

EDS241: Assignment 2

Briana Barajas

02/07/2024

Reminders: Make sure to read through the setup in markdown. Remember to write out interpretations and report your results in writing (and table/plot etc) forms.

1 Part 1 Treatment Ignorability Assumption and Applying Matching Estimators (19 points):

The goal is to estimate the causal effect of maternal smoking during pregnancy on infant birth weight using the treatment ignorability assumptions. The data are taken from the National Natality Detail Files, and the extract “SMOKING_EDS241.csv” is a random sample of all births in Pennsylvania during 1989-1991. Each observation is a mother-infant pair. The key variables are:

The outcome and treatment variables are:

- `birthwgt` = birth weight of infant in grams
- `tobacco` = indicator for maternal smoking

The control variables are:

continuous:

- `mage` (mother’s age),
- `meduc` (mother’s education),

categorical:

- `mblack` (=1 if mother identifies as Black)
- `alcohol` (=1 if consumed alcohol during pregnancy),
- `first` (=1 if first child)
- `diabete` (=1 if mother diabetic)
- `anemia` (=1 if mother anemic)

Table 1: Subtracted mean birthweight for smoking mothers from mean birthweight for non-smoking mothers to keep the mean difference positive.

tobacco	mean_birthwgt	mean_diff_bw
0	3430.286	244.5394
1	3185.747	NA

```
# Load data for Part 1
birth_weight <- read_csv(paste0(data_wd, "birthweight_simple.csv")) %>%
  clean_names()
```

Question (a) Mean Differences, Assumptions, and Covariates (3 pts)

- a) What is the mean difference in birth weight of infants with smoking and non-smoking mothers? [1 pt]

```
# calculate the mean difference in birthweight
birth_weight %>%
  group_by(tobacco) %>%
  summarise(mean_birthwgt = mean(birthwgt)) %>%
  mutate(mean_diff_bw = mean_birthwgt - dplyr::lead(mean_birthwgt)) %>%
  kbl(caption = "Subtracted mean birthweight for smoking mothers from mean birthweight for
    non-smoking mothers to keep the mean difference positive.") %>%
  kable_minimal()
```

```
# calculate statistical significance of difference in mean birthweights
t.test(birthwgt~tobacco, data = birth_weight)
```

```
##
## Welch Two Sample t-test
##
## data: birthwgt by tobacco
## t = 58.932, df = 26945, 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:
## 236.4060 252.6727
## sample estimates:
## mean in group 0 mean in group 1
## 3430.286 3185.747
```

ANS: As illustrated in **Table 1**, the mean difference in birth weight of infants with smoking versus non-smoking mothers is 244.54 grams. The mean values for smoking (`tobacco = 1`) and non-smoking (`tobacco = 0`) mothers, as well as the difference between these means (`mean_bw_diff`) is summarized in Table 1. The difference in birth weight is statistically significant ($p < 0.05$), so we can reject the null hypothesis that the difference in means is equal to 0.

Under what assumption does this correspond to the average treatment effect of maternal smoking during pregnancy on infant birth weight? [0.5 pt]

ANS: This mean difference corresponds with the average treatment effect (ATE) of maternal smoking, assuming that the treatment effect is constant across the entire population, and that the control variables (mother's age, mother's education, etc.) have no influence on maternal smoking or birth weight. The latter point refers to the ignorability assumption.

