

EDS241: FINAL

Briana Barajas

02/23/2024

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

1 Part 1: RCTs, treatment ignorability (selection on observables), propensity scores (*15 points total*)

Setup

This exercise is inspired by Costello et al. 2008 article in science “Can Catch Shares Prevent Fisheries Collapse”, which we also discussed in class (lecture 5). “Inspired” means that the data `final_fisheries_data.csv` are synthetically generated to simplify things for our purposes. It contains the variables on 11,135 fisheries (only cross sectional, no time observations): These fisheries were either regulated by an Individual Transferable Quota (ITQ) for all years between 1990 and 2012 or in none of those years. Variables in the dataset include:

The outcome and treatment variables are:

`COLL_SHARE` = share of years a fishery is collapsed between 1990 and 2012 (collapse defined as harvest being more than 10% below maximum recorded harvest).

`ITQ` = dummy variable indicating ‘treatment’ with an ITQ (equal to 1 if the fishery has been regulated by an ITQ and 0 otherwise).

The control variables are:

`MET1`, `MET2`, ... `MET6` = Dummy variables indicating to which Marine Ecosystem Type (MET) the fishery belongs to (coral reefs, kelp forests, seagrass meadows, open ocean, deep sea, mangrove forests). This type does not change over the relevant time period and does not depend on human influence.

`IND_SR` = Index of species richness in 1980 with values between 0 and 100 indicating the biodiversity with respect to species in the fishery. Bounds of 0 and 100 are the lowest and highest observed values of species diversity across all fisheries in 1980, respectively.

`COMM_VAL` = Commercial value of fisheries in 1980 in million US-\$

The basic question of interest is “**What is the average treatment effect of implementing an ITQ in the time period from 1990 to 2012 on the share of years with a collapse.**” It is likely that the probability a fishery is selected for an ITQ depends on the pre-treatment characteristics given. It is also quite likely that the pre-treatment characteristics have an effect on the share of collapse for each fishery, i.e. our outcome variable of interest.

```
## Load Data
fish <- read_csv(here("final", "data", "final_fisheries_data.csv")) %>%
  clean_names()
```

```
# pivot data
fish_pivot <- fish %>%
  pivot_longer(cols = 1:6, names_to = "ecosystem_type", values_to = "value") %>%
  mutate(ecosystem_type = as.factor(ecosystem_type)) %>%
  filter(value == 1)
```

Question (a) Pretreatment Ecosystem Characteristic Comparison, Visual (3 pts)

- (a) Compare the distributions of pre-treatment ecosystem characteristics (i.e. MET1, MET2, „, MET6) between the treated and the control groups by drawing back to back histograms [2 pts]. Write one sentence discussing the (dis)similarity between the two groups [1pt].

```
## Histograms comparing covariates
```

```
## Remember to include histograms in final product
```

```
# ATTEMPT 1 - Calculate prop scores and use histback like in assignment 2
```

```
# calculate propensity scores
```

```
prop_scores_fish <- glm(itq ~ met1 + met2 + met3 + met4 + met5 + met6, data = fish_raw, family = binomial)
```

```
# add ps-value column
```

```
fish_raw <- fish_raw %>%
```

```
  mutate(psvalue = predict(prop_scores_fish, type = "response"))
```

```
# plot histogram to compare covariates
```

```
histbackback(split(fish_raw$psvalue, fish_raw$itq))
```

```
# ATTEMPT 5 - Lengthen df
```

```
fish_long_prop <- glm(itq ~ ecosystem_type, data = fish, family = binomial())
```

```
fish$psvalue = predict(fish_long_prop, type = "response")
```

```
histbackback(split(fish$psvalue, fish$itq))
```

```
# ATTEMPT 4 - Pivot longer and create subsets to create mirrored plot of met count
```

```
fish_treat <- fish_pivot %>% filter(itq == 1)
```

```
fish_control <- fish_pivot %>% filter(itq == 0)
```

```
ggplot() +
```

```
  geom_col(data = fish_treat, aes(x = ecosystem_type, y = value), fill = "#2a9d8f") +
```

```
  geom_col(data = fish_control, aes(x = ecosystem_type, y = -value), fill = "#264653") +
```

