

# Lecture 7: Regression

Naijia Liu

Feb. 13, 2024

# Questions?

P-set I due on Thursday 11:59pm.

Utilize OHs and Slack channel.

Tidyverse and base R both fine!

We only accept compiled rmd files and pdf!

# OLS Regression

- Review of OLS linear regression.

# OLS Regression

- Review of OLS linear regression.
- In the view of DID, IV and Matching.

# OLS Regression

- Review of OLS linear regression.
- In the view of DID, IV and Matching.
- Variable selection using penalization.

# OLS Regression

- Review of OLS linear regression.
- In the view of DID, IV and Matching.
- Variable selection using penalization.
- Heterogenous treatment effect using penalized regression.

# Notations

- $Y_i$ : Outcome Variable / Dependent Variable
- $X_i$ : Independent Variables
- $\beta$ : Coefficient for IVs
- $\beta_0$ : Coefficient for Intercept
- $\epsilon_i$ : Error term

# Linear Regression: A Model for the Mean

Assume a model for an observed simple random sample  $Y_i$ :

$$\underbrace{Y_i}_{\text{Outcome}} = \underbrace{\beta_0}_{\text{Mean or intercept}} + \underbrace{\epsilon_i}_{\text{Residual, error term}}$$



# Linear Regression: A Model for the Mean

Assume a model for an observed simple random sample  $Y_i$ :

$$\underbrace{Y_i}_{\text{Outcome}} = \underbrace{\beta_0}_{\text{Mean or intercept}} + \underbrace{\epsilon_i}_{\text{Residual, error term}}$$

How to choose  $\hat{\beta}_0$ ?

$$\hat{\beta}_0 = \underset{\tilde{\beta}_0}{\operatorname{argmin}} \mathcal{L}(\tilde{\beta}_0)$$

We want  $\hat{\beta}_0$  to be the one to minimize some type of error, in predicting  $Y_i$ .

And  $\mathcal{L}(\tilde{\beta}_0)$  is a **loss** function.

# Commonly Encountered Loss Functions

- Criterion of Least Squares (OLS): **We want to minimize the sum of squared error between true data and our predictions.**

$$\hat{\beta}_0 = \underset{\tilde{\beta}_0}{\operatorname{argmin}} \sum_{i=1}^N \left( Y_i - \tilde{\beta}_0 \right)^2$$

# Commonly Encountered Loss Functions

- Criterion of Least Squares (OLS): **We want to minimize the sum of squared error between true data and our predictions.**

$$\hat{\beta}_0 = \underset{\tilde{\beta}_0}{\operatorname{argmin}} \sum_{i=1}^N \left( Y_i - \tilde{\beta}_0 \right)^2$$

- Criterion of Least Absolute Deviation: **We want to minimize the sum of absolute error between true data and our predictions.**

$$\hat{\beta}_0 = \underset{\tilde{\beta}_0}{\operatorname{argmin}} \sum_{i=1}^N \left| Y_i - \tilde{\beta}_0 \right|$$

# Commonly Encountered Loss Functions

- Criterion of Least Squares (OLS): **We want to minimize the sum of squared error between true data and our predictions.**

$$\hat{\beta}_0 = \underset{\tilde{\beta}_0}{\operatorname{argmin}} \sum_{i=1}^N \left( Y_i - \tilde{\beta}_0 \right)^2$$

- Criterion of Least Absolute Deviation: **We want to minimize the sum of absolute error between true data and our predictions.**

$$\hat{\beta}_0 = \underset{\tilde{\beta}_0}{\operatorname{argmin}} \sum_{i=1}^N |Y_i - \tilde{\beta}_0|$$

- Penalized Least Squares: **We want to minimize the sum of squared error between true data and our predictions, plus something else (later).**

$$\hat{\beta}_0 = \underset{\tilde{\beta}_0}{\operatorname{argmin}} \sum_{i=1}^N \left( Y_i - \tilde{\beta}_0 \right)^2 + \lambda \tilde{\beta}_0^2$$