Calculate and create a table demonstrating the differences in the mean proportions/values of covariates observed in smokers and non-smokers (remember to report whether differences are statistically significant) and discuss whether this provides empirical evidence for or against this assumption.

```
## =====
##      Prep subsets for testing      ----
## =====

# create df with tobacco and continuous variable
continuous <- birth_weight %>%
  select(tobacco, mage, meduc)

# create df with tobacco and categorical variables
binary <- birth_weight %>%
  select(tobacco, anemia, diabete, tobacco, alcohol, mblack, first)

# create list of variable names
binary_names <- binary %>% names()
continuous_names <- names(continuous)[2:3]

# create empty data frame to store results
prop_test_results <- data.frame()
t_test_results <- data.frame()

## =====
##      prop test for binary vars      ----
## =====

for (i in binary_names) {

  # split data into treated and untreated
  treated <- binary %>% filter(tobacco == 1) %>% pull(!!sym(i))
  untreated <- binary %>% filter(tobacco == 0) %>% pull(!!sym(i))

  # perform the prop test
  prop_test_result <- prop.test(x = c(sum(treated),
                                     sum(untreated)),
                               n = c(length(treated),
                                     length(untreated)), correct = FALSE)

  prop_test_result_tidy <- broom::tidy(prop_test_result)
  prop_test_result_tidy$Variable <- i
  prop_test_results <- rbind(prop_test_results, prop_test_result_tidy)
}

## =====
```

```

##          t-test for continuous vars      ----
## =====
for (i in continuous_names) {
  # Dynamically creating the formula for the t-test
  formula <- as.formula(paste(i, "~ tobacco"))

  # Performing the t-test
  t_test_result <- t.test(formula, data = continuous)

  # Storing the tidy results of the t-test in the data frame
  t_test_result_tidy <- broom::tidy(t_test_result)
  t_test_result_tidy$Variable <- i
  t_test_results <- rbind(t_test_results, t_test_result_tidy)
}

## =====
##          display results                ----
## =====

# combine results to single df
combine_results <- bind_rows(
  prop_test_results %>% select(Variable, estimate1, estimate2, p.value),
  t_test_results %>% select(Variable, estimate1, estimate2, p.value)
)

# create output table
combined_results_table <- kable(combine_results, format = "latex",
                                col.names = c("Variable",
                                                "Proportion or Mean Treated",
                                                "Proportion or Mean Control", "P-Value"),
                                caption = "Treated and Untreated Pre-treatment Proportion and T-Test Results",
                                kable_styling(font_size = 7, latex_options = "hold_position") %>%
                                kable_minimal())

# print table
combined_results_table

```

Table 2: Treated and Untreated Pre-treatment Proportion and T-Test Results

Variable	Proportion or Mean Treated	Proportion or Mean Control	P-Value
tobacco	1.0000000	0.0000000	0.0000000
anemia	0.0141031	0.0078005	0.0000000
diabete	0.0175187	0.0173636	0.8858005
alcohol	0.0441825	0.0071033	0.0000000
mblack	0.1354121	0.1086279	0.0000000
first	0.3645879	0.4360900	0.0000000
mage	27.4530853	25.5385632	0.0000000
meduc	13.2394207	11.9209454	0.0000000

ANS: [INPUT TABLE INTERPRETATION]

Remember that this is observational data. What other quantitative empirical evidence or test could help you assess the former assumption? [1.5 pt: 0.5 pt table, 1 pt discussion]

```
# calculate the statistical significance of the difference in mean

## Calculate mean difference. Remember to calculate a measure of statistical significance

## For continuous variables you can use the t-test
#t.test()

## For binary variables you should use the proportions test
#prop.test()

## Covariate Calculations and Tables (feel free to use code from Assignment 1 key)
```

Question (b) ATE and Covariate Balance (3 pts)

- b) Assume that maternal smoking is randomly assigned conditional on the observable covariates listed above. Estimate the effect of maternal smoking on birth weight using an OLS regression with NO linear controls for the covariates. [0.5 pts] Perform the same estimate including the control variables [0.5 pts]. Next, compute indices of covariate imbalance between the treated and non-treated regarding these covariates (see example file from class). Present your results in a table.[1 pts] What do you find and what does it say regarding whether the assumption you mentioned responding to a) is fulfilled? [1 pts]

```
## =====
##           run regressions          ----
## =====

# ATE Regression univariate
univariate <- lm(birthwgt~tobacco, data = birth_weight)

# ATE with covariates
multiple_reg <- lm(birthwgt~ anemia + diabete + tobacco +
                  alcohol + mblack + first + mage + meduc,
                  data = birth_weight)

## =====
##           plot results              ----
## =====

# tidy results
univariate <- broom::tidy(univariate)
multiple_reg <- broom::tidy(multiple_reg)

# print results
univariate %>% kbl(caption = "Univariate Linear Regression") %>%
  kable_minimal()

multiple_reg %>% kbl(caption = "Univariate Linear Regression") %>%
  kable_minimal()
```

