$

Social Pressure Turnout Experiment (Book Chapter 2.4.2)

$$Y = \alpha + \beta_1 \text{Control} + \beta_2 \text{Hawthorne} + \beta_3 \text{Neighbors} + \epsilon$$

```
lm(primary2008 ~ Control + Hawthorne + Neighbors, data = social)
```

```
## Coefficients:
```

```
##          (Intercept)      messagesControl      messagesHawthorne
```

```
##          0.314538          -0.017899          0.007837
```

```
## messagesNeighbors
```

```
##          0.063411
```

- Average outcome for control condition

$$\hat{\alpha} + \hat{\beta}_1 = 0.315 + (-0.018) = 0.297$$

- Average outcome for neighbors condition

$$\hat{\alpha} + \hat{\beta}_3 = 0.315 + 0.063 = 0.378$$

- Average effect of the Neighbors treatment (relative to the control)

$$\hat{\alpha} + \hat{\beta}_3 - (\hat{\alpha} + \hat{\beta}_1) = \hat{\beta}_3 - \hat{\beta}_1 = 0.063 - (-0.018) = 0.081$$

Heterogeneous Treatment Effects

- ▶ Is average aggregate treatment effect enough?
 - ▶ may depend on unit-specific traits
 - ▶ e.g. negative campaign strategy → electorates by supporting party
 - ▶ e.g. minimum wage → employers by employer wealth
 - ▶ Statistical approach for identification
 - ▶ Linear regression model with an interaction term
$$Y = \alpha + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \epsilon$$
 - ▶ **Hierarchy principle**: all low level effect terms (X_1, X_2) must be included
 - ▶ Social pressure experiment
 - ▶ Neighbors message has different effects to
 - ▶ Previous primary voters and non-voters
- primary2004 whether the voter voted in the 2004 primary election
- ▶ Using neighbors and control groups only with a modified model
$$Y = \alpha + \beta_1 \text{primary2004} + \beta_2 \text{Neighbors} + \beta_3 (\text{primary2004} \times \text{Neighbors}) + \epsilon$$

Heterogeneous Treatment Effects

$$Y = \alpha + \beta_1 \text{primary2004} + \beta_2 \text{Neighbors} + \beta_3 (\text{primary2004} \times \text{Neighbors}) + \epsilon$$

- ▶ Average treatment effect for voters in previous election ($\text{primary2004}_i = 1$)
 - ▶ $\hat{Y}(\text{Neighbors} = 1, \text{primary2004} = 1) - \hat{Y}(\text{Neighbors} = 0, \text{primary2004} = 1)$
- ▶ Average treatment effect for non-voters in previous election ($\text{primary2004}_i = 0$)
 - ▶ $\hat{Y}(\text{Neighbors} = 1, \text{primary2004} = 0) - \hat{Y}(\text{Neighbors} = 0, \text{primary2004} = 0)$
- ▶ Difference in the estimated average treatment effect between voters/non-voters

```
## lm(formula = primary2008 ~ primary2004 + messages + primary2004:messages,  
##      data = social.neighbor)  
##  
## Coefficients:  
##              (Intercept)  
##              0.23711  
##              primary2004  
##              0.14870  
##              messagesNeighbors  
##              0.06930  
## primary2004:messagesNeighbors  
##              0.02723
```


Heterogeneous Treatment Effects

- Interaction model with **linear age** effect

$$Y = \alpha + \beta_1 \text{age} + \beta_2 \text{Neighbors} + \beta_3 (\text{age} \times \text{Neighbors}) + \epsilon$$

- Average treatment effect difference when age increases by 1 year)

- Neighbor treatment effect for **age x population**

$$(\hat{\alpha} + \hat{\beta}_1 x + \hat{\beta}_2 + \hat{\beta}_3 x) - (\hat{\alpha} + \hat{\beta}_1 x) = \hat{\beta}_2 + \hat{\beta}_3 x$$