# Review on Taking Derivatives

For a function:

$$f(x) = (x + a)^2$$

The first derivative of it is:

$$f'(x) = \frac{d}{dx}f(x) = 2(x + a)$$

# Review on Taking Derivatives

For a function:

$$f(x) = (x + a)^2$$

The first derivative of it is:

$$f'(x) = \frac{d}{dx}f(x) = 2(x + a)$$

The second derivative of it is:

$$f''(x) = 2$$

# Review on Taking Derivatives

For a function:

$$f(x) = (x + a)^2$$

The first derivative of it is:

$$f'(x) = \frac{d}{dx}f(x) = 2(x + a)$$

The second derivative of it is:

$$f''(x) = 2$$

If we want to find the point that minimizes the function, we want to set first derivative to 0.

In this case we have:

$$x = -a$$

# OLS Solution

$$\hat{\beta}_0 = \underset{\tilde{\beta}_0}{\operatorname{argmin}} \sum_{i=1}^N \left( Y_i - \tilde{\beta}_0 \right)^2$$



# OLS Solution

$$\hat{\beta}_0 = \underset{\tilde{\beta}_0}{\operatorname{argmin}} \sum_{i=1}^N \left( Y_i - \tilde{\beta}_0 \right)^2$$

$$\Rightarrow \left. \frac{\partial \mathcal{L}(\tilde{\beta}_0)}{\partial \tilde{\beta}_0} \right|_{\tilde{\beta}_0 = \hat{\beta}_0} = 0$$

$$\frac{\partial}{\partial \tilde{\beta}_0} \sum_{i=1}^N \left( Y_i - \tilde{\beta}_0 \right)^2 = 0$$

$$\sum_{i=1}^N \frac{\partial}{\partial \tilde{\beta}_0} \left( Y_i - \tilde{\beta}_0 \right)^2 = 0$$

$$\sum_{i=1}^N -2 \cdot \left( Y_i - \hat{\beta}_0 \right) = 0$$

# OLS Solution

$$\hat{\beta}_0 = \underset{\tilde{\beta}_0}{\operatorname{argmin}} \sum_{i=1}^N \left( Y_i - \tilde{\beta}_0 \right)^2$$

$$\Rightarrow \left. \frac{\partial \mathcal{L}(\tilde{\beta}_0)}{\partial \tilde{\beta}_0} \right|_{\tilde{\beta}_0 = \hat{\beta}_0} = 0$$

$$\frac{\partial}{\partial \tilde{\beta}_0} \sum_{i=1}^N \left( Y_i - \tilde{\beta}_0 \right)^2 = 0$$

$$\sum_{i=1}^N \frac{\partial}{\partial \tilde{\beta}_0} \left( Y_i - \tilde{\beta}_0 \right)^2 = 0$$

$$\sum_{i=1}^N -2 \cdot \left( Y_i - \hat{\beta}_0 \right) = 0$$

$$\sum_{i=1}^N \left( Y_i - \hat{\beta}_0 \right) = 0$$

$$\sum_{i=1}^N Y_i = \sum_{i=1}^N \hat{\beta}_0$$

$$\sum_{i=1}^N Y_i = N \hat{\beta}_0$$

$$\frac{1}{N} \sum_{i=1}^N Y_i = \bar{Y}_i = \hat{\beta}_0$$

# Least Absolute Deviation Solution (optional)

$$\hat{\beta}_0 = \underset{\tilde{\beta}_0}{\operatorname{argmin}} \sum_{i=1}^N |Y_i - \tilde{\beta}_0|$$

$$\Rightarrow \left. \frac{\partial \mathcal{L}(\tilde{\beta}_0)}{\partial \tilde{\beta}_0} \right|_{\tilde{\beta}_0 = \hat{\beta}_0} = 0$$

$$\frac{\partial}{\partial \tilde{\beta}_0} \left\{ \sum_{i=1}^N |Y_i - \tilde{\beta}_0| \right\} = 0$$