Table 3: Univariate Linear Regression

term	estimate	std.error	statistic	p.value
(Intercept)	3430.2863	1.790782	1915.52379	0
tobacco	-244.5394	4.078907	-59.95219	0

Table 4: Univariate Linear Regression

term	estimate	std.error	statistic	p.value
(Intercept)	3362.2582445	11.9272609	281.8969309	0.0000000
anemia	-4.7963916	16.7544421	-0.2862758	0.7746675
diabete	73.2275309	12.1035155	6.0501043	0.0000000
tobacco	-228.0730765	4.1774982	-54.5956133	0.0000000
alcohol	-77.3497487	13.4653594	-5.7443508	0.0000000
mblack	-240.0303000	5.1062333	-47.0073116	0.0000000
first	-96.9441154	3.4466349	-28.1271784	0.0000000
mage	-0.6940244	0.3565637	-1.9464246	0.0516067
meduc	11.6883416	0.8604935	13.5832997	0.0000000

```
## =====
##               covariate balance      ----
## =====

# compute balance for multiple linear regression
covariate_balance <- xBalance(birthwgt~ anemia + diabete + tobacco +
                             alcohol + mblack + first + mage + meduc,
                             data = birth_weight,
                             report = c("std.diffs","chisquare.test", "p.values"))

# print results in neat balance table
covariate_balance %>% kbl() %>% kable_minimal()
```

Question (c) Propensity Score Estimation (3 pts)

- c) Next, estimate propensity scores (i.e. probability of being treated) for the sample, using the provided covariates. Create a regression table reporting the results of the regression and discuss what the covariate coefficients indicate and interpret one coefficient [1.5 pts]. Create histograms of the propensity scores comparing the distributions of propensity scores for smokers ('treated') and non-smokers

	std.diff.unstrat	p.unstrat		chisquare	df	p.value
anemia	-0.0000293	0.000006	unstrat	6752.456	8	0
diabete	0.0000378	0.000000				
tobacco	-0.0003811	0.000000				
alcohol	-0.0000940	0.000000				
mblack	-0.0003211	0.000000				
first	-0.0001533	0.000000				
mage	0.0001937	0.000000				
meduc	0.0002057	0.000000				

('control'), discuss the overlap and what it means [1.5 pts].

```
## Propensity Scores
ps <- glm(treat ~ age + educ + nodegree + re74 + re75,
          data = lalonde, family = binomial())
summary(ps)

##
## Call:
## glm(formula = treat ~ age + educ + nodegree + re74 + re75, family = binomial(),
##      data = lalonde)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -2.69393404  0.79887518  -3.372  0.000746 ***
## age          0.00246356  0.01024772   0.240  0.810019
## educ         0.15693194  0.05298613   2.962  0.003059 **
## nodegree     0.85020779  0.28125268   3.023  0.002503 **
## re74         -0.00012251  0.00002576  -4.756 0.00000198 ***
## re75          0.00002574  0.00003955   0.651  0.515252
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 751.49  on 613  degrees of freedom
## Residual deviance: 692.88  on 608  degrees of freedom
## AIC: 704.88
##
## Number of Fisher Scoring iterations: 5
```

```
# calculate propensity score for covariates
propensity_scores <- glm(birthwgt ~ anemia + diabete + tobacco + alcohol +
                          mblack + first + mage + meduc,
                          data = birth_weight)
summary(propensity_scores)
```

```
##
## Call:
## glm(formula = birthwgt ~ anemia + diabete + tobacco + alcohol +
##      mblack + first + mage + meduc, data = birth_weight)
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 3362.2582    11.9273 281.897 < 0.0000000000000002 ***
## anemia      -4.7964     16.7544  -0.286      0.7747
## diabete      73.2275     12.1035   6.050    0.00000000145 ***
## tobacco    -228.0731      4.1775 -54.596 < 0.0000000000000002 ***
## alcohol    -77.3497     13.4654  -5.744    0.00000000926 ***
## mblack     -240.0303      5.1062 -47.007 < 0.0000000000000002 ***
## first      -96.9441      3.4466 -28.127 < 0.0000000000000002 ***
## mage        -0.6940      0.3566  -1.946     0.0516 .
```

```
## meduc          11.6883      0.8605  13.583 < 0.0000000000000002 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 234966.3)
##
##      Null deviance: 23834376571  on 94172  degrees of freedom
## Residual deviance: 22125369801  on 94164  degrees of freedom
## AIC: 1431918
##
## Number of Fisher Scoring iterations: 2
```

```
## PS Histogram Unmatched
```

