

Tools for Planning and Analyzing Randomized Controlled Trials and A/B Tests

Johann Gagnon-Bartsch, Adam Sales, Duy Pham,
Charlotte Mann, and Jaylin Lowe

Department of Statistics, University of Michigan &
Department of Mathematical Sciences, Worcester Polytechnic Institute

Educational Data Mining

July 14, 2024



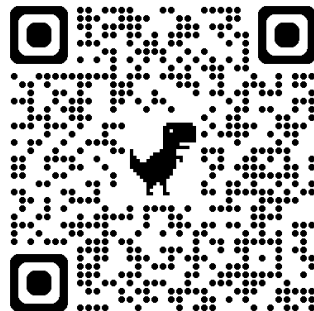
While we are getting settled...

Follow the steps at

<https://tinyurl.com/edmrct-setup>



Workshop Resources



Github site

to get everything ready to follow along in *RStudio*!

Today's Plan



- 1:30-1:45 – Part I: Conceptual Overview
- 1:45-3:00 – Part II: Estimating Effects with RCT Data
- 3:00-3:30 – Part III: Incorporating Auxiliary Data
- 3:30-3:45 – Break 15 min
- 3:45-4:15 – Part IV: Treatment Effect Heterogeneity
- 4:15-5:00 – Part V: Planning Experiments

What You Need



- Tutorial website: `https://tinyurl.com/edmrct`
- RStudio
- Clone repo from Github:
`https://github.com/manncz/edm-rct-tutorial/`

Tutorial Structure

We will be alternating between:

- Conceptual descriptions of the methods
- Detailed walk-throughs of the software
- Opportunities for you to run analyses yourself, with our help



Please feel free to ask questions at any time!

- Calling out (unmute yourself if on Zoom)
- Zoom chat
- Any other way you can think of to get our attention

Conceptual Overview

Estimating Effects with RCT Data

Incorporating Auxiliary Data

Break

Treatment Effect Heterogeneity

Planning Experiments

Experiments in Education Research

“Experiment” = “RCT” = “Randomized Controlled Trial”



- Randomize subjects (students? teachers? schools?) between condition
- Expose subjects to their randomized conditions
- Measure outcome(s) of interest

Experiments in Education Research

“Experiment” = “RCT” = “Randomized Controlled Trial”



- Randomize subjects (students? teachers? schools?) between condition
- Expose subjects to their randomized conditions
- Measure outcome(s) of interest
- Associations between condition and outcomes are causal

Experiments in Education Research

“Experiment” = “RCT” = “Randomized Controlled Trial”



- Randomize subjects (students? teachers? schools?) between condition
- Expose subjects to their randomized conditions
- Measure outcome(s) of interest
- Associations between condition and outcomes are causal
- Typical examples:
 - A/B tests in online learning
 - Field trials of (say) new curriculum vs. business as usual

Example 1: ASSISTments ETrials

ASSISTments TestBed Introduction

```
graph LR; 1[1. Start With Your Research Idea] --> 2[2. Create Your Problem Set]; 2 --> 3[3. Deliver to Teachers and then Students]; 3 --> 4[4. Analyze Data]; 4 --> 5[5. Write]; 5 --> RP((Your Research Paper))
```

- 1. Start With Your Research Idea**
Develop an intervention to study.
To use the WPI Subject Pool Submit your idea to WPI.
- 2. Create Your Problem Set**
Create your Problem Set in ASSISTments.
One problem set for the whole study is preferable.
- 3. Deliver to Teachers and then Students**
For the WPI subject pool the problem set will be approved and made available.
There are ways to deliver using LMS and personal links
- 4. Analyze Data**
Approved studies will get a weekly e-mail when there are more students with the data.
Or use the Data Request Form.
- 5. Write**
Write up your results and submit it for publication.

Your Research Paper

2:30 / 3:00

YouTube



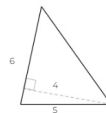
Integrated A/B test platform

Example 1: ASSISTments E-trials

- Question: Text or video hints?
- Outcome: Complete skill builder?
- $n = 683$ middle school students

Problem 2 ⓘ

What is the area of the triangle?



(Images not to scale)



Example 1: ASSISTments E-trials

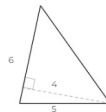
- Question: Text or video hints?
- Outcome: Complete skill builder?
- $n = 683$ middle school students

Results,

- Video: 205/337 (61%) completed
- Text: 193/346 (56%) completed

Problem 2

What is the area of the triangle?

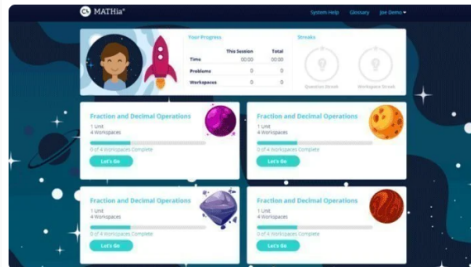


(Images not to scale)



Example II: Cognitive Tutor Effectiveness Trial

- 73 High Schools & 74 Middle Schools in 7 states
- Similar schools paired
- In each pair, one school randomized to treatment, one to control
- Algebra I students in Trt school used CTAI, Control school used business as usual
- All students took a posttest at the end of the year



Example II: Cognitive Tutor Effectiveness Trial



Results

	Average Posttest			
	Middle		High	
	Year 1	Year 2	Year 1	Year 2
Control	17.4	16.9	10.3	9.7
Treatment	14.3	15.2	10.1	10.6

Scientific Goals

1. What is the average effect of [intervention] on [outcome]?



2. How Does the effect vary?

Scientific Goals

1. What is the average effect of [intervention] on [outcome]?
 - “Intervention” AKA “Treatment” (the thing you’re randomizing)
 - Contrast between 2+ conditions
 - E.g. access to ChatGPT hint vs teacher-written hint vs no hint
 - For today: focus on 2 conditions, “Treatment” vs “Control”
 - (those labels may be arbitrary)
2. How Does the effect vary?



Scientific Goals

1. What is the average effect of [intervention] on [outcome]?
 - “Intervention” AKA “Treatment” (the thing you’re randomizing)
 - Contrast between 2+ conditions
 - E.g. access to ChatGPT hint vs teacher-written hint vs no hint
 - For today: focus on 2 conditions, “Treatment” vs “Control”
 - (those labels may be arbitrary)
 - “Outcome”
 - Scalar quantity that the intervention might affect
 - E.g. student correctness on the next problem (0 or 1)
2. How Does the effect vary?



Scientific Goals

1. What is the average effect of [intervention] on [outcome]?
 - “Intervention” AKA “Treatment” (the thing you’re randomizing)
 - Contrast between 2+ conditions
 - E.g. access to ChatGPT hint vs teacher-written hint vs no hint
 - For today: focus on 2 conditions, “Treatment” vs “Control”
 - (those labels may be arbitrary)
 - “Outcome”
 - Scalar quantity that the intervention might affect
 - E.g. student correctness on the next problem (0 or 1)
 - “Average Effect” ...to be defined soon!
2. How Does the effect vary?