- Neighbor treatment effect for **age x+1 population**

$$\{\hat{\alpha} + \hat{\beta}_1(x + 1) + \hat{\beta}_2 + \hat{\beta}_3(x + 1)\} - \{\hat{\alpha} + \hat{\beta}_1(x + 1)\} = \hat{\beta}_2 + \hat{\beta}_3(x + 1)$$

- Average treatment effect difference when age increases by 1 year

$$\hat{\beta}_3 = \{\hat{\beta}_2 + \hat{\beta}_3(x + 1)\} - (\hat{\beta}_2 + \hat{\beta}_3 x)$$

```
## lm(formula = primary2008 ~ age * messages, data = social.neighbor)
##
## Coefficients:
##              (Intercept)                  age
##              0.0894768                0.0039982
##      messagesNeighbors  age:messagesNeighbors
##              0.0485728                0.0006283
```

Heterogeneous Treatment Effects

- Interaction model with **linear age** and **quadratic age** effects

$$Y = \alpha + \beta_1 \text{age} + \beta_2 \text{age}^2 + \beta_3 \text{Neighbors} + \beta_4 (\text{age} \times \text{Neighbors}) + \beta_5 (\text{age}^2 \times \text{Neighbors}) + \epsilon.$$

- Average treatment effect difference when age increases by 1 year)

- Neighbor treatment effect for **age x population**

- Neighbor treatment effect for **age x+1 population**

- Average treatment effect difference when age increases by 1 year

```
##               (Intercept)                age
##               -9.700e-02                1.172e-02
##               I (age^2)                messagesNeighbors
##               -7.389e-05                -5.275e-02
##   age:messagesNeighbors  I (age^2) :messagesNeighbors
##               4.804e-03                -3.961e-05
```

