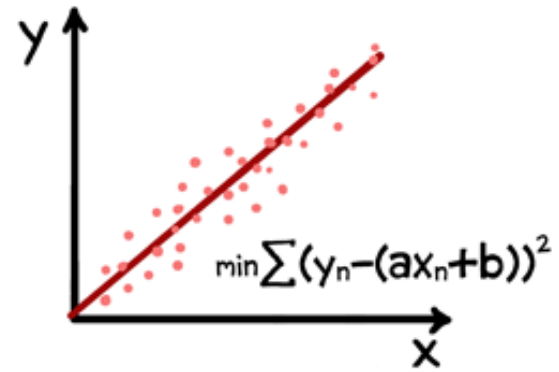
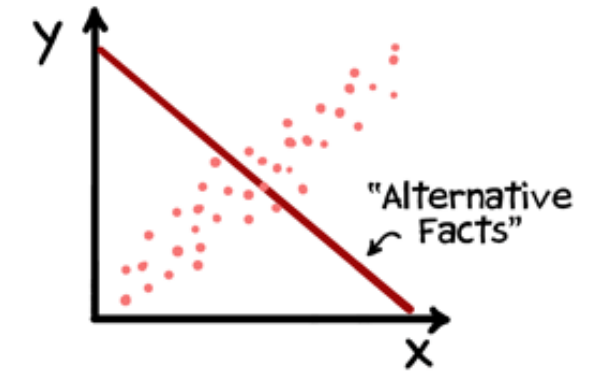


## Linear Regression



JORGE CHAM © 2016

## Societal Regression



WWW.PHDCOMICS.COM

INFO 251: Applied Machine Learning

# Regression and Impact Evaluation

# Announcements

- PS2 posted
- Tentative alternate quiz time: 7am (on day of quiz)
  - If you need/want to take the quiz at this alternate time, email me and GSI (Qutub) with an explanation
- Please have a pen and paper handy for today's lecture

# Course Outline

- Causal Inference and Research Design
  - **Experimental methods**
  - Non-experiment methods
- Machine Learning
  - Design of Machine Learning Experiments
  - Linear Models and Gradient Descent
  - Non-linear models
  - Neural models
  - Unsupervised Learning
  - Practicalities, Fairness, Bias
- Special topics

# Key Concepts (last lecture)

- Impact Evaluation
- Counterfactuals
- Identifying assumptions
- Single Difference research design
- Pre vs. Post research design
- Difference-in-Difference (Double Difference) research design
- Differential Trends
- Progresas

# Outline

- Regression recap
- Regression and Impact Evaluation
- Heterogeneous treatment effects
- Double-Difference via Regression
- Fixed effects and Normalization

# Key Concepts (today's lecture)

- Regression and causal effect estimation
- Dummy variables, “one-hot” vectors
- Heterogeneous treatment effects
- Parallel trends assumption
- Interaction variables
- Cross-sectional vs. panel data
- Between vs. within variation
- Difference regressions
- Normalization
- Fixed effects

