

The Experimental Ideal

EC 425/525, Set 2

Edward Rubin
03 April 2019

Prologue

Schedule

Last time

Research basics, our class, and R

Today

Admin: Canvas exists. Updated class [website](#).

Material: The Rubin causal model (not mine), [Chapter 2 MHE](#). R basics.

Schedule

Last time

Research basics, our class, and R

Today

Admin: Canvas exists. Updated class [website](#).

Material: The Rubin causal model (not mine), [Chapter 2 MHE](#). R basics.

Assignment₁ Install R and [RStudio](#) on your computer.

Schedule

Last time

Research basics, our class, and R

Today

Admin: Canvas exists. Updated class [website](#).

Material: The Rubin causal model (not mine), [Chapter 2 MHE](#). R basics.

Assignment₁ Install R and [RStudio](#) on your computer.

Assignment₂ Take 15 minutes to quietly think about your interests.

Schedule

Last time

Research basics, our class, and R

Today

Admin: Canvas exists. Updated class [website](#).

Material: The Rubin causal model (not mine), [Chapter 2 MHE](#). R basics.

Assignment₁ Install R and [RStudio](#) on your computer.

Assignment₂ Take 15 minutes to quietly think about your interests.

Future

Lab: Matrix work, regression, functions, simulation

Long run: Deepen understandings/intuitions for causality and inference.

Review

Research fundamentals

Review

Research fundamentals

Angrist and Pischke provide four **fundamental questions for research**:

1. What is the **causal relationship of interest**?
2. How would an **ideal experiment** capture this causal effect of interest?
3. What is your **identification strategy**?
4. What is your **mode of inference**?

Review

Research fundamentals

Angrist and Pischke provide four **fundamental questions for research**:

1. What is the **causal relationship of interest**?
2. How would an **ideal experiment** capture this causal effect of interest?
3. What is your **identification strategy**?
4. What is your **mode of inference**?

Seemingly straightforward questions can be fundamentally unanswerable.

Review

General research recommendations

More unsolicited advice:

- Be curious.
- Ask questions.
- Attend seminars.
- Meet faculty (UO + visitors).
- Focus on learning—especially intuition.[†]
- **Be kind and constructive.**

[†] Learning is not always the same as getting good grades.

The experimental ideal

The experimental ideal

What's so great about experiments?

Science widely regards **experiments as the gold standard** for research.

But why? The costs can be substantial.

Costs

- slow and expensive
- heavily regulated by review boards
- can abstract away from the actual question/setting

Benefits

So the benefits need to be pretty large, right?

The experimental ideal

Example: Hospitals and health

Imagine we want to know the **causal effect of hospitals on health**.

The experimental ideal

Example: Hospitals and health

Imagine we want to know the **causal effect of hospitals on health**.

Research question

Within the population of poor, elderly individuals, does visiting the emergency room for primary care improve health?

The experimental ideal

Example: Hospitals and health

Imagine we want to know the **causal effect of hospitals on health**.

Research question

Within the population of poor, elderly individuals, does visiting the emergency room for primary care improve health?

Empirical exercise

1. Collect data on *health status* and *hospital visits*.
2. Summarize health status by hospital-visit group.

The experimental ideal

Example: Hospitals and health

Our empirical exercise from the 2005 National Health Interview Survey:

Group	Sample Size	Mean Health Status	Std. Error
Hospital	7,774	3.21	0.014
No hospital	90,049	3.93	0.003

The experimental ideal

Example: Hospitals and health

Our empirical exercise from the 2005 National Health Interview Survey:

Group	Sample Size	Mean Health Status	Std. Error
Hospital	7,774	3.21	0.014
No hospital	90,049	3.93	0.003

We get a t statistic of 58.9 when testing a difference in groups' means (0.72).

The experimental ideal

Example: Hospitals and health

Our empirical exercise from the 2005 National Health Interview Survey:

Group	Sample Size	Mean Health Status	Std. Error
Hospital	7,774	3.21	0.014
No hospital	90,049	3.93	0.003

We get a t statistic of 58.9 when testing a difference in groups' means (0.72).

Conclusion? Hospitals make folks worse. Hospitals make sick people sicker.

The experimental ideal

Example: Hospitals and health

Our empirical exercise from the 2005 National Health Interview Survey:

Group	Sample Size	Mean Health Status	Std. Error
Hospital	7,774	3.21	0.014
No hospital	90,049	3.93	0.003

We get a t statistic of 58.9 when testing a difference in groups' means (0.72).

Conclusion? Hospitals make folks worse. Hospitals make sick people sicker.

Alternative conclusion: Perhaps we're making a mistake in our analysis...

The experimental ideal

Example: Hospitals and health

Our empirical exercise from the 2005 National Health Interview Survey:

Group	Sample Size	Mean Health Status	Std. Error
Hospital	7,774	3.21	0.014
No hospital	90,049	3.93	0.003

We get a t statistic of 58.9 when testing a difference in groups' means (0.72).

Conclusion? Hospitals make folks worse. Hospitals make sick people sicker.

Alternative conclusion: Perhaps we're making a mistake in our analysis... maybe sick people go to hospitals?

The experimental ideal

Potential outcomes framework

Let's develop a framework to better discuss the problem here.

The experimental ideal

Potential outcomes framework

Let's develop a framework to better discuss the problem here.

- Binary treatment variable (e.g., hospitalized): $D_i = 0, 1$
- Outcome for individual i (e.g., health): Y_i

This framework has a few names...

- Neyman potential outcomes framework
- Rubin causal model
- Neyman-Rubin "potential outcome" | "causal" "framework" | "model"

The experimental ideal

Potential outcomes framework

Research question: Does D_i affect Y_i ?

The experimental ideal

Potential outcomes framework

Research question: Does D_i affect Y_i ?

For each individual i , there are two **potential outcomes** (w/ binary D_i)

The experimental ideal

Potential outcomes framework

Research question: Does D_i affect Y_i ?

For each individual i , there are two **potential outcomes** (w/ binary D_i)

1. Y_{1i} if $D_i = 1$
 i 's health outcome if she went to the hospital

The experimental ideal

Potential outcomes framework

Research question: Does D_i affect Y_i ?

For each individual i , there are two **potential outcomes** (w/ binary D_i)

1. Y_{1i} if $D_i = 1$

i 's health outcome if she went to the hospital

2. Y_{0i} if $D_i = 0$

i 's health outcome if she did not go to the hospital

The experimental ideal

Potential outcomes framework

Research question: Does D_i affect Y_i ?

