

## 2. Introduction to causality

---

Frank Edwards

9/9/2020

## Causality

---

# The key question in causal inference

- Does treatment  $x$  affect outcome  $y$

## The key question in causal inference

- Does treatment  $x$  affect outcome  $y$
- In medicine: does a treatment affect a patient

## The key question in causal inference

- Does treatment  $x$  affect outcome  $y$
- In medicine: does a treatment affect a patient
- Typically designed by randomly assigning patients to treatment and control groups, where treatment groups are exposed to  $x$ , and control groups are not

# The fundamental problem of causal inference

*How much did the treatment matter?*

# The fundamental problem of causal inference

*How much did the treatment matter?*

We answer this question with counterfactuals:

*What would have happened if treated units were untreated? What would have happened if untreated units were treated?*

# The fundamental problem of causal inference

*How much did the treatment matter?*

We answer this question with counterfactuals:

*What would have happened if treated units were untreated? What would have happened if untreated units were treated?*

For an observation  $i$ , where  $Y_i(1)$  indicates treatment and  $Y_i(0)$  indicates no treatment, the causal effect of the treatment is defined as

$$Y_i(1) - Y_i(0)$$



# The fundamental problem of causal inference

*How much did the treatment matter?*

We answer this question with counterfactuals:

*What would have happened if treated units were untreated? What would have happened if untreated units were treated?*

For an observation  $i$ , where  $Y_i(1)$  indicates treatment and  $Y_i(0)$  indicates no treatment, the causal effect of the treatment is defined as

$$Y_i(1) - Y_i(0)$$

Why is this a problematic definition?

## Causal questions in social science

- Does race impact hiring decisions?
  - A Black candidate applied for a job, but did not get it.
  - Would a Black candidate have been offered a job if they were white?

## Causal questions in social science

- Does race impact hiring decisions?
  - A Black candidate applied for a job, but did not get it.
  - Would a Black candidate have been offered a job if they were white?
- Does the minimum wage increase unemployment?
  - Unemployment went up in a city after the minimum wage increased
  - Would unemployment have gone up were there not an increase in the minimum wage?

## Causal questions in social science

- Does race impact hiring decisions?
  - A Black candidate applied for a job, but did not get it.
  - Would a Black candidate have been offered a job if they were white?
- Does the minimum wage increase unemployment?
  - Unemployment went up in a city after the minimum wage increased
  - Would unemployment have gone up were there not an increase in the minimum wage?
- Does community policing decrease crime?
  - A police department implemented community policing in certain neighborhoods, and reported crime went down
  - Would reported crime have gone down without community policing?

*Evaluates how treatments causally effect outcomes by assigning different levels of treatment to different observations, then measuring the corresponding values of the outcome*

# Using an experiment to estimate the effects of a criminal record on employment

Pager, Devah. "The mark of a criminal record." American journal of sociology 108.5 (2003): 937-975.

*With over 2 million individuals currently incarcerated, and over half a million prisoners released each year, the large and growing number of men being processed through the criminal justice system raises important questions about the consequences of this massive institutional intervention. This article focuses on the consequences of incarceration for the employment outcomes of black and white job seekers. The present study adopts an experimental audit approach—in which matched pairs of individuals applied for real entry - level jobs - to formally test the degree to which a criminal record affects subsequent employment opportunities. The findings of this study reveal an important, and much underrecognized, mechanism of stratification. A criminal record presents a major barrier to employment, with important implications for racial disparities.*

1. Do employers use criminal histories to make hiring decisions?
2. Is racial discrimination a major barrier to employment?
3. Does the effect of a criminal record differ for white and Black applicants?

## What counterfactuals are needed for each question?

1. Do employers use criminal histories to make hiring decisions?



## What counterfactuals are needed for each question?

1. Do employers use criminal histories to make hiring decisions?
2. Is racial discrimination a major barrier to employment?

## What counterfactuals are needed for each question?

1. Do employers use criminal histories to make hiring decisions?
2. Is racial discrimination a major barrier to employment?
3. Does the effect of a criminal record differ for white and Black applicants?

## Before we begin

- grab the data from Slack

```
dat<-read.csv("../data/criminalrecord.csv")
```

## Variables in the data

**jobid** Job ID number

**callback** 1 if tester received a callback, 0 if the tester did not receive a callback.

**black** 1 if the tester is black, 0 if the tester is white.

**crimrec** 1 if the tester has a criminal record, 0 if the tester does not.

**interact** 1 if tester interacted with employer during the job application, 0 if tester does not interact with employer.