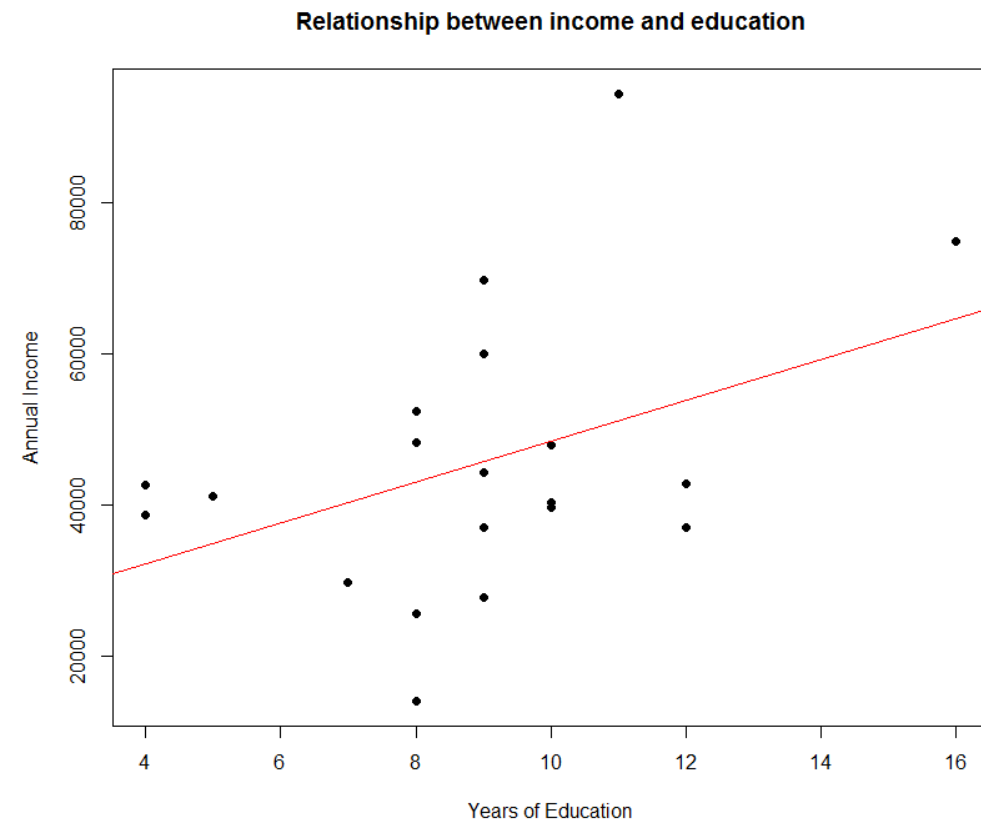
# Regression and Impact Evaluation

- Thus far we've used cross-tabs to “eyeball” the impact of a treatment
- Advantages
  - Very simple to compute
  - Easily interpretable
- Disadvantages
  - How to measure statistical precision?
  - How to deal with known confounds (e.g., differential trends)?
  - May be unbiased, but may not be precise - controlling for additional factors may increase precision

# Regression: Quick Recap

- Linear regression offers a concise summary of the mean of one variable as a function of the other variable through two parameters: the slope and the intercept of the regression line
  - Causality often *implied*, rarely *justified*

	Education	Age	Income
[1,]	8	35	30942.35
[2,]	8	23	37323.89
[3,]	8	58	49381.84
[4,]	5	41	31680.86
[5,]	13	35	81147.84
[6,]	9	43	38682.86
[7,]	8	35	34632.30
[8,]	7	56	14394.98
[9,]	11	62	22243.85
[10,]	14	24	51831.79
[11,]	12	25	23963.90
[12,]	12	32	66780.27
[13,]	4	41	26979.73
[14,]	8	49	38837.48
[15,]	10	21	40726.37
[16,]	8	33	40269.51
[17,]	4	36	34293.32
[18,]	10	38	61158.98
[19,]	11	36	64329.59
[20,]	9	48	51069.77





# Regression: Quick Recap

- Simple bivariate (linear) regression

- The regression model

$$wages_i = \alpha + \beta * education_i + error_i$$

- The fitted model

$$wages_i = 12409 + 3310 * education_i + error_i$$

- Intuition check

- What does  $\beta$  tell us?
- What is 12409?
- What are the expected wages be for someone with 14 years of education?
  - $12409 + 14 * 3310 = 58,749$

# Regression: Quick Recap

- Regression with binary predictor/independent variables

- The regression model

$$wages_i = \alpha + \beta * isForeign_i + error_i$$

- The fitted model

$$wages_i = 54212 - 2710 * isForeign_i + error_i$$

- Multiple (linear) regression

$$wages_i = \alpha + \beta * education_i + \gamma * isForeign_i + error_i$$

# Regression: Categorical variables

- What if our control variables are categorical?
  - Example: We want to study the relationship between wages and education, controlling for country
  - $Wages_i = \alpha + \beta Education_i + \gamma Country_i + \epsilon_i$
  