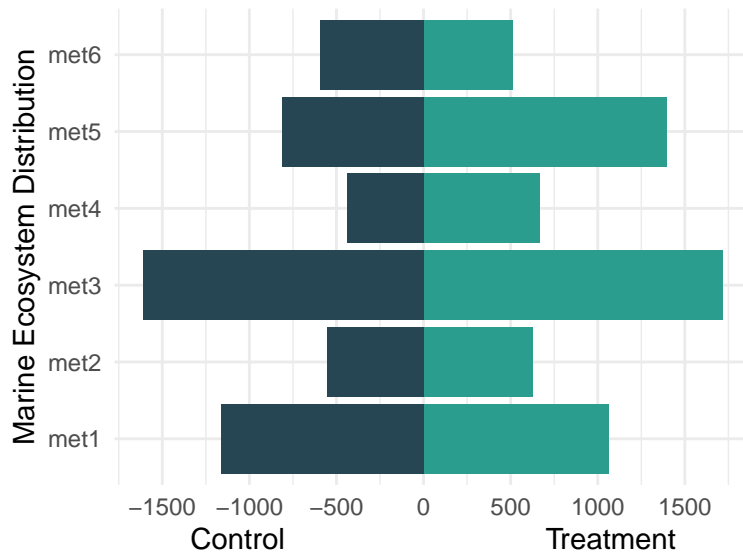
```
  scale_y_continuous(n.breaks = 10) +
```

```
  coord_flip() +
```

```
  theme_minimal() +
```

```
  labs(y = "Control", y2 = "Treatment",
```

```
       x = "Marine Ecosystem Distribution")
```



ANS: I decided to pivot the data frame so marine ecosystem type would be a single column, making it easier to compare the treatment and control groups. The distribution for ecosystem types in the control and treatment groups (where treatment is `itq`) are not equal. For example, in `met5` and `met4` there appear to be more observations for the treatment group.

Question (b) Pretreatment Ecosystem Characteristic Comparison, Mean differences *3 pts*)

- (b) Do a test on mean differences between the treated and control groups for the species richness index (`IND_SR`) and commercial value (`COMM_VAL`) variables. Interpret the results (estimated difference and significance) [2 pts] and make a conclusion regarding the similarity between the groups [1pt].

```
## Mean Differences (remember to use prop.test or t.test when applicable)

# calculate differences in mean for sp richness
mean_diff_ind <- broom::tidy(
  t.test(ind_sr ~ itq, data = fish)) %>%
  select(estimate1, estimate2, p.value) %>%
  mutate(Variable = "ind_sr")

# calculate differences in mean for commercial value
mean_diff_comm <- broom::tidy(
  t.test(comm_val ~ itq, data = fish)) %>%
  select(estimate1, estimate2, p.value) %>%
  mutate(Variable = "comm_val")

# bind result tables
t_test_results <- rbind(mean_diff_ind, mean_diff_comm) %>%
  relocate(Variable)

# print results
t_test_results %>%
  kbl(col.names = c("Variable", "Mean Treated", "Mean Control", "P-Value"),
```

Table 1: Mean differences of Treatment vs. Control (1b)

Variable	Mean Treated	Mean Control	P-Value
ind_sr	57.38515	48.55968	0
comm_val	117.22839	84.87908	0

```
caption = "Mean differences of Treatment vs. Control (1b)" %>%
kable_paper()
```

ANS: A t-test was used to compare the difference in the means for the continuous variables, and both tests resulted in a p-value that was less than 0.05. Given a significance level (α) of 0.05, we can reject the null hypothesis that the difference in means between the control group and treatment group equals zero. This is true for both `ind_sr` and `comm_val`.

Question (c) Treatment Ignorability (1 pt)

- (c) Based on your results from (a) and (b), do you see a problem with just comparing the outcome variable means between treated and untreated fisheries?

