

Gov 50: 13. Regression and Causality

Matthew Blackwell

Harvard University

Fall 2018

1. Today's agenda
2. Randomized experiments with regression
3. Categorical variables
4. Interaction terms

1/ Today's agenda

Where are we? Where are going?

- Past two weeks:
 - ▶ Predicting with past values
 - ▶ Predicting with another variable (linear regression)
- Today:
 - ▶ Analyzing experiments with regression
 - ▶ Interactions for estimating varying treatment effects
- HW3 due Thursday night.

2/ Randomized experiments with regression

Political effects of gov't programs

- Around 2000, Mexico implemented a conditional cash transfer program (CCT) called *Progresa*
 - ▶ Welfare payments given if children are enrolled in schools, get regular check-ups, etc.
- Do these programs have political effects?
 - ▶ Program had support from most parties.
 - ▶ Was implemented in a nonpartisan fashion.
 - ▶ Would the incumbent presidential party be rewarded?
- Randomized roll-out of the CCT program:
 - ▶ treatment: receive CCT 21 months before 2000 election
 - ▶ control: receive CCT 6 months before 2000 election
- Hypothesis: having CCT longer would mobilize voters for incumbent PRI party.

The data

Name	Description
<code>treatment</code>	early Progresa (1) or late Progresa (0)
<code>pri2000s</code>	PRI votes in the 2000 election as a share of adults in precinct
<code>t2000</code>	turnout in the 2000 election as share of adults in precinct

```
cct <- read.csv("data/progresa.csv")
```

Difference in means estimates

- Does CCT affect turnout?

```
mean(cct$t2000[cct$treatment == 1]) -  
  mean(cct$t2000[cct$treatment == 0])
```

```
## [1] 4.27
```

- Does CCT affect PRI (incumbent) votes?

```
mean(cct$pri2000s[cct$treatment == 1]) -  
  mean(cct$pri2000s[cct$treatment == 0])
```

```
## [1] 3.62
```


Binary independent variables

$$Y_i = \alpha + \beta X_i + \varepsilon_i$$

- When independent variable X_i is **binary**:
 - ▶ Intercept $\hat{\alpha}$ is the average outcome in the $X = 0$ group.
 - ▶ Slope $\hat{\beta}$ is the difference-in-means of Y between $X = 1$ group and $X = 0$ group.
- If there are other independent variables, this becomes the difference-in-means controlling for those covariates.

Linear regression for experiments

- Allows us to estimate the ATE with regression (as long as we have randomization!):

```
mean(cct$pri2000s[cct$treatment == 1]) -  
mean(cct$pri2000s[cct$treatment == 0])
```

```
## [1] 3.62
```

```
lm(pri2000s ~ treatment, data = cct)
```

```
##
```

```
## Call:
```

```
## lm(formula = pri2000s ~ treatment, data = cct)
```

```
##
```

```
## Coefficients:
```

```
## (Intercept)      treatment
```

```
##          34.49          3.62
```

3/ Categorical variables

Categorical variables in regression

- We often have **categorical variables**:
 - ▶ Race/ethnicity: white, black, Latino, Asian.
 - ▶ Partisanship: Democrat, Republican, Independent
- Strategy for including in a regression: create a **series of binary variables**

Unit	Party	Democrat	Republican	Independent
1	Democrat	1	0	0
2	Democrat	1	0	0
3	Independent	0	0	1
4	Republican	0	1	0
⋮	⋮	⋮	⋮	⋮

- Then include **all but one** of these categorical variables:

$$\text{turnout}_i = \alpha + \beta_1 \text{Republican}_i + \beta_2 \text{Independent}_i + \varepsilon_i$$

Interpreting categorical variables

$$\text{turnout}_i = \alpha + \beta_1 \text{Republican}_i + \beta_2 \text{Independent}_i + \varepsilon_i$$

- $\hat{\alpha}$: average outcome in the **omitted group** (Democrats).
- Other coefficients: difference-in-means between that group and the omitted group.
 - ▶ $\hat{\beta}_1$: average difference in turnout rates between Republicans and Democrats
 - ▶ $\hat{\beta}_2$: average difference in turnout rates between Independents and Democrats

Social pressure experiment

- Experimental study where each household for 2006 MI primary was randomly assigned to one of 4 conditions:
 - ▶ Control: no mailer
 - ▶ Civic Duty: mailer saying voting is your civic duty.
 - ▶ Hawthorne: a “we’re watching you” message.
 - ▶ Neighbors: naming-and-shaming social pressure mailer.
- Outcome: whether household members voted or not.