- How to deal with a categorical predictor?
  - Convert to a single binary variable:
    - $Wages_i = \alpha + \beta Education_i + \gamma USA_i + \epsilon_i$
    - $USA_i = 1$  iff worker  $i$  is from USA,  $USA_i = 0$  otherwise
    - Makes sense if we care about the effect of one category relative to others
  - Convert to a set of binary variables:
    - $Wages_i = \beta Education_i + \gamma_1 USA_i + \gamma_2 CHINA_i + \dots + \gamma_M Country_m + \epsilon_i$
    - $Wages_i = \beta Education_i + \sum_{c=1}^M \gamma_c \mathbf{1}(Country_i = c) + \epsilon_i$
    - $Wages_i = \beta Education_i + Country_i + \epsilon_i \leftarrow$  this is an abuse of notation, but it is very common

# Regression: “Dummy” variables

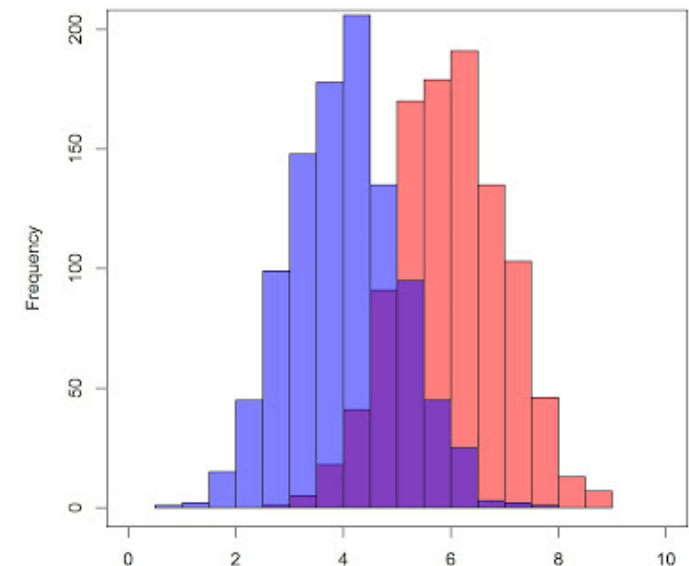
- Converting levels to a series of binary variables
  - $Y_i = \alpha + \beta Education_i + Country_i + \epsilon_i$
  - $Country_i$  is a “dummy variable” or “one-hot vector” or “fixed effect”
- Interpretation
  - Equivalent to creating a country-specific intercept
  - Each of the  $Country$  coefficients indicates average wage for workers from that particular country (when  $Education_i = 0$ ), *relative to the reference country*
  - With M countries, we have (M-1) coefficients and an intercept  $\alpha$
  - Alternatively, estimate with no intercept and M coefficients
    - $Y_i = \beta Education_i + Country_i + \epsilon_i$
    - Intuition check: Will the coefficients for “country” dummies be the same in both cases?

# Outline

- Regression recap
- **Regression and Impact Evaluation**
- Heterogeneous treatment effects
- Double-Difference via Regression
- Fixed effects and Normalization

# Regression and Impact: Basics

- How to measure the effect of treatment  $T$  on outcome  $Y$  in a regression?
  - The regression equation:
$$Y_i = \alpha + \beta T_i + \epsilon_i$$
  - Example: We estimate the effect of eating a cookie on happiness on a scale of 1-10. We estimate  $\hat{\alpha} = 4.1, \hat{\beta} = 1.3$ . What does this mean?
- If  $T$  is randomly assigned,  $\hat{\beta}$  is an estimate of the **causal impact** of  $T$  on  $Y$



# Revisiting Bertrand et al.

“Applied for loan”