Scientific Goals

1. What is the average effect of [intervention] on [outcome]?

- “Intervention” AKA “Treatment” (the thing you’re randomizing)
 - Contrast between 2+ conditions
 - E.g. access to ChatGPT hint vs teacher-written hint vs no hint
 - For today: focus on 2 conditions, “Treatment” vs “Control”
 - (those labels may be arbitrary)
- “Outcome”
 - Scalar quantity that the intervention might affect
 - E.g. student correctness on the next problem (0 or 1)
- “Average Effect” ...to be defined soon!

2. How Does the effect vary?

- From one (type of) student to the next
- From one context to the next



Statistical Goals

1. Get the most out of your data: more data \rightarrow better estimation!!
2. ...Without making unnecessary assumptions
3. Easily
4. Design better experiments to start with



Statistical Goals

1. Get the most out of your data: more data \rightarrow better estimation!!
 - Baseline covariate data
 - Historical user data
2. ...Without making unnecessary assumptions
3. Easily
4. Design better experiments to start with



Statistical Goals

1. Get the most out of your data: more data \rightarrow better estimation!!
 - Baseline covariate data
 - Historical user data
2. ...Without making unnecessary assumptions
 - “Design-based” methods
 - NO assumptions about confounding, models, etc. etc.
3. Easily
4. Design better experiments to start with





1. Get the most out of your data: more data → better estimation!!
 - Baseline covariate data
 - Historical user data
2. ...Without making unnecessary assumptions
 - “Design-based” methods
 - NO assumptions about confounding, models, etc. etc.
3. Easily
 - i.e. without a PhD in statistics
 - Use our software package :)
4. Design better experiments to start with

Types of Variables: Baseline Covariates

- Fixed at baseline
- Unaffected by treatment



Types of Variables: Baseline Covariates

- Fixed at baseline
- Unaffected by treatment

Uses:

- More precise estimates
- Explore effect variation



Example: Covariates in ASSISTments

- Log data. For each previous skillbuilder,
 - Completed skill builder?
 - # problems attempted / completed?
 - Time to mastery
 - ...
- Demographic data



Example: Covariates in ASSISTments

- Log data. For each previous skillbuilder,
 - Completed skill builder?
 - # problems attempted / completed?
 - Time to mastery
 - ...
- Demographic data

Don't use post-treatment variables!



Auxiliary Data

- Covariate and outcome data from *other* subjects
- Often: historical data



- Covariate and outcome data from *other* subjects
- Often: historical data
- Requirements
 - Separate sample from RCT
 - (some of the) same covariate data as for RCT subjects
 - similar outcome data as RCT



- Covariate and outcome data from *other* subjects
- Often: historical data
- Requirements
 - Separate sample from RCT
 - (some of the) same covariate data as for RCT subjects
 - similar outcome data as RCT

Uses:

- More precise estimates
- Planning experiments





Estimate treatment effects



Estimate treatment effects
Using all our data



Estimate treatment effects

Using all our data

- Covariates (even high-dimensional)
- Auxiliary/historical data



Estimate treatment effects

Using all our data

- Covariates (even high-dimensional)
- Auxiliary/historical data

Without bias or extra assumptions

Conceptual Overview

Estimating Effects with RCT Data

Incorporating Auxiliary Data

Break

Treatment Effect Heterogeneity

Planning Experiments

Potential Outcomes (Neyman-Rubin)

Consider a randomized experiment with:

- N participants
- One treatment group, one control group



Potential Outcomes (Neyman-Rubin)

- The outcome depends on treatment.



Potential Outcomes (Neyman-Rubin)



- The outcome depends on treatment.
If the coin had landed the other way, the outcome may have been different.

Potential Outcomes (Neyman-Rubin)



- The outcome depends on treatment.
If the coin had landed the other way, the outcome may have been different.
- Each subject has two **potential outcomes**.

Potential Outcomes (Neyman-Rubin)



- The outcome depends on treatment.
If the coin had landed the other way, the outcome may have been different.
- Each subject has two **potential outcomes**.
One for treatment, one for control.

Potential Outcomes (Neyman-Rubin)



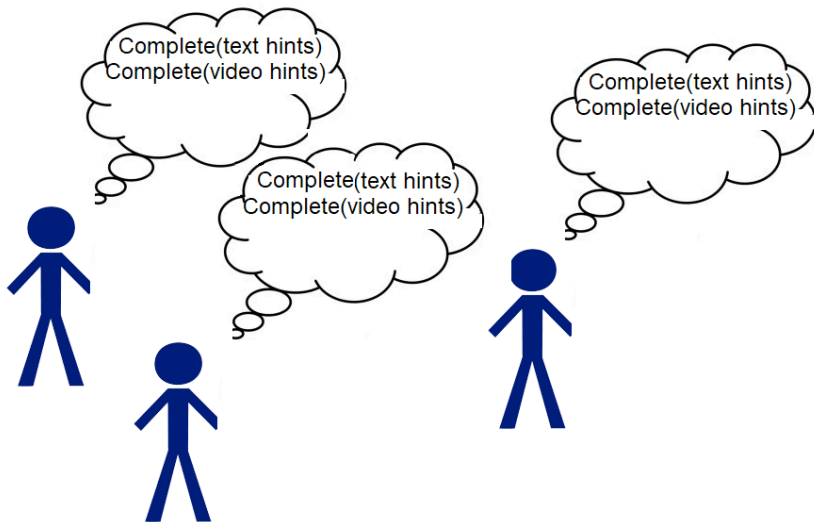
- The outcome depends on treatment.
If the coin had landed the other way, the outcome may have been different.
- Each subject has two **potential outcomes**.
One for treatment, one for control.
- We only ever observe **one** potential outcome.

Potential Outcomes (Neyman-Rubin)

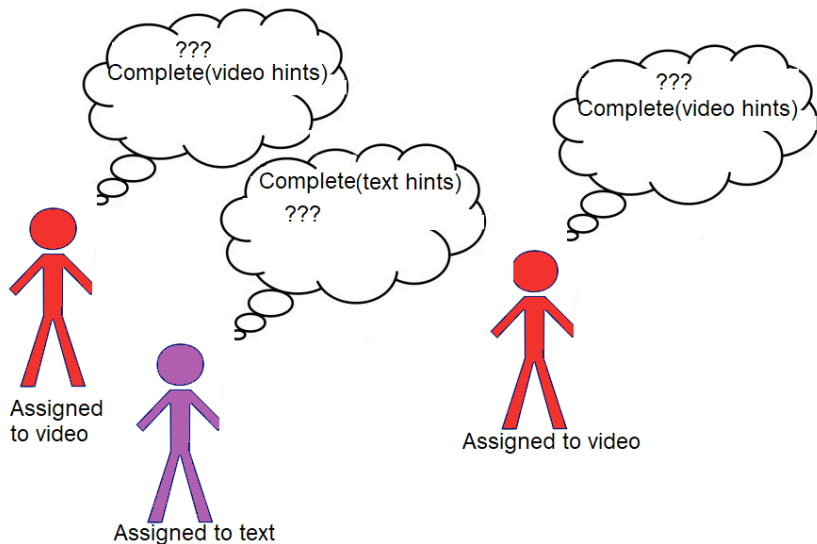


- The outcome depends on treatment.
If the coin had landed the other way, the outcome may have been different.
- Each subject has two **potential outcomes**.
One for treatment, one for control.
- We only ever observe **one** potential outcome.
The other is a counterfactual.