For each individual i , there are two **potential outcomes** (w/ binary D_i)

1. Y_{1i} if $D_i = 1$
 i 's health outcome if she went to the hospital
2. Y_{0i} if $D_i = 0$
 i 's health outcome if she did not go to the hospital

The difference between these two outcomes gives us the **causal effect of hospital treatment**, i.e.,

$$\tau_i = Y_{1i} - Y_{0i}$$

The experimental ideal

#problems

This simple equation

$$\tau_i = \textcolor{red}{Y}_{1i} - \textcolor{blue}{Y}_{0i}$$

leads us to ***the fundamental problem of causal inference.***

The experimental ideal

#problems

This simple equation

$$\tau_i = \textcolor{red}{Y}_{1i} - \textcolor{blue}{Y}_{0i}$$

leads us to ***the fundamental problem of causal inference.***

We can never simultaneously observe $\textcolor{red}{Y}_{1i}$ and $\textcolor{blue}{Y}_{0i}$.

The experimental ideal

#problems

This simple equation

$$\tau_i = \textcolor{orange}{Y}_{1i} - \textcolor{blue}{Y}_{0i}$$

leads us to ***the fundamental problem of causal inference.***

We can never simultaneously observe $\textcolor{red}{Y}_{1i}$ and $\textcolor{blue}{Y}_{0i}$.

Most of applied econometrics focuses on addressing this simple problem.

The experimental ideal

#problems

This simple equation

$$\tau_i = \textcolor{orange}{Y}_{1i} - \textcolor{blue}{Y}_{0i}$$

leads us to ***the fundamental problem of causal inference.***

We can never simultaneously observe $\textcolor{red}{Y}_{1i}$ and $\textcolor{blue}{Y}_{0i}$.

Most of applied econometrics focuses on addressing this simple problem.

Accordingly, our methods try to address the related question

For each $\textcolor{red}{Y}_{1i}$, what is a (reasonably) good counterfactual?

The experimental ideal

Solutions?

Problem We cannot directly calculate $\tau_i = Y_{1i} - Y_{0i}$.

The experimental ideal

Solutions?

Problem We cannot directly calculate $\tau_i = Y_{1i} - Y_{0i}$.

Proposed solution

Compare outcomes for people who visited the hospital ($Y_{1i} | D_i = 1$) to outcomes for people who did not visit the hospital ($Y_{0j} | D_i = 0$).

The experimental ideal

Solutions?

Problem We cannot directly calculate $\tau_i = Y_{1i} - Y_{0i}$.

Proposed solution

Compare outcomes for people who visited the hospital ($Y_{1i} | D_i = 1$) to outcomes for people who did not visit the hospital ($Y_{0j} | D_i = 0$).

$$E[Y_i | D_i = 1] - E[Y_i | D_i = 0]$$

which gives us the *observed difference in health outcomes*.

The experimental ideal

Solutions?

Problem We cannot directly calculate $\tau_i = \text{Y}_{1i} - \text{Y}_{0i}$.

Proposed solution

Compare outcomes for people who visited the hospital ($\text{Y}_{1i} | D_i = 1$) to outcomes for people who did not visit the hospital ($\text{Y}_{0j} | D_i = 0$).

$$E[\text{Y}_i | D_i = 1] - E[\text{Y}_i | D_i = 0]$$

which gives us the *observed difference in health outcomes*.

Q This comparison will return an answer, but is it *the* answer we want?

The experimental ideal

Selection

Q What does $E[Y_i | D_i = 1] - E[Y_i | D_i = 0]$ actually tell us?

The experimental ideal

Selection

Q What does $E[Y_i | D_i = 1] - E[Y_i | D_i = 0]$ actually tell us?

A First notice that we can write i 's outcome Y_i as

$$Y_i = Y_{0i} + D_i \underbrace{(Y_{1i} - Y_{0i})}_{\tau_i}$$

The experimental ideal

Selection

Q What does $E[Y_i | D_i = 1] - E[Y_i | D_i = 0]$ actually tell us?

A First notice that we can write i 's outcome Y_i as

$$Y_i = Y_{0i} + D_i \underbrace{(Y_{1i} - Y_{0i})}_{\tau_i}$$

Now write out our expectation, apply this definition, do creative math.

$$E[Y_i | D_i = 1] - E[Y_i | D_i = 0]$$

The experimental ideal

Selection

Q What does $E[Y_i | D_i = 1] - E[Y_i | D_i = 0]$ actually tell us?

A First notice that we can write i 's outcome Y_i as

$$Y_i = Y_{0i} + D_i \underbrace{(Y_{1i} - Y_{0i})}_{\tau_i}$$

Now write out our expectation, apply this definition, do creative math.

$$\begin{aligned} &E[Y_i | D_i = 1] - E[Y_i | D_i = 0] \\ &= E[Y_{1i} | D_i = 1] - E[Y_{0i} | D_i = 0] \end{aligned}$$

The experimental ideal

Selection

Q What does $E[Y_i | D_i = 1] - E[Y_i | D_i = 0]$ actually tell us?

A First notice that we can write i 's outcome Y_i as

$$Y_i = Y_{0i} + D_i \underbrace{(Y_{1i} - Y_{0i})}_{\tau_i}$$

Now write out our expectation, apply this definition, do creative math.

$$\begin{aligned} E[Y_i | D_i = 1] - E[Y_i | D_i = 0] &= E[Y_{1i} | D_i = 1] - E[Y_{0i} | D_i = 0] \\ &= E[Y_{1i} | D_i = 1] - E[Y_{0i} | D_i = 1] + E[Y_{0i} | D_i = 1] - E[Y_{0i} | D_i = 0] \end{aligned}$$

The experimental ideal

Selection

Q What does $E[Y_i | D_i = 1] - E[Y_i | D_i = 0]$ actually tell us?

A First notice that we can write i 's outcome Y_i as

$$Y_i = Y_{0i} + D_i \underbrace{(Y_{1i} - Y_{0i})}_{\tau_i}$$

Now write out our expectation, apply this definition, do creative math.

$$\begin{aligned} & E[Y_i | D_i = 1] - E[Y_i | D_i = 0] \\ &= E[Y_{1i} | D_i = 1] - E[Y_{0i} | D_i = 0] \\ &= \underbrace{E[Y_{1i} | D_i = 1] - E[Y_{0i} | D_i = 1]}_{\text{Average treatment effect on the treated}} + E[Y_{0i} | D_i = 1] - E[Y_{0i} | D_i = 0] \end{aligned}$$

Average treatment effect on the treated 😊

The experimental ideal

Selection

Q What does $E[Y_i | D_i = 1] - E[Y_i | D_i = 0]$ actually tell us?