“Female Photo”

$$Y_i = \alpha + \beta T_i + \epsilon_i$$

TABLE III

EFFECTS OF ADVERTISING CONTENT ON BORROWER BEHAVIOR

Dependent variable	Applied for loan before mailer deadline	Applied for loan before mailer deadline	Applied for loan before mailer deadline	Obtained loan before mailer deadline	Loan amount obtained before mailer deadline	Loan in collection status	Borrowed from other lender
Sample	Full	Males	Females	Full	Full	Obtained	Full
Estimator	Probit	Probit	Probit	Probit	OLS	Probit	Probit
Mean (dependent variable)	0.0850	0.0824	0.0879	0.0741	110.4363	0.1207	0.2183
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
Monthly interest rate in percentage point units (e.g., 8.2)	-0.0029*** (0.0005)	-0.0025*** (0.0007)	-0.0034*** (0.0008)	-0.0026*** (0.0005)	-4.7712*** (0.8238)	0.0071*** (0.0022)	0.0009 (0.0008)
1 = no photo	0.0013 (0.0040)	-0.0050 (0.0048)	0.0021 (0.0055)	0.0029 (0.0037)	3.9316 (7.6763)	0.0013 (0.0166)	-0.0024 (0.0060)
1 = female photo (System I: affective response)	0.0057** (0.0026)	0.0079** (0.0034)	0.0032 (0.0038)	0.0056** (0.0024)	8.3292 (5.0897)	-0.0076 (0.0107)	-0.0047 (0.0040)
1 = photo gender matches client's (System I: affinity/similarity)	0.0026 (0.0026)			-0.0033 (0.0024)	-7.1773 (5.0850)	-0.0059 (0.0107)	0.0041 (0.0040)
1 = photo race matches client's (System I: affinity/similarity)	-0.0056 (0.0048)	-0.0014 (0.0064)	-0.0099 (0.0070)	-0.0035 (0.0044)	9.0638 (10.4079)	0.0181 (0.0176)	-0.0018 (0.0072)
1 = one example loan shown (System I: avoid choice overload)	0.0068** (0.0028)	0.0099*** (0.0038)	0.0031 (0.0040)	0.0075*** (0.0026)	2.4394 (4.8383)	0.0073 (0.0117)	-0.0043 (0.0042)
1 = interest rate shown (System I: several, potentially offsetting, channels)	0.0025 (0.0030)	-0.0017 (0.0042)	0.0073 (0.0044)	0.0043 (0.0028)	2.8879 (6.7231)	0.0140 (0.0123)	0.0007 (0.0049)

# Regression: “Control” variables

- How to simultaneously measure the effect of a treatment  $T$  and a non-experimental control variable  $X$  on an outcome  $Y$  in a regression setting?

$$Y_i = \alpha_1 + \beta_1 T_i + \gamma X_i + e_i$$

- How is this different from a version without control variables?

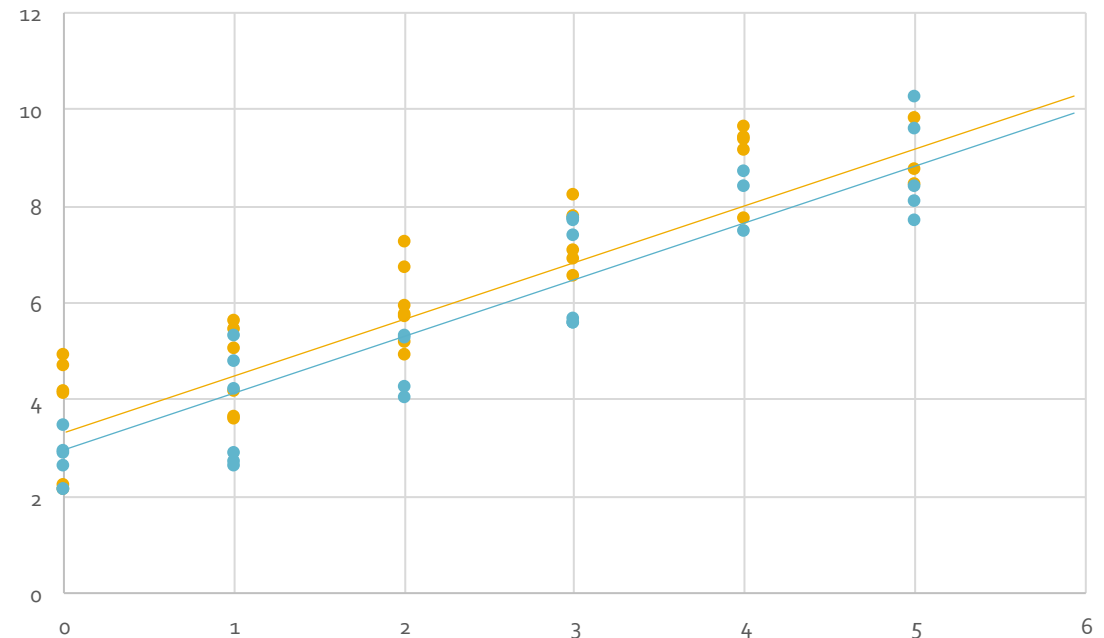
$$Y_i = \alpha_2 + \beta_2 T_i + \epsilon_i$$

