

Statistical Tools for Causal Inference

The SKY Community

2019-04-17

Contents

Introduction	5
I The Two Fundamental Problems of Inference	7
1 The Fundamental Problem of Causal Inference	11
1.1 Rubin Causal Model	11
1.2 Treatment effects	17
1.3 Fundamental problem of causal inference	20
1.4 Intuitive estimators, confounding factors and selection bias	21

Introduction

Tools of causal inference are the basic statistical building block behind most scientific results. It is thus extremely useful to have an open source collectively agreed upon resource presenting and assessing them, as well as listing the current unresolved issues. The content of this book covers the basic theoretical knowledge and technical skills required for implementing statistical methods of causal inference. This means:

- Understanding of the basic language to encode causality,
- Knowledge of the fundamental problems of inference and the biases of intuitive estimators,
- Understanding of how econometric methods recover treatment effects,
- Ability to compute these estimators along with an estimate of their precision using the statistical software R combined with latex using Rmarkdown.

All the notions and estimators are introduced using a numerical example and simulations.

All the code behind this book is written in Rmarkdown and is publically available on GitHub. Feel free to propose corrections and updates.

Part I

The Two Fundamental Problems of Inference

When trying to estimate the effect of a program on an outcome, we face two very important and difficult problems: the Fundamental Problem of Causal Inference (FPCI) and the Fundamental Problem of Statistical Inference (FPSI).

In its most basic form, the FPCI states that our causal parameter of interest (TT , short for Treatment on the Treated, that we will define shortly) is fundamentally unobservable, even when the sample size is infinite. The main reason for that is that one component of TT , the outcome of the treated had they not received the program, remains unobservable. We call this outcome a counterfactual outcome. The FPCI is a very dispiriting result, and is actually the basis for all of the statistical methods of causal inference. All of these methods try to find ways to estimate the counterfactual by using observable quantities that hopefully approximate it as well as possible. Most people, including us but also policymakers, generally rely on intuitive quantities in order to generate the counterfactual (the individuals without the program or the individuals before the program was implemented). Unfortunately, these approximations are generally very crude, and the resulting estimators of TT are generally biased, sometimes severely.

The Fundamental Problem of Statistical Inference (FPSI) states that, even if we have an estimator E that identifies TT in the population, we cannot observe E because we only have access to a finite sample of the population. The only thing that we can form from the sample is a sample equivalent \hat{E} to the population quantity E , and $\hat{E} \neq E$. Why is $\hat{E} \neq E$? Because a finite sample is never perfectly representative of the population. What can we do to deal with the FPSI? I am going to argue that there are mainly two things that we might want to do: estimating the extent of sampling noise and decreasing sampling noise.

Chapter 1

The Fundamental Problem of Causal Inference

In order to state the FPCI, we are going to describe the basic language to encode causality set up by Rubin, and named Rubin Causal Model (RCM). RCM being about partly observed random variables, it is hard to make these notions concrete with real data. That's why we are going to use simulations from a simple model in order to make it clear how these variables are generated. The second virtue of this model is that it is going to make it clear the source of selection into the treatment. This is going to be useful when understanding biases of intuitive comparisons, but also to discuss the methods of causal inference. A third virtue of this approach is that it makes clear the connexion between the treatment effects literature and models. Finally, a fourth reason that it is useful is that it is going to give us a source of sampling variation that we are going to use to visualize and explore the properties of our estimators.

I use X_i to denote random variable X all along the notes. I assume that we have access to a sample of N observations indexed by $i \in \{1, \dots, N\}$. “ i ” will denote the basic sampling units when we are in a sample, and a basic element of the probability space when we are in populations. Introducing rigorous measure-theoretic notations for the population is feasible but is not necessary for comprehension.

When the sample size is infinite, we say that we have a population. A population is a very useful fiction for two reasons. First, in a population, there is no sampling noise: we observe an infinite amount of observations, and our estimators are infinitely precise. This is useful to study phenomena independently of sampling noise. For example, it is in general easier to prove that an estimator is equal to TT under some conditions in the population. Second, we are most of the time much more interested in estimating the values of parameters in the population rather than in the sample. The population parameter, independent of sampling noise, gives a much better idea of the causal parameter for the population of interest than the parameter in the sample. In general, the estimator for both quantities will be the same, but the estimators for the effect of sampling noise on these estimators will differ. Sampling noise for the population parameter will generally be larger, since it is affected by another source of variability (sample choice).

1.1 Rubin Causal Model

The RCM is made of three distinct building blocks: a treatment allocation rule, that decides who receives the treatment; potential outcomes, that measure how each individual reacts to the treatment; the switching equation that relates potential outcomes to observed outcomes through the allocation rule.

1.1.1 Treatment allocation rule

The first building block of the RCM is the treatment allocation rule. Throughout this class, we are going to be interested in inferring the causal effect of only one treatment with respect to a control condition. Extensions to multi-valued treatments are in general self-explanatory.

In the RCM, treatment allocation is captured by the variable D_i . $D_i = 1$ if unit i receives the treatment and $D_i = 0$ if unit i does not receive the treatment and thus remains in the control condition.

The treatment allocation rule is critical for several reasons. First, because it switches the treatment on or off for each unit, it is going to be at the source of the FPCI. Second, the specific properties of the treatment allocation rule are going to matter for the feasibility and bias of the various econometric methods that we are going to study.

Let's take a few examples of allocation rules. These allocation rules are just examples. They do not cover the space of all possible allocation rules. They are especially useful as concrete devices to understand the sources of biases and the nature of the allocation rule. In reality, there exists even more complex allocation rules (awareness, eligibility, application, acceptance, active participation). Awareness seems especially important for program participation and has only been tackled recently by economists.

First, some notation. Let's imagine a treatment that is given to individuals. Whether each individual receives the treatment partly depends on the level of her outcome before receiving the treatment. Let's denote this variable Y_i^B , with B standing for "Before". It can be the health status assessed by a professional before deciding to give a drug to a patient. It can be the poverty level of a household used to assess its eligibility to a cash transfer program.

1.1.1.1 Sharp cutoff rule

The sharp cutoff rule means that everyone below some threshold \bar{Y} is going to receive the treatment. Everyone whose outcome before the treatment lies above \bar{Y} does not receive the treatment. Such rules can be found in reality in a lot of situations. They might be generated by administrative rules. One very simple way to model this rule is as follows:

$$D_i = \mathbb{1}[Y_i^B \leq \bar{Y}], \quad (1.1)$$

where $\mathbb{1}[A]$ is the indicator function, taking value 1 when A is true and 0 otherwise.