**city** 1 if job is located in the city center, 0 if job is located in the suburbs.

**distance** Job's average distance to downtown.

**custserv** 1 if job is in the customer service sector, 0 if it is not.

**manualskill** 1 if job requires manual skills, 0 if it does not.

# Take a look at the data

```
head(dat)
```

```
##   jobid callback black crimrec interact city distance custserv manualskill
## 1   108         1     0         1         1     0         15         1         0
## 2   113         0     0         0         1     0         20         0         1
## 3   101         1     0         0         0     0         15         1         0
## 4    64         1     0         0         0     1          7         1         0
## 5    33         0     0         1         0     1          5         1         0
## 6    73         0     0         1         0     1         10         0         1
```

## Exploring the data: univariate crosstabs

```
table(race = dat$black)
```

```
## race  
##    0    1  
## 300 396
```

```
table(crimrec = dat$crimrec)
```

```
## crimrec  
##    0    1  
## 349 347
```

## Exploring the data: bivariate crosstabs

```
table(race = dat$black, crimrec = dat$crimrec)
```

```
##      crimrec
## race    0    1
##    0 150 150
##    1 199 197
```

```
table(Black = dat$black == 1, callback = dat$callback)
```

```
##          callback
## Black      0    1
##  FALSE 224  76
##   TRUE 358  38
```

What was the callback rate for subjects assigned a criminal record?

```
crim_rec<-table(crimrec = dat$crimrec, callback = dat$callback)  
crim_rec
```

```
##          callback  
## crimrec    0    1  
##          0 270  79  
##          1 312  35
```



## Using crosstabs

What was the callback rate for subjects assigned a criminal record?

```
crim_rec
```

```
##          callback
## crimrec    0    1
##          0 270  79
##          1 312  35
```

```
## Divide those with a criminal record and callback
```

```
## By all those with a criminal record
```

```
crim_rec[2, 2] / sum(crim_rec[2,])
```

```
## [1] 0.1008646
```

## Using crosstabs

What was the callback rate for subjects *not* assigned a criminal record?

```
crim_rec
```

```
##          callback
## crimrec    0    1
##          0 270  79
##          1 312  35
```

```
## Divide those with a criminal record and callback
```

```
## By all those with a criminal record
```

```
crim_rec[1, 2] / sum(crim_rec[1,])
```

```
## [1] 0.226361
```

## Subsetting and an aside on logicals

---

```
temp<-c(TRUE, FALSE, TRUE)  
str(temp)
```

```
##  logi [1:3] TRUE FALSE TRUE
```

```
temp<-c(TRUE, FALSE, TRUE)  
sum(temp)
```

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

```
temp<-c(TRUE, FALSE, TRUE)  
mean(temp)
```

```
## [1] 0.6666667
```

```
## AND: &  
TRUE & FALSE
```

```
## [1] FALSE
```

```
## OR: |  
TRUE | FALSE
```

```
## [1] TRUE
```



```
## NOT: !
```

```
!TRUE
```

```
## [1] FALSE
```

## We often use logicals in conjunction with comparisons

- $<$  and  $>$  less than and greater than

## We often use logicals in conjunction with comparisons

- `<` and `>` less than and greater than
- `<=` and `>=` less/greater than or equal to

## We often use logicals in conjunction with comparisons

- `<` and `>` less than and greater than
- `<=` and `>=` less/greater than or equal to
- `==` equal to

## We often use logicals in conjunction with comparisons

- < and > less than and greater than
- <= and >= less/greater than or equal to
- == equal to
- != not equal to

## We often use logicals in conjunction with comparisons

- `<` and `>` less than and greater than
- `<=` and `>=` less/greater than or equal to
- `==` equal to
- `!=` not equal to
- `%in%` element in vector

That's neat (but kinda useless?)