- In a perfectly randomized experiment...
  - What, if anything, can we say about  $Cor(T_i, X_i)$ ?
  - What, if anything, can we say about our estimates of  $\beta_1$  and  $\beta_2$ ?
  - What, if anything, can we say about  $\gamma$ ?



# Control variables: Example

- Example: We are estimating effect of eating a cookie on happiness on a scale of 1-10, while controlling for years in grad school
- Regression equation?
  - $Y_i = \alpha + \beta T_i + \gamma X_i + \epsilon_i$
- Coefficient estimates
  - $\hat{\alpha} = 3.4, \hat{\beta} = -0.5, \hat{\gamma} = 1.2$
  - What do these results mean?



# Outline

- Regression recap
- Regression and Impact Evaluation
- **Heterogeneous treatment effects**
- Double-Difference via Regression
- Fixed effects and Normalization

# Treatment effect heterogeneity

- How to simultaneously measure the effect of a treatment  $T$  **and** a non-experimental control variable  $X$  **and** a differential effect of treatment by control variable on an outcome  $Y$  in a regression setting?
  - When might this happen in practice?
    - The treatment effect is different for different types of people (where “type” is measured with  $X_i$ )
    - For instance: Tall students may respond more to cookies than short students
  - We can estimate this with a regression:

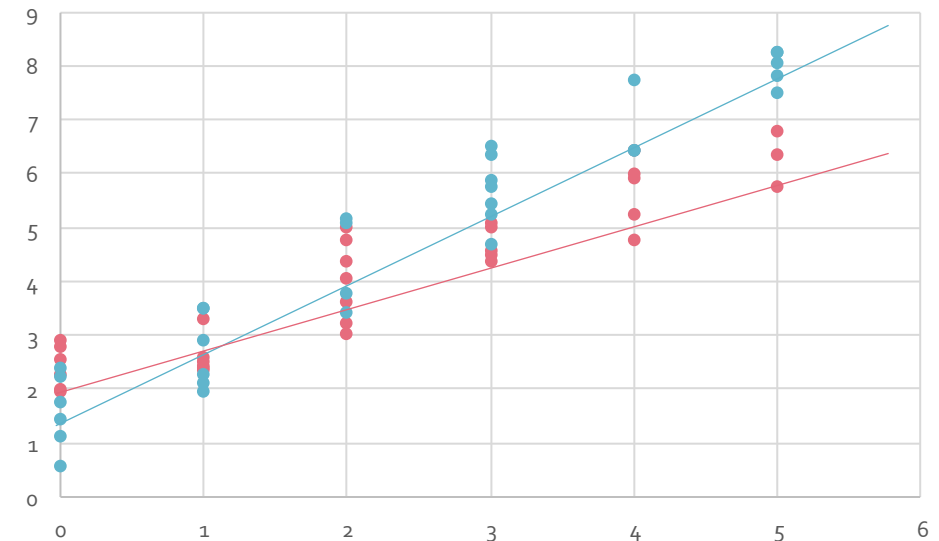
$$Y_i = \alpha + \beta T_i + \gamma X_i + \delta(T_i * X_i) + \epsilon_i$$