Example 1.1 (Sharp cutoff rule). Imagine that $Y_i^B = \exp(y_i^B)$, with $y_i^B = \mu_i + U_i^B$, $\mu_i \sim \mathcal{N}(\bar{\mu}, \sigma_\mu^2)$ and $U_i^B \sim \mathcal{N}(0, \sigma_U^2)$. Now, let's choose some values for these parameters so that we can generate a sample of individuals and allocate the treatment among them. I'm going to switch to R for that.

```
param <- c(8,.5,.28,1500)
names(param) <- c("barmu","sigma2mu","sigma2U","barY")
param
```

```
##      barmu sigma2mu  sigma2U      barY
##      8.00      0.50      0.28 1500.00
```

Now, I have chosen values for the parameters in my model. For example, $\bar{\mu} = 8$ and $\bar{Y} = 1500$. What remains to be done is to generate Y_i^B and then D_i . For this, I have to choose a sample size ($N = 1000$) and then generate the shocks from a normal.

```
# for reproducibility, I choose a seed that will give me the same random sample each time I run the pro
set.seed(1234)
N <- 1000
mu <- rnorm(N,param["barmu"],sqrt(param["sigma2mu"]))
```

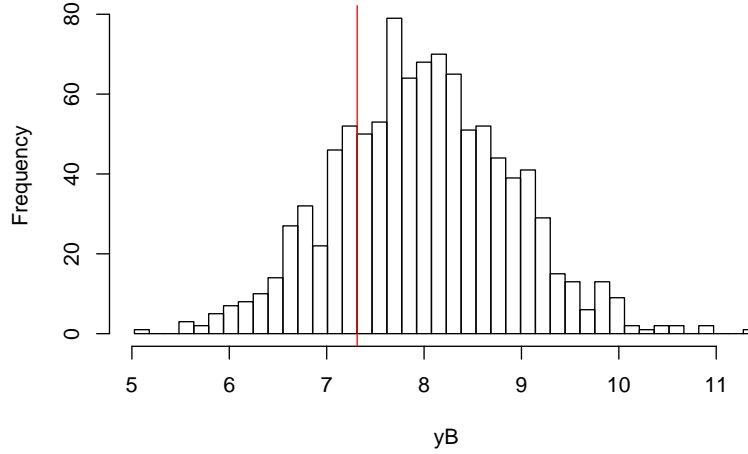
Figure 1.1: Histogram of y_B

Table 1.1: Treatment allocation with sharp cutoff rule

0	771
1	229

```
UB <- rnorm(N,0,sqrt(param["sigma2U"]))
yB <- mu + UB
YB <- exp(yB)
Ds <- ifelse(YB<=param["barY"],1,0)
```

Let's now build a histogram of the data that we have just generated.

```
# building histogram of yB with cutoff point at ybar
# Number of steps
Nsteps.1 <- 15
#step width
step.1 <- (log(param["barY"])-min(yB[Ds==1]))/Nsteps.1
Nsteps.0 <- (-log(param["barY"])+max(yB[Ds==0]))/step.1
breaks <- cumsum(c(min(yB[Ds==1]),c(rep(step.1,Nsteps.1+Nsteps.0+1))))
hist(yB,breaks=breaks,main="")
abline(v=log(param["barY"]),col="red")
```

You can see on Figure 1.1 a histogram of y_i^B with the red line indicating the cutoff point: $\bar{y} = \ln(\bar{Y}) = 7.3$. All the observations below the red line are treated according to the sharp rule while all the one located above are not. In order to see how many observations eventually receive the treatment with this allocation rule, let's build a contingency table.

```
table.D.sharp <- as.matrix(table(Ds))
knitr::kable(table.D.sharp,caption='Treatment allocation with sharp cutoff rule',booktabs=TRUE)
```

We can see on Table 1.1 that there are 229 treated observations.

1.1.1.2 Fuzzy cutoff rule

This rule is less sharp than the sharp cutoff rule. Here, other criteria than Y_i^B enter into the decision to allocate the treatment. The doctor might measure the health status of a patient following official guidelines, but he might also measure other factors that will also influence his decision of giving the drug to the patient.

The officials administering a program might measure the official income level of a household, but they might also consider other features of the household situation when deciding to enroll the household into the program or not. If these additional criteria are unobserved to the econometrician, then we have the fuzzy cutoff rule. A very simple way to model this rule is as follows:

$$D_i = \mathbb{1}[Y_i^B + V_i \leq \bar{Y}], \quad (1.2)$$

where V_i is a random variable unobserved to the econometrician and standing for the other influences that might drive the allocation of the treatment. V_i is distributed according to a, for the moment, unspecified cumulative distribution function F_V . When V_i is degenerate (*i.e.* it has only one point of support: it is a constant), the fuzzy cutoff rule becomes the sharp cutoff rule.

1.1.1.3 Eligibility + self-selection rule

It is also possible that households, once they have been made eligible to the treatment, can decide whether they want to receive it or not. A patient might be able to refuse the drug that the doctor suggests she should take. A household might refuse to participate in a cash transfer program to which it has been made eligible. Not all programs have this feature, but most of them have some room for decisions by the agents themselves of whether they want to receive the treatment or not. One simple way to model this rule is as follows:

$$D_i = \mathbb{1}[D_i^* \geq 0]E_i, \quad (1.3)$$

where D_i^* is individual i 's valuation of the treatment and E_i is whether or not she is deemed eligible for the treatment. E_i might be chosen according to the sharp cutoff rule or to the fuzzy cutoff rule, or to any other eligibility rule. We will be more explicit about D_i^* in what follows.

SIMULATIONS ARE MISSING FOR THESE LAST TWO RULES

1.1.2 Potential outcomes

The second main building block of the RCM are potential outcomes. Let's say that we are interested in the effect of a treatment on an outcome Y . Each unit i can thus be in two potential states: treated or non treated. Before the allocation of the treatment is decided, both of these states are feasible for each unit.

Definition 1.1 (Potential outcomes). For each unit i , we define two potential outcomes:

- Y_i^1 : the outcome that unit i is going to have if it receives the treatment,
- Y_i^0 : the outcome that unit i is going to have if it **does not** receive the treatment.