---

It is very useful!

```
2<3
```

```
## [1] TRUE
```



It is very useful!

```
2<3 & 2>3
```

```
## [1] FALSE
```

It is very useful!

```
2<3 | 2>3
```

```
## [1] TRUE
```

It is very useful!

```
!(2<3)
```

```
## [1] FALSE
```

## Vectorized comparisons

```
temp<-c(2,3,4,5)
```

```
3<temp
```

```
## [1] FALSE FALSE  TRUE  TRUE
```

## Vectorized comparisons

```
temp<-c(2,3,4,5)
```

```
3==temp
```

```
## [1] FALSE TRUE FALSE FALSE
```

## Vectorized comparisons

```
temp<-c(2,3,4,5)
```

```
3%in%temp
```

```
## [1] TRUE
```

## Let's use these to subset

```
## Note that recoding here is not needed
## subset for all rows where Black is equal to 1
dat_blk<-dat[dat$black == 1, ]
head(dat_blk)
```

```
##      jobid callback black crimrec interact city distance custserv manualskill
## 301  1179         0     1         1         0     0         12         1         0
## 302  1180         0     1         1         0     1          1         1         0
## 303  1136         0     1         0         0     0         15         0         1
## 304  1095         0     1         0         1     0         14         1         1
## 305  1076         0     1         1         0     1          5         0         1
## 306  1143         0     1         1         0     1          2         0         1
```

Use this variable to subset the data into Black/white applicants

```
dat_blk<-dat[dat$black == 1, ]
```



Use this variable to subset the data into Black/white applicants

```
dat_wht<-dat[dat$black == 0, ]
```

## Use this variable to subset the data into Black/white applicants

```
nrow(dat_blk)
```

```
## [1] 396
```

```
nrow(dat_wht)
```

```
## [1] 300
```

```
nrow(dat)
```

```
## [1] 696
```

## Let's subset into race/crimrec datasets

```
dat_blk_crim<-dat_blk[dat_blk$crimrec==1,]  
## OR dat_blk_crim<-dat[dat$black==1 & dat$crimrec==1, ]
```

## Let's subset into race/crimrec datasets

```
dat_wht_crim<-dat_wht[dat_wht$crimrec==1,]
```

## Let's subset into race/crimrec datasets

```
## check number of cases
```

```
nrow(dat_blk_crim)
```

```
## [1] 197
```

```
nrow(dat_wht_crim)
```

```
## [1] 150
```

Questions on logicals and filters?

---

## Recoding and conditionals

Let's make distance categorical, with cuts at the 25th, 50th, and 75th quantile

```
summary(dat$distance)
```

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
##	0.00	8.00	12.00	11.96	16.00	25.00	2

```
## NA???
```

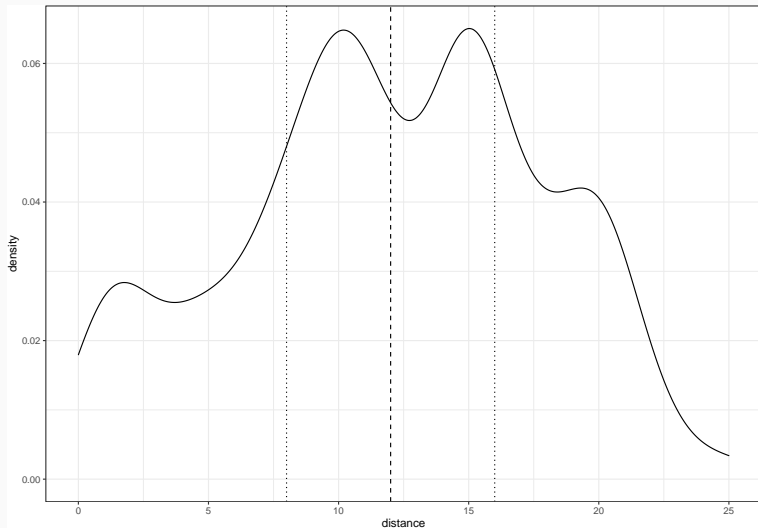
## Subsetting to remove missing values

```
## remove pesky NA values
dat_clean<-dat[!(is.na(dat$distance)),]
### wait, what did you do there???!
### also works, but more aggressive: dat_clean<-na.omit(dat)
summary(dat_clean$distance)
```

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	0.00	8.00	12.00	11.96	16.00	25.00



## Visualizing quantiles: remember area under the curve?



## Making a recode with one condition

Make a new variable for distance, with value “near” if below the median, and “far” if above

```
dat_clean$distance_far <- dat_clean$distance >  
  median(dat_clean$distance)  
  
table(dat_clean$distance_far)
```

```
##  
## FALSE  TRUE  
##   374   320
```

## Making a recode with one condition: ifelse()

Make a new variable for distance, with value “near” if below the median, and “far” if above

```
dat_clean$distance_binary <-ifelse(  
  dat_clean$distance<median(dat_clean$distance), # CONDITION  
  "near", # IF TRUE  
  "far") # IF FALSE  
table(dat_clean$distance_binary)
```

```
##  
## far near  
## 370 324
```