Question (d) Matching Balance (3 pts)

- (d) Next, match treated/control mothers using your estimated propensity scores and nearest neighbor matching. Compare the balancing of pretreatment characteristics (covariates) between treated and non-treated units in the original dataset (from c) with the matched dataset (think about comparing histograms/regressions) [2 pts]. Make sure to report and discuss the balance statistics [1 pts].

```
## Nearest-neighbor Matching
```

```
## Covariate Imbalance post matching:
```

```
## Histogram of PS after matching
```

Question (e) ATE with Nearest Neighbor (3 pts)

- (e) Estimate the ATT using the matched dataset. Report and interpret your result (Note: no standard error or significance test is required here)

```
## Nearest Neighbor
```

```
## ATT
```

Question (f) ATE with WLS Matching (3 pts)

- f) Last, use the original dataset and perform the weighted least squares estimation of the ATE using the propensity scores (including controls). Report and interpret your results, here include both size and precision of estimate in reporting and interpretation.

```
## Weighted least Squares (WLS) estimator Preparation
```

```
## Weighted least Squares (WLS) Estimates
```

```
## Present Results
```


Question (g) Differences in Estimates (1 pts)

- g) Explain why it was to be expected given your analysis above that there is a difference between your estimates in e) and f)?

2 Part 2 Panel model and fixed effects (6 points)

****We will use the progres data (progesa.csv) from last time as well as a new dataset, progres_pre.csv. In the original dataset, treatment households had been receiving the transfer for a year. Now, you get an additional dataset with information on the same households from before the program was implemented, establishing a baseline study (year 1997). **Note: You will need to install the packages plm and dplyr (included in template preamble). Again, you can find a description of the variables at the bottom of PDF and HERE.***

Question (a) Estimating Effect with First Difference (3 pts)

Load the new baseline data (pre-program) and the follow-up data (post-program, from Assignment 1) into R. Create a time denoting variable (with the same name) in BOTH datasets with a value of 0 for the pre-program dataset and 1 for the other one. Create a panel dataset by appending the data (i.e. binding the dataset row-wise together creating a single dataset). We want to examine the same outcome variable as before, value of animal holdings (vani)=. Estimate a standard difference-in-differences (DiD) regression and interpret the results.

```
rm(list=ls()) # clean environment

## Load the datasets
# progres_pre <- read.csv() insert your filepath etc
# progres_post <- read.csv()

## Append post to pre dataset
#progres <- rbind(progres_pre, progres_post)
```

- a) Estimate a first-difference (FD) regression manually, interpret the results briefly (size of coefficient and precision!) ****Note: Calculate the difference between pre- and post-program for each individual and for each variable used (i.e the outcome and the independent variables).[3 pts]** To do that, follow these steps and the code given in the R-template:*

```
### Code included to help get you started
## i. Sort the panel data in the order in which you want to take differences, i.e. by household and time

## Create first differences of variables
# progres <- progres %>%
#   arrange(hhid, year) %>%
#   group_by(hhid)

## ii. Calculate the first difference using the lag function from the dplyr package.
#   mutate(vani_fd = vani - dplyr::lag(vani))

## iii. Estimate manual first-difference regression (Estimate the regression using the newly created variable)
# fd_manual <- lm(vani_fd ~ ...)
```

Question (b) Fixed Effects Estimates (2 pts)

- b) Now also run a fixed effects (FE or ‘within’) regression and compare the results. Interpret the estimated treatment effects briefly (size of coefficient and precision!)

Fixed Effects Regression

Present Regression Results

Question (c) First Difference and Fixed Effects and Omitted Variable Problems (1 pts)

- c) Explain briefly how the FD and FE estimator solves a specific omitted variable problem? Look at the example on beer tax and traffic fatalities from class to start thinking about omitted variables. Give an example of a potential omitted variable for the example we are working with here that might confound our results? For that omitted variable, is a FE or FD estimator better? One example is enough.