Heterogeneous Treatment Effects

- ▶ Interaction model with **linear age** and **quadratic age** effects

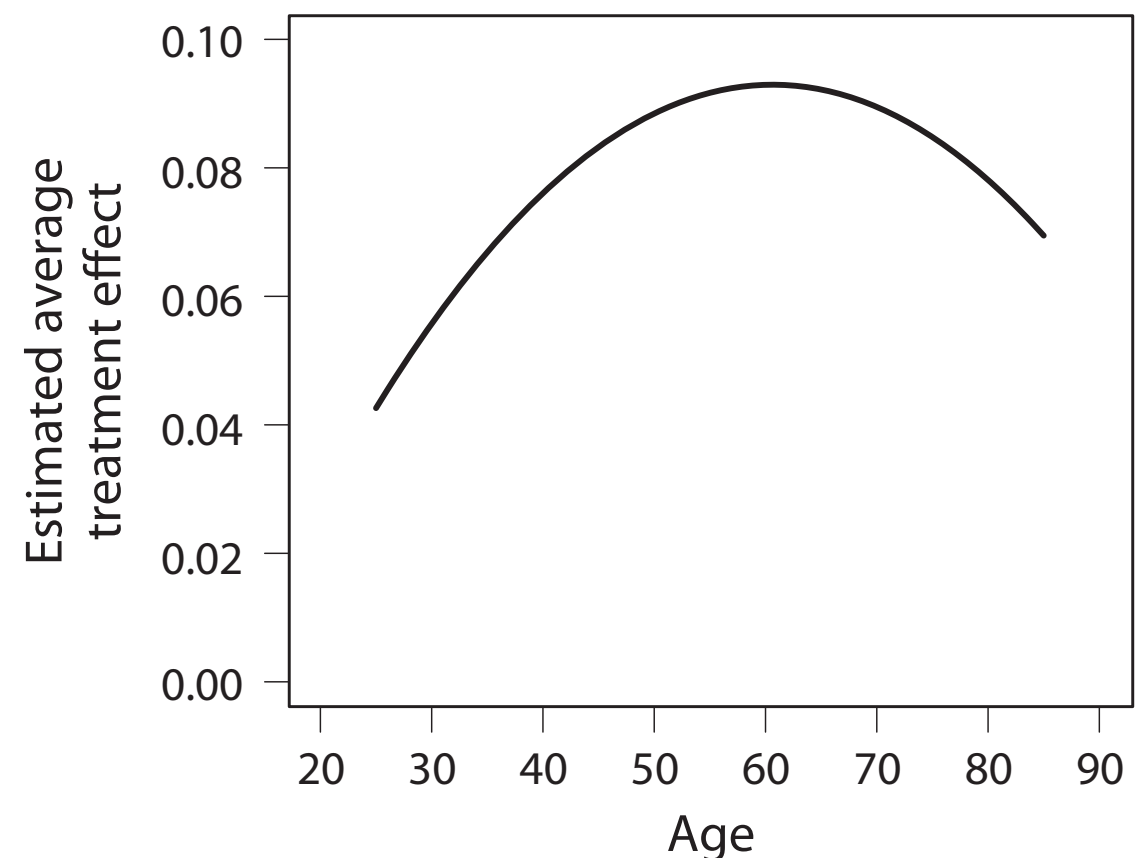
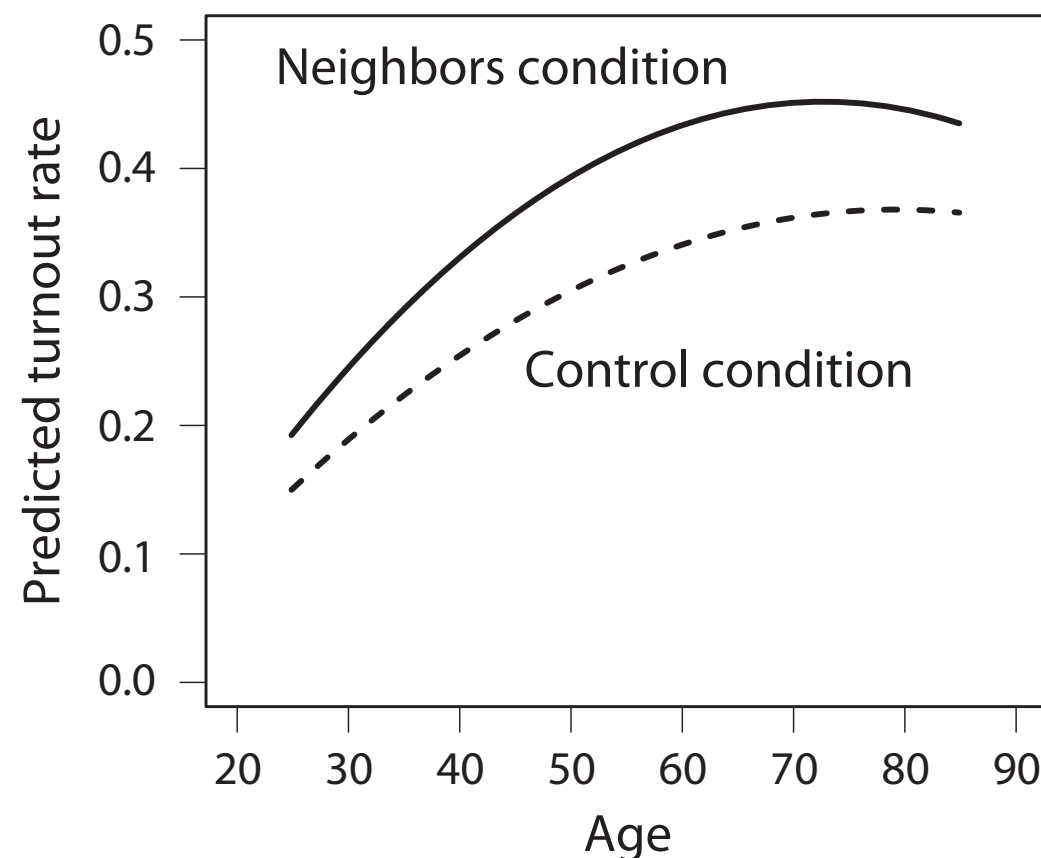
$$Y = \alpha + \beta_1 \text{age} + \beta_2 \text{age}^2 + \beta_3 \text{Neighbors} + \beta_4 (\text{age} \times \text{Neighbors}) + \beta_5 (\text{age}^2 \times \text{Neighbors}) + \epsilon.$$

