

# Problem Set 2

Your Name - Net ID - Section Number

Due Nov 10, 2023

This homework must be turned in on Brightspace by Nov. 10, 2023. It must be your own work, and your own work only – you must not copy anyone’s work, or allow anyone to copy yours. This extends to writing code. You may consult with others, but when you write up, you must do so alone.

Your homework submission must be written and submitted using Rmarkdown. No handwritten solutions will be accepted. **No zip files will be accepted. Make sure we can read each line of code in the pdf document.** You should submit the following:

1. A compiled PDF file named yourNetID\_solutions.pdf containing your solutions to the problems.
2. A .Rmd file containing the code and text used to produce your compiled pdf named your-NetID\_solutions.Rmd.

Note that math can be typeset in Rmarkdown in the same way as Latex. Please make sure your answers are clearly structured in the Rmarkdown file:

1. Label each question part
2. Do not include written answers as code comments.
3. The code used to obtain the answer for each question part should accompany the written answer. Comment your code!

## Question 1 (Total: 50)

In new democracies and post-conflict settings, Truth and Reconciliation Commissions (TRCs) are often tasked with investigating and reporting about wrongdoing in previous governments. Depending on the context, institutions such as TRCs are expected to reduce hostilities (e.g. racial hostilities) and promote peace.

In 1995, South Africa's new government formed a national TRC in the aftermath of apartheid. [Gibson 2004] uses survey data collected from 2000-2001 to examine whether this TRC promoted inter-racial reconciliation. The outcome of interest is respondent racial attitudes (as measured by the level of agreement with the prompt: "I find it difficult to understand the customs and ways of [the opposite racial group]"). The treatment is "exposure to the TRC" as measured by the individual's level of self-reported knowledge about the TRC.

You will need to use the `trc_data.dta` file for this question. The relevant variables are:

- RUSTAND - Outcome: respondent's racial attitudes (higher values indicate greater agreement)
- TRCKNOW - Treatment dummy (1 = if knows about the TRC, 0 = otherwise)
- age - Respondent age (in 2001)
- female - Respondent gender
- wealth - Measure of wealth constructed based on asset ownership (assets are fridge, floor polisher, vacuum cleaner, microwave oven, hi-fi, washing machine, telephone, TV, car)
- religiosity - Self-reported religiosity (7 point scale)
- ethsalience - Self-reported ethnic identification (4 point scale)
- rcbblack - Respondent is black
- rcwhite - Respondent is white
- rccol - Respondent is coloured (distinct multiracial ethnic group)
- EDUC - Level of education (9 point scale)

### Part a (15 points)

Estimate the average treatment effect of TRC exposure on respondents' racial attitudes under the assumption that TRC exposure is ignorable. Report a 95% confidence interval for your estimate and interpret your results. (Use robust standard errors throughout.)

```
library(tidyverse)
library(haven)
library(estimatr) # for lm with robust se : ?lm_robust()

# Load in the TRC data (it's a STATA .dta so we use the haven package)
TRC_data <- haven::read_dta("trc_data.dta")

linreg <- lm_robust(RUSTAND ~ TRCKNOW, data = TRC_data)

tidy(linreg)
```

```
##           term    estimate  std.error statistic      p.value  conf.low  conf.high
## 1 (Intercept)  2.5311438  0.02805761  90.212369  0.000000e+00  2.476131  2.5861565
## 2      TRCKNOW -0.2177317  0.04433111  -4.911488  9.491614e-07 -0.304652 -0.1308115
##      df outcome
## 1 3203 RUSTAND
## 2 3203 RUSTAND
```

The estimate for the average effect of knowing about TRC is -0.2177 on respondents' racial attitudes. The 95% confidence interval is equal to [-0.3047, -0.1308]. 0 is noticeably not within this interval, therefore we reject the null hypothesis of no effect since the effect is statistically significant at 95% confidence level.

## Part b (15 points)

Examine whether exposed and nonexposed respondents differ on the full set of observed covariates using a series of balance tests. Briefly discuss, in which ways do exposed and nonexposed respondents differ?

```
TRC_data <- TRC_data %>%
  mutate(age_std = age / sd(age),
         female_std = female / sd(female),
         wealth_std = wealth / sd(wealth),
         religiosity_std = religiosity / sd(religiosity),
         ethsalience_std = ethsalience / sd(ethsalience),
         rcblack_std = rcblack / sd(rcblack),
         rcwhite_std = rcwhite / sd(rcwhite),
         rccol_std = rccol / sd(rccol),
         EDUC_std = EDUC / sd(EDUC))

means <- TRC_data %>%
  group_by(TRCKNOW) %>%
  summarize(age = mean(age_std),
         female = mean(female_std),
         wealth = mean(wealth_std),
         religiosity = mean(religiosity_std),
         ethsalience = mean(ethsalience_std),
         rcblack = mean(rcblack_std),
         rcwhite = mean(rcwhite_std),
         rccol = mean(rccol_std),
         EDUC = mean(EDUC_std))

sds <- TRC_data %>%
  group_by(TRCKNOW) %>%
  summarize(age = sd(age_std),
         female = sd(female_std),
         wealth = sd(wealth_std),
         religiosity = sd(religiosity_std),
         ethsalience = sd(ethsalience_std),
         rcblack = sd(rcblack_std),
         rcwhite = sd(rcwhite_std),
         rccol = sd(rccol_std),
         EDUC = sd(EDUC_std))

pval = function(x) {
  return(lm_robust(x ~ TRCKNOW, data = TRC_data)$p.value[2])
}

pvals <- TRC_data %>%
  summarize(age = pval(age),
         female = pval(female),
         wealth = pval(wealth),
         religiosity = pval(religiosity),
         ethsalience = pval(ethsalience),
         rcblack = pval(rcblack),
         rcwhite = pval(rcwhite),
         rccol = pval(rccol),
         EDUC = pval(EDUC))
```

```
# means
meansoutput <- as.data.frame(means[-c(1)])
rownames(meansoutput) <- c('Control Means', 'Treatment Means')

# sds
sdsoutput <- as.data.frame(sds[-c(1)])
rownames(sdsoutput) <- c('Control SDs', 'Treatment SDs')

# pvals
pvalsoutput <- c('p value')
```

```
## Warning: Setting row names on a tibble is deprecated.
```

```
balancetest <- rbind(meansoutput, sdsoutput, pvals) %>%
  mutate(across(everything(), round, 6))
```

```
## Warning: There was 1 warning in `mutate()`.
## i In argument: `across(everything(), round, 6)`.
```