A First notice that we can write i 's outcome Y_i as

$$Y_i = Y_{0i} + D_i \underbrace{(Y_{1i} - Y_{0i})}_{\tau_i}$$

Now write out our expectation, apply this definition, do creative math.

$$E[Y_i | D_i = 1] - E[Y_i | D_i = 0]$$

$$= E[Y_{1i} | D_i = 1] - E[Y_{0i} | D_i = 0]$$

$$= \underbrace{E[Y_{1i} | D_i = 1] - E[Y_{0i} | D_i = 1]}_{\text{Average treatment effect on the treated } \smiley} + \underbrace{E[Y_{0i} | D_i = 1] - E[Y_{0i} | D_i = 0]}_{\text{Selection bias } \frowny}$$

The experimental ideal

Selection

The **first term** is *good variation*—essentially the answer that we want.

$$E[\mathbf{Y}_{1i} \mid \mathbf{D}_i = 1] - E[\mathbf{Y}_{0i} \mid \mathbf{D}_i = 1]$$

The experimental ideal

Selection

The **first term** is *good variation*—essentially the answer that we want.

$$\begin{aligned} & E[\text{Y}_{1i} \mid D_i = 1] - E[\text{Y}_{0i} \mid D_i = 1] \\ &= E[\text{Y}_{1i} - \text{Y}_{0i} \mid D_i = 1] \end{aligned}$$

The experimental ideal

Selection

The **first term** is *good variation*—essentially the answer that we want.

$$\begin{aligned} E[Y_{1i} \mid D_i = 1] - E[Y_{0i} \mid D_i = 1] \\ = E[Y_{1i} - Y_{0i} \mid D_i = 1] \\ = E[\tau_i \mid D_i = 1] \end{aligned}$$

The experimental ideal

Selection

The **first term** is *good variation*—essentially the answer that we want.

$$\begin{aligned} & E[Y_{1i} \mid D_i = 1] - E[Y_{0i} \mid D_i = 1] \\ &= E[Y_{1i} - Y_{0i} \mid D_i = 1] \\ &= E[\tau_i \mid D_i = 1] \end{aligned}$$

The **average causal effect** of hospitalization for hospitalized individuals.

The experimental ideal

Selection

The **first term** is *good variation*—essentially the answer that we want.

$$\begin{aligned} E[Y_{1i} \mid D_i = 1] - E[Y_{0i} \mid D_i = 1] \\ = E[Y_{1i} - Y_{0i} \mid D_i = 1] \\ = E[\tau_i \mid D_i = 1] \end{aligned}$$

The **average causal effect** of hospitalization for hospitalized individuals.

The **second term** is *bad variation*—preventing us from knowing the answer.

$$E[Y_{0i} \mid D_i = 1] - E[Y_{0i} \mid D_i = 0]$$

The experimental ideal

Selection

The **first term** is *good variation*—essentially the answer that we want.

$$\begin{aligned} E[Y_{1i} \mid D_i = 1] - E[Y_{0i} \mid D_i = 1] \\ = E[Y_{1i} - Y_{0i} \mid D_i = 1] \\ = E[\tau_i \mid D_i = 1] \end{aligned}$$

The **average causal effect** of hospitalization for hospitalized individuals.

The **second term** is *bad variation*—preventing us from knowing the answer.

$$E[Y_{0i} \mid D_i = 1] - E[Y_{0i} \mid D_i = 0]$$

The difference in the average untreated outcome between the treatment and control groups.

The experimental ideal

Selection

The **first term** is *good variation*—essentially the answer that we want.

$$\begin{aligned} E[Y_{1i} \mid D_i = 1] - E[Y_{0i} \mid D_i = 1] \\ = E[Y_{1i} - Y_{0i} \mid D_i = 1] \\ = E[\tau_i \mid D_i = 1] \end{aligned}$$

The **average causal effect** of hospitalization for hospitalized individuals.

The **second term** is *bad variation*—preventing us from knowing the answer.

$$E[Y_{0i} \mid D_i = 1] - E[Y_{0i} \mid D_i = 0]$$

The difference in the average untreated outcome between the treatment and control groups.

Selection bias The extent to which the "control group" provides a bad counterfactual for the treated individuals.

The experimental ideal

Selection

Angrist and Pischke (MHE, p. 15),

The goal of most empirical economic research is to overcome selection bias, and therefore to say something about the causal effect of a variable like D_i .

The experimental ideal

Selection

Angrist and Pischke (MHE, p. 15),

The goal of most empirical economic research is to overcome selection bias, and therefore to say something about the causal effect of a variable like D_i .

Q So how do experiments—the gold standard of empirical economic (and scientific) research—accomplish this goal and overcome selection bias?

The experimental ideal

Back to experiments

Q How do experiments overcome selection bias?

The experimental ideal

Back to experiments

Q How do experiments overcome selection bias?

A Experiments break the link between potential outcomes and treatment.

In other words: Randomly assigning D_i makes D_i independent of which outcome we observe (meaning Y_{1i} or Y_{0i}).

The experimental ideal

Back to experiments

Q How do experiments overcome selection bias?

A Experiments break the link between potential outcomes and treatment.

In other words: Randomly assigning D_i makes D_i independent of which outcome we observe (meaning Y_{1i} or Y_{0i}).

Difference in means with random assignment of D_i

$$E[Y_i | D_i = 1] - E[Y_i | D_i = 0]$$

The experimental ideal

Back to experiments

Q How do experiments overcome selection bias?

A Experiments break the link between potential outcomes and treatment.

In other words: Randomly assigning D_i makes D_i independent of which outcome we observe (meaning Y_{1i} or Y_{0i}).

Difference in means with random assignment of D_i

$$\begin{aligned} E[Y_i | D_i = 1] - E[Y_i | D_i = 0] \\ = E[Y_{1i} | D_i = 1] - E[Y_{0i} | D_i = 0] \end{aligned}$$

The experimental ideal

Back to experiments

Q How do experiments overcome selection bias?

A Experiments break the link between potential outcomes and treatment.

In other words: Randomly assigning D_i makes D_i independent of which outcome we observe (meaning Y_{1i} or Y_{0i}).

Difference in means with random assignment of D_i

$$\begin{aligned} & E[Y_i | D_i = 1] - E[Y_i | D_i = 0] \\ &= E[Y_{1i} | D_i = 1] - E[Y_{0i} | D_i = 0] \\ &= E[Y_{1i} | D_i = 1] - E[Y_{0i} | D_i = 1] \end{aligned}$$

The experimental ideal

Back to experiments

Q How do experiments overcome selection bias?

A Experiments break the link between potential outcomes and treatment.

In other words: Randomly assigning D_i makes D_i independent of which outcome we observe (meaning Y_{1i} or Y_{0i}).

