

EDS241: Assignment 1

Amanda Herbst

01/24/2024

1 Part 1

(NOTE: Uses the RCT.R code provided with lecture to generate data) DO NOT CHANGE ANYTHING BELOW UNTIL IT SAYS EXPLICITLY

1.1 BELOW YOU CAN (AND HAVE TO) CHANGE AND ADD CODE TO DO ASSIGNMENT

Part 1: Use the small program above that generates synthetic potential outcomes without treatment, Y_{i0} , and with treatment, Y_{i1} . When reporting findings, report them using statistical terminology (i.e. more than y/n.) Please do the following and answer the respective questions (briefly).

- Create equally sized treatment and control groups by creating a binary random variable D_i where the units with the *1's" are chosen randomly.

```
# by replicating the vector (0,1) to the length N, then random sampling without replacement, it guarant
Di <- sample(rep(c(0,1),length.out=N),N, replace = FALSE)

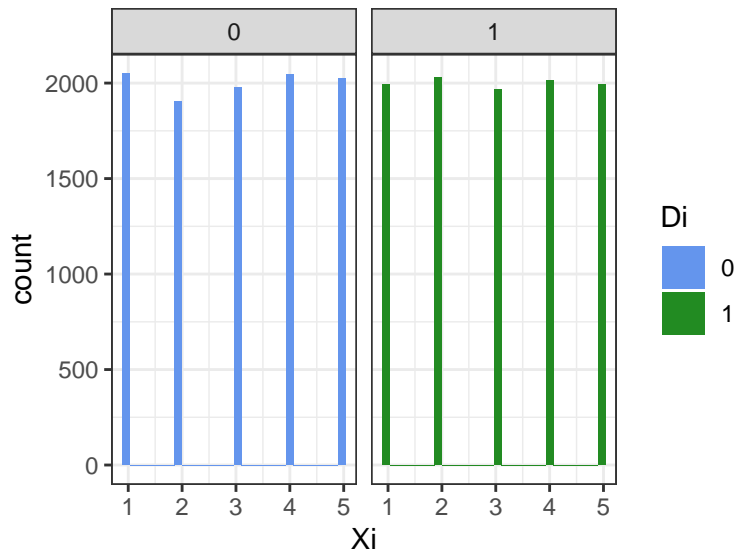
# add treatment variable to dataframe
df <- df %>%
  cbind(Di)

# # check for equal groups
# sum(dtest==1)
```

- Make two separate histograms of X_i for the treatment and control group. What do you see and does it comply with your expectations, explain why or why not?

The distributions of X_i for treatment and control groups are approximately the same. There is around 2000 observations of each characteristic for both treatment and control groups. This complies with my expectations because the treatment and control groups are of equal size and the characteristics were randomly sampled in a way that they would be approximately equal sizes

```
# create histogram of  $X_i$  for treatment and control groups
df %>%
  mutate(Di = as.factor(Di)) %>%
  ggplot(aes(x =  $X_i$ , fill = Di)) +
  scale_fill_manual(values = c("cornflowerblue", "forestgreen")) +
  geom_histogram() +
  # facet wrap by  $D_i$  (0 = control, 1 = treatment)
  facet_wrap(~Di) +
  theme_bw()
```



c) Test whether Di is uncorrelated with the pre-treatment characteristic Xi and report your finding.

****We fail to reject the null hypothesis**

```
cor_Xi <- cor.test(df$Xi, df$Di)
sjPlot::tab_model(cor_Xi,
                  title = "Di Correlation with Xi")
```

Di Correlation with Xi

Dependent variable

estimate

CI

p

-0.00

-0.02 – 0.01

0.615

d) Test whether Di is uncorrelated with the potential outcomes Yi_0 and Yi_1 and report your finding (only possible for this synthetic dataset where we know all potential outcomes).

We fail to reject the null hypothesis for the correlation test between both Di and Yi_0 and Di and Yi_1. Therefore, Di is uncorrelated with the potential outcomes Yi_0 and Yi_1

```
cor_Yi_0 <- cor.test(df$Yi_0, df$Di)
sjPlot::tab_model(cor_Yi_0,
                  title = "Di Correlation with Yi_0")
```

Di Correlation with Yi_0

Dependent variable

estimate

CI

p
-0.01
-0.02 – 0.01
0.270

```
cor_Yi_1 <- cor.test(df$Yi_1, df$Di)
sjPlot::tab_model(cor_Yi_1,
                  title = "Di Correlation with Yi_1")
```

Di Correlation with Yi_1

Dependent variable

estimate

CI

p

-0.01
-0.02 – 0.01
0.396

e) Estimate the ATE by comparing mean outcomes for treatment and control group. Test for mean difference between the groups and report your findings.