ANS: Based on the histogram from *Question 1a*, there appears to be differences in the marine ecosystem type (`met`) of the control group versus the treatment. This would cause bias when statistically estimating whether or not regulation by an Individual Transferable Quota (ITQ) altered the share of years a fishery collapsed (`coll_share`). Similarly the t-test results above demonstrate that the difference in the mean species richness index (`ind_sr`) and commercial value (`comm_val`) are statistically significant. This adds to the potential for bias, since pre-treatment characteristics vary for the treated and untreated.

Question (d) Propensity Scores (2 pts)

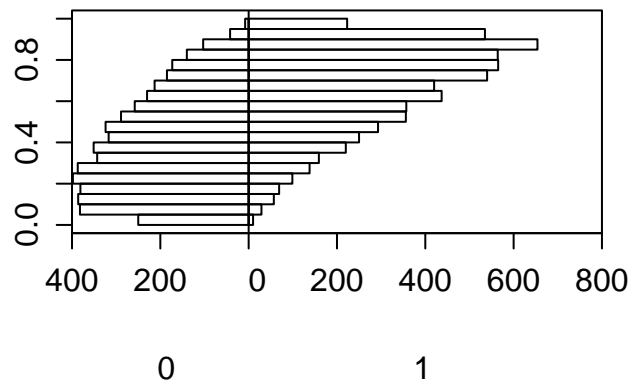
- (d) Estimate the propensity scores (probability of being treated) using a logit model, assume that all covariates are relevant and should be included in the estimation [0.5 pt]. Draw separate histograms (back to back) of the propensity scores for the treated and the untreated group [0.5 pt]. Comment on the overlap, do you have any concerns? Why/why not? [1]

```
## Propensity Score Estimates
```

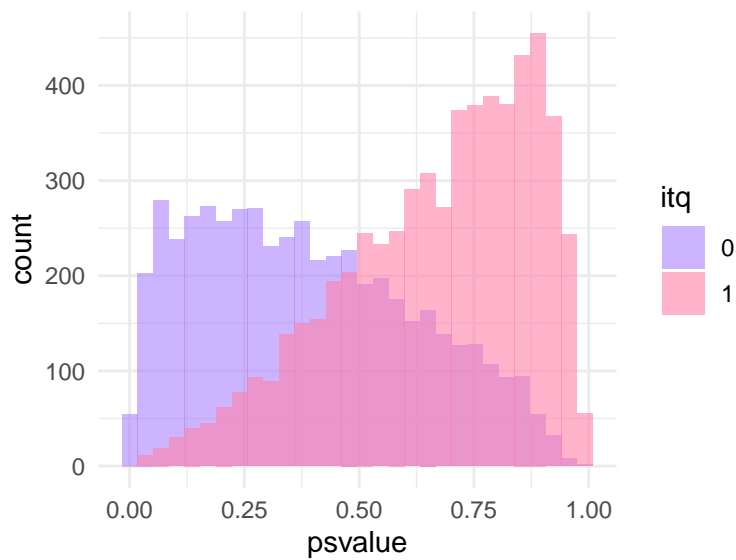
```
propensity_scores <- glm(itq ~ met1 + met2 + met3 + met4 + met5 + met6 +
                        ind_sr + comm_val, data = fish, family = binomial())

fish$psvalue <- predict(propensity_scores, type = "response")

histbackback(split(fish$psvalue, fish$itq))
```



```
fish %>%
  mutate(itq = as.factor(itq)) %>%
  ggplot(aes(x = psvalue, fill = itq)) +
  geom_histogram(alpha = 0.6, position = "identity") +
  scale_fill_manual(values = c("mediumpurple1", "palevioletred1")) +
  theme_minimal()
```



ANS: There is a decent amount of overlap for the center values, so matching should work. I am a bit concerned about the observations on the more extreme ends. For example, there are many observations in the treated group with a value around 1, but there are not many counterfactuals.