Difference in means with random assignment of D_i

$$\begin{aligned} & E[Y_i \mid D_i = 1] - E[Y_i \mid D_i = 0] \\ &= E[Y_{1i} \mid D_i = 1] - E[Y_{0i} \mid D_i = 0] \\ &= E[Y_{1i} \mid D_i = 1] - E[Y_{0i} \mid D_i = 1] \quad \text{from random assignment of } D_i \end{aligned}$$

The experimental ideal

Back to experiments

Q How do experiments overcome selection bias?

A Experiments break the link between potential outcomes and treatment.

In other words: Randomly assigning D_i makes D_i independent of which outcome we observe (meaning Y_{1i} or Y_{0i}).

Difference in means with random assignment of D_i

$$\begin{aligned} & E[Y_i \mid D_i = 1] - E[Y_i \mid D_i = 0] \\ &= E[Y_{1i} \mid D_i = 1] - E[Y_{0i} \mid D_i = 0] \\ &= E[Y_{1i} \mid D_i = 1] - E[Y_{0i} \mid D_i = 1] \quad \text{from random assignment of } D_i \\ &= E[Y_{1i} - Y_{0i} \mid D_i = 1] \end{aligned}$$

The experimental ideal

Back to experiments

Q How do experiments overcome selection bias?

A Experiments break the link between potential outcomes and treatment.

In other words: Randomly assigning D_i makes D_i independent of which outcome we observe (meaning Y_{1i} or Y_{0i}).

Difference in means with random assignment of D_i

$$\begin{aligned} & E[Y_i \mid D_i = 1] - E[Y_i \mid D_i = 0] \\ &= E[Y_{1i} \mid D_i = 1] - E[Y_{0i} \mid D_i = 0] \\ &= E[Y_{1i} \mid D_i = 1] - E[Y_{0i} \mid D_i = 1] \quad \text{from random assignment of } D_i \\ &= E[Y_{1i} - Y_{0i} \mid D_i = 1] \\ &= E[\tau_i \mid D_i = 1] \end{aligned}$$

The experimental ideal

Back to experiments

Q How do experiments overcome selection bias?

A Experiments break the link between potential outcomes and treatment.

In other words: Randomly assigning D_i makes D_i independent of which outcome we observe (meaning Y_{1i} or Y_{0i}).

Difference in means with random assignment of D_i

$$\begin{aligned} & E[Y_i \mid D_i = 1] - E[Y_i \mid D_i = 0] \\ &= E[Y_{1i} \mid D_i = 1] - E[Y_{0i} \mid D_i = 0] \\ &= E[Y_{1i} \mid D_i = 1] - E[Y_{0i} \mid D_i = 1] \quad \text{from random assignment of } D_i \\ &= E[Y_{1i} - Y_{0i} \mid D_i = 1] \\ &= E[\tau_i \mid D_i = 1] \\ &= E[\tau_i] \end{aligned}$$

The experimental ideal

Back to experiments

Q How do experiments overcome selection bias?

A Experiments break the link between potential outcomes and treatment.

In other words: Randomly assigning D_i makes D_i independent of which outcome we observe (meaning Y_{1i} or Y_{0i}).

Difference in means with random assignment of D_i

$$E[Y_i | D_i = 1] - E[Y_i | D_i = 0]$$

$$= E[Y_{1i} | D_i = 1] - E[Y_{0i} | D_i = 0]$$

$$= E[Y_{1i} | D_i = 1] - E[Y_{0i} | D_i = 1] \quad \text{from random assignment of } D_i$$

$$= E[Y_{1i} - Y_{0i} | D_i = 1]$$

$$= E[\tau_i | D_i = 1]$$

$$= E[\tau_i] \quad \text{Random assignment of } D_i \text{ breaks selection bias.}$$

The experimental ideal

Example: Training programs

Governments subsidize training programs to assist disadvantaged workers.

The experimental ideal

Example: Training programs

Governments subsidize training programs to assist disadvantaged workers.

Q Do these programs have the desired effects (*i.e.*, increase wages)?

The experimental ideal

Example: Training programs

Governments subsidize training programs to assist disadvantaged workers.

Q Do these programs have the desired effects (*i.e.*, increase wages)?

A Observational studies—comparing wage data from participants and non-participants—often find that people who complete these programs actually make *lower wages*.

The experimental ideal

Example: Training programs

Governments subsidize training programs to assist disadvantaged workers.

Q Do these programs have the desired effects (*i.e.*, increase wages)?

A Observational studies—comparing wage data from participants and non-participants—often find that people who complete these programs actually make *lower wages*.

Challenges Participants self select. + Programs target lower-wage workers.

The experimental ideal

Example: Training programs

How do we formalize these concerns in our framework?

The experimental ideal

Example: Training programs

How do we formalize these concerns in our framework?

Observational program evaluations

$$E[\text{Wage}_i \mid \text{Program}_i = 1] - E[\text{Wage}_i \mid \text{Program}_i = 0] =$$

$$\underbrace{E[\text{Wage}_{1i} \mid \text{Program}_i = 1] - E[\text{Wage}_{0i} \mid \text{Program}_i = 1]}_{\text{Average causal effect of training program on wages for participants, } i.e., \bar{\tau}_1} +$$

Average causal effect of training program on wages for participants, *i.e.*, $\bar{\tau}_1$

$$\underbrace{E[\text{Wage}_{0i} \mid \text{Program}_i = 1] - E[\text{Wage}_{0i} \mid \text{Program}_i = 0]}_{\text{Selection bias}}$$

Selection bias

The experimental ideal

Example: Training programs

How do we formalize these concerns in our framework?

Observational program evaluations

$$E[\text{Wage}_i \mid \text{Program}_i = 1] - E[\text{Wage}_i \mid \text{Program}_i = 0] =$$

$$\underbrace{E[\text{Wage}_{1i} \mid \text{Program}_i = 1] - E[\text{Wage}_{0i} \mid \text{Program}_i = 1]}_{\text{Average causal effect of training program on wages for participants, } i.e., \bar{\tau}_1} +$$

Average causal effect of training program on wages for participants, *i.e.*, $\bar{\tau}_1$

$$\underbrace{E[\text{Wage}_{0i} \mid \text{Program}_i = 1] - E[\text{Wage}_{0i} \mid \text{Program}_i = 0]}_{\text{Selection bias}}$$

Selection bias

If the program attracts/selects individuals who, on average, have lower wages without the program (sort of the point of the program), then we have negative selection bias.