Example 1.2. Let's choose functional forms for our potential outcomes. For simplicity, all lower case letters will denote log outcomes. $y_i^0 = \mu_i + \delta + U_i^0$, with δ a time shock common to all the observations and $U_i^0 = \rho U_i^B + \epsilon_i$, with $|\rho| < 1$. In the absence of the treatment, part of the shocks U_i^B that the individuals experienced in the previous period persist, while some part vanish. $y_i^1 = y_i^0 + \bar{\alpha} + \theta\mu_i + \eta_i$. In order to generate the potential outcomes, one has to define the laws for the shocks and to choose parameter values. Let's assume that $\epsilon_i \sim \mathcal{N}(0, \sigma_\epsilon^2)$ and $\eta_i \sim \mathcal{N}(0, \sigma_\eta^2)$. Now let's choose some parameter values:

```
1 <- length(param)
param <- c(param,0.9,0.01,0.05,0.05,0.05,0.1)
names(param)[(1+1):length(param)] <- c("rho","theta","sigma2epsilon","sigma2eta","delta","baralpha")
param
```

##	barmu	sigma2mu	sigma2U	barY	rho
##	8.00	0.50	0.28	1500.00	0.90

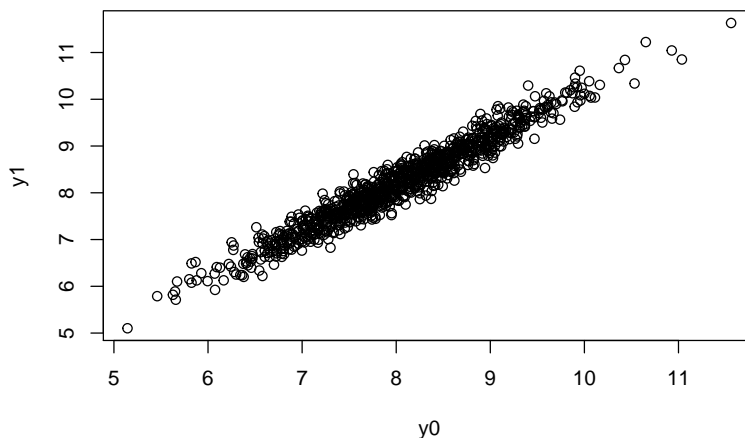


Figure 1.2: Potential outcomes

```
##      theta sigma2epsilon      sigma2eta      delta      baralpha
##      0.01      0.05      0.05      0.05      0.10
```

We can finally generate the potential outcomes;

```
epsilon <- rnorm(N,0,sqrt(param["sigma2epsilon"]))
eta <- rnorm(N,0,sqrt(param["sigma2eta"]))
U0 <- param["rho"]*UB + epsilon
y0 <- mu + U0 + param["delta"]
alpha <- param["baralpha"]+ param["theta"]*mu + eta
y1 <- y0+alpha
Y0 <- exp(y0)
Y1 <- exp(y1)
```

Now, I would like to visualize my potential outcomes:

```
plot(y0,y1)
```

You can see on the resulting Figure 1.2 that both potential outcomes are positively correlated. Those with a large potential outcome when untreated (*e.g.* in good health without the treatment) also have a positive health with the treatment. It is also true that individuals with bad health in the absence of the treatment also have bad health with the treatment.

1.1.3 Switching equation

The last building block of the RCM is the switching equation. It links the observed outcome to the potential outcomes through the allocation rule:

$$\begin{aligned}
 Y_i &= \begin{cases} Y_i^1 & \text{if } D_i = 1 \\ Y_i^0 & \text{if } D_i = 0 \end{cases} \\
 &= Y_i^1 D_i + Y_i^0 (1 - D_i)
 \end{aligned} \tag{1.4}$$

Example 1.3. In order to generate observed outcomes in our numerical example, we simply have to enforce the switching equation:

```
y <- y1*Ds+y0*(1-Ds)
Y <- Y1*Ds+Y0*(1-Ds)
```

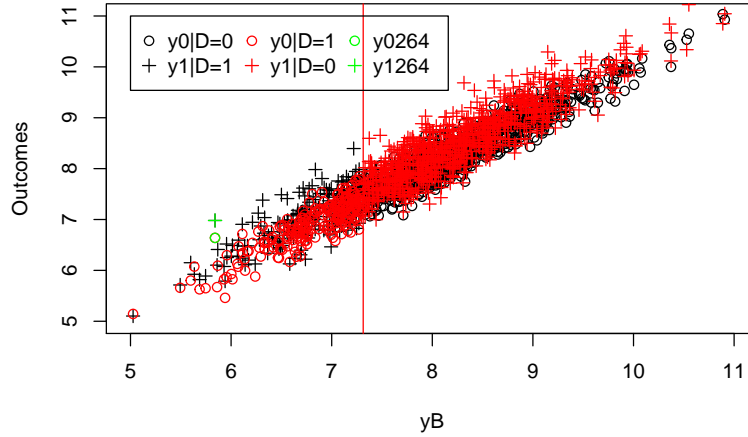


Figure 1.3: Potential outcomes

What the switching equation (1.4) means is that, for each individual i , we get to observe only one of the two potential outcomes. When individual i belongs to the treatment group (*i.e.* $D_i = 1$), we get to observe Y_i^1 . When individual i belongs to the control group (*i.e.* $D_i = 0$), we get to observe Y_i^0 . Because the same individual cannot be at the same time in both groups, we can NEVER see both potential outcomes for the same individual at the same time.

For each of the individuals, one of the two potential outcomes is unobserved. We say that it is a **counterfactual**. A counterfactual quantity is a quantity that is, according to Hume's definition, contrary to the observed facts. A counterfactual cannot be observed, but it can be conceived by an effort of reason: it is the consequence of what would have happened had some action not been taken.

Remark. One very nice way of visualising the switching equation has been proposed by Jerzy Neyman in a 1923 prescient paper. Neyman proposes to imagine two urns, each one filled with N balls. One urn is the treatment urn and contains balls with the id of the unit and the value of its potential outcome Y_i^1 . The other urn is the control urn, and it contains balls with the value of the potential outcome Y_i^0 for each unit i . Following the allocation rule D_i , we decide whether unit i is in the treatment or control group. When unit i is in the treatment group, we take the corresponding ball from the first urn and observe the potential outcome on it. But, at the same time, the urns are connected so that the corresponding ball with the potential outcome of unit i in the control urn disappears as soon as we draw ball i from the treatment urn.