Neighbors mailer

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	_____
9995 JENNIFER KAY SMITH		Voted	_____
9997 RICHARD B JACKSON		Voted	_____
9999 KATHY MARIE JACKSON		Voted	_____

Social pressure data

```
social <- read.csv("data/social.csv")
head(social[, c("messages", "control", "civic",
               "hawthorne", "neighbors", "primary2006")])
```

```
##      messages control civic hawthorne neighbors
## 1 Civic Duty      0      1          0          0
## 2 Civic Duty      0      1          0          0
## 3 Hawthorne      0      0          1          0
## 4 Hawthorne      0      0          1          0
## 5 Hawthorne      0      0          1          0
## 6   Control      1      0          0          0
##      primary2006
## 1              0
## 2              0
## 3              1
## 4              1
## 5              1
## 6              0
```


Categorical variables in R

```
lm(primary2006 ~ civic + hawthorne + neighbors, data = social)
```

```
##
```

```
## Call:
```

```
## lm(formula = primary2006 ~ civic + hawthorne + neighbors, data = s
```

```
##
```

```
## Coefficients:
```

```
## (Intercept)      civic    hawthorne    neighbors
##      0.2966      0.0179      0.0257      0.0813
```

- (Intercept): average turnout when all independent vars = 0
 - ▶ \rightsquigarrow ~30% turnout rate in the “Control” condition
- neighbors: difference in turnout rates between “Civic Duty” condition and “Control” condition.
 - ▶ \rightsquigarrow social pressure mailer leads to 8pp increase in turnout rates.

Factor variables in lm()

- Including a **factor variable** in `lm()` will automatically create binary variables and exclude one group:

```
lm(primary2006 ~ messages, data = social)
```

```
##  
## Call:  
## lm(formula = primary2006 ~ messages, data = social)  
##  
## Coefficients:  
##      (Intercept)      messagesControl  
##      0.31454      -0.01790  
## messagesHawthorne  messagesNeighbors  
##      0.00784      0.06341
```

- Omitted group is “Civic Duty” \rightsquigarrow not ideal!

Changing the factor reference level

- To see what group will be the reference, check the `levels()` function:

```
levels(social$messages)
```

```
## [1] "Civic Duty" "Control"      "Hawthorne"  "Neighbors"
```

- Can change the omitted group using `relevel()`:

```
social$messages <- relevel(social$messages, ref = "Control")  
levels(social$messages)
```

```
## [1] "Control"      "Civic Duty"  "Hawthorne"  "Neighbors"
```

Comparing the results

```
coef(lm(primary2006 ~ civic + hawthorne + neighbors, data = social))
```

```
## (Intercept)      civic    hawthorne    neighbors
##      0.2966      0.0179      0.0257      0.0813
```

```
coef(lm(primary2006 ~ messages, data = social))
```

```
##      (Intercept) messagesCivic Duty
##      0.2966      0.0179
## messagesHawthorne messagesNeighbors
##      0.0257      0.0813
```

```
mean(social$primary2006[social$neighbors == 1]) -
  mean(social$primary2006[social$control == 1])
```

```
## [1] 0.0813
```

4/ Interaction terms

Heterogeneous treatment effects

- **Heterogeneous treatment effects:** the treatment effect varies across groups.
 - ▶ Average effect of a drug is 0, but positive for men and negative for women.
 - ▶ Massively important questions for determining who should receive treatment.
- Social pressure experiment:
 - ▶ `primary2004` measures whether the person voted in 2004, before the experiment.
 - ▶ Do 2004 voters respond differently to social pressure mailer than non-voters?
- Two approaches:
 - ▶ Subsets, subsets, subsets.
 - ▶ Interaction terms in regression.

Subset approach

- Easy way to estimate heterogeneous effects: our old friend, `subset()`.
- First, estimate the ATE for the voters:

```
voters <- subset(social, primary2004 == 1)
ate.v <- mean(voters$primary2006[voters$neighbors == 1]) -
  mean(voters$primary2006[voters$control == 1])
ate.v
```

```
## [1] 0.0965
```

- Now, estimate the ATE for the nonvoters:

```
nonvoters <- subset(social, primary2004 == 0)
ate.nv <- mean(nonvoters$primary2006[nonvoters$neighbors == 1]) -
  mean(nonvoters$primary2006[nonvoters$control == 1])
ate.nv
```

```
## [1] 0.0693
```

Difference in effects

- How much does the estimated treatment effect differ between groups?

```
ate.v - ate.nv
```

```
## [1] 0.0272
```

- Any easier way to allow for different effects of treatment by groups?