The experimental ideal

Example: Training programs

$$E[\text{Wage}_i \mid \text{Program}_i = 1] - E[\text{Wage}_i \mid \text{Program}_i = 0] =$$

$$E[\text{Wage}_{1i} \mid \text{Program}_i = 1] - E[\text{Wage}_{0i} \mid \text{Program}_i = 1] +$$

$$E[\text{Wage}_{0i} \mid \text{Program}_i = 1] - E[\text{Wage}_{0i} \mid \text{Program}_i = 0]$$

So even if the program, on average, has a positive wage effect (in the participant group), *i.e.*, $\bar{\tau}_1 > 0$, we will detect a lower effect due to the negative selection bias.

The experimental ideal

Example: Training programs

$$E[\text{Wage}_i \mid \text{Program}_i = 1] - E[\text{Wage}_i \mid \text{Program}_i = 0] =$$

$$E[\text{Wage}_{1i} \mid \text{Program}_i = 1] - E[\text{Wage}_{0i} \mid \text{Program}_i = 1] +$$

$$E[\text{Wage}_{0i} \mid \text{Program}_i = 1] - E[\text{Wage}_{0i} \mid \text{Program}_i = 0]$$

So even if the program, on average, has a positive wage effect (in the participant group), *i.e.*, $\bar{\tau}_1 > 0$, we will detect a lower effect due to the negative selection bias.

If the bias is sufficiently large (relative to the treatment effect), our estimate will even get the sign of the effect wrong.

The experimental ideal

Example: Training programs

$$E[\text{Wage}_i \mid \text{Program}_i = 1] - E[\text{Wage}_i \mid \text{Program}_i = 0] =$$

$$E[\text{Wage}_{1i} \mid \text{Program}_i = 1] - E[\text{Wage}_{0i} \mid \text{Program}_i = 1] +$$

$$E[\text{Wage}_{0i} \mid \text{Program}_i = 1] - E[\text{Wage}_{0i} \mid \text{Program}_i = 0]$$

So even if the program, on average, has a positive wage effect (in the participant group), *i.e.*, $\bar{\tau}_1 > 0$, we will detect a lower effect due to the negative selection bias.

If the bias is sufficiently large (relative to the treatment effect), our estimate will even get the sign of the effect wrong.

Related While observational studies typically found negative program effects, several experiments found positive program effects.

The experimental ideal

Example: The STAR experiment

The Tennessee STAR experiment is a famous/popular example of an experiment that allows us to answer an important social/policy question.

Research question Do classroom resources affect student performance?

The experimental ideal

Example: The STAR experiment

The Tennessee STAR experiment is a famous/popular example of an experiment that allows us to answer an important social/policy question.

Research question Do classroom resources affect student performance?

- Statewide(-ish) in Tennessee for the 1985–1986 kindergarten cohort
- Ran for 4 years with ~11,600 children. Cost ~\$12 million.

The experimental ideal

Example: The STAR experiment

The Tennessee STAR experiment is a famous/popular example of an experiment that allows us to answer an important social/policy question.

Research question Do classroom resources affect student performance?

- Statewide(-ish) in Tennessee for the 1985–1986 kindergarten cohort
- Ran for 4 years with ~11,600 children. Cost ~\$12 million.

Treatments

1. *Small* classes (13–17 students)
2. *Regular* classes (22–35 students) plus part-time teacher's aide
3. *Regular* classes (22–35 students) plus full-time teacher's aide

The experimental ideal

Example: The STAR experiment

First question Did the randomization balance participants' characteristics across the treatment groups?

The experimental ideal

Example: The STAR experiment

First question Did the randomization balance participants' characteristics across the treatment groups?

Ideally, we would have pre-experiment data on outcome variable.

Unfortunately, we only have a few demographic attributes.

Table 2.2.1, MHE

Variable	Treatment: Class Size			
	Small	Regular	Regular + Aide	P-value
<i>Free lunch</i>	0.47	0.48	0.50	0.09
<i>White/Asian</i>	0.68	0.67	0.66	0.26
<i>Age in 1985</i>	5.44	5.43	5.42	0.32
<i>Attrition rate</i>	0.49	0.52	0.53	0.02
<i>K. class size</i>	15.10	22.40	22.80	0.00
<i>K. test percentile</i>	54.70	48.90	50.00	0.00

Table 2.2.1, MHE

Variable	Treatment: Class Size			P-value
	Small	Regular	Regular + Aide	
<i>Free lunch</i>	0.47	0.48	0.50	0.09
<i>White/Asian</i>	0.68	0.67	0.66	0.26
<i>Age in 1985</i>	5.44	5.43	5.42	0.32
Attrition rate	0.49	0.52	0.53	0.02
K. class size	15.10	22.40	22.80	0.00
K. test percentile	54.70	48.90	50.00	0.00

Demographics appear balanced across the three treatment groups.

Table 2.2.1, MHE

Variable	Treatment: Class Size			P-value
	Small	Regular	Regular + Aide	
<i>Free lunch</i>	0.47	0.48	0.50	0.09
<i>White/Asian</i>	0.68	0.67	0.66	0.26
<i>Age in 1985</i>	5.44	5.43	5.42	0.32
Attrition rate	0.49	0.52	0.53	0.02
<i>K. class size</i>	15.10	22.40	22.80	0.00
<i>K. test percentile</i>	54.70	48.90	50.00	0.00

The three groups differ significantly on attrition rate.

Table 2.2.1, MHE

Variable	Treatment: Class Size			
	Small	Regular	Regular + Aide	P-value
<i>Free lunch</i>	0.47	0.48	0.50	0.09
<i>White/Asian</i>	0.68	0.67	0.66	0.26
<i>Age in 1985</i>	5.44	5.43	5.42	0.32
<i>Attrition rate</i>	0.49	0.52	0.53	0.02
K. class size	15.10	22.40	22.80	0.00
<i>K. test percentile</i>	54.70	48.90	50.00	0.00

The randomization generated variation in the treatment.

Table 2.2.1, MHE

Variable	Treatment: Class Size			
	Small	Regular	Regular + Aide	P-value
<i>Free lunch</i>	0.47	0.48	0.50	0.09
<i>White/Asian</i>	0.68	0.67	0.66	0.26
<i>Age in 1985</i>	5.44	5.43	5.42	0.32
<i>Attrition rate</i>	0.49	0.52	0.53	0.02
K. class size	15.10	22.40	22.80	0.00
K. test percentile	54.70	48.90	50.00	0.00

The small-class treatment significantly increased test scores.

The experimental ideal

The STAR experiment

The previous table estimated/compared the treatment effects using simple differences in means.

We can make the same comparisons using regressions.

Specifically, we regress our outcome (test percentile) on dummy variables (binary indicator variables) for each treatment group.