# Binary heterogeneity

- Example: We are estimating the effect of eating a cookie (T) on happiness (Y) on a scale of 1-10, while controlling for height (X) **and** the interaction (T\*X)
  - $Y_i = \alpha + \beta T_i + \gamma X_i + \delta(T_i * X_i) + \epsilon_i$
  - $\hat{\alpha} = 1.4, \hat{\beta} = -1.0, \hat{\gamma} = 0.8, \hat{\delta} = 0.6$
  - What do these results mean?
  - What is the expected happiness for:
    - A short student with a cookie
    - A tall student with a cookie
    - A short student without a cookie
    - A tall student without a cookie

# Treatment effect heterogeneity

- Example: We are estimating the effect of eating a cookie on happiness on a scale of 1-10, while controlling for grad school years **and** the interaction effect. We estimate:
  - $Y_i = \alpha + \beta T_i + \gamma X_i + \delta(T_i * X_i) + \epsilon_i$
  - $\hat{\alpha} = 1.4, \hat{\beta} = -1.0, \hat{\gamma} = 0.8, \hat{\delta} = 0.6$
  - What do these results mean?
  - How to visualize?
  - Draw regression lines!



# Outline

- Regression recap
- Regression and Impact Evaluation
- Heterogeneous treatment effects
- **Double-Difference via Regression**
- Fixed effects and Normalization

# Regression and Impact: Diff-in-Diff

- In a typical double-difference setting, data is collected at baseline (pre-treatment) and at endline (post-treatment)
- Very similar to a situation where you have a treatment  $T$  **and** a non-experimental binary control variable  $X$  **and** a differential effect of treatment by a binary control variable
- In this case, the control variable is time!
  - Instead of a (continuous)  $X_i$  we have a binary  $Post_i$

# Regression and Impact: Diff-in-Diff

- For example, our data look like this:

ID	Period	Treatment	Y	...
1	PRE	Treat	12	...
1	POST	Treat	14	...
2	PRE	Control	11	...
2	POST	Control	11	...
3	Pre	Control	24	...
...	...	...	...	...
12419	...	...	...	...

	Control	Treat
Pre	A	B
Post	C	D

- Write down the regression equation:

$$Y_i = \alpha + \beta T_i + \gamma Post_i + \delta(T_i * Post_i) + \epsilon_i$$

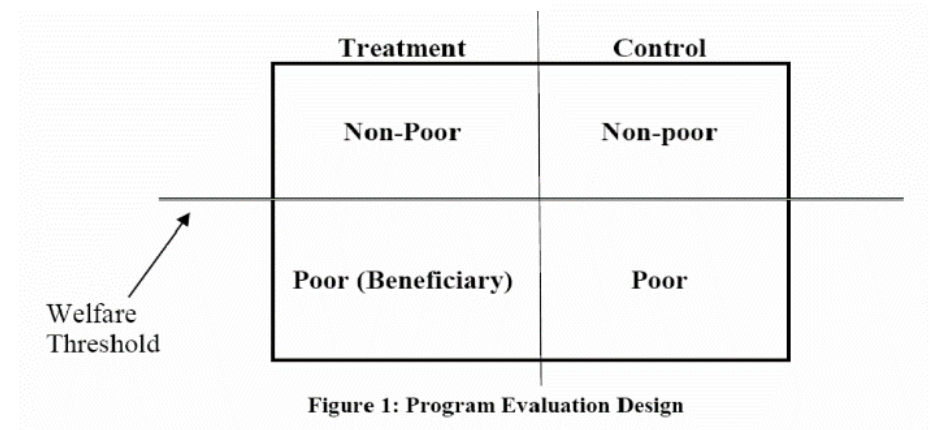


# Regression and Impact: Diff-in-Diff

- Let's return to Progresa (Shultz, 2004: page 209)

$$S_i = \alpha_0 + \alpha_1 P_i + \alpha_2 E_i + \alpha_3 P_i E_i + \sum_{k=1}^K \gamma_{ki} C_{ki} + \sum_{j=1}^J \beta_j X_{ji} + e_i \quad i = 1, 2, \dots, n \quad (1)$$

- $P_i$  = Progresa village ("T")
- $E_i$  = Eligible (poor) household ("X")
- $P_i E_i$  = Difference-in-difference estimator
- $C_{ki}$  = Dummy variable for grade of child
  - Controls for fact that enrollment rates vary across grades
- $X_{ji}$  = Other control variables (I think?)