# Making a recode with multiple conditions

```
### define quartile cut points
q1<-quantile(dat_clean$distance, 0.25)
q2<-quantile(dat_clean$distance, 0.5)
q3<-quantile(dat_clean$distance, 0.75)
q1; q2; q3
```

```
## 25%
##    8
```

```
## 50%
##   12
```

```
## 75%
##   16
```

```
table(dat_clean$distance_quartile)
```

```
## < table of extent 0 >
```

# Making a recode with multiple conditions

```
### make factor variable
dat_clean$distance_quartile<-NA

dat_clean$distance_quartile[dat_clean$distance<q1]<-"1st"

dat_clean$distance_quartile[dat_clean$distance>=q1 &
                             dat_clean$distance<q2]<-"2nd"

dat_clean$distance_quartile[dat_clean$distance>=q2 &
                             dat_clean$distance<q3]<-"3rd"

dat_clean$distance_quartile[dat_clean$distance>q3]<-"4th"

table(dat_clean$distance_quartile)

##
## 1st 2nd 3rd 4th
## 148 176 184 158
```

## Returning to Pager's experiment

---

# The counterfactual and potential outcomes

```
c_fact<-data.frame(callback = dat$callback,  
                  crimrec = dat$crimrec)  
  
### create explicit counterfactual  
  
c_fact$callback_crimT<-ifelse(c_fact$crimrec==1, c_fact$callback, NA)  
c_fact$callback_crimF<-ifelse(c_fact$crimrec==0, c_fact$callback, NA)  
  
head(c_fact)
```

```
##   callback crimrec callback_crimT callback_crimF  
## 1      1      1      1      NA  
## 2      0      0      NA      0  
## 3      1      0      NA      1  
## 4      1      0      NA      1  
## 5      0      1      0      NA  
## 6      0      1      0      NA
```

For observation  $i$ , the sample average treatment effect (SATE) is equal to:

$$\text{callback\_crimTRUE}_i - \text{callback\_crimFALSE}_i$$



## What is the causal effect for rows 1 - 6

For observation  $i$ , the treatment effect is equal to:

$\text{callback\_crimTRUE}_i - \text{callback\_crimFALSE}_i$

```
head(c_fact)
```

```
##      callback crimrec callback_crimT callback_crimF
## 1           1         1              1             NA
## 2           0         0             NA              0
## 3           1         0             NA              1
## 4           1         0             NA              1
## 5           0         1              0             NA
## 6           0         1              0             NA
```

## What is the causal effect for rows 1 - 6

For observation  $i$ , the treatment effect is equal to:

$\text{callback\_crimTRUE}_i - \text{callback\_crimFALSE}_i$

```
head(c_fact)
```

```
##      callback crimrec callback_crimT callback_crimF
## 1           1         1              1             NA
## 2           0         0             NA              0
## 3           1         0             NA              1
## 4           1         0             NA              1
## 5           0         1              0             NA
## 6           0         1              0             NA
```

*The fundamental problem of causal inference is that we only observe one of these outcomes*

- By randomizing assignment to treatment, we can treat units as equivalent

## Randomized experiments (or RCTs)

- By randomizing assignment to treatment, we can treat units as equivalent
- If units are equivalent, we can estimate the average treatment effect as a difference in means on the outcome between the treatment and control group

## Randomized experiments (or RCTs)

- By randomizing assignment to treatment, we can treat units as equivalent
- If units are equivalent, we can estimate the average treatment effect as a difference in means on the outcome between the treatment and control group
- If we don't randomize, we have no assurance that the treated and control groups are equivalent, meaning we can't argue that we've observed the counterfactual

## The SATE for Pager's experiment

We assume that we can estimate the counterfactual for people with criminal records (i.e. no criminal record), by using the mean value of the callback outcome for people assigned to have no criminal record.

```
### obtain the mean callback rate of those with a criminal record
dat_crimrecT <- mean(dat[dat$crimrec==1, "callback"])
### and those without
dat_crimrecF <- mean(dat[dat$crimrec==0, "callback"])
### the mean callback rate for the treatment group and the control
dat_crimrecT

## [1] 0.1008646

dat_crimrecF

## [1] 0.226361
```

- Homework: More work with Pager's data
- Causality, part 2. Observational studies
- Measuring characteristics of the distribution of a variable