The switching equation works a lot like Schrodinger's cat paradox. Schrodinger's cat is placed in a sealed box and receives a dose of poison when an atom emits a radiation. As long as the box is sealed, there is no way we can know whether the cat is dead or alive. When we open the box, we observe either a dead cat or a living cat, but we cannot observe the cat both alive and dead at the same time. The switching equation is like opening the box, it collapses the observed outcome into one of the two potential ones.

Example 1.4. One way to visualize the inner workings of the switching equation is to plot the potential outcomes along with the criteria driving the allocation rule. In our simple example, it simply amounts to plotting observed (y_i) and potential outcomes (y_i^1 and y_i^0) along y_i^B .

```
plot(yB[Ds==0], y0[Ds==0], pch=1, xlim=c(5, 11), ylim=c(5, 11), xlab="yB", ylab="Outcomes")
points(yB[Ds==1], y1[Ds==1], pch=3)
points(yB[Ds==0], y1[Ds==0], pch=3, col='red')
points(yB[Ds==1], y0[Ds==1], pch=1, col='red')
test <- 5.8
i.test <- which(abs(yB-test)==min(abs(yB-test)))
points(yB[abs(yB-test)==min(abs(yB-test))], y1[abs(yB-test)==min(abs(yB-test))], col='green', pch=3)
points(yB[abs(yB-test)==min(abs(yB-test))], y0[abs(yB-test)==min(abs(yB-test))], col='green')
abline(v=log(param["barY"]), col="red")
legend(5, 11, c('y0|D=0', 'y1|D=1', 'y0|D=1', 'y1|D=0', paste('y0', i.test, sep=''), paste('y1', i.test, sep='')),
```

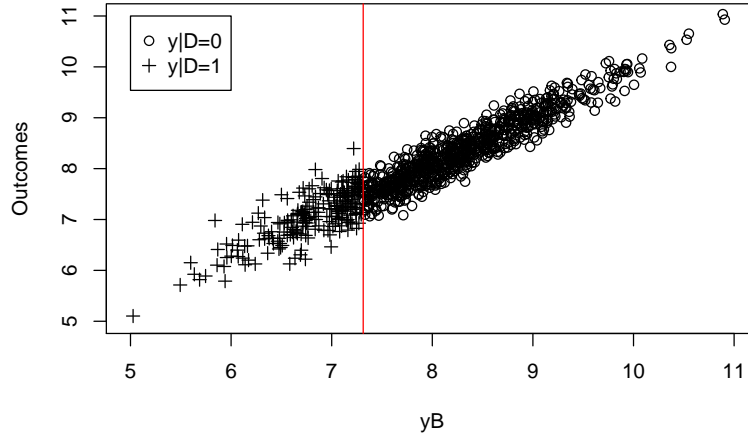



Figure 1.4: Observed outcomes

```
plot(yB[Ds==0], y0[Ds==0], pch=1, xlim=c(5, 11), ylim=c(5, 11), xlab="yB", ylab="Outcomes")
points(yB[Ds==1], y1[Ds==1], pch=3)
legend(5, 11, c('y|D=0', 'y|D=1'), pch=c(1, 3))
abline(v=log(param["barY"]), col="red")
```

Figure 1.4 plots the observed outcomes y_i that results from applying the switching equation. Figure 1.4 shows that each individual in the sample is endowed with two potential outcomes, represented by a circle and a cross. Figure 1.3 plots the observed outcomes y_i along with the unobserved potential outcomes. Only one of the two potential outcomes is observed (the cross for the treated group and the circle for the untreated group) and the other is not. The observed sample in Figure 1.3 only shows observed outcomes, and is thus silent on the values of the missing potential outcomes.

1.2 Treatment effects

The RCM enables the definition of causal effects at the individual level. In practice though, we generally focus on a summary measure: the effect of the treatment on the treated.

1.2.1 Individual level treatment effects

Potential outcomes enable us to define the central notion of causal inference: the causal effect, also labelled the treatment effect, which is the difference between the two potential outcomes.

Definition 1.2 (Individual level treatment effect). For each unit i , the causal effect of the treatment on outcome Y is: $\Delta_i^Y = Y_i^1 - Y_i^0$.

Example 1.5. The individual level causal effect in log terms is: $\Delta_i^y = \alpha_i = \bar{\alpha} + \theta\mu_i + \eta_i$. The effect is the sum of a part common to all individuals, a part correlated with μ_i : the treatment might have a larger or a smaller effect depending on the unobserved permanent ability or health status of individuals, and a random shock. It is possible to make the effect of the treatment to depend on U_i^B also, but it would complicate the model.

In Figure 1.4, the individual level treatment effects are the differences between each cross and its corresponding circle. For example, for observation 264, the two potential outcomes appear in green in Figure 1.4. The effect of the treatment on unit 264 is equal to:

$$\Delta_{264}^y = y_{264}^1 - y_{264}^0 = 6.98 - 6.64 = 0.34.$$

Since observation 264 belongs to the treatment group, we can only observe the potential outcome in the presence of the treatment, y_{264}^1 .

The RCM allows for heterogeneity of treatment effects. The treatment has a large effect on some units and a much smaller effect on other units. We can even have some units that benefit from the treatment and some units that are harmed by the treatment. The individual level effect of the treatment is itself a random variable (and not a fixed parameter). It has a distribution, F_{Δ^Y} .

Heterogeneity of treatment effects seems very natural: the treatment might interact with individuals' different backgrounds. The effect of a drug might depend on the genetic background of an individual. An education program might only work for children that already have sufficient non-cognitive skills, and thus might depend in turn on family background. An environmental regulation or a behavioral intervention might only trigger reactions by already environmentally aware individuals. A CCT might have a larger effect when individuals are credit-constrained or face shocks.

Example 1.6. In our numerical example, the distribution of $\Delta_i^Y = \alpha_i$ is a normal: $\alpha_i \sim \mathcal{N}(\bar{\alpha} + \theta\bar{\mu}, \theta^2\sigma_\mu^2 + \sigma_\eta^2)$. We would like to visualize treatment effect heterogeneity. For that, we can build a histogram of the individual level causal effect.