# Regression and Impact: Summary

- Double Difference

$$Y_i = \alpha + \beta T_i + \gamma Post_i + \delta(T_i * Post_i) + \epsilon_i$$

- Simple difference

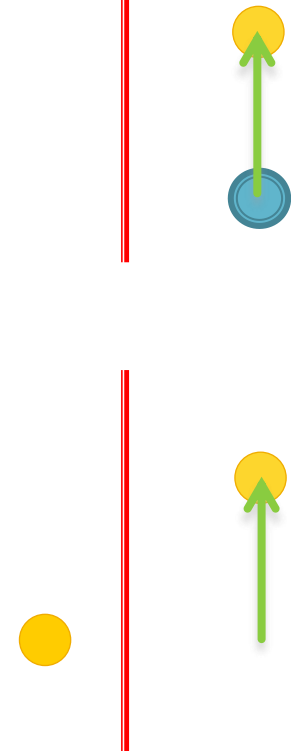
$$Y_i = \alpha + \beta T_i + \epsilon_i$$

- Estimated only in Post period

- Pre vs. Post

$$Y_i = \alpha + \gamma Post_i + \epsilon_i$$

- Estimated only on treatment group



# Outline

- Regression recap
- Regression and Impact Evaluation
- Heterogeneous treatment effects
- Double-Difference via Regression
- **Fixed effects and Normalization**

# Fixed Effects and Normalization

- **Cross-sectional data:** Data from several units at a single point in time
  - Schooling outcomes in 1998
  - Happiness of students on Jan 20
  - Number of logins per person in March
  
- **Time series (panel) data:** Multiple observations for each unit over time
  - Schooling outcomes in 1997 and 1998
  - Happiness of students on Jan 20 and Jun 20
  - Number of logins per person in each month of 2010

# Between and Within Variation

- **Cross-sectional data:**

- Variation is between/across units

Location	Year	Price	Quantity sold (per capita)
Chicago	2003	\$75	2.0
Seattle	2003	\$50	1.0
Milwaukee	2003	\$60	1.5
Madison	2003	\$55	0.8

- What do you notice about relationship between price and quantity?
  - Across the four cities, price and quantity are positively correlated
  - This is not what we would expect – omitted variables likely matter

# Between and Within Variation

- **Panel data:**
  - Variation is between/across units *and* **within units over time**
- What do you notice about relationship between price and quantity?
  - Within each of the four cities, price and quantity are inversely correlated (as expected with downward sloping demand)

Location	Year	Price	Quantity
Chicago	2003	\$75	2.0
Chicago	2004	\$85	1.8
Seattle	2003	\$50	1.0
Seattle	2004	\$48	1.1
Milwaukee	2003	\$60	1.5
Milwaukee	2004	\$65	1.4
Madison	2003	\$55	0.8
Madison	2004	\$60	0.7

# Exploiting Within-Unit Variation

- How to isolate changes in outcomes correlated with changes *within* a unit ?
  - e.g., Changes in demand caused by changes in price *within* a given city (over time)

# Exploiting Within-Unit Variation

- One (familiar?) approach: “difference” regressions
  - Isolate differences over time in  $X$  and  $Y$
  - Instead of:  $Y_{it} = \alpha + \beta X_{it} + \epsilon_{it}$
  - $(Y_{it} - Y_{i(t-1)}) = \alpha + \beta (X_{it} - X_{i(t-1)}) + \epsilon_{it}$ 
    - Are *changes* in  $Y$  related to *changes* in  $X$ ?
- Another approach: normalization
  - Instead of:  $Y_{it} = \alpha + \beta X_{it} + \epsilon_{it}$
  - $(Y_{it} - \bar{Y}_i) = \alpha + \beta (X_{it} - \bar{X}_i) + \epsilon_{it}$



