

Exercise: Analyzing Field Experiments

Following on from our simulation of omitted variable bias last week, we will now simulate an experimental analysis (you can adapt your script from last week or start again; the first 4 steps are the same):

1. Generate data on a population of 1,000 people. Specifically, create a variable (a vector) that randomly assigns these people to be male or female (50:50). *Hint: In R, try `rbinom` and in Stata, try `rbino`.*

```
N <- 1000
x <- rbinom(N,1,0.5)
```

2. Now we are going to simulate the potential outcomes - a measure of attitudes - *under control* (y_0) for our population. Create another variable of random normally-distributed values with mean of 5 and standard deviation of 1.

```
y0 <- rnorm(N,5,1)
```

3. One problem with observational data is that potential outcomes are often correlated with other variables such as gender. Adjust your value of y_0 to add 1 (one) for all units who are male.

```
y0 <- y0 + x
```

4. Now simulate potential outcomes *under treatment* (y_1) for all units. Define a *constant* treatment effect of $c = 2$ and create another vector $y_1 = y_0 + c$.

```
c <- 2
y1 <- y0 + c
```

5. Next, let's assume a specific Treatment Assignment Mechanism – random assignment. Create a treatment variable D that gives each unit a 50% chance of binary treatment (like a coin flip).

```
data <- tibble(x,y0,y1) %>%
  mutate(D=rbinom(N,1,0.5))
```

6. Are gender and treatment correlated (they were strongly last week)? Calculate the correlation between x and D .

```
cor(data$x,data$D)
```

```
## [1] 0.01317426
```

7. In practice, we only observe one value: y_{obs} . Create a new variable y_{obs} which equals y_1 if $D = 1$ but which equals y_0 if $D = 0$.

```
data <- data %>% mutate(y_obs=case_when(D==1~y1,
                                         D==0~y0))
```

8. Compare the balance of observed pre-treatment covariates (gender) in the treatment and control groups using a difference-in-means test.

```
data %>% t.test(x ~ D, data=.)
```

```
##
## Welch Two Sample t-test
##
## data: x by D
## t = -0.41622, df = 990.88, p-value = 0.6773
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
```

```
## -0.07533930 0.04897279
## sample estimates:
## mean in group 0 mean in group 1
##      0.5031315      0.5163148
```

9. Compare the balance of observed potential outcomes (y_0 , y_1) in the treatment and control groups using a difference-in-means test.

```
data %>% t.test(y0 ~ D, data=.)
```

```
##
## Welch Two Sample t-test
##
## data: y0 by D
## t = -0.040568, df = 997.99, p-value = 0.9676
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -0.1391681 0.1335306
## sample estimates:
## mean in group 0 mean in group 1
##      5.480087      5.482906
```

```
data %>% t.test(y1 ~ D, data=.)
```

```
##
## Welch Two Sample t-test
##
## data: y1 by D
## t = -0.040568, df = 997.99, p-value = 0.9676
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -0.1391681 0.1335306
## sample estimates:
## mean in group 0 mean in group 1
##      7.480087      7.482906
```

10. Based on the observable data, analyze the results of the experiment using a difference-in-means test of outcomes (y_{obs}) by treatment status (D). Interpret the result. Is this an accurate estimate of the treatment effect that we created at the start?

```
data %>% t.test(y_obs ~ D, data=.)
```

```
##
## Welch Two Sample t-test
##
## data: y_obs by D
## t = -28.825, df = 997.99, p-value < 2.2e-16
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -2.139168 -1.866469
## sample estimates:
## mean in group 0 mean in group 1
##      5.480087      7.482906
```

11. As an alternative, run the basic regression of observable outcomes (y_{obs}) on treatment (D). How does this compare to the difference-in-means test above?

```
data %>% lm(y_obs ~ D, data=.) %>%
  stargazer(header=F, keep.stat=c("n"))
```

Table 1:

<i>Dependent variable:</i>	
	y_obs
D	2.003*** (0.070)
Constant	5.480*** (0.050)
Observations	1,000
<i>Note:</i> *p<0.1; **p<0.05; ***p<0.01	

12. Now add the control variable x to the regression. How do the results change?

```
data %>% lm(y_obs ~ D + x, data=.) %>%
  stargazer(header=F, keep.stat=c("n"))
```

Table 2:

<i>Dependent variable:</i>	
	y_obs
D	1.990*** (0.062)
x	0.997*** (0.062)
Constant	4.978*** (0.055)
Observations	1,000
<i>Note:</i> *p<0.1; **p<0.05; ***p<0.01	

13. Now let's try to understand how things change when treatment is clustered:

- In one of your datasets, add an additional column which transforms ('bins') the y_0 variable into a categorical variable with 20 categories, where category '1' represents the lowest values of y_0 and category '20' contains the highest values of y_0 . This reflects the fact that people in the same cluster have more similar characteristics.

```
data <- data %>% mutate(cluster=as.factor(ntile(y0,20)))
```

- Add a new binary treatment variable, $D_cluster$, that randomly assigns each cluster to treatment or control (so treatment is at the level of the cluster).

```
cluster_treat <- tibble(cluster=factor(1:20),
  D_cluster=rbinom(20,1,0.5))
```

	Model 1
(Intercept)	5.37*** (0.04)
D_cluster	2.32*** (0.07)
R ²	0.51
Adj. R ²	0.51
Num. obs.	1000
RMSE	1.09

*** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$

Table 3: Statistical models

	Model 1
(Intercept)	5.37*** (0.33)
D_cluster	2.32*** (0.52)
R ²	0.51
Adj. R ²	0.51
Num. obs.	1000
RMSE	1.09

*** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$

Table 4: Statistical models

```
data <- data %>% left_join(cluster_treat, by="cluster")
```

c. Recalculate y_{obs} based on the new treatment variable, $D_cluster$.

```
data <- data %>% mutate(y_obs_cluster=case_when(D_cluster==1~y1,
                                                D_cluster==0~y0))
```

d. Now run the regression of observable outcomes (y_{obs}) on clustered treatment ($D_cluster$). What is the confidence interval around the size of the effect $D_cluster$? Does this makes sense?

```
data %>% lm(y_obs_cluster ~ D_cluster, data=.) %>%
  texreg(include.ci=F)
```

e. Now run the regression of observable outcomes (y_{obs}) on clustered treatment ($D_cluster$), and cluster standard errors at the cluster level. How does the confidence interval on $D_cluster$ differ from the previous analysis without clustered standard errors?

```
data %>% lm_robust(y_obs_cluster ~ D_cluster, data=., clusters=cluster) %>%
  texreg(include.ci=F)
```

14. Repeat the experiment and regression from Question 11 one hundred times with the same x , y_0 and y_1 every time, but re-randomize treatment assignment D each time. For each experiment, perform the basic regression and calculate the 95% confidence interval on the treatment variable. Finally, calculate how many of the 100 confidence intervals cover the real treatment effect ($c = 2$).

```
out <- list()
for (i in 1:100) {
  out[[i]] <- data %>% mutate(D=rbinom(N,1,0.5),
                             y_obs=case_when(D==1~y1,
                                                D==0~y0)) %>%
```

```

lm(y_obs ~ D + x, data=.) %>%
tidy() %>%
filter(term=="D") %>%
mutate(conf.lo=estimate-1.96*std.error,
        conf.hi=estimate+1.96*std.error) %>%
select(conf.lo, conf.hi)
}

out %>% bind_rows() %>%
mutate(covers_2=ifelse(conf.lo>2|conf.hi<2,0,1)) %>%
tally(covers_2)

## # A tibble: 1 x 1
##       n
##   <dbl>
## 1     95

```