We reject the null hypothesis so the difference is true

```
# create new "observed" Yi variable based on the Di
# when Di = 1, Yi_obs = Yi_1
# when Di = 0, Yi_obs = Yi_0
df <- df %>%
  mutate(Yi_obs = Di*Yi_1+(1-Di)*Yi_0)

# mean difference
df %>%
  group_by(Di) %>%
  summarize(mean_Yi = mean(Yi_obs))
```

```
## # A tibble: 2 x 2
##       Di mean_Yi
##   <dbl>   <dbl>
## 1     0     1.52
## 2     1     3.00
```

```
# test for mean difference = t.test
t.test(df$Yi_obs ~ df$Di)
```

```
##
## Welch Two Sample t-test
##
## data: df$Yi_obs by df$Di
## t = -70.206, df = 18005, p-value < 0.00000000000000022
## alternative hypothesis: true difference in means between group 0 and group 1 is not equal to 0
## 95 percent confidence interval:
## -1.526022 -1.443126
```

```
## sample estimates:
## mean in group 0 mean in group 1
##      1.515060      2.999633
```

f) Estimate the ATE using a simple regression of (i) Y_i on D_i and (ii) Y_i on D_i and X_i and report your findings.

(i) The average treatment effect (ATE) is estimated to be 1.51 (ii) The (ATE) is estimated to be 1.50 holding the pre-treatment characteristics constant

```
# create new "observed" Yi variable based on the Di
```

```
# when Di = 1, Yi_obs = Yi_1
```

```
# when Di = 0, Yi_obs = Yi_0
```

```
df <- df %>%
  mutate(Yi_obs = Di*Yi_1+(1-Di)*Yi_0)
```

```
# (i) Yi on Di
```

```
summary(lm(Yi_obs ~ Di, data = df))
```

```
##
## Call:
## lm(formula = Yi_obs ~ Di, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -5.1883 -1.0355 -0.0053  1.0391  5.3289
##
## Coefficients:
##              Estimate Std. Error t value      Pr(>|t|)
## (Intercept)   1.51506    0.01495  101.33 <0.0000000000000002 ***
## Di            1.48457    0.02115   70.21 <0.0000000000000002 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.495 on 19998 degrees of freedom
## Multiple R-squared:  0.1977, Adjusted R-squared:  0.1977
## F-statistic: 4929 on 1 and 19998 DF, p-value: < 0.00000000000000022
```

```
# (ii) Yi on Di and Xi
```

```
summary(lm(Yi_obs ~ Di + Xi, data = df))
```

```
##
## Call:
## lm(formula = Yi_obs ~ Di + Xi, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.6434 -0.7226  0.0010  0.7250  4.1843
##
## Coefficients:
##              Estimate Std. Error t value      Pr(>|t|)
## (Intercept) -0.728853    0.019054  -38.25 <0.0000000000000002 ***
## Di           1.492105    0.014947   99.83 <0.0000000000000002 ***
## Xi           0.745684    0.005269  141.53 <0.0000000000000002 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##  
## Residual standard error: 1.057 on 19997 degrees of freedom  
## Multiple R-squared:  0.5992, Adjusted R-squared:  0.5992  
## F-statistic: 1.495e+04 on 2 and 19997 DF,  p-value: < 0.00000000000000022
```

2 Part 2

Part 2 is based on Gertler, Martinez, and Rubio-Codina (2012) (article provided on canvas) and covers impact evaluation of the Mexican conditional cash transfer Progresa (later called Oportunidades, now Prospera). Basically, families with low-incomes received cash benefits if they complied to certain conditions, such as regular school attendance for children and regular healthcare visits. You can read more about the program in the Boxes 2.1 (p.10) & 3.1 (p.40) of the Handbook on impact evaluation: quantitative methods and practices by Khandker, B. Koolwal, and Samad (2010). The program followed a randomized phase-in design. You have data on households (hh) from 1999, when treatment hh have been receiving benefits for a year and control hh have not yet received any benefits. You can find a description of the variables at the end of the assignment. Again, briefly report what you find or respond to the questions.