The experimental ideal

Example of our three treatment dummies.

i	y_i	Trt_{1i}	Trt_{2i}	Trt_{3i}
1	y_1	1	0	0
2	y_2	1	0	0
\vdots	\vdots	\vdots	\vdots	\vdots
ℓ	y_ℓ	1	0	0
$\ell + 1$	$y_{\ell+1}$	0	1	0
\vdots	\vdots	\vdots	\vdots	\vdots
p	y_p	0	1	0
$p + 1$	y_{p+1}	0	0	1
\vdots	\vdots	\vdots	\vdots	\vdots
N	y_N	0	0	1

The experimental ideal

Regression analysis

Assume for the moment that the treatment effect is constant[†], i.e.,

$$\mathbf{Y}_{1i} - \mathbf{Y}_{0i} = \rho \quad \forall i$$

You'll often hear econometricians say "homogeneous" (vs. "heterogeneous").

The experimental ideal

Regression analysis

Assume for the moment that the treatment effect is constant[†], i.e.,

$$\text{Y}_{1i} - \text{Y}_{0i} = \rho \quad \forall i$$

then we can rewrite

$$\text{Y}_i = \text{Y}_{0i} + D_i (\text{Y}_{1i} - \text{Y}_{0i})$$

You'll often hear econometricians say "homogeneous" (vs. "heterogeneous").

The experimental ideal

Regression analysis

Assume for the moment that the treatment effect is constant[†], i.e.,

$$\mathbf{Y}_{1i} - \mathbf{Y}_{0i} = \rho \quad \forall i$$

then we can rewrite

$$\mathbf{Y}_i = \mathbf{Y}_{0i} + D_i (\mathbf{Y}_{1i} - \mathbf{Y}_{0i})$$

as

$$\mathbf{Y}_i = \underbrace{\alpha}_{=E[\mathbf{Y}_{0i}]} + D_i \underbrace{\rho}_{\mathbf{Y}_{1i} - \mathbf{Y}_{0i}} + \underbrace{\eta_i}_{\mathbf{Y}_{0i} - E[\mathbf{Y}_{0i}]}$$

You'll often hear econometricians say "homogeneous" (vs. "heterogeneous").

The experimental ideal

Regression analysis

$$Y_i = \alpha + D_i \rho + \eta_i$$

Now write out the conditional expectation of Y_i for both levels of D_i

The experimental ideal

Regression analysis

$$Y_i = \alpha + D_i \rho + \eta_i$$

Now write out the conditional expectation of Y_i for both levels of D_i

$$E[Y_i \mid D_i = 1] =$$

The experimental ideal

Regression analysis

$$Y_i = \alpha + D_i \rho + \eta_i$$

Now write out the conditional expectation of Y_i for both levels of D_i

$$E[Y_i \mid D_i = 1] = E[\alpha + \rho + \eta_i \mid D_i = 1]$$

The experimental ideal

Regression analysis

$$Y_i = \alpha + D_i \rho + \eta_i$$

Now write out the conditional expectation of Y_i for both levels of D_i

$$E[Y_i \mid D_i = 1] = E[\alpha + \rho + \eta_i \mid D_i = 1] = \alpha + \rho + E[\eta_i \mid D_i = 1]$$

The experimental ideal

Regression analysis

$$Y_i = \alpha + D_i \rho + \eta_i$$

Now write out the conditional expectation of Y_i for both levels of D_i

$$E[Y_i \mid D_i = 1] = E[\alpha + \rho + \eta_i \mid D_i = 1] = \alpha + \rho + E[\eta_i \mid D_i = 1]$$

$$E[Y_i \mid D_i = 0]$$

The experimental ideal

Regression analysis

$$Y_i = \alpha + D_i \rho + \eta_i$$

Now write out the conditional expectation of Y_i for both levels of D_i

$$E[Y_i \mid D_i = 1] = E[\alpha + \rho + \eta_i \mid D_i = 1] = \alpha + \rho + E[\eta_i \mid D_i = 1]$$

$$E[Y_i \mid D_i = 0] = E[\alpha + \eta_i \mid D_i = 0]$$

The experimental ideal

Regression analysis

$$Y_i = \alpha + D_i \rho + \eta_i$$

Now write out the conditional expectation of Y_i for both levels of D_i

$$E[Y_i \mid D_i = 1] = E[\alpha + \rho + \eta_i \mid D_i = 1] = \alpha + \rho + E[\eta_i \mid D_i = 1]$$

$$E[Y_i \mid D_i = 0] = E[\alpha + \eta_i \mid D_i = 0] = \alpha + E[\eta_i \mid D_i = 0]$$

The experimental ideal

Regression analysis

$$Y_i = \alpha + D_i \rho + \eta_i$$

Now write out the conditional expectation of Y_i for both levels of D_i

$$E[Y_i \mid D_i = 1] = E[\alpha + \rho + \eta_i \mid D_i = 1] = \alpha + \rho + E[\eta_i \mid D_i = 1]$$

$$E[Y_i \mid D_i = 0] = E[\alpha + \eta_i \mid D_i = 0] = \alpha + E[\eta_i \mid D_i = 0]$$

Take the difference...

$$E[Y_i \mid D_i = 1] - E[Y_i \mid D_i = 0]$$

The experimental ideal

Regression analysis

$$Y_i = \alpha + D_i \rho + \eta_i$$

Now write out the conditional expectation of Y_i for both levels of D_i

$$E[Y_i | D_i = 1] = E[\alpha + \rho + \eta_i | D_i = 1] = \alpha + \rho + E[\eta_i | D_i = 1]$$

$$E[Y_i | D_i = 0] = E[\alpha + \eta_i | D_i = 0] = \alpha + E[\eta_i | D_i = 0]$$

Take the difference...

$$\begin{aligned} & E[Y_i | D_i = 1] - E[Y_i | D_i = 0] \\ &= \rho + \underbrace{E[\eta_i | D_i = 1] - E[\eta_i | D_i = 0]}_{\text{Selection bias}} \end{aligned}$$

The experimental ideal

Regression analysis

$$E[Y_i \mid D_i = 1] - E[Y_i \mid D_i = 0] = E[\eta_i | D_i = 1] - E[\eta_i | D_i = 0]$$

Again, our estimate of the **treatment effect** (ρ) is only going to be as good as our ability to shut down the **selection bias**.

Selection bias in regression model: $E[\eta_i | D_i = 1] - E[\eta_i | D_i = 0]$

Selection bias here should remind you a lot of

The experimental ideal

Regression analysis

$$E[Y_i | D_i = 1] - E[Y_i | D_i = 0] = E[\eta_i | D_i = 1] - E[\eta_i | D_i = 0]$$

Again, our estimate of the **treatment effect** (ρ) is only going to be as good as our ability to shut down the **selection bias**.

