

Problem Set 1

Alex, Micah, Scott, and David

12/01/2021

Contents

1	Potential Outcomes Notation	2
1.1	Explain the notation $Y_i(1)$.	2
1.2	Explain the notation $Y_1(1)$.	2
1.3	Explain the notation $E[Y_i(1) d_i = 0]$.	2
1.4	Explain the difference between the notation $E[Y_i(1)]$ and $E[Y_i(1) d_i = 1]$	2
2	Potential Outcomes and Treatment Effects	3
2.1	Illustration	3
2.2	Data Possibilities	3
3	Visual Acuity	4
3.1	Treatment effect	4
3.2	Story time	4
3.3	True ATE	5
3.4	Even-Odd split	5
3.5	Biased or Unbiased?	6
3.6	How many splits are possible?	6
3.7	Observational study	7
3.8	Observational ATE	7
4	Randomization and Experiments	8
4.1	Define your terms	8
4.2	Does a random, iid sample produce an unbiased treatment effect estimate?	8
4.3	What if an official agency produces the iid sample?	8
4.4	What if someone else randomly assigns	8
5	Moral Panic	9
5.1	Explain the statements	9
5.2	Can you believe it	9

1 Potential Outcomes Notation

1.1 Explain the notation $Y_i(1)$.

Answer: it is the potential outcome if the i th subject were treated.

1.2 Explain the notation $Y_1(1)$.

Answer: potential outcome if subject “1” were treated.

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

Answer: the expectation of treated potential outcome for a subject that does not receive treatment.

1.4 Explain the difference between the notation $E[Y_i(1)]$ and $E[Y_i(1)|d_i = 1]$

Answer: ... While $Y_i(1)$ is the potential outcome if the i th subject were treated from the entire dataset, $E[Y_i(1)|d_i = 1]$ is the expectation of treated potential outcome for a subject that does receive treatment.

2 Potential Outcomes and Treatment Effects

2.1 Illustration

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

```
left_side_of_equation <- table[, mean(y_1) - mean(y_0)]
```

```
right_side_of_equation <- table[, mean(y_1 - y_0)]
```

```
left_side_of_equation
```

```
## [1] 2
```

```
right_side_of_equation
```

```
## [1] 2
```

Answer: I demonstrated above in the code that $E[Y_i(1)] - E[Y_i(0)] = E[Y_i(1) - Y_i(0)] = 2$ in our table above.

2.2 Data Possibilities

Is it possible to collect all necessary values and construct a table like the one provided in real life? Explain why or why not?

Answer: No, in real life we cannot obtain all the values of y_1 and y_0 for all data points. This is because in reality our data points are in either y_0 group or y_1 group, but not both. If in y_0 then y_1 remains counterfactual and vice versa.

3 Visual Acuity

3.1 Treatment effect

Compute the individual treatment effect for each of the ten children.

```
d[, treatment_effect := y_1 - y_0]  
d
```

```
##      child y_0 y_1 treatment_effect  
## 1:      1 1.2 1.2              0.0  
## 2:      2 0.1 0.7              0.6  
## 3:      3 0.5 0.5              0.0  
## 4:      4 0.8 0.8              0.0  
## 5:      5 1.5 0.6             -0.9  
## 6:      6 2.0 2.0              0.0  
## 7:      7 1.3 1.3              0.0  
## 8:      8 0.7 0.7              0.0  
## 9:      9 1.1 1.1              0.0  
## 10:     10 1.4 1.4              0.0
```

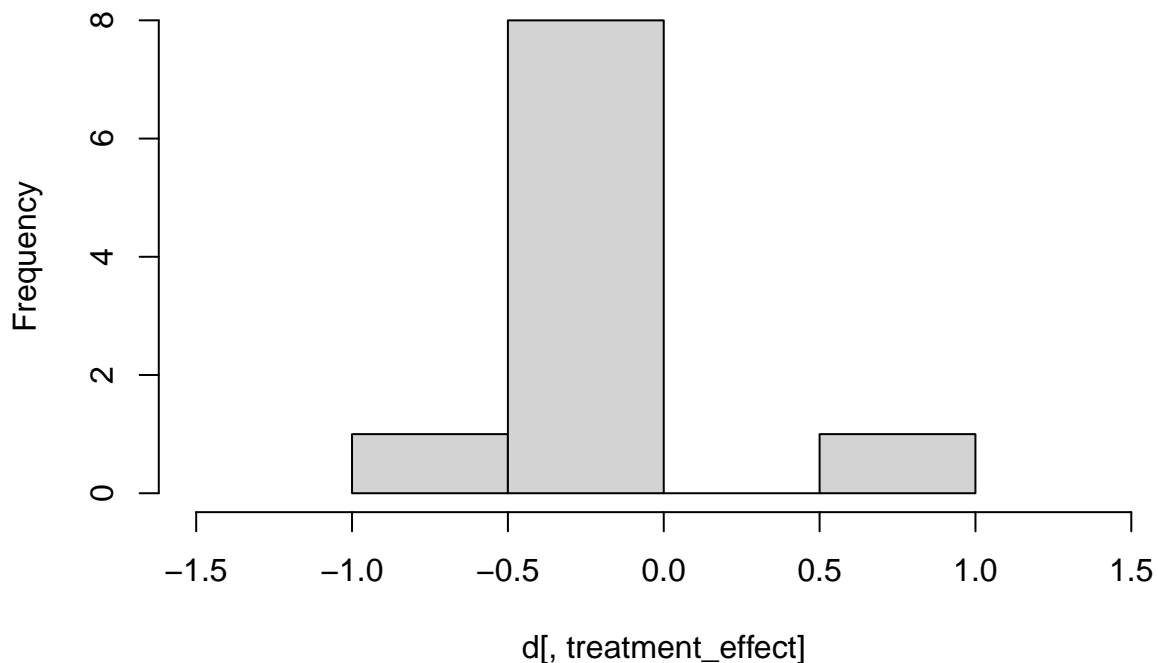
Answer: ...