Question (e) ATT with Nearest Neighbor Matching (3 pts: 2 pt estimate, 1 pt interpretation)

- (e) Use the propensity scores from (d) to estimate the Average Treatment Effect on the Treated (ATT) with a nearest neighbor matching estimator. Interpret the result (just the size of the estimate)

```
## Nearest Neighbor Matching
```

```
## Estimate ATT
```

```
# find matches using nearest-neighbors method
match_itq <- matchit(itq ~ met1 + met2 + met3 + met4 + met5 + met6 +
                    ind_sr + comm_val, data = fish,
                    method = "nearest", ratio = 1)

# store results
match_itq_df <- match.data(match_itq)

# calculate the difference by subgroup
subgroup_diff <- match_itq_df %>%
  group_by(subclass) %>%
  mutate(diff = mean(coll_share[itq == 1]) - mean(coll_share[itq == 0]),
         .groups = "drop")

# calculate the mean for all differences of aTT
ATT <- mean(subgroup_diff$diff)

paste("Estimate of ATT using matched dataset:", ATT)
```

```
## [1] "Estimate of ATT using matched dataset: -0.0713262272262593"
```

ANS: The average treatment effect, -0.0713, indicates the difference in the proportion of years that a fishery was collapsed (between 1990-2012). On average, fisheries in the treated group were collapsed for fewer years than the control group.

Question (f) ATE with WLS (3 pts: 1 pt estimate, 1 pt interpretation)

- (f) Estimate the Average Treatment Effect (ATE) using the weighted least squares on the full sample. Interpret the estimated size and conclude if it is significantly different from zero from a statistical perspective.

```
## WLS Matching
```

```
## Estimate ATE
```

2 Part 2 Difference in Difference Estimation (10 points total + 3pts extra credit)

Here we return for a final time to the dataset from Gertler, Martinez, and Rubio-Codina (2012) and use a different way of estimating the effect of the Mexican conditional cash transfer on the value of animal

holdings of recipients. We'll use the panel data from assignment 2, where you have both the pre-program and post-program observations. See Template for dataset preparation instructions.

****Data Preparation****

*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](#).

Prepare Data: Load the new data (`progres_a_pre_1997.csv`) and the follow-up data (`progres_a_post_1999.csv`) into R. Note that we created a time denoting variable (with the same name, 'year') in BOTH datasets. Again, you will 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 family animal holdings (`vani`). You will use the full dataset for each estimate. NOTE: you should not change any NAs from the TREATED column in your analysis, as we expect that spillover was likely in this program. NAs will be excluded from your calculations/estimations.

Question (a) DiD Estimator, ATE (5 pts: 3 pts estimate, 2 pts interpretation)

- (a) Calculate the DiD estimator of the treatment effect (ATE) of the program on the value of animal holdings (`vani`) "manually" i.e. based on group mean values without running a regression. Report and interpret the result (Note: no significance test or standard errors is possible, so you do not need to report these values).

```
## Estimate ATE with DiD estimator manually.  
# You will need to calculate various means to get this estimate  
  
## Compute the Difference-in-Differences
```

Question (b) Difference in Difference using OLS (5 pts)

- (b) Now set up an OLS-regression using group mean values to estimate the same ATE. Interpret the estimated treatment effect [3 pts]. Also interpret the coefficients on the time dummy and the group dummy variable (see interpretation done in class in lecture 9) [2 pts].

****Hints:**** You will need to create a new dataframe with a variety of dummy variables to do this. The R example provided with the DiD module (and/or the excel file) should help.

```
## Create a new data frame for OLS regression  
  
## Run the OLS regression w/dummies  
  
## Report OLS Model results Print the summary of the OLS model
```

3 Extra Credit: ATE with OLS using full dataset (3 pts: 2 pts estimate, 1 pt interpretation)

- (c) Estimate the ATE with an OLS-regression based on the original units as observations (i.e. not with group mean values, you will need to use the entire dataset). Even though the specification is the same as in the regression with the group mean values above, you'll need to create new indicator variables for the treatment group and the post treatment time period as well as their interaction term. Verify that

you get the same result as above. Now report also on the precision of the estimation and test whether the estimated coefficient is different from zero.

```
## Create the dummy variables (you'll need 3)
```

```
## OLS regression
```

```
# Present Regressions in Table
```