On top of the histogram, we can also draw the theoretical distribution of the treatment effect: a normal with mean 0.18 and variance 0.05.

```
hist(alpha,main="",prob=TRUE)
curve(dnorm(x, mean=(param["baralpha"]+param["theta"]*param["barmu"]), sd=sqrt(param["theta"]^2*param["
```

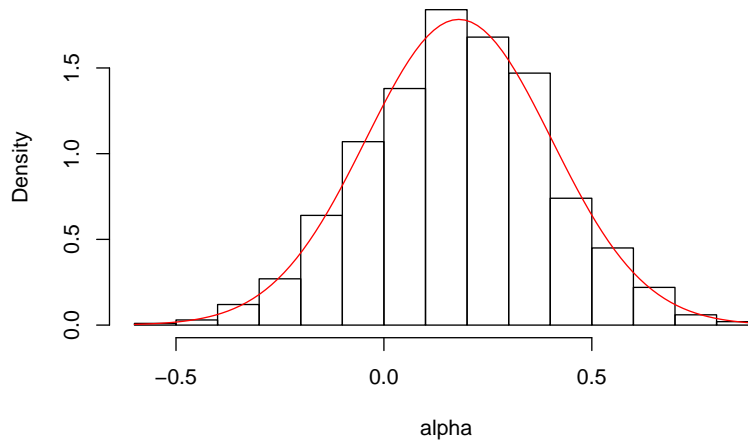


Figure 1.5: Histogram of Δ^Y

The first thing that we can see on Figure 1.5 is that the theoretical and the empirical distributions nicely align with each other. We also see that the majority of the observations lies to the right of zero: most people experience a positive effect of the treatment. But there are some individuals that do not benefit from the treatment: the effect of the treatment on them is negative.

1.2.2 Average treatment effect on the treated

We do not generally estimate individual-level treatment effects. We generally look for summary statistics of the effect of the treatment. By far the most widely reported causal parameter is the Treatment on the Treated parameter (TT). It can be defined in the sample at hand or in the population.

Definition 1.3 (Average and expected treatment effects on the treated). The Treatment on the Treated parameters for outcome Y are:

- The average Treatment effect on the Treated in the sample:

$$\Delta_{TT_s}^Y = \frac{1}{\sum_{i=1}^N D_i} \sum_{i=1}^N (Y_i^1 - Y_i^0) D_i,$$

- The expected Treatment effect on the Treated in the population:

$$\Delta_{TT}^Y = \mathbb{E}[Y_i^1 - Y_i^0 | D_i = 1].$$

The TT parameters measure the average effect of the treatment on those who actually take it, either in the sample at hand or in the population. It is generally considered to be the most policy-relevant parameter since it measures the effect of the treatment as it has actually been allocated. For example, the expected causal effect on the overall population is only relevant if policymakers are considering implementing the treatment even on those who have not been selected to receive it. For a drug or an anti-poverty program, it would mean giving the treatment to healthy or rich people, which would make little sense.

TT does not say anything about how the effect of the treatment is distributed in the population or in the sample. TT does not account for the heterogeneity of treatment effects. In Lecture 7, we will look at other parameters of interest that look more closely into how the effect of the treatment is distributed.

Example 1.7. The value of TT in our sample is:

$$\Delta_{TT_s}^y = 0.168.$$

Computing the population value of TT is slightly more involved: we have to use the formula for the conditional expectation of a censored bivariate normal random variable:

$$\begin{aligned} \Delta_{TT}^y &= \mathbb{E}[\alpha_i | D_i = 1] \\ &= \bar{\alpha} + \theta \mathbb{E}[\mu_i | \mu_i + U_i^B \leq \bar{y}] \\ &= \bar{\alpha} + \theta \left(\bar{\mu} - \frac{\sigma_\mu^2}{\sqrt{\sigma_\mu^2 + \sigma_U^2}} \frac{\phi\left(\frac{\bar{y} - \bar{\mu}}{\sqrt{\sigma_\mu^2 + \sigma_U^2}}\right)}{\Phi\left(\frac{\bar{y} - \bar{\mu}}{\sqrt{\sigma_\mu^2 + \sigma_U^2}}\right)} \right) \\ &= \bar{\alpha} + \theta \bar{\mu} - \theta \left(\frac{\sigma_\mu^2}{\sqrt{\sigma_\mu^2 + \sigma_U^2}} \frac{\phi\left(\frac{\bar{y} - \bar{\mu}}{\sqrt{\sigma_\mu^2 + \sigma_U^2}}\right)}{\Phi\left(\frac{\bar{y} - \bar{\mu}}{\sqrt{\sigma_\mu^2 + \sigma_U^2}}\right)} \right), \end{aligned}$$

where ϕ and Φ are respectively the density and the cumulative distribution functions of the standard normal. The second equality follows from the definition of α_i and D_i and from the fact that η_i is independent from μ_i and U_i^B . The third equality comes from the formula for the expectation of a censored bivariate normal random variable. In order to compute the population value of TT easily for different sets of parameter values, let's write a function in R:

```
delta.y.tt <- function(param){return(param["baralpha"]+param["theta"]*param["barmu"]
                                     -param["theta"]*(param["sigma2mu"]*dnorm((log(param["barY"])-param
                                     /sqrt(param["sigma2mu"]+param["sigma2U"])))
                                     *pnorm((log(param["barY"])-param["barmu"])/(sqrt
```

The population value of TT computed using this function is: $\Delta_{TT}^y = \text{'r round(delta.y.tt(param),3)}$. We can see that the values of TT in the sample and in the population differ slightly. This is because of sampling noise: the units in the sample are not perfectly representative of the units in the population.

1.3 Fundamental problem of causal inference

At least in this lecture, causal inference is about trying to infer TT, either in the sample or in the population. The FPCI states that it is impossible to directly observe TT because one part of it remains fundamentally unobserved.

Theorem 1.1 (Fundamental problem of causal inference). *It is impossible to observe TT, either in the population or in the sample.*

Proof. The proof of the FPCI is rather straightforward. Let me start with the sample TT:

$$\begin{aligned}\Delta_{TT_s}^Y &= \frac{1}{\sum_{i=1}^N D_i} \sum_{i=1}^N (Y_i^1 - Y_i^0) D_i \\ &= \frac{1}{\sum_{i=1}^N D_i} \sum_{i=1}^N Y_i^1 D_i - \frac{1}{\sum_{i=1}^N D_i} \sum_{i=1}^N Y_i^0 D_i \\ &= \frac{1}{\sum_{i=1}^N D_i} \sum_{i=1}^N Y_i D_i - \frac{1}{\sum_{i=1}^N D_i} \sum_{i=1}^N Y_i^0 D_i.\end{aligned}$$

Since Y_i^0 is unobserved whenever $D_i = 1$, $\frac{1}{\sum_{i=1}^N D_i} \sum_{i=1}^N Y_i^0 D_i$ is unobserved, and so is $\Delta_{TT_s}^Y$. The same is true for the population TT:

$$\begin{aligned}\Delta_{TT}^Y &= \mathbb{E}[Y_i^1 - Y_i^0 | D_i = 1] \\ &= \mathbb{E}[Y_i^1 | D_i = 1] - \mathbb{E}[Y_i^0 | D_i = 1] \\ &= \mathbb{E}[Y_i | D_i = 1] - \mathbb{E}[Y_i^0 | D_i = 1].\end{aligned}$$

$\mathbb{E}[Y_i^0 | D_i = 1]$ is unobserved, and so is Δ_{TT}^Y . □

The key insight in order to understand the FPCI is to see that the outcomes of the treated units had they not been treated are unobservable, and so is their average or expectation. We say that they are counterfactual, contrary to what has happened.

