

R Lab 0 - R Programming Refresher

Advanced Topics in Causal Inference

Goals for this lab:

1. Review of R programming and good programming practice.
2. Create a basic function to generate data with $O = (W, A, Y)$ structure.
3. Implement g-computation, IPTW, and TMLE estimators inspired by the average treatment effect (ATE) of a point treatment intervention.

Next lab:

1. Explore different data structures, inspired by real-world experiments.
2. Simulate data generating process that gives rise to the data we observe.
3. Think of studies in which these data generating systems may occur.

1 Introduction and Motivation

Welcome to Advanced Topics in Causal Inference R Labs! This semester we'll be working in R to implement concepts we've learned in lectures and discussion sections. In this lab we'll provide a refresher of R programming.

Section 2 will give you a refresher on R basics, including objects, subsetting, making plots, etc. Section 3 will walk you through how to generate data, and Section 4 will implement the 3 estimators we learned in Causal I to estimate the average treatment effect. You won't be turning in this lab, so feel free to complete as many or few sections as you feel will help you get a solid foundation before we start labs (to turn in) next week.

In future labs, please turn in relevant R code, text, and figures. Consider turning in your homework as a compiled R Markdown (RMD) file, which “knits” your code, text, and figures into a single PDF, Word, or html document. A .Rmd template with further instructions can be found on bCourses.



Figure 1: Welcome to Causal II labs!

2 R review and refresher

2.1 Help files

Help files in R are extremely useful. Coding in R isn't so much about memorizing functions; it's more useful to know how to look up commands. Knowing how to read the help documentation will make you a much more independent and creative programmer.

For example, type `?rnorm` in your console in the command line. This pulls up the help page for the `rnorm()` function. What does the function `rnorm()` do? The help file gives you a description of all things having to do with the normal distribution (you will see a few other functions listed, such as `pnorm()`, `dnorm()`, and `qnorm()`).

The *Usage* section lets you know what arguments go into which functions. If there's an equal sign (i.e., `=`), that means that argument has a default value, so you don't need to specify that argument. You only need to specify if you want something different than the default.

The *Arguments* section describes what each of the arguments are.

So, for example, `rnorm(n = 3, mean = 4, sd = 1)` generates 3 random numbers from a normal distribution with mean 4 and standard deviation 1. And `rnorm(n = 3)` generates 3 random numbers from a normal distribution with mean 0 and standard deviation 1. Notice we don't need to specify the `mean` and `sd` arguments in `rnorm()` if we want to generate numbers from a $Normal(\mu = 0, \sigma = 1)$ distribution because these are the default parameters (i.e., in the Usage section of the help file, the arguments `mean` and `sd` are equal to 0 and 1, respectively). Conversely, we always need to specify what `n` is because there's no default (i.e., no '=' sign after `n`).

2.2 Assigning values to R objects

We can assign results of computations to variables. Type each of the following lines either in an R script (and run each line), or in your console:

```
> x = (6-2)^.5/3^2
```

Here we are assigning a complicated arithmetic computation to the variable `x`. Notice what comes up in the workspace window when you assign the variable `x` to that value. Notice what comes up in the console when you run `x`:

```
> x
```

```
[1] 0.2222222
```

`identical()` is a function that outputs a logical vector telling us whether two variables are identical:

```
> X = 3
> identical(x, X)
```

```
[1] FALSE
```

As you can see here: R is case sensitive! Here we're catenating `x` and `X` into a vector called `y`:

```
> y = c(x, X)
> y
```

```
[1] 0.2222222 3.0000000
```

On your own (remember, use `?` to figure out what a function does and what arguments that go into a function):

1. Use the `rep()` function to make a vector that repeats the number “4” 5 times (i.e., 4 4 4 4 4). Assign that vector to the variable `z`.
2. Use the `seq()` function to make a vector that shows odd numbers from 1 to 10 (i.e., 1 3 5 7 9). Assign that vector to the variable `y`. Notice how we’ve overwritten the variable `y` we had before.
3. Use the `sample()` function to sample 17 of elements from `y` with replacement. Assign that vector to the variable `w`.
4. Add `z + y`.
5. Add `y + w`. What happened? Why? (Hint: use the `length()` function on `y` and `w`).
6. Remove `w` using the `rm()` function.