# Exploiting Within-Unit Variation

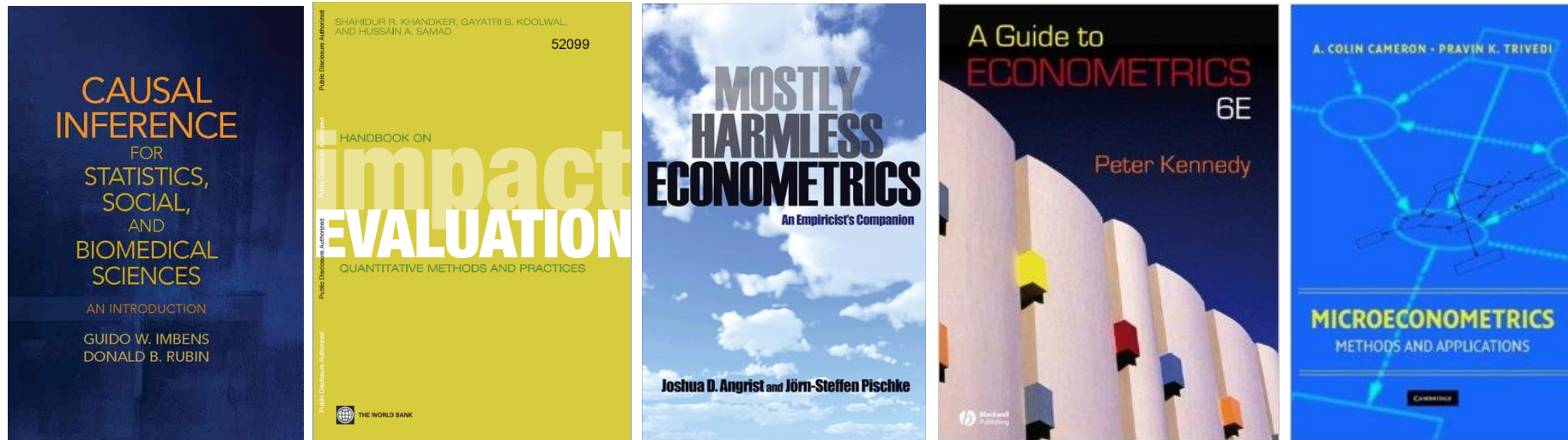
- A third approach: **Fixed effects**
- Basic idea (refer to lecture notes for details)
  - Instead of:  $(Y_{it} - \bar{Y}_i) = \alpha + \beta(X_{it} - \bar{X}_i) + \epsilon_{it}$
  - Add “dummy” variables for each  $i$ :  $Y_{it} = \alpha + \beta X_{it} + (\mu_{i=2} + \dots + \mu_{i=N}) + \epsilon_{it}$ 
    - Equivalent to adding an intercept for each  $i$
    - Conceptually the same as the “country fixed effect” example from slide 12
  - Shorthand:  $Y_{it} = \alpha + \beta X_{it} + \mu_i + \epsilon_{it}$
- Note: we can also “normalize” for time FE's
- $Y_{it} = \alpha + \beta X_{it} + \mu_i + \pi_t + \epsilon_{it}$

# Fixed Effects: Digging Deeper

- Key advantage of Fixed Effects
  - Fixed effects control for unobserved heterogeneity
  - They remove the effect of time-invariant characteristics to assess the net effect of the predictors on the outcome
- Extensions
  - Time trends
  - Region-specific slopes

## Additional Resources

Beginner —————> Advanced



# For Next Class:

- Read Chapters 6 and 7:
- Get started on problem set 2!

## Handbook on Impact Evaluation

Quantitative Methods and Practices

Shahidur R. Khandker  
Gayatri B. Koolwal  
Hussain A. Samad