Selection bias in regression model: $E[\eta_i | D_i = 1] - E[\eta_i | D_i = 0]$

Selection bias here should remind you a lot of **omitted-variable bias**.

There is something in our disturbance η_i that is affecting Y_i and is also correlated with D_i .

The experimental ideal

Regression analysis

$$E[Y_i | D_i = 1] - E[Y_i | D_i = 0] = E[\eta_i | D_i = 1] - E[\eta_i | D_i = 0]$$

Again, our estimate of the **treatment effect** (ρ) is only going to be as good as our ability to shut down the **selection bias**.

Selection bias in regression model: $E[\eta_i | D_i = 1] - E[\eta_i | D_i = 0]$

Selection bias here should remind you a lot of **omitted-variable bias**.

There is something in our disturbance η_i that is affecting Y_i and is also correlated with D_i .

In other metrics-y words: Our treatment D_i is endogenous.

The experimental ideal

Solutions and covariates

Selection bias in regression model: $E[\eta_i | \mathbf{D}_i = 1] - E[\eta_i | \mathbf{D}_i = 0]$

The experimental ideal

Solutions and covariates

Selection bias in regression model: $E[\eta_i | \mathbf{D}_i = 1] - E[\eta_i | \mathbf{D}_i = 0]$

As before, if we randomly assign D_i , then selection bias disappears.

The experimental ideal

Solutions and covariates

Selection bias in regression model: $E[\eta_i | \mathbf{D}_i = 1] - E[\eta_i | \mathbf{D}_i = 0]$

As before, if we randomly assign D_i , then selection bias disappears.

Another potential route to identification is to condition on covariates in the hopes that they "take care of" the relationship between \mathbf{D}_i and whatever is in our disturbance η_i .

The experimental ideal

Solutions and covariates

Selection bias in regression model: $E[\eta_i | \mathbf{D}_i = 1] - E[\eta_i | \mathbf{D}_i = 0]$

As before, if we randomly assign D_i , then selection bias disappears.

Another potential route to identification is to condition on covariates in the hopes that they "take care of" the relationship between \mathbf{D}_i and whatever is in our disturbance η_i .

Without very clear reasons explaining how you know you've controlled for the "bad variation", clean and convincing identification on this path is going to be challenging.

The experimental ideal

Covariates

That said, covariates can help with two things:

1. Even experiments may need **conditioning/controls**: The STAR experiment was random *within school*—not across schools.
2. Covariates can soak up unexplained variation—**increasing precision.**

The experimental ideal

Covariates

That said, covariates can help with two things:

1. Even experiments may need **conditioning/controls**: The STAR experiment was random *within school*—not across schools.
2. Covariates can soak up unexplained variation—**increasing precision**.

Now that we've seen regression can analyze experiments, let's estimate the STAR example...

Table 2.2.2, MHE

Explanatory variable	1	2	3
<i>Small class</i>	4.82 (2.19)	5.37 (1.26)	5.36 (1.21)
<i>Regular + aide</i>	0.12 (2.23)	0.29 (1.13)	0.53 (1.09)
<i>White/Asian</i>			8.35 (1.35)
<i>Female</i>			4.48 (0.63)
<i>Free lunch</i>			-13.15 (0.77)
<i>School F.E.</i>	F	T	T

The omitted level is *Regular* (with part-time aide).

Table 2.2.2, MHE

Explanatory variable	1	2	3
<i>Small class</i>	4.82 (2.19)	5.37 (1.26)	5.36 (1.21)
<i>Regular + aide</i>	0.12 (2.23)	0.29 (1.13)	0.53 (1.09)
<i>White/Asian</i>			8.35 (1.35)
<i>Female</i>			4.48 (0.63)
<i>Free lunch</i>			-13.15 (0.77)
<i>School F.E.</i>	F	T	T

Results without other controls are very similar to the difference in means.

Table 2.2.2, MHE

Explanatory variable	1	2	3
<i>Small class</i>	4.82 (2.19)	5.37 (1.26)	5.36 (1.21)
<i>Regular + aide</i>	0.12 (2.23)	0.29 (1.13)	0.53 (1.09)
<i>White/Asian</i>			8.35 (1.35)
<i>Female</i>			4.48 (0.63)
<i>Free lunch</i>			-13.15 (0.77)
<i>School F.E.</i>	F	T	T

School FEs enforce the experiment's design and increase precision.

Table 2.2.2, MHE

Explanatory variable	1	2	3
<i>Small class</i>	4.82 (2.19)	5.37 (1.26)	5.36 (1.21)
<i>Regular + aide</i>	0.12 (2.23)	0.29 (1.13)	0.53 (1.09)
<i>White/Asian</i>			8.35 (1.35)
<i>Female</i>			4.48 (0.63)
<i>Free lunch</i>			-13.15 (0.77)
<i>School F.E.</i>	F	T	T

Additional controls slightly increase precision.

Object types/classes

As we discussed last class, R revolves around objects, e.g., `test ← 123.`

Object types/classes

As we discussed last class, R revolves around objects, e.g., `test ← 123.`

Objects have types/classes.

Object types/classes

As we discussed last class, R revolves around objects, e.g., `test ← 123.`

Objects have types/classes.

- `1`, `2/3`, and are `numeric`.

Object types/classes

As we discussed last class, R revolves around objects, e.g., `test ← 123.`

Objects have types/classes.

- `1`, `2/3`, and are `numeric`.
- `"Hello"` and `'cruel world'` are both `character`.

Object types/classes

As we discussed last class, R revolves around objects, e.g., `test ← 123`.

Objects have types/classes.

- `1`, `2/3`, and are `numeric`.
- `"Hello"` and `'cruel world'` are both `character`.
- `TRUE`, `T`, `FALSE`, and `F` are `logical` (as is the result of `3 > 2`).

Object types/classes

As we discussed last class, R revolves around objects, e.g., `test ← 123`.

Objects have types/classes.

- `1`, `2/3`, and are `numeric`.
- `"Hello"` and `'cruel world'` are both `character`.
- `TRUE`, `T`, `FALSE`, and `F` are `logical` (as is the result of `3 > 2`).

The `class(x)` function tells you the class of object `x`.

Structure

In addition to having types/classes, objects have some type of structure.

- `1:3`, `c(1, 2)`, and `seq(2, 8, 2)` each produce a `numeric`-class vector.

Structure

In addition to having types/classes, objects have some type of structure.

- `1:3`, `c(1, 2)`, and `seq(2, 8, 2)` each produce a `numeric`-class `vector`.
- `c("Alright", "already")` produces a `vector` of `character` class.