2.3 Data types and subsetting

Subsetting just means picking out elements of a vector (or other objects). This will be useful for picking out certain variables or observations of a dataset. Let’s reset `z` and `y` to the following (i.e., run these two lines):

```
> z = c(2, 4, 6, 7, 10)
> y = c(1, 3, 5, 7, 9)
```

Catenate `y` and `z`, and assign to `x`.

```
> x = c(y, z)
```

2.3.1 Subsetting

- Subsetting by *position*

```
> x[3]
```

```
[1] 5
```

Here, we are getting the 3rd element of `x`. This is subsetting by *position*.

- Subsetting by *exclusion*

```
> x[-3]
```

```
[1] 1 3 7 9 2 4 6 7 10
```

Here, we are getting all but the 3rd element of `x`. This is subsetting by *exclusion*.

- Subsetting by *name*

Write the following code. There’s a “built-in” vector called `letters` in R that has all the letters in the alphabet:

```
> names(x) = letters[1:length(x)]
> x
```

```

a b c d e f g h i j
1 3 5 7 9 2 4 6 7 10

```

Above, we're taking the first 10 (since that's the length of `x`) letters of the alphabet, and assigning those letters as names for the element in `x`. Below, we're getting all elements of `x` with the name 'c':

```

> x['c']

c
5

```

This is subsetting by *name*.

- Subsetting by *logical*

The following returns a vector of “logicals” telling us which elements of `y` are equal to 3:

```

> y==3

[1] FALSE TRUE FALSE FALSE FALSE

```

The second element returned is `TRUE` because the second element of `y` is 3. Notice our use of the double equals sign (`==`). You can also do `>` for greater than, `<` for less than, `>=` for greater than or equal to, `<=` for less than or equal to. Below, we're subsetting the elements of `z` based on the elements of `y` are equal to 3:

```

> z[y==3]

[1] 4

```

This is subsetting by *logical*.

- Subsetting by *everything*:

```

> x[]

a b c d e f g h i j
1 3 5 7 9 2 4 6 7 10

```

Here, we are getting everything from `x`.

On your own:

1. Get every other element of `x` (hint: use one of the functions you implemented above!)
2. Get elements with the name “a”, “c”, and “f”
3. Make every element of `x` equal to 4.

2.4 Matrices and dataframes

2.4.1 Matrices

Matrices are “shaped” vectors, in which all items have to be of the same type. Let's reset `x` to be a vector with numbers from 1 to 10:

```

> x = 1:10

```

Now let's construct a matrix, `mat`, which uses the elements of `x`, and makes a "shape" out of them so that it has 2 columns and 5 rows:

```
> mat = matrix(x, nrow = 5)
```

What are the dimensions of the matrix, `mat`?

```
> dim(mat)
```

```
[1] 5 2
```

Let's subset out the second row of `mat`:

```
> mat[2,]
```

```
[1] 2 7
```

Now let's subset out the first column of `mat`. Notice how it returns a vector.

```
> mat[,1]
```

```
[1] 1 2 3 4 5
```

On your own:

1. Get the element that's the third row and first column of the `mat` matrix.
2. Assign names to the `mat` matrix so that it looks like this:

```
b c
a 1 4
b 3 4
c 5 4
d 7 4
e 9 4
```

Hint: Use `rownames(mat)` and `colnames(mat)`.

2.4.2 Dataframes

A data frame is a collection of columns that can be of different types. This is what you typically think of as a dataset. In your group projects, you may need to load a data file read/manipulate it in R. If you haven't done so already, download the `CAquakes.csv` file from the Lab 0 folder in bCourses. Notice which folder on your computer you've saved it to, or move it to a folder that makes sense for this assignment. First, set your working directory:

```
> setwd("/Users/linamontoya/Dropbox/252E 2017/RLabs/RLab0")
```

This is where CAQuakes.csv lives on my computer. Change the filepath according to where CAQuakes.csv lives on your computer and make sure to put quotes around the filepath. Read in the CAQuakes .csv file:

```
> CAQuakes = read.csv(file = "CAQuakes.csv")
```

Note, the above will only work if the working directory is set to where CAQuakes.csv lives on your computer. Also notice how the dataframe CAQuakes has appeared in your workspace. Let's get the data type of CAQuakes:

```
> class(CAQuakes)
```

```
[1] "data.frame"
```

We see R read it in as a data frame. Type View(CAQuakes) in your console. Inspect your data:

```
> dim(CAQuakes) # dimensions of your data frame
```

```
[1] 383  4
```

```
> head(CAQuakes) # this gets the first 6 elements of CAQuakes
```