3.2 Story time

Tell a “story” that could explain this distribution of treatment effects. In particular, discuss what might cause some children to have different treatment effects than others.

```
hist(d[, treatment_effect], xlim = c(-1.5, 1.5), )
```

Histogram of d[, treatment_effect]



Answer: From the table above, we can see that only two subjects had non-zero treatment effects (0.6 an increase in acuity and -0.9 a decrease in acuity). This could be due to confounding variables that were not accounted for. Moreover, since the non-treatment effects are at the edges of the average treatment effect value, this could show that there may have been other sorts of errors in data entry or data collection. Since the average treatment effect is close to 0, excluding the two outliers, we fail to show a causal inference that playing outside more than 10 hours/week between ages of 3-6 increases acuity.

3.3 True ATE

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

```
mean(d[,treatment_effect])
```

```
## [1] -0.03
```

Answer: ATE is -0.03

3.4 Even-Odd split

Suppose we are able to do an experiment in which we can control the amount of time 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? (Please describe your work.)

```
#creating treatment column based on the criteria described above.
```

```
d[, treatment := child %% 2 == 1,] #odd
d
```

```
##      child y_0 y_1 treatment_effect treatment
## 1:      1 1.2 1.2           0.0      TRUE
## 2:      2 0.1 0.7           0.6     FALSE
## 3:      3 0.5 0.5           0.0      TRUE
## 4:      4 0.8 0.8           0.0     FALSE
## 5:      5 1.5 0.6          -0.9      TRUE
## 6:      6 2.0 2.0           0.0     FALSE
## 7:      7 1.3 1.3           0.0      TRUE
## 8:      8 0.7 0.7           0.0     FALSE
## 9:      9 1.1 1.1           0.0      TRUE
## 10:     10 1.4 1.4           0.0     FALSE
```

```
#creating a new column name acuity_level based on treatment column created above.
```

```
d[, acuity_level := ifelse(test = treatment == TRUE, yes = y_1, no = y_0)]
d
```

```
##      child y_0 y_1 treatment_effect treatment acuity_level
## 1:      1 1.2 1.2           0.0      TRUE      1.2
## 2:      2 0.1 0.7           0.6     FALSE      0.1
## 3:      3 0.5 0.5           0.0      TRUE      0.5
## 4:      4 0.8 0.8           0.0     FALSE      0.8
## 5:      5 1.5 0.6          -0.9      TRUE      0.6
```

```
## 6:      6 2.0 2.0      0.0    FALSE      2.0
## 7:      7 1.3 1.3      0.0     TRUE      1.3
## 8:      8 0.7 0.7      0.0    FALSE      0.7
## 9:      9 1.1 1.1      0.0     TRUE      1.1
## 10:     10 1.4 1.4      0.0    FALSE      1.4
```

```
#taking the mean of only control subjects
mean_control_acuity_level <- d[treatment==FALSE, mean(acuity_level)]

#taking the mean of only treatment subjects
mean_treatment_acuity_level <- d[treatment==TRUE, mean(acuity_level)]

#finding ATE
ATE <- mean_treatment_acuity_level - mean_control_acuity_level
ATE
```

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

Answer: ATE is -0.06

3.5 Biased or Unbiased?

How different is the estimate from the truth? In your own words, why is there a difference? Does this mean that the estimator is a biased or an unbiased estimator? Does this mean that the estimate is biased or unbiased?

```
# Use this code chunk to show your code work
```

Answer: The actual ATE is 50% smaller than the estimate ATE. This difference could be due to the fact that we have such small sample and it could be that there is some noise/error. A statistical t-test might help us determine if the difference is statistically significant.