- ▶ Average treatment effect difference when age increases by 1 year)
 - ▶ Neighbor treatment effect for **age x population**
 - ▶ Neighbor treatment effect for **age x+c population**
 - ▶ Average treatment effect difference when age increases by c years

Heterogeneous Treatment Effects

- Interaction model with **linear age** and **quadratic age** effects

$$Y = \alpha + \beta_1 \text{age} + \beta_2 \text{age}^2 + \beta_3 \text{Neighbors} + \beta_4 (\text{age} \times \text{Neighbors}) + \beta_5 (\text{age}^2 \times \text{Neighbors}) + \epsilon.$$



Summary

- ▶ Using linear regression for causal inference
 - ▶ How can you associate coefficients with treatment effects?
 - ▶ Example: Randomized natural experiment in India
- ▶ Regression with multiple predictors
 - ▶ Multiple regression
- ▶ Heterogeneous treatment effects
 - ▶ Linear regression with interaction terms

Statistics Is the First Step for Becoming

Top-paying entry-level jobs of 2019

RANK	JOB TITLE	MEDIAN STARTING SALARY
1	Data Scientist	\$95,000
2	Software Engineer	\$90,000
3	Product Manager	\$89,000
4	Investment Banking Analyst	\$85,000
5	Product Designer	\$85,000
6	UX Designer	\$73,000
7	Implementation Consultant	\$72,000
8	Java Developer	\$72,000
9	Systems Engineer	\$70,000
10	Software Developer	\$68,600

Social Science Majors Have Every Required Skill

- ▶ If and only if you learn statistics, econometrics and programming

**Harvard
Business
Review**

Latest

Magazine

Popular

Topics

Podcasts

Video

Store

The Big Idea

Visu

ECONOMICS

Why Tech Companies Hire So Many Economists

by [Susan Athey](#) and [Michael Luca](#)

FEBRUARY 12, 2019

<https://hbr.org/2019/02/why-tech-companies-hire-so-many-economists?fbclid=IwAR3rUFlqEnQJp6GhXdDc6bYI3N1wbh9Fj-Gc-7p-yPlvdXsd7Ax2lqgl5Jk>

Social Science Majors Have Every Required Skill

- ▶ If and only if you learn statistics, econometrics and programming



Markets Tech Media Success Perspectives Video

International Edition +

Amazon gets an edge with its secret squad of PhD economists

By Lydia DePillis, CNN Business

Updated 1039 GMT (1839 HKT) March 13, 2019

<https://edition.cnn.com/2019/03/13/tech/amazon-economists/index.html>

Exams

- ▶ R quiz
 - ▶ will be during the lecture class on paper (May 28th 4:30~6:00PM)
 - ▶ Sorry — examination platform (Qualtrics) unstable
 - ▶ 30 multiple choice questions
 - ▶ All questions from weeks 1-6 lab materials (except addendum)
 - ▶ No questions from those not covered during lab sessions
 - ▶ I will go over the lecture contents during the lab sessions next week.
 - ▶ No digital device (including calculator) allowed
- ▶ Final exam
 - ▶ June 4th 4:30~6:00PM (no lab session next week)
 - ▶ No digital device (including calculator) allowed

Final Exam

- The definitions of the following measures will be provided

$$\text{correlation}(x, y) =$$

$$\frac{1}{n-1} \sum_{i=1}^n (\text{z-score of } x_i \times \text{z-score of } y_i)$$

$$\hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X}$$

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

$$\text{RMSE} = \sqrt{\frac{1}{n} \text{SSR}} = \sqrt{\frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2}$$

Good luck!