	Date	Latitude	Longitude	M
1	18001011	36.8	-121.5	5.5
2	18001122	32.9	-117.8	6.3
3	18030000	34.2	-118.1	5.5
4	18030525	32.8	-117.1	5.5
5	18060325	34.4	-119.7	5.5
6	18080621	37.8	-122.6	5.5

```
> summary(CAQuakes)
```

	Date	Latitude	Longitude	M
Min.	:18001011	Min. :31.50	Min. : -125.9	Min. :5.500
1st Qu.	:18890676	1st Qu.:34.15	1st Qu.: -121.8	1st Qu.:5.600
Median	:19261024	Median :36.50	Median : -119.6	Median :5.900
Mean	:19235315	Mean :36.37	Mean : -119.7	Mean :5.992
3rd Qu.	:19585801	3rd Qu.:38.30	3rd Qu.: -117.2	3rd Qu.:6.200
Max.	:19991016	Max. :42.31	Max. : -114.5	Max. :7.900

The summary() function gives you information on your object, depending on what kind of object it is. What kind of information does the summary() function give for a dataframe?

Here are 3 ways of extracting the first 6 elements of Latitude column vector:

```
> head(CAQuakes$Latitude)
```

```
[1] 36.8 32.9 34.2 32.8 34.4 37.8
```

```
> head(CAQuakes[["Latitude"]])
```

```
[1] 36.8 32.9 34.2 32.8 34.4 37.8
```

```
> head(CAquakes[,2])
```

```
[1] 36.8 32.9 34.2 32.8 34.4 37.8
```

On your own:

1. Using the `CAquakes` data, get the mean Longitude.
2. Look up the function `hist()` and plot the magnitudes (`M`) of the earthquakes. Try changing the main title and x-axis title to something more informative than the default.
3. Use a scatterplot to visualize how earthquake magnitude behaves as latitude increases (*Hint*: look up and use the `plot()` function). Again, try changing the titles.

3 Creating functions and generating data

There are many built-in functions in R (for example, above we saw the `rep()`, `seq()`, and `rnorm()` functions, to name a few). However, it's also possible to write your own functions. We will often create functions in these labs, in particular, to generate many copies of variables that comprise a dataset. The general syntax for a function is as follows:

```
> examplefun = function(x) {
+
+   y = x + 3 # what the function does
+
+   return(y) # what we want the function to return back (output)
+ }
```

`examplefun` is the name of the function (we make that up), and `x` is the argument (input). `examplefun()` takes in as input `x`, and outputs `y`, which is the input plus 3. For example, if input 4 for `x`, we anticipate `examplefun()` to spit out 7 (i.e., $4 + 3$):

```
> examplefun(x = 4)
```

```
[1] 7
```

On your own:

1. Set the seed to 252.
Note: Setting the seed to a value (generally) ensures that running the same chunk of code twice will return the exact same output, even if the code involves “random” numbers. This will also help your GSIs help you with your labs/homework assignments/projects.
2. Create a function, `Uniform()`, that generates n i.i.d. random variables uniformly distributed between 0 and x .
Hint: the function should take in 2 arguments as input: 1) n , the number of uniform variables to generate, and 2) x , the maximum value the uniform variable can take. *Extra hint*: look up the `runif()` function to randomly generate values from the uniform distribution.

3. Create a function, `generate_data()`, that generates a dataframe with n i.i.d. random variables that follow these distributions:

$$\begin{aligned}U_{W1} &\sim \text{Uniform}(0, 1) \\U_{W2} &\sim \text{Normal}(\mu = 1, \sigma^2 = 2^2) \\U_A &\sim \text{Uniform}(0, 1) \\U_Y &\sim \text{Uniform}(0, 1)\end{aligned}$$

Hint: the skeleton of your function should look like this (fill in the commented out parts with your own code):

```
> generate_data = function(n) {
+
+   # generate n Uniform(min = 0, max = 1) variables here, set equal to U.W1
+   # generate n Normal(mean = 0, max = 1) variables here, set equal to U.W2
+   # generate n Uniform(min = 0, max = 1) variables here, set equal to U.A
+   # generate n Uniform(min = 0, max = 1) variables here, set equal to U.Y
+
+ }
```

where `generate_data` takes in the single argument `n`, the number of observations. Try using the `Uniform()` function you created in the previous step to generate the uniform random variables.