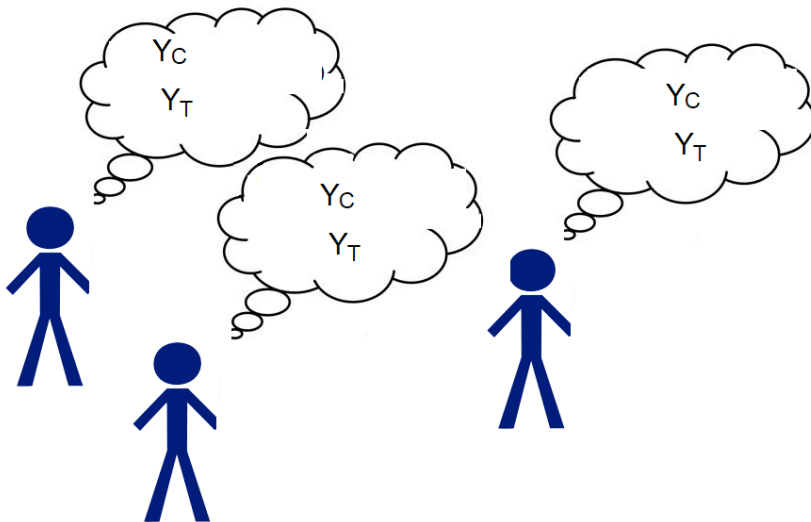
Potential Outcomes (Neyman-Rubin)



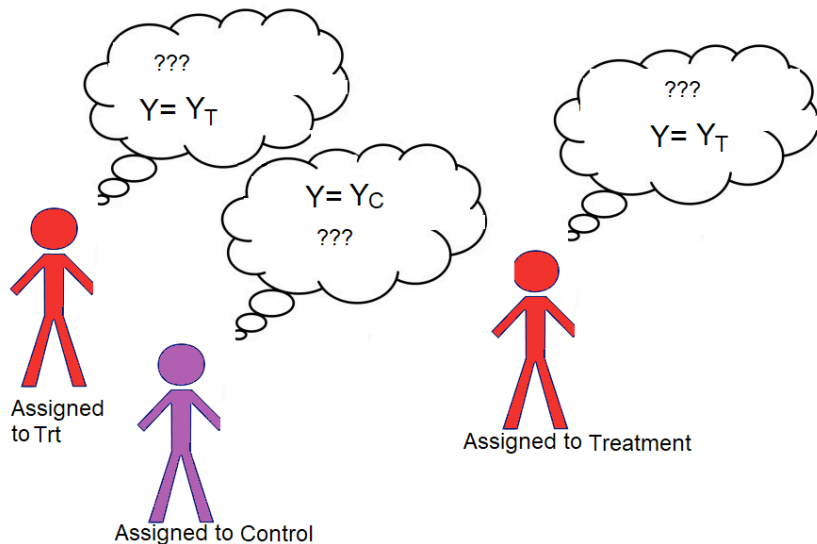
Potential Outcomes (Neyman-Rubin)



Potential Outcomes (Neyman-Rubin)



Potential Outcomes (Neyman-Rubin)



Potential Outcomes (Neyman-Rubin)

- For each participant i there are two potential outcomes, y_i^t and y_i^c



Potential Outcomes (Neyman-Rubin)

- For each participant i there are two potential outcomes, y_i^t and y_i^c
- Potential outcomes are **fixed** values, not random



Potential Outcomes (Neyman-Rubin)

- For each participant i there are two potential outcomes, y_i^t and y_i^c
- Potential outcomes are **fixed** values, not random
- Let T_i be the treatment assignment of unit i

$$T_i = \begin{cases} 1, & \text{Unit } i \text{ is assigned to treatment} \\ 0, & \text{Unit } i \text{ is assigned to control} \end{cases}$$



Potential Outcomes (Neyman-Rubin)

- For each participant i there are two potential outcomes, y_i^t and y_i^c
- Potential outcomes are **fixed** values, not random
- Let T_i be the treatment assignment of unit i

$$T_i = \begin{cases} 1, & \text{Unit } i \text{ is assigned to treatment} \\ 0, & \text{Unit } i \text{ is assigned to control} \end{cases}$$

- Let Y_i be the observed outcome for unit i . If unit i is assigned to treatment, we observe y_i^t ; otherwise, we observe y_i^c :

$$Y_i = \begin{cases} y_i^c & \text{if } T_i = 0 \\ y_i^t & \text{if } T_i = 1 \end{cases}$$



Individual and Average Treatment Effects

- The individual treatment effect is

$$\tau_i = y_i^t - y_i^c$$



Individual and Average Treatment Effects

- The individual treatment effect is

$$\tau_i = y_i^t - y_i^c$$

- The individual treatment effect is **never observed**.



Individual and Average Treatment Effects

- The individual treatment effect is

$$\tau_i = y_i^t - y_i^c$$

- The individual treatment effect is **never observed**.
- The average treatment effect (ATE) is

$$\bar{\tau} = \frac{1}{N} \sum_{i=1}^N \tau_i$$



Individual and Average Treatment Effects

- The individual treatment effect is

$$\tau_i = y_i^t - y_i^c$$

- The individual treatment effect is **never observed**.
- The average treatment effect (ATE) is

$$\bar{\tau} = \frac{1}{N} \sum_{i=1}^N \tau_i$$

- The average treatment effect can be estimated.