Structure

In addition to having types/classes, objects have some type of structure.

- `1:3`, `c(1, 2)`, and `seq(2, 8, 2)` each produce a `numeric`-class `vector`.
- `c("Alright", "already")` produces a `vector` of `character` class.
- `c(1, 3, T, "Hello")` produces a `vector` of `character` class.

Structure

In addition to having types/classes, objects have some type of structure.

- `1:3`, `c(1, 2)`, and `seq(2, 8, 2)` each produce a `numeric`-class `vector`.
- `c("Alright", "already")` produces a `vector` of `character` class.
- `c(1, 3, T, "Hello")` produces a `vector` of `character` class.
- `matrix(data = 1:15, ncol = 5)` creates a `matrix` with class from `data`.

Structure

In addition to having types/classes, objects have some type of structure.

- `1:3`, `c(1, 2)`, and `seq(2, 8, 2)` each produce a `numeric`-class `vector`.
- `c("Alright", "already")` produces a `vector` of `character` class.
- `c(1, 3, T, "Hello")` produces a `vector` of `character` class.
- `matrix(data = 1:15, ncol = 5)` creates a `matrix` with class from `data`.
- `data.frame(x = 1:2, y = c("a", "b"), z = T)` produces a `data.frame` with three columns and two rows. The first column (`x`) is `numeric`; the second column (`y`) is `character`, and the third column (`z`) is `logical`.

R

Our matrix

```
matrix(data = 1:15, ncol = 5)
```

```
#>      [,1] [,2] [,3] [,4] [,5]
#> [1,]     1     4     7    10    13
#> [2,]     2     5     8    11    14
#> [3,]     3     6     9    12    15
```

R

Our matrix

```
matrix(data = 1:15, ncol = 5)
```

```
#>      [,1] [,2] [,3] [,4] [,5]
#> [1,]     1    4    7   10   13
#> [2,]     2    5    8   11   14
#> [3,]     3    6    9   12   15
```

Our first data.frame!

```
data.frame(x = 1:2, y = c("a", "b"),
```

```
#>   x y   z
#> 1 1 a TRUE
#> 2 2 b TRUE
```

R

Our matrix

```
matrix(data = 1:15, ncol = 5)
```

```
#>      [,1] [,2] [,3] [,4] [,5]
#> [1,]     1    4    7   10   13
#> [2,]     2    5    8   11   14
#> [3,]     3    6    9   12   15
```

Our first data.frame!

```
data.frame(x = 1:2, y = c("a", "b"),
```

```
#>   x y   z
#> 1 1 a TRUE
#> 2 2 b TRUE
```

Notice how R helps 'fill' out the columns when lengths don't match.

Packages

Straight out of the box, R has a ton of useful features, but it really gets its power from the additional packages (libraries) that users create.

- **Open-source greatness** Users find needs and create amazing solutions.
- **Caveat utilitor** There are a lot of packages, each with a lot of functions. Mistakes can happen.
- **Open-source greatness₂** Again, R is open source: Check the code!

Packages

Straight out of the box, R has a ton of useful features, but it really gets its power from the additional packages (libraries) that users create.

- **Open-source greatness** Users find needs and create amazing solutions.
- **Caveat utilitor** There are a lot of packages, each with a lot of functions. Mistakes can happen.
- **Open-source greatness₂** Again, R is open source: Check the code! (Maybe. Sometimes it's very hard.)

Packages

Straight out of the box, R has a ton of useful features, but it really gets its power from the additional packages (libraries) that users create.

- **Open-source greatness** Users find needs and create amazing solutions.
- **Caveat utilitor** There are a lot of packages, each with a lot of functions. Mistakes can happen.
- **Open-source greatness₂** Again, R is open source: Check the code! (Maybe. Sometimes it's very hard.)

Examples `ggplot2` (plotting), `dplyr` (data work that can link with SQL), `sf` and `raster` (geospatial work), `lfe` (high-dimensional fixed-effect regression), `data.table` (fast and efficient data work)

Installing packages

Once you find a function/package that you need to install,[†] you'll typically install it via `install.packages("newAmazingPackage")`.^{††}

We'll use the package `dplyr` throughout the course. Let's install it.

```
# Install 'dplyr' package  
install.packages("dplyr")
```

Aside Notice the comment above the actual code (R uses `#` for comments).

[†] Tool #1: Google. ^{††} The quotation marks are important.

Installing packages

Once you find a function/package that you need to install,[†] you'll typically install it via `install.packages("newAmazingPackage")`.^{††}

We'll use the package `dplyr` throughout the course. Let's install it.

```
# Install 'dplyr' package  
install.packages("dplyr")
```

Aside Notice the comment above the actual code (R uses `#` for comments). While not necessary for R to work, comments are necessary for research.

[†] Tool #1: Google. ^{††} The quotation marks are important.

Using packages

Once you install a package, it is on your machine.

You don't need to install it again—though you probably should update them from time to time.

Using packages

Once you install a package, it is on your machine.

You don't need to install it again—though you probably should update them from time to time.

To **load a package**, use the `library(package)` function[†], e.g., to load `dplyr`

```
# Load 'dplyr'  
library(dplyr)
```

[†] Notice `library()` doesn't *need* quotation marks. I know...

Using packages

Once you install a package, it is on your machine.

You don't need to install it again—though you probably should update them from time to time.

To **load a package**, use the `library(package)` function[†], e.g., to load `dplyr`

```
# Load 'dplyr'  
library(dplyr)
```

Now all functions contained in `dplyr` are available (until you close R).

[†] Notice `library()` doesn't *need* quotation marks. I know...

Package management

All of this installing, loading, updating, checking-for-existance-and-then-loading can get old.

As can typing `library(pacakge1)`, `library(package2)`, ...

Package management

All of this installing, loading, updating, checking-for-existance-and-then-loading can get old.

As can typing `library(pacakge1)`, `library(package2)`, ...

[Enter] The `pacman` package... for package management, of course.

Package management

All of this installing, loading, updating, checking-for-existance-and-then-loading can get old.

As can typing `library(package1)`, `library(package2)`, ...

[Enter] The `pacman` package... for package management, of course.

After installing (`install.packages("pacman")`), you can

- Install and load packages via `p_load(package1, ..., packageN)`
- Update packages via `p_update()`

The `p_load` paradigm is especially helpful for collaborations or projects across multiple machines.

Table of contents

Admin

1. Schedule
2. Review

R

1. Object types/classes
2. Packages

Experimental ideal

1. Neyman/Rubin framework
2. Selection
3. Experiments
4. Example: Training programs
5. Example: STAR experiment
 - Mean differences
 - Dummy variables
 - Regression analysis
 - Covariates