Interaction terms

- Can allow for different slopes/coefficients/effects of a variable by including an **interaction term**:

$$\text{turnout}_i = \alpha + \beta_1 \text{primary2004}_i + \beta_2 \text{neighbors}_i + \beta_3 (\text{primary2004}_i \times \text{neighbors}_i) + \varepsilon_i$$

- Literally a new variable that the primary 2004 variable multiplied by the neighbors variable.
- Equal to 1 if voted in 2004 ($\text{primary2004} == 1$) and received neighbors mailer ($\text{neighbors} == 1$)
- Logic comes through when considering the predicted values from the regression.

Predicted values from non-interacted model

- Let $X_i = \text{primary2004}_i$ and $Z_i = \text{neighbors}_i$:

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta}_1 X_i + \hat{\beta}_2 Z_i$$

	Control ($Z_i = 0$)	Neighbors ($Z_i = 1$)
non-voter ($X_i = 0$)	$\hat{\alpha} + \hat{\beta}_1 0 + \hat{\beta}_2 0 = \hat{\alpha}$	$\hat{\alpha} + \hat{\beta}_1 0 + \hat{\beta}_2 1 = \hat{\alpha} + \hat{\beta}_2$
voter ($X_i = 1$)	$\hat{\alpha} + \hat{\beta}_1$	$\hat{\alpha} + \hat{\beta}_1 + \hat{\beta}_2$

- Effect of Neighbors for non-voters: $(\hat{\alpha} + \hat{\beta}_2) - \hat{\alpha} = \hat{\beta}_2$
- Effect of Neighbors for voters: $(\hat{\alpha} + \hat{\beta}_1 + \hat{\beta}_2) - (\hat{\alpha} + \hat{\beta}_1) = \hat{\beta}_2$

Predicted from interacted model

- Now for the interacted model:

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta}_1 X_i + \hat{\beta}_2 Z_i + \hat{\beta}_3 X_i Z_i$$

	Control ($Z_i = 0$)	Neighbors ($Z_i = 1$)
non-voter ($X_i = 0$)	$\hat{\alpha} + \hat{\beta}_1 0 + \hat{\beta}_2 0 + \hat{\beta}_3 0 \cdot 0 = \hat{\alpha}$	$\hat{\alpha} + \hat{\beta}_1 0 + \hat{\beta}_2 1 + \hat{\beta}_3 0 \cdot 1 = \hat{\alpha} + \hat{\beta}_2$
voter ($X_i = 1$)	$\hat{\alpha} + \hat{\beta}_1$	$\hat{\alpha} + \hat{\beta}_1 + \hat{\beta}_2 + \hat{\beta}_3$

- Effect of Neighbors for non-voters: $(\hat{\alpha} + \hat{\beta}_2) - \hat{\alpha} = \hat{\beta}_2$
- Effect of Neighbors for voters:
 $(\hat{\alpha} + \hat{\beta}_1 + \hat{\beta}_2 + \hat{\beta}_3) - (\hat{\alpha} + \hat{\beta}_1) = \hat{\beta}_2 + \hat{\beta}_3$

Interpreting coefficients

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta}_1 \text{primary2004}_i + \hat{\beta}_2 \text{neighbors}_i + \hat{\beta}_3 (\text{primary2004}_i \times \text{neighbors}_i)$$

	Control Group	Neighbors Group
2004 primary non-voter	$\hat{\alpha}$	$\hat{\alpha} + \hat{\beta}_2$
2004 primary voter	$\hat{\alpha} + \hat{\beta}_1$	$\hat{\alpha} + \hat{\beta}_1 + \hat{\beta}_2 + \hat{\beta}_3$

- $\hat{\alpha}$: turnout rate for 2004 non-voters in control group.
- $\hat{\beta}_1$: difference between turnout rates between 2004 voters and non-voters.
- $\hat{\beta}_2$: effect of neighbors for 2004 non-voters.
- $\hat{\beta}_3$: difference in the effect of neighbors mailer between 2004 voters and 2004 non-voters.

Interactions in R

- You can include an interaction with `var1:var2`:

```
social.neighbor <- subset(social, neighbors == 1 | control == 1)
fit <- lm(primary2006 ~ primary2004 + neighbors + primary2004:neighbors,
          data = social.neighbor)
coef(fit)
```

```
##           (Intercept)           primary2004
##           0.2371           0.1487
##           neighbors primary2004:neighbors
##           0.0693           0.0272
```

- Compare coefficients to subset approach:

```
ate.nv
```

```
## [1] 0.0693
```

```
ate.v - ate.nv
```

```
## [1] 0.0272
```

On deck

- More interactions.
- Non-linear relationships in regression
- Next week: start with more statistical theory.