Definition 1.4 (Counterfactual). Both $\frac{1}{\sum_{i=1}^N D_i} \sum_{i=1}^N Y_i^0 D_i$ and $\mathbb{E}[Y_i^0 | D_i = 1]$ are counterfactual quantities that we will never get to observe.

Example 1.8. The average counterfactual outcome of the treated is the mean of the red circles in the y axis on Figure 1.4:

$$\frac{1}{\sum_{i=1}^N D_i} \sum_{i=1}^N y_i^0 D_i = 6.91.$$

Remember that we can estimate this quantity only because we have generated the data ourselves. In real life, this quantity is hopelessly unobserved.

$\mathbb{E}[y_i^0 | D_i = 1]$ can be computed using the formula for the expectation of a censored normal random variable:

$$\begin{aligned}
\mathbb{E}[y_i^0 | D_i = 1] &= \mathbb{E}[\mu_i + \delta + U_i^0 | D_i = 1] \\
&= \mathbb{E}[\mu_i + \delta + \rho U_i^B + \epsilon_i | D_i = 1] \\
&= \delta + \mathbb{E}[\mu_i + \rho U_i^B | y_i^B \leq \bar{y}] \\
&= \delta + \bar{\mu} - \frac{\sigma_\mu^2 + \rho\sigma_U^2}{\sqrt{\sigma_\mu^2 + \sigma_U^2}} \frac{\phi\left(\frac{\bar{y} - \bar{\mu}}{\sqrt{\sigma_\mu^2 + \sigma_U^2}}\right)}{\Phi\left(\frac{\bar{y} - \bar{\mu}}{\sqrt{\sigma_\mu^2 + \sigma_U^2}}\right)}.
\end{aligned}$$

We can write a function in R to compute this value:

```
esp.y0.D1 <- function(param){
  return(param["delta"]+param["barmu"]
    -((param["sigma2mu"]+param["rho"]*param["sigma2U"])
      *dnorm((log(param["barY"])-param["barmu"])/(sqrt(param["sigma2mu"]+param["sigma2U"]))))
    /(sqrt(param["sigma2mu"]+param["sigma2U"])*pnorm((log(param["barY"])-param["barmu"])/
      (sqrt(param["sigma2mu"]+param["sigma2U"]))))))
}
```

The population value of TT computed using this function is: $\mathbb{E}[y_i^0 | D_i = 1] = 6.9$.

1.4 Intuitive estimators, confounding factors and selection bias

In this section, we are going to examine the properties of two intuitive comparisons that laypeople, policymakers but also ourselves make in order to estimate causal effects: the with/without comparison (WW) and the before/after comparison (BA). WW compares the average outcomes of the treated individuals with those of the untreated individuals. BA compares the average outcomes of the treated after taking the treatment to their average outcomes before they took the treatment. These comparisons try to proxy for the expected counterfactual outcome in the treated group by using an observed quantity. WW uses the expected outcome of the untreated individuals as a proxy. BA uses the expected outcome of the treated before they take the treatment as a proxy.

Unfortunately, both of these proxies are generally poor and provide biased estimates of TT . The reason that these proxies are poor is that the treatment is not the only factor that differentiates the treated group from the groups used to form the proxy. The intuitive comparisons are biased because factors, other than the treatment, are correlated to its allocation. The factors that bias the intuitive comparisons are generally called confounding factors or confounders.

The treatment effect measures the effect of a ceteris paribus change in treatment status, while the intuitive comparisons capture both the effect of this change and that of other correlated changes that spuriously contaminate the comparison. Intuitive comparisons measure correlations while treatment effects measure causality. The old motto “correlation is not causation” applies vehemently here.

Remark. A funny anecdote about this expression “correlation is not causation”. This expression is due to Karl Pearson, the father of modern statistics. He coined the phrase in his famous book “The Grammar of Science.” Pearson is famous for inventing the correlation coefficient. He actually thought that correlation was a much superior, much more rigorous term, than causation. In his book, he actually used the sentence to argue in favor of abandoning causation altogether and focusing on the much better-defined and measurable concept of correlation. Interesting turn of events that his sentence is now used to mean that correlation is weaker than causation, totally reverting the original intended meaning.

In this section, we are going to define both comparisons, study their biases and state the conditions under which they identify TT . This will prove to be a very useful introduction to the notion of identification. It is

also very important to be able to understand the sources of bias of comparisons that we use every day and that come very naturally to policy makers and lay people.

Remark. In this section, we state the definitions and formulae in the population. This is for two reasons. First, it is simpler, and lighter in terms of notation. Second, it emphasizes that the problems with intuitive comparisons are independent of sampling noise. Most of the results stated here for the population extend to the sample, replacing the expectation operator by the average operator. I will nevertheless give examples in the sample, since it is so much simpler to compute. I will denote sample equivalents of population estimators with a hat.

1.4.1 The with/without comparison, selection bias and cross-sectional confounders

The with/without comparison (*WW*) is very intuitive: just compare the outcomes of the treated and untreated individuals in order to estimate the causal effect. This approach is nevertheless generally biased. We call the bias of *WW* selection bias (*SB*). Selection bias is due to unobserved confounders that are distributed differently in the treatment and control group and that generate differences in outcomes even in the absence of the treatment. In this section, I define the *WW* estimator, derives its bias, introduces the confounders and states conditions under which it is unbiased.

1.4.1.1 The with/without comparison

The with/without comparison (*WW*) is very intuitive: just compare the outcomes of the treated and untreated individuals in order to estimate the causal effect.