Individual and Average Treatment Effects

- The individual treatment effect is

$$\tau_i = y_i^t - y_i^c$$

- The individual treatment effect is **never observed**.
- The average treatment effect (ATE) is

$$\bar{\tau} = \frac{1}{N} \sum_{i=1}^N \tau_i$$

- The average treatment effect can be estimated.
- Also: average effects for subgroups of subjects (more later)



Estimating Average Treatment Effects

“The fundamental problem of causal inference”

- We only observe one potential outcome for each subject
 - For treatment subjects y^t
 - For control, y^c



Estimating Average Treatment Effects

“The fundamental problem of causal inference”

- We only observe one potential outcome for each subject
 - For treatment subjects y^t
 - For control, y^c
- One potential outcome is always missing
- We need to *impute* the missing potential outcome



Estimating Average Treatment Effects

“The fundamental problem of causal inference”



- We only observe one potential outcome for each subject
 - For treatment subjects y^t
 - For control, y^c
- One potential outcome is always missing
- We need to *impute* the missing potential outcome
- Two approaches to imputation:
 1. Use randomization: unbiased, but imprecise

“The fundamental problem of causal inference”



- We only observe one potential outcome for each subject
 - For treatment subjects y^t
 - For control, y^c
- One potential outcome is always missing
- We need to *impute* the missing potential outcome
- Two approaches to imputation:
 1. Use randomization: unbiased, but imprecise
 2. Use covariates & model: biased, but precise

“The fundamental problem of causal inference”



- We only observe one potential outcome for each subject
 - For treatment subjects y^t
 - For control, y^c
- One potential outcome is always missing
- We need to *impute* the missing potential outcome
- Two approaches to imputation:
 1. Use randomization: unbiased, but imprecise
 2. Use covariates & model: biased, but precise
 3. Our approach: use both!

Our Method, in a Nutshell

Step 1:

Train algorithms to predict y^c, y^t as a function of covariates

$$f^c : \mathbf{X} \rightarrow y^c \text{ (use data from ctl group)}$$

$$f^t : \mathbf{X} \rightarrow y^t \text{ (use data from trt group)}$$



Our Method, in a Nutshell

Step 1:

Train algorithms to predict y^c, y^t as a function of covariates

$$f^c : \mathbf{X} \rightarrow y^c \text{ (use data from ctl group)}$$

$$f^t : \mathbf{X} \rightarrow y^t \text{ (use data from trt group)}$$

Step 2:

Use algorithms to get imputations:

$$\hat{y}_i^c = f^c(X_i)$$

$$\hat{y}_i^t = f^t(X_i)$$



Our Method, in a Nutshell

Step 1:

Train algorithms to predict y^c, y^t as a function of covariates

$$f^c : \mathbf{X} \rightarrow y^c \text{ (use data from ctl group)}$$

$$f^t : \mathbf{X} \rightarrow y^t \text{ (use data from trt group)}$$

Step 2:

Use algorithms to get imputations:

$$\hat{y}_i^c = f^c(X_i)$$

$$\hat{y}_i^t = f^t(X_i)$$

Step 3: Calculate \hat{m}_i : weighted average of \hat{y}_i^c and \hat{y}_i^t



Our Method, in a Nutshell

Step 1:

Train algorithms to predict y^c, y^t as a function of covariates

$$f^c : \mathbf{X} \rightarrow y^c \text{ (use data from ctl group)}$$

$$f^t : \mathbf{X} \rightarrow y^t \text{ (use data from trt group)}$$

Step 2:

Use algorithms to get imputations:

$$\hat{y}_i^c = f^c(X_i)$$

$$\hat{y}_i^t = f^t(X_i)$$

Step 3: Calculate \hat{m}_i : weighted average of \hat{y}_i^c and \hat{y}_i^t

Step 4:

Use randomization-based method to estimate effects on $Y - \hat{m}$ instead of Y



Important Caveat

For this to be *strictly* unbiased, we need:

\hat{m}_i independent of T_i



Important Caveat

For this to be *strictly* unbiased, we need:

$$\hat{m}_i \text{ independent of } T_i$$

Since Y_i is a function of T_i , that means we need:

$$\hat{y}^c \text{ and } \hat{y}^t \text{ independent of } Y_i$$



Important Caveat

For this to be *strictly* unbiased, we need:

$$\hat{m}_i \text{ independent of } T_i$$

Since Y_i is a function of T_i , that means we need:

$$\hat{y}^c \text{ and } \hat{y}^t \text{ independent of } Y_i$$

We can't use i 's data to train f^c and f^t !



Important Caveat

For this to be *strictly* unbiased, we need:

$$\hat{m}_i \text{ independent of } T_i$$

Since Y_i is a function of T_i , that means we need:

$$\hat{y}^c \text{ and } \hat{y}^t \text{ independent of } Y_i$$

We can't use i 's data to train f^c and f^t !

Solution: re-train f^c and f^t for each subject i , leaving out i 's data



Important Caveat

For this to be *strictly* unbiased, we need:

$$\hat{m}_i \text{ independent of } T_i$$

Since Y_i is a function of T_i , that means we need:

$$\hat{y}^c \text{ and } \hat{y}^t \text{ independent of } Y_i$$

We can't use i 's data to train f^c and f^t !

Solution: re-train f^c and f^t for each subject i , leaving out i 's data

“Leave-One-Out Potential Outcomes” or “LOOP”



Sticks and Stones May Break my Bones, but Bad Models Won't Hurt Me

- What if f^c and f^t are totally wrong and bad??



Sticks and Stones May Break my Bones, but Bad Models Won't Hurt Me

- What if f^c and f^t are totally wrong and bad??
- Estimate will still be unbiased!



Sticks and Stones May Break my Bones, but Bad Models Won't Hurt Me



- What if f^c and f^t are totally wrong and bad??
- Estimate will still be unbiased!
- Standard errors, p-values, and confidence intervals will still be valid!

Sticks and Stones May Break my Bones, but Bad Models Won't Hurt Me



- What if f^c and f^t are totally wrong and bad??
- Estimate will still be unbiased!
- Standard errors, p-values, and confidence intervals will still be valid!
- (core of inference is based on randomization)

Sticks and Stones May Break my Bones, but Bad Models Won't Hurt Me



- What if f^c and f^t are totally wrong and bad??
- Estimate will still be unbiased!
- Standard errors, p-values, and confidence intervals will still be valid!
- (core of inference is based on randomization)
- Covariate adjustment won't help much

Sticks and Stones May Break my Bones, but Bad Models Won't Hurt Me



- What if f^c and f^t are totally wrong and bad??
- Estimate will still be unbiased!
- Standard errors, p-values, and confidence intervals will still be valid!
- (core of inference is based on randomization)
- Covariate adjustment won't help much
- In moderate/large samples, it won't hurt either!

Digression: What about old-fashioned regression?

Regression method:

Fit model:

$$Y_i = \beta_0 + \beta_1 T_i + \beta_2 X_{1i} + \beta_3 X_{2i} + \dots$$

Estimated effect: $\hat{\beta}_2$



Digression: What about old-fashioned regression?

Regression method:

Fit model:

$$Y_i = \beta_0 + \beta_1 T_i + \beta_2 X_{1i} + \beta_3 X_{2i} + \dots$$

Estimated effect: $\hat{\beta}_2$

Problem: What if the model is false?

- E.g. Y isn't linear in covariates
- E.g. What if there should be interactions?



Digression: What about old-fashioned regression?