## Caused by warning:

```
## ! The `...` argument of `across()` is deprecated as of dplyr 1.1.0.
## Supply arguments directly to `.fns` through an anonymous function instead.
##
## # Previously
##   across(a:b, mean, na.rm = TRUE)
##
## # Now
##   across(a:b, \(x) mean(x, na.rm = TRUE))
```

```
balancetest
```

```
##           age  female  wealth religiosity ethsalience  rcblack
## Control Means  2.618827 0.865799 0.773647      2.147880    4.691181 1.027792
## Treatment Means 2.520788 1.076452 0.927554      2.106486    4.727592 1.105420
## Control SDs     1.031005 0.991807 0.980711      0.987703    1.020535 1.001644
## Treatment SDs   0.958177 0.998127 1.017117      1.014767    0.974185 0.996662
## p value         0.005398 0.000000 0.000015      0.245140    0.303006 0.028719
##           rcwhite  rccol  EDUC
## Control Means  0.576580 0.458660 3.317230
## Treatment Means 0.614206 0.321938 3.701274
## Control SDs     0.990721 1.061433 0.940997
## Treatment SDs   1.011231 0.913756 1.029533
## p value         0.290459 0.000091 0.000000
```

### Part c (10 points)

Now assume that TRC exposure is conditionally ignorable given the set of observed covariates:

1. Use a logistic regression model to estimate the propensity score for each observation. (For purposes of this question, do not include any interactions.)
2. With this model, construct inverse propensity of treatment weights (IPTW) for each observation using the unstabilized weights.
3. Use the propensity score to construct an IPW estimator and report the point estimate for the ATE.

Use the following covariates: age, female, wealth, religiosity, ethsalience, rcblack, rcwhite, rccol, EDUC, far

```
logreg <- glm(TRCKNOW ~ age + female + religiosity + ethsalience + rcblack + rcwhite + rccol + EDUC, far
```

```
TRC_data$scores <- predict(logreg,)
```

#### Part d (10 points)

Using the bootstrap method (resampling individual rows of the data with replacement), obtain an estimate for the standard error of your IPTW estimator for the ATE. Compute a 95% confidence interval and interpret your findings. (You should report estimate, standard error, 95% CI lower, 95% CI upper, for interpretation, compare your results in Part C/D to your estimate from Part A and briefly discuss your findings.)

```
# Set random seed
set.seed(123)
```

### Question 2 (Total: 50 points)

Use the same data set as in Question 1.

#### Part a (15 points)

Estimate the ATT of TRC exposure on respondents' racial attitudes using the MatchIt approach. You can use the matchit function from MatchIt package in R. Implement the nearest neighbor matching algorithm and estimate the ATT. Report the 95% confidence interval of your estimate.

```
library(MatchIt)
```

```
## Warning: package 'MatchIt' was built under R version 4.3.2
```

```
# Read the help file first! Check out the default settings
# ?matchit()
```

#### Part b (15 points)

Estimate the ATT of TRC exposure on respondents' racial attitudes using the MatchIt approach. You can use the matchit function from MatchIt package in R. Implement the exact matching algorithm and estimate the ATT. Report the 95% confidence interval of your estimate.

#### Part c (10 points)

Estimate the ATT of TRC exposure on respondents' racial attitudes using the MatchIt approach. You can use the matchit function from MatchIt package in R. Implement the **coarsened exact matching** algorithm and estimate the ATT. Report the 95% confidence interval of your estimate.

#### part d (10 points)

Compare and contrast the three different matching algorithms. Provide evidence and an argument about which one we should use.

### BONUS ONLY: Question 3 (Total: Up to +12)

Question 3 is for bonus points. (See forthcoming lecture on Nov. 7th)

#### part a (+4 points)

Using the regression method to predict potential outcomes for all individuals in the dataset and calculate the ATE with bootstrapped standard errors. Report and interpret your results. (Hint: Start by fitting the treatment and control model with subsets of the data.)

```

## Fit a model among TRCKNOW == 1 to get  $E[Y_i(1) | X]$ 

## Fit a model among TRCKNOW == 0 to get  $E[Y_i(0) | X]$ 

## Predict the potential outcome under treatment for all units

## Predict the potential outcome under control for all units

## Average of the differences

### Bootstrap for SEs
set.seed(123)

```

**part b (+4 points)**

Using the regression method to predict potential outcomes for all individuals and calculate the ATT with bootstrapped standard errors. Report and interpret your results.

```

## Fit a model among TRCKNOW == 1 to get  $E[Y_i(1) | X]$ 

## Fit a model among TRCKNOW == 0 to get  $E[Y_i(0) | X]$ 

## Predict the potential outcome under treatment for all units

## Predict the potential outcome under control for all units

## Average of the differences

### Bootstrap for SEs
set.seed(123)

```

**part c (+4 points)**

Compare and contrast the ATE and ATT from the regression approach.