- a) Some variables in the dataset were collected in 1997 before treatment began. Use these variables to test whether there are systematic differences between the control and the treatment group before the cash transfer began (i.e. test for systematic differences on all 1997 variables). Describe your results. Does it matter whether there are systematic differences? Why or why not? Would it be a mistake to do the same test with these variables if they were collected after treatment began and if so why? Note: If your variable is a proportion (e.g. binary variables), you should use a proportions test, otherwise you can use a t-test.

no difference: hhsize (5.5/5.8), dirt floor (0.33/0.29), bathroom (0.42/0.44), home ownership (0.065/0.059) difference: electricity (0.28/0.38) Systematic differences matter because they affect our ability to create counterfactuals

```
# variables collected in 1997
progres97 <- progres97 %>%
  select(dirtfloor97, bathroom97, electricity97, homeown97, hhsize97, treatment)

## For continuous variables you can use the t-test
#t.test()
# household size
t.test(progres97$hhsize97 ~ progres97$treatment)

## For binary variables you should use the proportions test
#prop.test()
## Dirt Floor (1 if dirt floor)
progres97 %>%
  select(treatment, dirtfloor97) %>%
  table() %>%
  prop.test()

## Bathroom (1 if bathroom for exclusive hh use)
progres97 %>%
  select(treatment, bathroom97) %>%
  table() %>%
  prop.test()

## Electricity (1 if electricity in dwelling)
progres97 %>%
  select(treatment, electricity97) %>%
  table() %>%
  prop.test()

## Home ownership (1 if home in ownership)
progres97 %>%
  select(treatment, homeown97) %>%
```

```
table() %>%
prop.test()
```

- b) Estimate the impact of program participation on the household's value of animal holdings (vani) using a simple univariate regression. Interpret the intercept and the coefficient. Is this an estimate of a treatment effect?

Intercept: The groups that were not treated owned an average of \$1,691.47 in animals in 1997 price. **coefficient:** The groups that were treated owned an average of \$50.21 more animals in 1997 price than those not treated This is an estimate of ATE but doesn't control for pre-treatment characteristics

```
model1 <- lm(vani ~ treatment, data = progresas)
summary(model1)
```

```
##
## Call:
## lm(formula = vani ~ treatment, data = progresas)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1742  -1691  -1313   -137   50495
##
## Coefficients:
##              Estimate Std. Error t value      Pr(>|t|)
## (Intercept)  1691.47      48.43   34.929 <0.0000000000000002 ***
## treatment      50.21      64.28    0.781      0.435
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3702 on 13512 degrees of freedom
## (862 observations deleted due to missingness)
## Multiple R-squared:  4.516e-05, Adjusted R-squared:  -2.884e-05
## F-statistic: 0.6103 on 1 and 13512 DF, p-value: 0.4347
```

- c) Now, include at least 6 independent control variables in your regression. How does the impact of program participation change? Choose one of your other control variables and interpret the coefficient.

The estimated ATE of the program participation increases from approximately \$50 to approximately \$56 (1997 price). The average effect of having electricity in 1997 in the control group, holding years of spouse's education, if a primary school is in the community, if a health center is nearby, and if the household is using any piece of land constant, is about \$600 increase in animal value (1997 price).

```
model2 <- lm(vani ~ treatment + educ_sp + primary + age_hh + healthcenter + land + electricity97, data = progresas)
summary(model2)
```

```
##
## Call:
## lm(formula = vani ~ treatment + educ_sp + primary + age_hh +
##      healthcenter + land + electricity97, data = progresas)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4266  -1718   -792    129   48117
##
## Coefficients:
```

```
##               Estimate Std. Error t value          Pr(>|t|)
## (Intercept)  -345.348    231.148  -1.494          0.1352
## treatment      56.316     63.338   0.889          0.3739
## educ_sp       35.215     14.024   2.511          0.0120 *
## primary       442.072    180.426   2.450          0.0143 *
## age_hh        24.315      2.229  10.908 <0.0000000000000002 ***
## healthcenter -970.197     82.281 -11.791 <0.0000000000000002 ***
## land          1394.163    62.958  22.144 <0.0000000000000002 ***
## electricity97  600.312    66.306   9.054 <0.0000000000000002 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3589 on 13506 degrees of freedom
## (862 observations deleted due to missingness)
## Multiple R-squared:  0.06035,    Adjusted R-squared:  0.05987
## F-statistic: 123.9 on 7 and 13506 DF,  p-value: < 0.00000000000000022
```

- d) The dataset also contains a variable `intention_to_treat`. This variable identifies eligible households in participating villages. Most of these households ended up in the treatment group receiving the cash transfer, but some did not. Test if the program has an effect on the value of animal holdings of these non-participants (spillover effects). Think of a reason why there might or might not be spillover effects.

there doesn't appear to be any spillover effects. could be because families that do get the cash are not sharing in the community...

Hint: Create a pseudo-treatment variable that is = 1 for individuals who were intended to get treatment but did not receive it, = 0 for the normal control group and excludes the normal treatment group.