Regression method:

Fit model:

$$Y_i = \beta_0 + \beta_1 T_i + \beta_2 X_{1i} + \beta_3 X_{2i} + \dots$$

Estimated effect: $\hat{\beta}_2$

Problem: What if the model is false?

- E.g. Y isn't linear in covariates
- E.g. What if there should be interactions?

Good news: $\hat{\beta}$ is *approximately unbiased in large samples*



Digression: What about old-fashioned regression?



Why our method?

1. *Exactly unbiased in any sample*
2. Use *any* algorithm for f^c, f^t
 - High dimensional covariates
 - Flexible for non-linearity, interactions

Digression: What about old-fashioned regression?



Why our method?

1. *Exactly unbiased in any sample*
2. Use *any* algorithm for f^c, f^t
 - High dimensional covariates
 - Flexible for non-linearity, interactions
 - \Rightarrow better imputations
 - \Rightarrow better effect estimates

Digression: What about old-fashioned regression?



Why our method?

1. *Exactly unbiased in any sample*
2. Use *any* algorithm for f^c, f^t
 - High dimensional covariates
 - Flexible for non-linearity, interactions
 - \Rightarrow better imputations
 - \Rightarrow better effect estimates
 - We recommend random forest

What You Need for Our Method



1. Randomized treatment variable
2. Outcome variable
3. Covariates
4. What is the experimental design?

This is not a promise.

What You Need for Our Method



1. Randomized treatment variable
2. Outcome variable
3. Covariates
4. What is the experimental design?

One last digression¹: experimental designs

¹This is not a promise.

The Two Questions of Experimental Design



1. Who or What is being randomized?
2. How are they being randomized?

Who/What is being randomized?



- Individual randomization
- Cluster or Group randomization

How are they being randomized?



- What's the probability each unit is assigned to treatment?
- How does one unit's assignment affect other units?

Examples We'll Cover



- Individual randomization
 - Bernoulli
 - Paired
- Cluster randomization
 - Paired

- ASSISTments E-Trials A/B test
 - Students are randomized individually
 - Students are randomized independently
 - \Rightarrow Bernoulli



- ASSISTments E-Trials A/B test
 - Students are randomized individually
 - Students are randomized independently
 - \Rightarrow Bernoulli
- Cognitive Tutor Effectiveness Study
 - *Schools* are randomized
 - Randomization is within pairs
 - (if your school is randomized to treatment, its pair *must* be randomized to control)
 - \Rightarrow paired cluster design



To be implemented (hopefully) soon:

- “Completely randomized design”
 - At the outset, fix # randomized to treatment, # randomized to control
 - Now T_i and T_j are dependent!
- Block-randomized design
 - e.g. a separate completely randomized experiment in each classroom
 - Paired designs are a special case



Other Designs

To be implemented (hopefully) soon:

- “Completely randomized design”
 - At the outset, fix # randomized to treatment, # randomized to control
 - Now T_i and T_j are dependent!
- Block-randomized design
 - e.g. a separate completely randomized experiment in each classroom
 - Paired designs are a special case

Probably won't get to for a while:

- Bandit designs
 - Probability i is assigned to treatment depends on previous subjects' outcomes





Estimating Effects in Practice

Installation:

- You will need to install the package from Github using the *devtools* package in R
- e.g. `install_github("mannncz/dRCT")`



Primary Functions:

`loop(Y, Tr, Z, pred, p, ...)`

`p_loop(Y, Tr, Z, pred, P, n, ...)`

Covariate Adjustment with Bernoulli Randomized Trails (LOOP)

loop(Y , Tr , Z , $pred$, p , ...)

- Y : outcome vector
- Tr : treatment assignment vector
- Z : matrix of covariates
- $pred$: interpolation algorithm
- p : probability of treatment
- ...: optional inputs for interpolation algorithm





pred

- *loop_rf*
- *loop_ols*
- *loop_glm*

Covariate Adjustment with Paired Trails (P-LOOP)

$p_loop(Y, Tr, Z, pred, P, n, \dots)$

- Y : outcome vector
- Tr : treatment assignment vector
- Z : matrix of covariates
- $pred$: interpolation algorithm
- P : vector of pair assignments
- n : optional vector of cluster sizes
- \dots : optional inputs for interpolation algorithm



pred

- *p_ols_po*
- *p_ols_v12*
- *p_ols_interp*
- *p_rf_po*
- *p_rf_v12*
- *p_rf_interp*



Real Data Example: Texas School Data

- AEIS: School-level data from Texas Education Agency from 2003-2011
- > 3,000 schools
- TAKS (standardized test) passing rates
- Thousands of additional possible covariates



Real Data Example: Synthetic School-Level RCT

- Inspired by the Cognitive Tutor Algebra I study (Pane et al. 2014)



Real Data Example: Synthetic School-Level RCT

- Inspired by the Cognitive Tutor Algebra I study (Pane et al. 2014)
- **RCT Sample:** 50 Texas middle schools
- **Treatment:** Alternative 8th grade mathematics curriculum
- **Design:** Schools randomly assigned to implement new curriculum or continue standard in the 2007/8 school year



Real Data Example: Synthetic School-Level RCT



- Inspired by the Cognitive Tutor Algebra I study (Pane et al. 2014)
- **RCT Sample:** 50 Texas middle schools
- **Treatment:** Alternative 8th grade mathematics curriculum
- **Design:** Schools randomly assigned to implement new curriculum or continue standard in the 2007/8 school year
- **Outcome:** 2008 8th grade math TAKS passing rate

Real Data Example: Synthetic School-Level RCT



- Inspired by the Cognitive Tutor Algebra I study (Pane et al. 2014)
- **RCT Sample:** 50 Texas middle schools
- **Treatment:** Alternative 8th grade mathematics curriculum
- **Design:** Schools randomly assigned to implement new curriculum or continue standard in the 2007/8 school year
- **Outcome:** 2008 8th grade math TAKS passing rate
- **Pretest:** 2007 8th grade math TAKS passing rate



1. Follow along while we talk through *01-explore-aeis-data.Rmd*
2. Work through *02-effect-est.Rmd*
 - Effect estimate for Bernoilli randomized trial
 - Effect estimate for paired randomed trial
 - Effect esitmate for paired cluster randomed trial
3. Flag any of us down as you have questions!

Conceptual Overview

Estimating Effects with RCT Data

Incorporating Auxiliary Data

Break

Treatment Effect Heterogeneity

Planning Experiments

Auxiliary Data

By “Auxiliary Data” we mean a dataset that meets these criteria:

1. Doesn't include data from RCT participants
2. Includes covariate data
3. Includes outcome data



Auxiliary Data

By “Auxiliary Data” we mean a dataset that meets these criteria:

1. Doesn't include data from RCT participants
2. Includes covariate data
3. Includes outcome data

Examples:

- A/B test: historical log data from users who worked on similar modules before the experiment started
- Field trial: Administrative (e.g. SLDS) data from students in schools that were not part of the RCT



Auxiliary Data

By “Auxiliary Data” we mean a dataset that meets these criteria:

1. Doesn't include data from RCT participants
2. Includes covariate data
3. Includes outcome data

Examples:

- A/B test: historical log data from users who worked on similar modules before the experiment started
- Field trial: Administrative (e.g. SLDS) data from students in schools that were not part of the RCT

Note: we have sometimes called this the “remnant”



What use is more data??