4. Add to the function `generate_data()` so that it also creates 4 new variables based on the following functions that relate them:

$$\begin{aligned}W1 &= \mathbb{I}[U_{W1} < 0.45] \\W2 &= 0.75 * U_{W2} \\A &= \mathbb{I}[U_A < \text{expit}(-1 + 2.6 * W1 + 0.9 * W2)] \\Y &= \mathbb{I}[U_Y < \text{expit}(-2 + A + 0.7 * W1)]\end{aligned}$$

Hint: Use the `as.numeric()` function for indicator variables.

Example: $\mathbb{I}[X > 0]$ can be coded as `as.numeric(X>0)`.

Extra hint: Recall that the *expit* function is the inverse of the logistic function:

$$\begin{aligned}\text{logit}(x) &= \log\left(\frac{x}{1-x}\right) \\ \text{expit}(x) &= \frac{1}{1 + e^{-x}}\end{aligned}$$

In R, the *expit* function is called `plogis()`.

5. Within `generate_data()`, create a dataframe using either the `data.frame()` function or the `cbind()` function consisting of only $W1$, $W2$, A , and Y and let that be your output (i.e., what you `return`) for the function `generate_data()`.
6. Create 1000 copies of $W1$, $W2$, A , and Y using `generate_data()` and set the output of the function equal to the object `ObsData`.
7. Show the first 6 lines and summary statistics of `ObsData`, using the `head()` and `summary()` functions, respectively.

4 Implementation of g-computation (simple substitution estimator), IPTW, and TMLE estimation of average treatment effect of a point treatment

Using `ObsData` we created in the previous section, we will implement g-computation, IPTW, and TMLE estimators using `SuperLearner` to generate estimates inspired by the average treatment effect (e.g., $\mathbb{E}[Y_{a=1} - Y_{a=0}]$). Recall from Causal I labs (note: not to worry if you don't remember much about these estimators – we will review them [including how to implement them] later on in this course):

1. Simple substitution estimator based on the G-Computation formula:

$$\hat{\Psi}_{SS}(\hat{\mathbb{P}}) = \frac{1}{n} \sum_{i=1}^n (\hat{Q}(1, W_i) - \hat{Q}(0, W_i))$$

where $\hat{\mathbb{P}}$ is the empirical distribution and $\hat{Q}(A, W)$ is the estimate of the conditional mean outcome given the exposure and baseline covariates $Q_0(A, W) \equiv \mathbb{E}_0(Y|A, W)$.

- Consistency of the simple (non-targeted) substitution estimator depends on consistent estimation of the conditional mean outcome $Q_0(A, W)$.

2. Standard (unstabilized) inverse probability weighted estimator (IPTW):

$$\hat{\Psi}_{IPTW}(\hat{\mathbb{P}}) = \frac{1}{n} \sum_{i=1}^n \left[\frac{\mathbb{I}(A_i = 1)}{\hat{g}(1|W_i)} - \frac{\mathbb{I}(A_i = 0)}{\hat{g}(0|W_i)} \right] Y_i$$

where $\hat{g}(1|W_i) = \hat{\mathbb{P}}(A_i = 1|W_i)$ is an estimate of the exposure mechanism.

- Consistency of IPTW estimators depends on consistent estimation of the exposure mechanism $g_0(1|W) \equiv \mathbb{P}_0(A = 1|W)$.

3. Targeted maximum likelihood estimation (TMLE):

$$\hat{\Psi}_{TMLE}(\hat{\mathbb{P}}) = \frac{1}{n} \sum_{i=1}^n (\bar{Q}_n^*(1, W_i) - \bar{Q}_n^*(0, W_i))$$

where $\bar{Q}_n^*(A, W)$ denotes the *targeted* estimate of the conditional mean outcome, given the exposure and baseline covariates $\bar{Q}_0(A, W)$.

- Implementation requires estimation of both the conditional mean function $\bar{Q}_0(A, W)$ and the exposure mechanism $g_0(A|W)$.

- Double robust estimators are consistent if either $\bar{Q}_0(A, W)$ or $g_0(A|W)$ is estimated consistently.

Reminder: An estimator is *consistent* if the point estimates converge (in probability) to the estimand as sample size $n \rightarrow \infty$.

Reminder 2: The notation of a subscript n or hat denotes that the value is an estimate (i.e., it comes from the empirical distribution \mathbb{P}_n or $\hat{\mathbb{P}}$). The notation of subscript 0 denotes that it is a true value (i.e., the value is comes from the true population distribution, \mathbb{P}_0). Again, we will review this later in the course :)

Warning: in practice you must go through the entire roadmap to infer any causality. Here, we are simply demonstrating how to implement the estimators.