An estimator is unbiased if the expected value of the estimates it produces is equal to our parameter of interest, in this case ATE. We saw that there was a difference between actual and estimated ATE given the random sampling of even and odd. If the difference is statistically significant, then we have a biased estimator, if the difference is statistically insignificant, then we have an unbiased estimator.

Assuming that the difference is caused by noise, then we can say that our estimator is unbiased and the generated result is unbiased estimates.

3.6 How many splits are possible?

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)?

```
# the formula to calculate this is as follows:

#formula: (2 to the power n) - k
#n is the number of observations
#k is the is the number of subsets, in our case treatment and control.
(2 ** 10) - 2
```

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

Answer: 1022

3.7 Observational study

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.

```
difference_in_mean <- mean(d$y_1[0:5]) - mean(d$y_0[6:10])
difference_in_mean
```

```
## [1] -0.54
```

Answer: -0.54

3.8 Observational ATE

Compare your answer in Problem 3.G to the true ATE. In your own words what causes the difference? Does this mean that the estimator is a biased or an unbiased estimator? Does this mean that the estimate is biased or unbiased?

```
# Use this code chunk to show your code work (if needed)
```

Answer: The estimate is an order of magnitude smaller than the actual ATE that we calculated above (i.e. the estimate is a much bigger negative number than actual ATE). The reason for it is that if you look at the data table, only the first 5 groups are getting a non-zero treatment effect, the last 5 (control) groups are getting zero treatment effects. Since the sample is small, the grouping has been exaggerated. That is, if the sample grows, the estimator ATE will reach the actual ATE.

In this case, the estimator is biased and the generated estimates would be biased.

4 Randomization and Experiments

The following questions can be a little bit challenging. This is because the argument that you are being asked to make is based on the rote application of a definition. To begin with, it is useful for you to define what you mean when you write about either *an experiment* or *an observational study*. Then, with these definitions on hand, use the definitions to answer the following questions.

4.1 Define your terms

- **An experiment is:** An experiment is an intervention that creates variations to teach us causal questions.
- **An experiment provides the following statistical guarantees:** An experiment is the only method that guarantees a causal inference. It allows researchers to create an intervention to manipulate the independent variable and control the confounding variables.
- **An observational study is:** An observation is not an intervention, instead, it is a measure or survey of members of a sample without trying to affect them.
- **An observational study provides the following statistical guarantees:** Observational studies do not guarantee causal inference, it does not control the dependent variables.

4.2 Does a random, iid sample produce an unbiased treatment effect estimate?

Assume that a researcher takes a random sample of elementary school children and compares the grades of those who were previously enrolled in an early childhood education program with the grades of those who were not enrolled in such a program. Is this an experiment, an observational study, or something in between?

Answer: This is an observational study.

4.3 What if an official agency produces the iid sample?

Assume that the researcher works together with an organization that provides early childhood education and offer free programs to certain children. However, which children that received this offer was not randomly selected by the researcher but rather chosen by the local government. (Assume that the government did not use random assignment but instead gives the offer to students who are deemed to need it the most) The research follows up a couple of years later by comparing the elementary school grades of students offered free early childhood education to those who were not. Is this an experiment, an observational study, or something in between? Explain!

Answer: This is a quasi experiment because the treatment and control groups were not randomly assigned.

4.4 What if someone else randomly assigns

If the government assigned students to treatment and control by “coin toss”, rather than simply sampling the population, would you say that the study is experimental or observational? Why? What, if any guarantees does this process provide?

Answer: This would be a experimental study, because coin toss has a probability of 50/50, meaning 50 percent of the sample would be assigned to treatment and 50% to control at random. This would guarantee a causal inference.

5 Moral Panic

5.1 Explain the statements

Explain the statement $E[Y_i(0)|D_i = 0] = E[Y_i(0)|D_i = 1]$ in words. First, state the rote English language translation. Second, tell us the *meaning* of this statement. A full points solution will use the term “potential outcomes” twice.

Answer: the expectation of control potential outcome for a subject that does not receive the treatment is equal to the expectation of control potential outcome for a subject that does receive the treatment.

This basically means the outcome(test score) when listening to death metal at least one time per week when the students actually listened to death metal at least one time per week would be the same as the outcome (test score) for students who listened to death metal less than one time per week.

5.2 Can you believe it

Do you expect that this circumstance actually matches with the meaning that you’ve just written down? Why or why not?

Answer: I believe that the circumstance actually matches with what I have written above. In other words, the conclusion that high school students who listen to death metal music at least once per week are more likely to perform badly on standardized test does not account for a variety of confounding variables. For example, maybe students who listen to death metal music at least once a week does not sleep well than those students who do not listen to death metal music at least once a week. In this case, sleep would be a confounding variable that may have caused low scores.