- Already imputing potential outcomes with f^c and f^t in LOOP
- f^c and f^t can be flexible, high dimensional
- They are fit to representative data



What use is more data??

- Already imputing potential outcomes with f^c and f^t in LOOP
- f^c and f^t can be flexible, high dimensional
- They are fit to representative data

Limits on f^c and f^t

- RCT sample size might be too small to fit *really* good models
- Human-adaptive modeling: no good!



Example 1: ASSISTments

Covariates:

- Log data. For each previous skillbuilder,
 - Completed skill builder?
 - # problems attempted / completed?
 - Time to mastery
- Demographic data



Example 1: ASSISTments

Covariates:

- Log data. For each previous skillbuilder,
 - Completed skill builder?
 - # problems attempted / completed?
 - Time to mastery
- Demographic data

Auxiliary Data:

- Observational
- Students who were *not* randomized
 - Previous users
 - Current users not assigned to that skillbuilder
- Same covariates available



Observational

RCT

Control

Treatment



Observational

Step 1:

Train Model $\hat{y}(\cdot) : \mathbf{x} \rightarrow Y$

With auxiliary data

RCT

Control

Treatment



Observational

Step 1:

Train Model $\hat{y}(\cdot) : \mathbf{x} \rightarrow Y$
With auxiliary data

Step 2:

Extrapolate
With fitted model & RCT
data

RCT

Control

$$\hat{y}(\mathbf{x}_i)$$

Treatment

$$\hat{y}(\mathbf{x}_j)$$



Observational

Step 1:

Train Model $\hat{y}(\cdot) : \mathbf{x} \rightarrow Y$
With auxiliary data

Step 2:

Extrapolate
With fitted model & RCT
data

Step 3:

Use $\hat{y}(\mathbf{x})$ as a
“super-covariate”

RCT

Control

$$\hat{y}(\mathbf{x}_i)$$

Treatment

$$\hat{y}(\mathbf{x}_j)$$



Basic idea: Use auxiliary-based predictions $\hat{y}(x_i)$ as a *covariate* in the RCT.





Basic idea: Use auxiliary-based predictions $\hat{y}(x_i)$ as a *covariate* in the RCT.

- The function $\hat{y}(\cdot)$ is fit on auxiliary data
- The covariates x are pre-treatment
- $\Rightarrow \hat{y}(x)$ is invariant to treatment assignment



Basic idea: Use auxiliary-based predictions $\hat{y}(x_i)$ as a *covariate* in the RCT.

- The function $\hat{y}(\cdot)$ is fit on auxiliary data
- The covariates x are pre-treatment
- $\Rightarrow \hat{y}(x)$ is invariant to treatment assignment
- $\hat{y}(x)$ might be an amazing covariate



Basic idea: Use auxiliary-based predictions $\hat{y}(x_i)$ as a *covariate* in the RCT.

- The function $\hat{y}(\cdot)$ is fit on auxiliary data
- The covariates x are pre-treatment
- $\Rightarrow \hat{y}(x)$ is invariant to treatment assignment
- $\hat{y}(x)$ might be an amazing covariate
- ...or it might not

Special Prediction Algorithm for LOOP

- If $\hat{y}(x)$ predicts Y really well, we would expect a linear relationship
 - \Rightarrow fit OLS models within LOOP



Special Prediction Algorithm for LOOP



- If $\hat{y}(x)$ predicts Y really well, we would expect a linear relationship
 - \Rightarrow fit OLS models within LOOP
- Maybe $\hat{y}(x)$ isn't that much better than other covariates (or, maybe it's useless)
 - \Rightarrow use random forest within LOOP

Special Prediction Algorithm for LOOP



- If $\hat{y}(x)$ predicts Y really well, we would expect a linear relationship
 - \Rightarrow fit OLS models within LOOP
- Maybe $\hat{y}(x)$ isn't that much better than other covariates (or, maybe it's useless)
 - \Rightarrow use random forest within LOOP
- Let the data decide!
 - *pred=reloop*



Incorporating Auxiliary Data in Practice

Incorporating Auxiliary Information

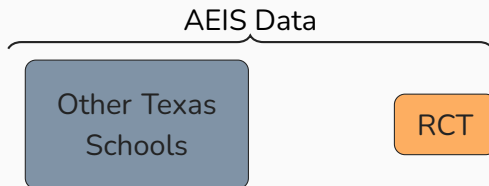
`loop(Y, Tr, Z, pred = reloop, p, yhat, ...)`

- *Y*: outcome vector
- *Tr*: treatment assignment vector
- *Z*: matrix of covariates
- *pred* = *reloop*: specify auxiliary data interpolation algorithm
- *p*: probability of treatment
- *yhat*: vector of auxiliary predictions
- ...: optional inputs for interpolation algorithm



Real Data Example: Texas Schools

- AEIS data includes thousands of schools not in our RCT
- A great setting for integrating auxiliary and RCT data





1. Work through *03-integrate-aux.Rmd*
 - We fit an auxiliary model and generate predictions to input as *yhat*
2. Apply what you learned in *04-effect-estABtest.Rmd*
3. Flag any of us down as you have questions!

Conceptual Overview

Estimating Effects with RCT Data

Incorporating Auxiliary Data

Break

Treatment Effect Heterogeneity

Planning Experiments



Take a 15 minute break!

See you back at 3:45 pm

Conceptual Overview

Estimating Effects with RCT Data

Incorporating Auxiliary Data

Break

Treatment Effect Heterogeneity

Planning Experiments

Heterogeneous Treatment Effects



- Traditionally, most problems in causal inference focus on the ATE as the estimand of interest.
- However, it is not unreasonable that **the same treatment might have different effects on different individuals.**
- As an average, the ATE **cannot** account for such potential variations.

Example



Example

$ITE = +5$ $ITE = +10$ $ITE = +5$ $ITE = +10$ $ITE = -15$ $ITE = +10$



$ITE = -15$ $ITE = +10$ $ITE = +10$ $ITE = +10$ $ITE = +10$ $ITE = +5$



$ITE = -15$ $ITE = +10$ $ITE = +5$ $ITE = -15$ $ITE = +10$ $ITE = +10$



Conditional Average Treatment Effect

- The conditional average treatment effect (CATE) is

$$\tau(x) = \mathbb{E}[\tau_i | X_i = x] = \mathbb{E}[y_i^t - y_i^c | X_i = x]$$





- The conditional average treatment effect (CATE) is

$$\tau(x) = \mathbb{E}[\tau_i | X_i = x] = \mathbb{E}[y_i^t - y_i^c | X_i = x]$$

- The expected treatment effect **conditional on having a specific set of covariate values**.



