

Problem Set #1

Experiments and Causality

January 19, 2019

1. Potential Outcomes Notation

- Explain the notation $Y_i(1)$.

It's the potential outcome for treatment of the i^{th} observation.

- Explain the notation $E[Y_i(1)|d_i = 0]$.

It's the expectation of the potential outcome for treatment of the i^{th} observation, given that the i^{th} observation is in the control group.

- Explain the difference between the notation $E[Y_i(1)]$ and the notation $E[Y_i(1)|d_i = 1]$. (Extra credit)

The first is the expectation of the potential outcome for treatment. The second is the expectation of the potential outcome for treatment given that the observation is in the treatment group. Given random assignment, these two numbers should be identical. However, selection bias can cause the two numbers to differ.

- Explain the difference between the notation $E[Y_i(1)|d_i = 1]$ and the notation $E[Y_i(1)|D_i = 1]$. Use exercise 2.7 from FE to give a concrete example of the difference.

2. Potential Outcomes Practice

Use the values in the following table to illustrate that $E[Y_i(1)] - E[Y_i(0)] = E[Y_i(1) - Y_i(0)]$.

	$Y_i(0)$	$Y_i(1)$	τ_i
Individual 1	5	6	1
Individual 2	3	8	5
Individual 3	10	12	2
Individual 4	5	5	0
Individual 5	10	8	-2

$$E[Y_i(1)] = \frac{5 + 3 + 10 + 5 + 10}{5} = 6.6$$

$$E[Y_i(0)] = \frac{6 + 8 + 12 + 5 + 8}{5} = 7.8$$

$$E[Y_i(1) - Y_i(0)] = \frac{1 + 5 + 2 + 0 - 2}{5} = 1.2$$

$$E[Y_i(1)] - E[Y_i(0)] = E[Y_i(1) - Y_i(0)]$$

3. Conditional Expectations

Consider the following table:

	$Y_i(0)$	$Y_i(1)$	τ_i
Individual 1	10	15	5
Individual 2	15	15	0
Individual 3	20	30	10
Individual 4	20	15	-5
Individual 5	10	20	10
Individual 6	15	15	0
Individual 7	15	30	15
Average	15	20	5

Use the values depicted in the table above to complete the table below.

$Y_i(0)$	15	20	30	Marginal $Y_i(0)$
10	n: 1 %: 0.14	n: 1 %: 0.14	n: 0 %: 0.00	0.29
15	n: 2 %: 0.29	n: 0 %: 0.00	n: 1 %: 0.14	0.43
20	n: 1 %: 0.14	n: 0 %: 0.00	n: 1 %: 0.14	0.29
Marginal $Y_i(1)$	0.57	0.14	0.29	1.0

- Fill in the number of observations in each of the nine cells;
- Indicate the percentage of all subjects that fall into each of the nine cells.
- At the bottom of the table, indicate the proportion of subjects falling into each category of $Y_i(1)$.
- At the right of the table, indicate the proportion of subjects falling into each category of $Y_i(0)$.
- Use the table to calculate the conditional expectation that $E[Y_i(0)|Y_i(1) > 15]$.

$$E[Y_i(0)|Y_i(1) > 15] = \frac{10 + 15 + 20}{3} = 15$$

- Use the table to calculate the conditional expectation that $E[Y_i(1)|Y_i(0) > 15]$.

$$E[Y_i(1)|Y_i(0) > 15] = \frac{15 + 30}{2} = 22.5$$

4. More Practice with Potential Outcomes

Suppose we are interested in the hypothesis that children playing outside leads them to have better eyesight.

Consider the following population of ten representative children whose visual acuity we can measure. (Visual acuity is the decimal version of the fraction given as output in standard eye exams. Someone with 20/20 vision has acuity 1.0, while someone with 20/40 vision has acuity 0.5. Numbers greater than 1.0 are possible for people with better than “normal” visual acuity.)

child	y0	y1
1	1.1	1.1
2	0.1	0.6
3	0.5	0.5
4	0.9	0.9
5	1.6	0.7
6	2.0	2.0
7	1.2	1.2
8	0.7	0.7
9	1.0	1.0
10	1.1	1.1

In the table, state $Y_i(1)$ means “playing outside an average of at least 10 hours per week from age 3 to age 6,” and state $Y_i(0)$ means “playing outside an average of less than 10 hours per week from age 3 to age 6.” Y_i represents visual acuity measured at age 6.

- Compute the individual treatment effect for each of the ten children. Note that this is only possible because we are working with hypothetical potential outcomes; we could never have this much information with real-world data. (We encourage the use of computing tools on all problems, but please describe your work so that we can determine whether you are using the correct values.)

Just take the difference between y1 and y2.

```
answer.P0a <- d$y1 - d$y0  
answer.P0a
```

```
## [1] 0.0 0.5 0.0 0.0 -0.9 0.0 0.0 0.0 0.0 0.0
```

- b. In a single paragraph, tell a story that could explain this distribution of treatment effects.

The treatment effect is zero for most of the children. By random chance, a few of the children's acuities are affected by playing outside.

- c. What might cause some children to have different treatment effects than others?

Each child responds differently to being outside. Most of them receive no effect to their acuities. A positive treatment effect might be one leading from less time in front a computer screen; a negative effect, from harmful sunlight.

- d. For this population, what is the true average treatment effect (ATE) of playing outside.