Definition 1.5 (With/without comparison). The with/without comparison is the difference between the expected outcomes of the treated and the expected outcomes of the untreated:

$$\Delta_{WW}^Y = \mathbb{E}[Y_i | D_i = 1] - \mathbb{E}[Y_i | D_i = 0].$$

Example 1.9. In the population, *WW* can be computed using the traditional formula for the expectation of a truncated normal distribution:

$$\begin{aligned} \Delta_{WW}^y &= \mathbb{E}[y_i | D_i = 1] - \mathbb{E}[y_i | D_i = 0] \\ &= \mathbb{E}[y_i^1 | D_i = 1] - \mathbb{E}[y_i^0 | D_i = 0] \\ &= \mathbb{E}[\alpha_i | D_i = 1] + \mathbb{E}[\mu_i + \rho U_i^B | \mu_i + U_i^B \leq \bar{y}] - \mathbb{E}[\mu_i + \rho U_i^B | \mu_i + U_i^B > \bar{y}] \\ &= \bar{\alpha} + \theta \left(\bar{\mu} - \frac{\sigma_\mu^2}{\sqrt{\sigma_\mu^2 + \sigma_U^2}} \frac{\phi\left(\frac{\bar{y} - \bar{\mu}}{\sqrt{\sigma_\mu^2 + \sigma_U^2}}\right)}{\Phi\left(\frac{\bar{y} - \bar{\mu}}{\sqrt{\sigma_\mu^2 + \sigma_U^2}}\right)} \right) - \frac{\sigma_\mu^2 + \rho\sigma_U^2}{\sqrt{\sigma_\mu^2 + \sigma_U^2}} \left(\frac{\phi\left(\frac{\bar{y} - \bar{\mu}}{\sqrt{\sigma_\mu^2 + \sigma_U^2}}\right)}{\Phi\left(\frac{\bar{y} - \bar{\mu}}{\sqrt{\sigma_\mu^2 + \sigma_U^2}}\right)} + \frac{\phi\left(\frac{\bar{y} - \bar{\mu}}{\sqrt{\sigma_\mu^2 + \sigma_U^2}}\right)}{1 - \Phi\left(\frac{\bar{y} - \bar{\mu}}{\sqrt{\sigma_\mu^2 + \sigma_U^2}}\right)} \right). \end{aligned}$$

In order to compute this parameter, we are going to set up a R function. For reasons that will become clearer later, we will define two separate functions to compute the first and second part of the formula. In the first part, you should have recognised *TT*, that we have already computed in Lecture 1. We are going to call the second part *SB*, for reasons that will become explicit in a bit.

```
delta.y.tt <- function(param){
  return(param["baralpha"]+param["theta"]*param["barmu"]-param["theta"]
    *((param["sigma2mu"]*dnorm((log(param["barY"])-param["barmu"])/
      (sqrt(param["sigma2mu"]+param["sigma2U"]))))
    /(sqrt(param["sigma2mu"]+param["sigma2U"]))*pnorm((log(param["barY"])-param["barmu"])
```

```

                                                    /(sqrt(param["sigma2mu"]+param["sigma2U"])))
}
delta.y.sb <- function(param){
  return(-(param["sigma2mu"]+param["rho"]*param["sigma2U"])/sqrt(param["sigma2mu"]+param["sigma2U"]))
  *dnorm((log(param["barY"])-param["barmu"])/(sqrt(param["sigma2mu"]+param["sigma2U"])))
  *(1/pnorm((log(param["barY"])-param["barmu"])/(sqrt(param["sigma2mu"]+param["sigma2U"])))
  +1/(1-pnorm((log(param["barY"])-param["barmu"])/(sqrt(param["sigma2mu"]+param["sigma2U"]))))))
}
delta.y.ww <- function(param){
  return(delta.y.tt(param)+delta.y.sb(param))
}

```

As a conclusion of all these derivations, WW in the population is equal to -1.298. Remember that the value of TT in the population is 0.172.

In order to compute the WW estimator in a sample, I'm going to generate a brand new sample and I'm going to choose a seed for the pseudo-random number generator so that we obtain the same result each time we run the code. I use `set.seed(1234)` in the code chunk below.

```

param <- c(8,.5,.28,1500)
names(param) <- c("barmu","sigma2mu","sigma2U","barY")
set.seed(1234)
N <- 1000
mu <- rnorm(N,param["barmu"],sqrt(param["sigma2mu"]))
UB <- rnorm(N,0,sqrt(param["sigma2U"]))
yB <- mu + UB
YB <- exp(yB)
Ds <- rep(0,N)
Ds[YB<=param["barY"]] <- 1
l <- length(param)
param <- c(param,0.9,0.01,0.05,0.05,0.05,0.1)
names(param)[(l+1):length(param)] <- c("rho","theta","sigma2epsilon","sigma2eta","delta","baralpha")
epsilon <- rnorm(N,0,sqrt(param["sigma2epsilon"]))
eta<- rnorm(N,0,sqrt(param["sigma2eta"]))
U0 <- param["rho"]*UB + epsilon
y0 <- mu + U0 + param["delta"]
alpha <- param["baralpha"]+ param["theta"]*mu + eta
y1 <- y0+alpha
Y0 <- exp(y0)
Y1 <- exp(y1)
y <- y1*Ds+y0*(1-Ds)
Y <- Y1*Ds+Y0*(1-Ds)

```

In this sample, the average outcome of the treated in the presence of the treatment is

$$\frac{1}{\sum_{i=1}^N D_i} \sum_{i=1}^N D_i y_i = 7.074.$$

It is materialized by a circle on Figure 1.6. The average outcome of the untreated is

$$\frac{1}{\sum_{i=1}^N (1 - D_i)} \sum_{i=1}^N (1 - D_i) y_i = 8.383.$$

It is materialized by a plus sign on Figure 1.6.