$$\sum_{i=1}^N \frac{\partial}{\partial \tilde{\beta}_0} |Y_i - \tilde{\beta}_0| = 0$$

$$\sum_{i=1}^N \operatorname{sgn}(Y_i - \hat{\beta}_0) = 0$$

$$\Rightarrow \tilde{Y} = \hat{\beta}_0$$

- If we rank all observations from small to large:

$$Y_{(1)} = \min(Y_i); Y_{(N)} = \max(Y_i)$$

# Least Absolute Deviation Solution (optional)

$$\hat{\beta}_0 = \underset{\tilde{\beta}_0}{\operatorname{argmin}} \sum_{i=1}^N |Y_i - \tilde{\beta}_0|$$

$$\Rightarrow \left. \frac{\partial \mathcal{L}(\tilde{\beta}_0)}{\partial \tilde{\beta}_0} \right|_{\tilde{\beta}_0 = \hat{\beta}_0} = 0$$

$$\frac{\partial}{\partial \tilde{\beta}_0} \left\{ \sum_{i=1}^N |Y_i - \tilde{\beta}_0| \right\} = 0$$

$$\sum_{i=1}^N \frac{\partial}{\partial \tilde{\beta}_0} |Y_i - \tilde{\beta}_0| = 0$$

$$\sum_{i=1}^N \operatorname{sgn}(Y_i - \tilde{\beta}_0) = 0$$

$$\Rightarrow \tilde{Y} = \hat{\beta}_0$$

- If we rank all observations from small to large:

$$Y_{(1)} = \min(Y_i); Y_{(N)} = \max(Y_i)$$

- Median value of  $Y_i$ :  $Y_{(N+1)/2}$ .

# Least Absolute Deviation Solution (optional)

$$\hat{\beta}_0 = \underset{\tilde{\beta}_0}{\operatorname{argmin}} \sum_{i=1}^N |Y_i - \tilde{\beta}_0|$$

$$\Rightarrow \left. \frac{\partial \mathcal{L}(\tilde{\beta}_0)}{\partial \tilde{\beta}_0} \right|_{\tilde{\beta}_0 = \hat{\beta}_0} = 0$$

$$\frac{\partial}{\partial \tilde{\beta}_0} \left\{ \sum_{i=1}^N |Y_i - \tilde{\beta}_0| \right\} = 0$$

$$\sum_{i=1}^N \frac{\partial}{\partial \tilde{\beta}_0} |Y_i - \tilde{\beta}_0| = 0$$

$$\sum_{i=1}^N \operatorname{sgn}(Y_i - \tilde{\beta}_0) = 0$$

$$\Rightarrow \tilde{Y} = \hat{\beta}_0$$

- If we rank all observations from small to large:

$$Y_{(1)} = \min(Y_i); Y_{(N)} = \max(Y_i)$$

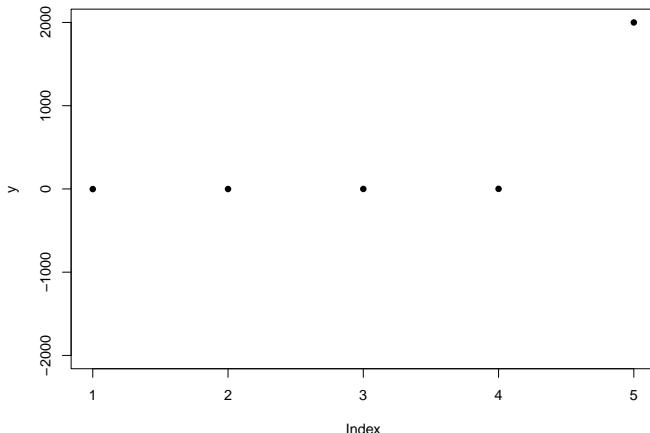
- Median value of  $Y_i$ :  $Y_{(N+1)/2}$ .
- Median more robust to extreme values than mean.