Just take the mean of the previous answer.

```
answer.P0d <- mean(answer.P0a)  
answer.P0d
```

```
## [1] -0.04
```

- e. Suppose we are able to do an experiment in which we can control the amount of time that these children play outside for three years. We happen to randomly assign the odd-numbered children to treatment and the even-numbered children to control. What is the estimate of the ATE you would reach under this assignment? (Again, please describe your work.)

Assign the children to treatment and control. Take the mean of y1 of the former minus the y0 of the latter.

```
library(data.table)  
dt = data.table(d)  
answer.P0e <- mean(dt[child %% 2 == 1, y1] - dt[child %% 2 == 0, y0])  
answer.P0e
```

```
## [1] -0.06
```

- f. How different is the estimate from the truth? Intuitively, why is there a difference?

The estimate is off by .02. This difference is just due to random chance.

- g. We just considered one way (odd-even) an experiment might split the children. How many different ways (every possible way) are there to split the children into a treatment versus a control group (assuming at least one person is always in the treatment group and at least one person is always in the control group)?

This is just the sum of the combinations, N choose i, for i from 1 to N-1.

$$\sum_{i=1}^{N-1} \binom{N}{i}$$

```
answer.P0g <- sum(choose(10, 1:9))  
answer.P0g
```

```
## [1] 1022
```

- h. Suppose that we decide it is too hard to control the behavior of the children, so we do an observational study instead. Children 1-5 choose to play an average of more than 10 hours per week from age 3 to

age 6, while Children 6-10 play less than 10 hours per week. Compute the difference in means from the resulting observational data.

```
answer.P0h <- mean(dt[child <= 5, y1] - dt[child >= 6, y0])
answer.P0h
```

```
## [1] -0.44
```

- i. Compare your answer in (h) to the true ATE. Intuitively, what causes the difference?

The children who chose to play an average of more than 10 hours might do so because they like to play outside, and vice versa. One could envision that children whose eyesights are damaged by being outside might choose to stay inside. The converse is probably true as well. Therefore we see much larger effect here.

5. Randomization and Experiments

Suppose that a reasearcher wants to investigate whether after-school math programs improve grades. The researcher randomly samples a group of students from an elementary school and then compare the grades between the group of students who are enrolled in an after-school math program to those who do not attend any such program. Is this an experiment or an observational study? Why?

This is an observational study. There was no intervention made on the students sampled from the school. The researcher had no say on whether those students went to an after-school program or not. It is easy to see that students who attend such programs voluntarily (or are forced by their parents) can have different grades than their peers, regardless of the effectiveness of the programs.

6. Lotteries

A researcher wants to know how winning large sums of money in a national lottery affect people's views about the estate tax. The researcher interviews a random sample of adults and compares the attitudes of those who report winning more than \$10,000 in the lottery to those who claim to have won little or nothing. The researcher reasons that the lottery choose winners at random, and therefore the amount that people report having won is random.

- a. Critically evaluate this assumption.

The assumption is reasonable if all the interviewed people played the lottery with roughly equal frequency, even those who did not win anything. The lottery randomly selects winners' sums within the subset of people who play, but does nothing to the ones who do not. It is possible that lottery players and non-players have different views on estate taxes but winning large sums of money has no effect.

- b. Suppose the researcher were to restrict the sample to people who had played the lottery at least once during the past year. Is it safe to assume that the potential outcomes of those who report winning more than \$10,000 are identical, in expectation, to those who report winning little or nothing?

This assumption is still questionable if the chances of winning larges sums is significantly increased by people who buy numerous tickets. In practice, it may be infeasible for anyone to buy enough lottery tickets to significantly boost thier chances of winning while maintaining positive gains. However, it is still worth controlling for the frequency at which people play.

Clarifications

1. Please think of the outcome variable as an individual's answer to the survey question "Are you in favor of raising the estate tax rate in the United States?"
2. The hint about potential outcomes could be rewritten as follows: Do you think those who won the lottery would have had the same views about the estate tax if they had actually not won it as those who actually did not win it? (That is, is $E[Y_i(0)|D = 1] = E[Y_i(0)|D = 0]$, comparing what would have happened to the actual winners, the $|D = 1$ part, if they had not won, the $Y_i(0)$ part, and what actually happened to those who did not win, the $Y_i(0)|D = 0$ part.) In general, it is just another way of asking, "are those who win the lottery and those who have not won the lottery comparable?"
3. Assume lottery winnings are always observed accurately and there are no concerns about under- or over-reporting.

7. Inmates and Reading

A researcher studying 1,000 prison inmates noticed that prisoners who spend at least 3 hours per day reading are less likely to have violent encounters with prison staff. The researcher recommends that all prisoners be required to spend at least three hours reading each day. Let d_i be 0 when prisoners read less than three hours each day and 1 when they read more than three hours each day. Let $Y_i(0)$ be each prisoner's PO of violent encounters with prison staff when reading less than three hours per day, and let $Y_i(1)$ be their PO of violent encounters when reading more than three hours per day.

In this study, nature has assigned a particular realization of d_i to each subject. When assessing this study, why might one be hesitant to assume that $E[Y_i(0)|D_i = 0] = E[Y_i(0)|D_i = 1]$ and $E[Y_i(1)|D_i = 0] = E[Y_i(1)|D_i = 1]$? In your answer, give some intuitive explanation in English for what the mathematical expressions mean.

The prisoners who were already reading 3 hours per day may have some other reason to be less violently inclined. This selection bias causes $E[Y_i(0)|D_i = 0] \neq E[Y_i(0)|D_i = 1]$