The estimate of the WW comparison in the sample is thus:

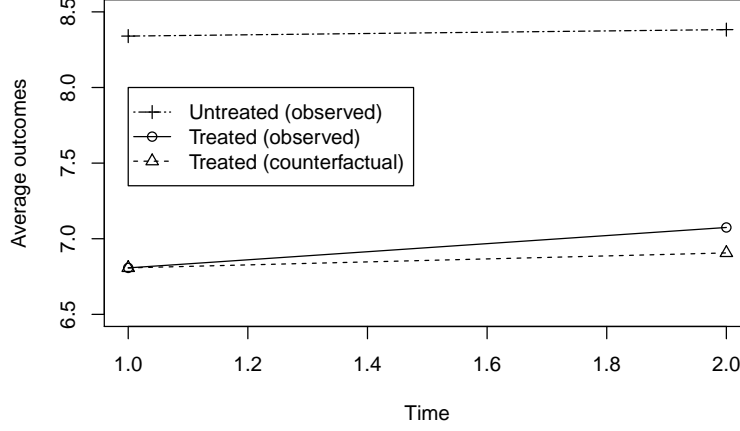


Figure 1.6: Evolution of average outcomes in the treated and control group before (Time =1) and after (Time=2) the treatment

$$\Delta_{WW}^{\hat{y}} = \frac{1}{\sum_{i=1}^N D_i} \sum_{i=1}^N Y_i D_i - \frac{1}{\sum_{i=1}^N (1 - D_i)} \sum_{i=1}^N Y_i (1 - D_i).$$

We have $\Delta_{WW}^{\hat{y}} = -1.308$. Remember that the value of TT in the sample is $\Delta_{TT_s}^y = 0.168$.

Overall, WW severely underestimates the effect of the treatment in our example. WW suggests that the treatment has a negative effect on outcomes whereas we know by construction that it has a positive one.

1.4.1.2 Selection bias

When we form the with/without comparison, we do not recover the TT parameter. Instead, we recover TT plus a bias term, called **selection bias**:

$$\Delta_{WW}^Y = \Delta_{TT}^Y + \Delta_{SB}^Y.$$

Definition 1.6 (Selection bias). Selection bias is the difference between the with/without comparison and the treatment on the treated parameter:

$$\Delta_{SB}^Y = \Delta_{WW}^Y - \Delta_{TT}^Y.$$

WW tries to approximate the counterfactual expected outcome in the treated group by using $\mathbb{E}[Y_i^0 | D_i = 0]$, the expected outcome in the untreated group. Selection bias appears because this proxy is generally poor. It is very easy to see that selection bias is indeed directly due to this bad proxy problem:

Theorem 1.2. *Selection bias is the difference between the counterfactual expected potential outcome in the absence of the treatment among the treated and the expected potential outcome in the absence of the treatment among the untreated.*

$$\Delta_{SB}^Y = \mathbb{E}[Y_i^0 | D_i = 1] - \mathbb{E}[Y_i^0 | D_i = 0].$$

Proof.

$$\begin{aligned}\Delta_{SB}^Y &= \Delta_{WW}^Y - \Delta_{TT}^Y \\ &= \mathbb{E}[Y_i|D_i = 1] - \mathbb{E}[Y_i|D_i = 0] - \mathbb{E}[Y_i^1 - Y_i^0|D_i = 1] \\ &= \mathbb{E}[Y_i^0|D_i = 1] - \mathbb{E}[Y_i^0|D_i = 0].\end{aligned}$$

The first and second equalities stem only from the definition of both parameters. The third equality stems from using the switching equation: $Y_i = Y_i^1 D_i + Y_i^0 (1 - D_i)$, so that $\mathbb{E}[Y_i|D_i = 1] = \mathbb{E}[Y_i^1|D_i = 1]$ and $\mathbb{E}[Y_i|D_i = 0] = \mathbb{E}[Y_i^0|D_i = 0]$. \square

Example 1.10. In the population, SB is equal to

$$\Delta_{SB}^y = \Delta_{WW}^y - \Delta_{TT}^y = -1.298 - 0.172 = \text{'round}(\text{delta.y.ww}(\text{param}) - \text{delta.y.tt}(\text{param}), 3)\text{'}$$

We could have computed SB directly using the formula from Theorem~??:

$$\begin{aligned}\Delta_{SB}^y &= \mathbb{E}[y_i^0|D_i = 1] - \mathbb{E}[y_i^0|D_i = 0] \\ &= -\frac{\sigma_\mu^2 + \rho\sigma_U^2}{\sqrt{\sigma_\mu^2 + \sigma_U^2}} \left(\frac{\phi\left(\frac{\bar{y}-\bar{\mu}}{\sqrt{\sigma_\mu^2 + \sigma_U^2}}\right)}{\Phi\left(\frac{\bar{y}-\bar{\mu}}{\sqrt{\sigma_\mu^2 + \sigma_U^2}}\right)} + \frac{\phi\left(\frac{\bar{y}-\bar{\mu}}{\sqrt{\sigma_\mu^2 + \sigma_U^2}}\right)}{1 - \Phi\left(\frac{\bar{y}-\bar{\mu}}{\sqrt{\sigma_\mu^2 + \sigma_U^2}}\right)} \right).\end{aligned}$$

When using the R function for SB that we have defined earlier, we indeed find: $\Delta_{SB}^y = -1.471$.

In the sample, $\Delta_{SB}^{\hat{y}} = -1.308 - 0.168 = -1.476$. Selection bias emerges because we are using a bad proxy for the counterfactual. The average outcome for the untreated is equal to $\frac{1}{\sum_{i=1}^N (1-D_i)} \sum_{i=1}^N (1-D_i)y_i = 8.383$ while the counterfactual average outcome for the treated is $\frac{1}{\sum_{i=1}^N D_i} \sum_{i=1}^N D_i y_i^0 = 6.906$. Their difference is

as expected equal to SB : $\Delta_{SB}^{\hat{y}} = 6.906 - 8.383 = -1.476$. The counterfactual average outcome of the treated is much smaller than the average outcome of the untreated. On Figure 1.6, this is materialized by the fact that the plus sign is located much above the triangle.

Remark. The concept of selection bias is related to but different from the concept of sample selection bias. With sample selection bias, we worry that selection into the sample might bias the estimated effect of a treatment on outcomes. With selection bias, we worry that selection into the treatment itself might bias the effect of the treatment on outcomes. Both biases are due to unobserved covariates, but they do not play out in the same way.

For example, estimating the effect of education on women's wages raises both selection bias and sample selection bias issues. Selection bias stems from the fact that more educated women are more likely to be more dynamic and thus to have higher earnings even when less educated. Selection bias would be positive in that case, overestimating the effect of education on earnings.

Sample selection bias stems from the fact that we can only use a sample of working women in order to estimate the effect of education on wages, since we do not observe the wages on non working women. But, selection into the labor force might generate sample selection bias. More educated women participate more in the labor market, while less educated women participate less. As a consequence, less educated women that work are different from the overall sample of less educated women. They might be more dynamic and work-focused. As a consequence, their wages are higher than the average wages of the less educated women. Comparing the wages of less educated women that work to those of more educated women that work might understate the effect of education on earnings. Sample selection bias would generate a negative bias on the education coefficient.