# A toy example

- Five observations:

$$Y_1 = -2, Y_2 = -1, Y_3 = 0, Y_4 = 1, Y_5 = 2000$$

- Mean as 399.6, median as 0

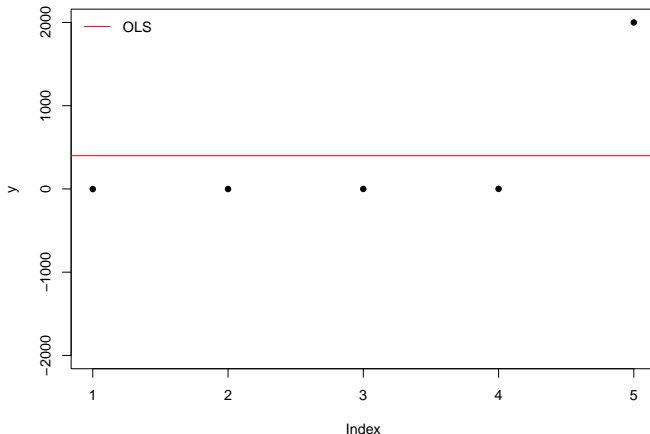


# A toy example

- Five observations:

$$Y_1 = -2, Y_2 = -1, Y_3 = 0, Y_4 = 1, Y_5 = 2000$$

- Mean as 399.6, median as 0

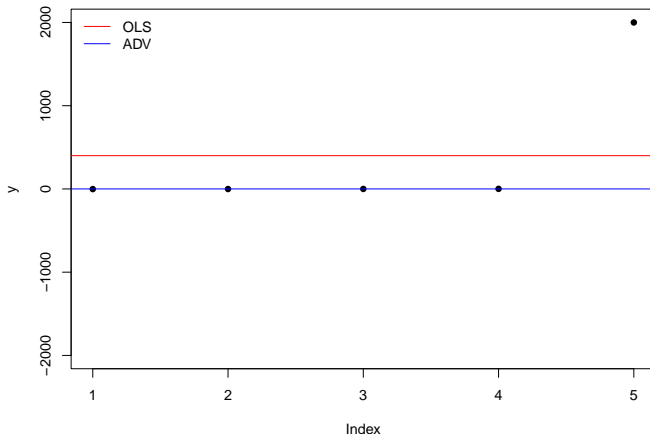


# A toy example

- Five observations:

$$Y_1 = -2, Y_2 = -1, Y_3 = 0, Y_4 = 1, Y_5 = 2000$$

- Mean as 399.6, median as 0





# Linear Regression Model

- A model for a **linear** relationship between two variables

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

- $X$ : Independent (explanatory) variable
- $Y$ : Dependent (outcome, response) variable
- $\epsilon$ : error (disturbance) term

# Linear Regression Model

- A model for a **linear** relationship between two variables

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

- $X$ : Independent (explanatory) variable
- $Y$ : Dependent (outcome, response) variable
- $\epsilon$ : error (disturbance) term
- Given a value of  $X$ , the model predicts the average of  $Y$
- Abuse of regression: extrapolation, causal misinterpretation

**Correlation is not causation!**

# Regression

Regression analysis answers:

1. What is the **best** line that describes an outcome variable (aka dependent variable) in terms of an independent variable?

# Regression

Regression analysis answers:

1. What is the **best** line that describes an outcome variable (aka dependent variable) in terms of an independent variable?
2. Given a value of the independent variable, what is my best guess for the dependent variable?

# Regression

Regression analysis answers:

1. What is the **best** line that describes an outcome variable (aka dependent variable) in terms of an independent variable?
2. Given a value of the independent variable, what is my best guess for the dependent variable?
3. How close is the line to the data?

**Loss function of choice**

# Estimating GDP

Given GDP growth rate in 2007, how can we estimate GDP growth rate in 2008?