1. Load the `SuperLearner` and `tmle` or `ltmle` packages by using the `library()` function, and set the seed to 252.

2. Set the object `n` to the number of observations in `ObsData`.

Hint: use the `nrow()` function.

3. Create a dataframe, `X`, that contains only `W1`, `W2`, and `A` of `ObsData`

Hint: use the `subset()` function.

```
> X = subset(ObsData, select = c(W1, W2, A))
```

4. Make new dataframes `X1` and `X0` that are copies of `X`. For the dataframe `X1`, set `A = 1`; for the dataframe `X0`, set `A = 0`. For example:

```
> X1$A = 1 # set A = 1 in X1 dataframe. Repeat for X0 but with A = 0.
```

5. Stack `X`, `X1`, and `X0` by using the `rbind()` function. Set the stacked data equal to the object `newdata`.

6. Set `SL.library = c('SL.glm', 'SL.mean', 'SL.glm.interaction')`. This is the SuperLearner library we'll use to estimate the outcome regression and treatment mechanism needed for the estimators.

7. **Call `SuperLearner()` to estimate $Q_0(A, W1, W2) = \mathbb{E}_0[Y|A, W1, W2]$, the outcome regression.** Specify the arguments `Y = ObsData$Y`, `X = X`, `newX = newdata`, `SL.library = SL.library`, and `family = 'gaussian'`. Set the SuperLearner object equal to the `Qbar`.

8. Obtain the predicted/expected `Y` given:

- (a) the observed `A` and covariates

```
> QbarAW = Qbar$SL.predict[1:n]
```

- (b) `A = 1` and the observed covariates

```
> Qbar1W = Qbar$SL.predict[(n+1):(2*n)]
```

- (c) `A = 0` and the observed covariates

```
> Qbar0W = Qbar$SL.predict[(2*n+1):(3*n)]
```

9. **Evaluate the simple substitution estimator** by taking the mean of `Qbar1W` and subtracting it from the mean of `Qbar0W`:

$$\begin{aligned}\hat{\Psi}_{SS}(\hat{\mathbb{P}}) &= \frac{1}{n} \sum_{i=1}^n \left[\hat{\mathbb{E}}(Y_i | A = 1, W1_i, W2_i) - \hat{\mathbb{E}}(Y_i | A = 0, W1_i, W2_i) \right] \\ &= \frac{1}{n} \sum_{i=1}^n \left[\bar{Q}_n(1, W1_i, W2_i) - \bar{Q}_n(0, W1_i, W2_i) \right]\end{aligned}$$

10. Create a dataframe with only the covariates. Call this dataframe `W`.

11. **Call `SuperLearner()` to estimate $g_0(A|W1, W2) = P_0(A|W1, W2)$, the treatment mechanism.** Specify the arguments `Y = ObsData$A`, `X = W`, `SL.library = SL.library`, and `family = 'binomial'`. Set the SuperLearner object equal to the `gHatSL`.

12. Generate the predicted probability of:

- (a) `A = 1` given covariates

```
> gHat1W = gHatSL$SL.predict
```

- (b) `A = 0` given covariates

```
> gHat0W = 1 - gHatSL$SL.predict
```

13. **Implement the IPTW estimator** by taking a weighted average of `Y`, with weights equal to $w_i = \frac{1}{g_n(A_i|W1_i, W2_i)}$:

$$\begin{aligned}\hat{\Psi}_{IPTW}(\hat{\mathbb{P}}) &= \frac{1}{n} \sum_{i=1}^n \frac{\mathbb{I}(A_i = 1)}{g_n(A_i|W1_i, W2_i)} Y_i - \frac{1}{n} \sum_{i=1}^n \frac{\mathbb{I}(A_i = 0)}{g_n(A_i|W1_i, W2_i)} Y_i \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{I}(A_i = 1)^* \hat{w}_i^* Y_i - \frac{1}{n} \sum_{i=1}^n \mathbb{I}(A_i = 0)^* \hat{w}_i^* Y_i\end{aligned}$$

14. **Implement the TMLE estimator** by calling `tmle()` or `ltmle()`.

- If you're using the `tmle()` package, specify the arguments `Y = ObsData$Y`, `A = ObsData$A`, `W = W`, `Q.SL.library = SL.library`, `g.SL.library = SL.library`.
- If you're using the `ltmle()` package, specify the arguments `data=ObsData`, `Anodes='A'`, `Ynodes='Y'`, `abar=list(1,0)`, `SL.library=SL.library`.