- The conditional average treatment effect (CATE) is

$$\tau(x) = \mathbb{E}[\tau_i | X_i = x] = \mathbb{E}[y_i^t - y_i^c | X_i = x]$$

- The expected treatment effect **conditional on having a specific set of covariate values**.
- Unlike the ATE, the CATE accounts for the different characteristics of individuals as reflected by the covariates ($X_i = x$).

Estimating the Conditional Average Treatment Effect

- Consider the following decomposition:

$$\begin{aligned}\tau(x) &= \mathbb{E}[\tau_i | X_i = x] = \mathbb{E}[y_i^t - y_i^c | X_i = x] \\ &= \mathbb{E}[y_i^t | X_i = x] - \mathbb{E}[y_i^c | X_i = x] \\ &= \mathbb{E}[Y_i | X_i = x, T_i = 1] - \mathbb{E}[Y_i | X_i = x, T_i = 0]\end{aligned}$$



Estimating the Conditional Average Treatment Effect

- Consider the following decomposition:

$$\begin{aligned}\tau(x) &= \mathbb{E}[\tau_i | X_i = x] = \mathbb{E}[y_i^t - y_i^c | X_i = x] \\ &= \mathbb{E}[y_i^t | X_i = x] - \mathbb{E}[y_i^c | X_i = x] \\ &= \mathbb{E}[Y_i | X_i = x, T_i = 1] - \mathbb{E}[Y_i | X_i = x, T_i = 0]\end{aligned}$$

- Regress observed outcomes on the covariates and treatment assignments.
- Estimate the CATE as the difference between estimated treatment ($T_i = 1$) and control ($T_i = 0$) outcomes.



Interaction Term(s) with Treatment Assignment

- Consider the following linear regression model:

$$\mu(x, t) = \mathbb{E}[Y_i | X_i = x, T_i = t] = \beta_0 + \beta_1 x + \beta_2 t + \beta_3(x \cdot t)$$



Interaction Term(s) with Treatment Assignment

- Consider the following linear regression model:

$$\mu(x, t) = \mathbb{E}[Y_i | X_i = x, T_i = t] = \beta_0 + \beta_1 x + \beta_2 t + \beta_3 (x \cdot t)$$



- Thus, we can estimate the CATE as:

$$\hat{\tau}(x) = \hat{\mu}(x, 1) - \hat{\mu}(x, 0) = (\hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 + \hat{\beta}_3 x) - (\hat{\beta}_0 + \hat{\beta}_1 x) = \hat{\beta}_2 + \hat{\beta}_3 x$$

Interaction Term(s) with Treatment Assignment

- Consider the following linear regression model:

$$\mu(x, t) = \mathbb{E}[Y_i | X_i = x, T_i = t] = \beta_0 + \beta_1 x + \beta_2 t + \beta_3 (x \cdot t)$$



- Thus, we can estimate the CATE as:

$$\hat{\tau}(x) = \hat{\mu}(x, 1) - \hat{\mu}(x, 0) = (\hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 + \hat{\beta}_3 x) - (\hat{\beta}_0 + \hat{\beta}_1 x) = \hat{\beta}_2 + \hat{\beta}_3 x$$

- Linear parametric model: **Ease of interpretation and statistical inference.**
- However, strict linearity is **also restrictive and thus potentially imprecise.**



- Broadly speaking, there are two categories:
 - **Meta-learners:** Methods that leverage off-the-shelf machine learning algorithms to **indirectly** estimate the CATE by learning its components. *Example:* S-, T-, & X-Learner (Künzel et al. 2019).
 - **Direct Estimators:** Methods specifically designed from the ground up to **directly** estimate the CATE. *Example:* Causal Forest (Wager and Athey 2018).



- Broadly speaking, there are two categories:
 - **Meta-learners:** Methods that leverage off-the-shelf machine learning algorithms to **indirectly** estimate the CATE by learning its components. *Example:* S-, T-, & X-Learner (Künzel et al. 2019).
 - **Direct Estimators:** Methods specifically designed from the ground up to **directly** estimate the CATE. *Example:* Causal Forest (Wager and Athey 2018).
- More flexible machine learning methods mean **potentially more precise estimates** but at the cost of **interpretability and inference**.

Leveraging Estimates of Individual Treatment Effects



- **Accuracy** is undeniably important. Knowing how much a treatment will help or harm someone to the best of our ability is always good...
- ...but so is **Interpretability**, especially since we are conditioning on and considering specific values of covariates.
- In other words, **how to have our cake and eat (a bit of) it** too?

Leveraging Estimates of Individual Treatment Effects

- Recall the definition of the CATE:

$$\tau(x) = \mathbb{E}[\tau_i | X_i = x] = \mathbb{E}[y_i^t - y_i^c | X_i = x]$$

- If we know the true individual treatment effects, we can fit a parametric model estimating τ_i with X_i to estimate the CATE.



Leveraging Estimates of Individual Treatment Effects

- Recall the definition of the CATE:

$$\tau(x) = \mathbb{E}[\tau_i | X_i = x] = \mathbb{E}[y_i^t - y_i^c | X_i = x]$$

- If we know the true individual treatment effects, we can fit a parametric model estimating τ_i with X_i to estimate the CATE.
- Unfortunately, τ_i is **never available** as we exclusively observe **either** y_i^t **or** y_i^c .



Leveraging Estimates of Individual Treatment Effects



- Recall the definition of the CATE:

$$\tau(x) = \mathbb{E}[\tau_i | X_i = x] = \mathbb{E}[y_i^t - y_i^c | X_i = x]$$

- If we know the true individual treatment effects, we can fit a parametric model estimating τ_i with X_i to estimate the CATE.
- Unfortunately, τ_i is **never available** as we exclusively observe **either** y_i^t **or** y_i^c .
- **However**, we have the estimate $\hat{\tau}_i$ from the LOOP estimator:

$$\hat{\tau}_i = \{Y_i - [(1-p)\hat{y}_i^t + p\hat{y}_i^c]\} \frac{T_i - p}{p(1-p)}$$

Leveraging Estimates of Individual Treatment Effects

- Given a Bernoulli randomization, $\hat{\tau}_i$ is an **unbiased** estimate of τ_i ($\mathbb{E}[\hat{\tau}_i] = \tau_i$).
In addition, if the imputations y_i^t and y_i^c are accurate, $\hat{\tau}_i$ will also be a **precise** estimate ($\hat{\tau}_i \approx \tau_i$) (Wu and Gagnon-Bartsch 2018).