- Assume: GDP growth in 2008 is GDP growth in 2007 times a constant plus an intercept

$$Y_i = \beta_0 + \beta_1 X_i$$

# Estimating GDP

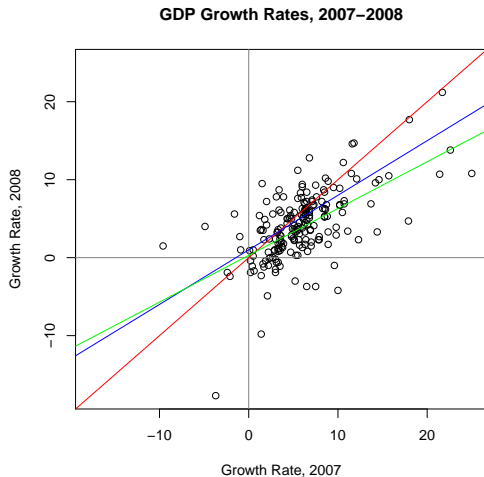
Given GDP growth rate in 2007, how can we estimate GDP growth rate in 2008?

- Assume: GDP growth in 2008 is GDP growth in 2007 times a constant plus an intercept

$$Y_i = \beta_0 + \beta_1 X_i$$

- Do we expect the coefficient estimate of GDP 2007 on GDP 2008 to be positive or negative?

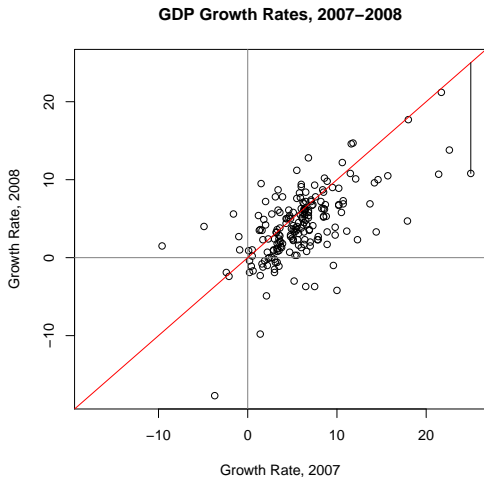
# Which Line to Choose





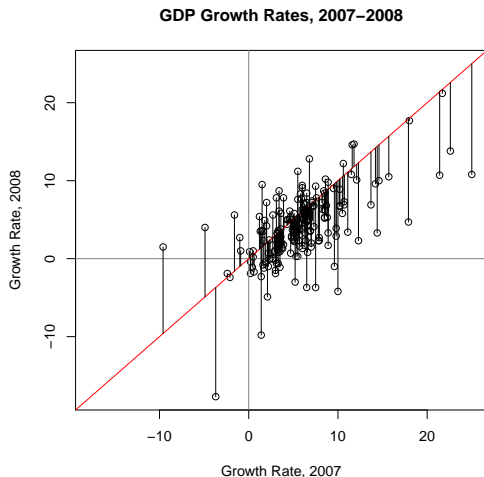
# How far is a point from the line?

The distance from one point to the line, called the residual



# How far are all of the points from the line?

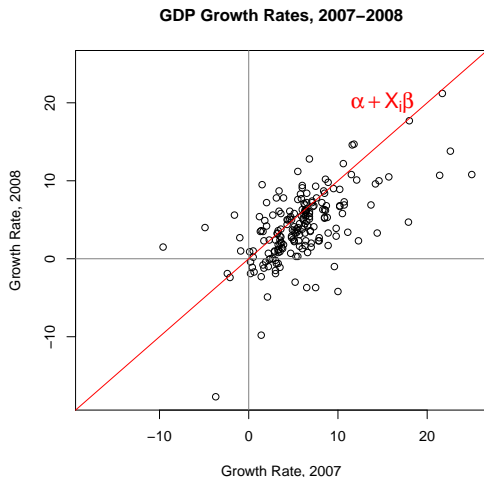
The total distances from the data to the line (residuals)



# Determining the line of best fit

Determining the line of best fit (aka the line of least squares)

