

POLS3004 Causal Analysis: Tutorial Exercise 5

We'll use data from LaLonde's evaluation of economic training programs in the United States. The study used observational data from a treatment group of people who took part in the National Supported Work Demonstration (NSW), a job training program in the mid 1970s. The comparison group is a sample of people who did not go through training, taken from the Current Population Survey (CPS). Treatment assignment was not random. What can we learn from the data?

We get our data from the MatchIt package:

```
install.packages('MatchIt')
library(MatchIt)
ldata <- MatchIt::lalonde
```

Also, install and load the `Matching` package.

The dataset contains:

- N - 614 observations (185 treated, 429 control)
- *treat* - 1 = treated, 0 = control
- *age* - age in years
- *educ* - number of years of education
- *black* - 1 = African American, 0 otherwise
- *hispan* - 1 = Hispanic origin, 0 otherwise
- *married* - 1 = married, 0 otherwise
- *nodegree* - 1 = no high school degree, 0 otherwise
- *re74* = income in 1974 (dollars)
- *re75* = income in 1975, in U.S. dollars (dollars)
- *re78* = income in 1978, in U.S. dollars (dollars) [the main outcome variable]

- a) Regress income in 1978 on the treatment, controlling for *hispan*, *educ* and *married*. What is the estimated Average Treatment Effect of training? Is it statistically significant?
- b) Now, carry out exact matching to find the effect of treatment on income in 1978, matching on *educ*, *hispan* and *married*. Report the ATT and its standard error. Does the result differ from part (a)?

Code Hints:

- Use the `Match()` function, and look up the help file to discover how to do exact matching
- Supply `Match()` with the arguments *Y* (the outcome variable), *TR* (the treatment variable) and *X* (a matrix of the variables used for matching, created using `cbind()`)
- `summary.Match()` will summarise the results

- c) Examine balance between treated and untreated units before and after matching, using t-tests for equality of means. How does balance differ before and after matching?

Code Hints:

- Use the `MatchBalance()` function. It requires two main arguments.
- The first is a formula in the format: `treatment ~ matchingvar1 + matchingvar2 + ...`
- The second is the name of your matching estimate in part (b)

- d) Now repeat parts (b) and (c) with the addition of *re74* and *re75* to the list of matching variables, using the Mahalanobis distance metric instead of exact matching. How good is balance after matching? Does the estimated ATT change?

Code Hint:

- Look at the help file to find out how to select the Mahalanobis distance metric

- e) Estimate propensity scores for all units using a logistic regression of treatment on all variables except *re78*, and add them to the dataset as a new column

Code Hints:

- Remember that a propensity score is the estimated probability of treatment. This can be calculated using the `fitted.value()` function

- f) Repeat part (b) and (c) again, this time using only your estimated propensity scores to match. How good is balance after matching? Does the estimated ATT change?

- g) An alternative to the matching estimator is to carry out matching, and then run a regression on the matched dataset. Do this using the following steps:

- i) Estimate the ATT from matching as in part (f), this time with the option `ties=FALSE`. Your estimate will be similar but not identical to (f)¹

¹This makes things less complicated for the sake of this example: without `ties=FALSE`, the matching estimator includes all ties and weights tied observations. To get a similar regression estimate, we would also have to weight our regression. With `ties=FALSE`, R instead breaks the ties at random, selecting only one of the tied observations, so our regression does not need to be weighted. This also means that everyone will get slightly different answers to this question, depending on how R randomly breaks the ties. In general it is more principled to use `ties=TRUE` since

ii) Create a dataset of only the observations used in matching in (i).

Code Hint: You can extract the rows used from the dataset by adding `$index.treated` and `$indexcontrol` to the name if your matching estimate

iii) Run a regression using only your dataset from (ii).

How similar is your regression in (iii) to the matching estimate in (i)?

we do not arbitrarily throw away data in that scenario. Note that it is also possible to reduce the number of ties using the `distance.tolerance` option