- $Y_i$ : 2008 GDP growth rate for country  $i$
- $X_i$ : 2007 GDP growth rate for country  $i$



# How far are all of the points from the line?

To allow for some difference between  $Y_i$  and  $\beta_0 + X_i\beta_1$ , we say

$$Y_i = \beta_0 + X_i\beta_1 + \epsilon_i$$

This is our **assumed model**

# How far are all of the points from the line?

To allow for some difference between  $Y_i$  and  $\beta_0 + X_i\beta_1$ , we say

$$Y_i = \beta_0 + X_i\beta_1 + \epsilon_i$$

This is our **assumed model**

After we see our data, we are going to estimate a model,

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

This is our **fitted model** or estimated model

# How far are all of the points from the line?

To allow for some difference between  $Y_i$  and  $\beta_0 + X_i\beta_1$ , we say

$$Y_i = \beta_0 + X_i\beta_1 + \epsilon_i$$

This is our **assumed model**

After we see our data, we are going to estimate a model,

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

This is our **fitted model** or estimated model

The fitted model varies from sample to sample (like in a survey).  
The assumed model does not necessarily.

# Assumptions

- Linearity among variables and error terms.  
**age and age square, income**

# Assumptions

- Linearity among variables and error terms.  
**age and age square, income**
- Error terms have a mean of zero



# Assumptions

- Linearity among variables and error terms.

**age and age square, income**

- Error terms have a mean of zero
- Error terms are uncorrelated with each other.

**One observation of the error term should not predict the next observation**

# Regression in Observational Causal Inference

## Regression and Difference in Difference

$$Y_i = \beta_0 + \beta_1 \cdot \text{period} + \beta_2 \cdot \text{treatment} + \beta_3 \cdot \text{period} \cdot \text{treatment} + \epsilon_i$$

# Regression in Observational Causal Inference

## Regression and Difference in Difference

$$Y_i = \beta_0 + \beta_1 \cdot \text{period} + \beta_2 \cdot \text{treatment} + \beta_3 \cdot \text{period} \cdot \text{treatment} + \epsilon_i$$

	Period= 0	Period= 1
Treatment = 0	$Y_i = \beta_0 + \epsilon_i$	$Y_i = \beta_0 + \beta_1 + \epsilon_i$
Treatment = 1	$Y_i = \beta_0 + \beta_2 + \epsilon_i$	$Y_i = \beta_0 + \beta_1 + \beta_2 + \beta_3 + \epsilon_i$

The DID estimator would be:

$$\underbrace{(\beta_0 + \beta_1 + \beta_2 + \beta_3 + \epsilon_i)}_{\text{Period 1, treated}} - \underbrace{(\beta_0 + \beta_2 + \epsilon_i)}_{\text{Period 0, treated}} - \left( \underbrace{(\beta_0 + \beta_1 + \epsilon_i)}_{\text{Period 1, control}} - \underbrace{(\beta_0 + \epsilon_i)}_{\text{Period 0, control}} \right)$$

Which yields:  $\beta_3$

# Regression in Observational Causal Inference

- Regression and Instrumental Variables

This video ([link](#)) explains why two stage least square regression gives us the IV estimator.

# Regression in Observational Causal Inference

- Regression and Instrumental Variables

This video ([link](#)) explains why two stage least square regression gives us the IV estimator.

- Regression vs Matching

# Regression in Observational Causal Inference

- Regression and Instrumental Variables

This video (link) explains why two stage least square regression gives us the IV estimator.

- Regression vs Matching

► Regression without matching:  $Y_i = \beta_0 + \beta_1 \cdot T_i + \beta \cdot X + \epsilon_i$

# Regression in Observational Causal Inference

- Regression and Instrumental Variables

This video (link) explains why two stage least square regression gives us the IV estimator.

- Regression vs Matching

- ▶ Regression without matching:  $Y_i = \beta_0 + \beta_1 \cdot T_i + \beta \cdot X + \epsilon_i$
- ▶ With matching