Leveraging Estimates of Individual Treatment Effects

- Given a Bernoulli randomization, $\hat{\tau}_i$ is an **unbiased** estimate of τ_i ($\mathbb{E}[\hat{\tau}_i] = \tau_i$).
In addition, if the imputations y_i^t and y_i^c are accurate, $\hat{\tau}_i$ will also be a **precise** estimate ($\hat{\tau}_i \approx \tau_i$) (Wu and Gagnon-Bartsch 2018).
- We first estimate the true individual treatment effects using the LOOP estimator.



Leveraging Estimates of Individual Treatment Effects

- Given a Bernoulli randomization, $\hat{\tau}_i$ is an **unbiased** estimate of τ_i ($\mathbb{E}[\hat{\tau}_i] = \tau_i$). In addition, if the imputations y_i^t and y_i^c are accurate, $\hat{\tau}_i$ will also be a **precise** estimate ($\hat{\tau}_i \approx \tau_i$) (Wu and Gagnon-Bartsch 2018).
- We first estimate the true individual treatment effects using the LOOP estimator.
- We can regress $\hat{\tau}_i$ on X_i and estimate the CATE as $\tau(x) = \mathbb{E}[\hat{\tau}_i | X_i = x]$, using the estimated effect $\hat{\tau}_i$ as a proxy for the true effect τ_i :

$$\hat{\tau}(x) = \hat{\alpha}_1 + \hat{\alpha}_2 x$$



Leveraging Estimates of Individual Treatment Effects

- Given a Bernoulli randomization, $\hat{\tau}_i$ is an **unbiased** estimate of τ_i ($\mathbb{E}[\hat{\tau}_i] = \tau_i$). In addition, if the imputations y_i^t and y_i^c are accurate, $\hat{\tau}_i$ will also be a **precise** estimate ($\hat{\tau}_i \approx \tau_i$) (Wu and Gagnon-Bartsch 2018).
- We first estimate the true individual treatment effects using the LOOP estimator.
- We can regress $\hat{\tau}_i$ on X_i and estimate the CATE as $\tau(x) = \mathbb{E}[\hat{\tau}_i | X_i = x]$, using the estimated effect $\hat{\tau}_i$ as a proxy for the true effect τ_i :

$$\hat{\tau}(x) = \hat{\alpha}_1 + \hat{\alpha}_2 x$$

- Notice that this formulation is (roughly) equivalent to the interaction model:

$$\hat{\tau}(x) = \hat{\mu}(x, 1) - \hat{\mu}(x, 0) = (\hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 + \hat{\beta}_3 x) - (\hat{\beta}_0 + \hat{\beta}_1 x) = \hat{\beta}_2 + \hat{\beta}_3 x$$

Leveraging Estimates of Individual Treatment Effects

- If the LOOP estimator's requirements are satisfied, the estimates will be unbiased – **even with poor-fitting models.**



Leveraging Estimates of Individual Treatment Effects

- If the LOOP estimator's requirements are satisfied, the estimates will be unbiased – **even with poor-fitting models**.
- Thus, we do not have to worry about bias from the regression in the second stage carrying over or, worse, amplifying bias from the first.



Leveraging Estimates of Individual Treatment Effects



- If the LOOP estimator's requirements are satisfied, the estimates will be unbiased – **even with poor-fitting models**.
- Thus, we do not have to worry about bias from the regression in the second stage carrying over or, worse, amplifying bias from the first.
- **Flexible model(s) to estimate the ITE in the first:** More **precise** than a strictly linear model with interaction(s).
- **Parametric model to estimate the CATE in the second:** More **interpretable** than a powerful but non-parametric model.

You can use the function `getITE` to retrieve the ITE estimates from a LOOP estimator.





You can use the function `getITE` to retrieve the ITE estimates from a LOOP estimator.

- This function will work for an estimator built with or without auxiliary data, which allows us to improve precision further.



You can use the function `getITE` to retrieve the ITE estimates from a LOOP estimator.

- This function will work for an estimator built with or without auxiliary data, which allows us to improve precision further.
- However, it is currently only for Bernoulli-randomized experiments.



You can use the function `getITE` to retrieve the ITE estimates from a LOOP estimator.

- This function will work for an estimator built with or without auxiliary data, which allows us to improve precision further.
- However, it is currently only for Bernoulli-randomized experiments.
- Once you have retrieve the estimates, choose your favorite model and do some regressing!



1. Work through *05-heterogeneousEffects.Rmd*
 - We fit retrieve ITE estimates from the model in *04-effect-estABtest.Rmd*.
 - We then estimate the CATE by regressing these estimates on the covariates.
2. Flag any of us down as you have questions!

Conceptual Overview

Estimating Effects with RCT Data

Incorporating Auxiliary Data

Break

Treatment Effect Heterogeneity

Planning Experiments

What You Need

- We'll be using the *dRCTpower* package to plan experiments
- Main function is *run_app*
- You can download the package in R using the following commands:

```
install.packages("devtools")  
devtools::install_github("jaylinlowe/dRCTpower")
```

- We will be using the *aux_dat_small.csv* file from the Github repo





How to choose a sample size for our experiment, particularly if auxiliary data will be incorporated?



- Incorporating auxiliary data in our analysis can improve precision, meaning we can have a smaller sample size with the same power



- Incorporating auxiliary data in our analysis can improve precision, meaning we can have a smaller sample size with the same power
- Gain in precision is determined by how predictive a model fit on the auxiliary data is for the RCT



- Incorporating auxiliary data in our analysis can improve precision, meaning we can have a smaller sample size with the same power
- Gain in precision is determined by how predictive a model fit on the auxiliary data is for the RCT
- But....we don't have the RCT data!



1. Break auxiliary dataset into subgroups



1. Break auxiliary dataset into subgroups
2. For each subgroup, treat it as the RCT and the rest of the data as the auxiliary data



1. Break auxiliary dataset into subgroups
2. For each subgroup, treat it as the RCT and the rest of the data as the auxiliary data
3. Calculate the required sample size under this framework



Large auxiliary dataset that:

- is substantially larger than the RCT will be



Large auxiliary dataset that:

- is substantially larger than the RCT will be
- has covariates



Large auxiliary dataset that:

- is substantially larger than the RCT will be
- has covariates
- has the same outcome of interest as the RCT



- Method is only plausible if it's reasonable to assume the RCT looks like some subgroup of the auxiliary data, even if we don't know what subgroup that is



- Method is only plausible if it's reasonable to assume the RCT looks like some subgroup of the auxiliary data, even if we don't know what subgroup that is
- Dangerous to assume RCT looks like any one subgroup



- Method is only plausible if it's reasonable to assume the RCT looks like some subgroup of the auxiliary data, even if we don't know what subgroup that is
- Dangerous to assume RCT looks like any one subgroup
- Dangerous to choose most optimistic option

General Power Calculations

$$n = 2\sigma^2 \frac{(\xi_{1-\alpha/2} + \xi_{1-\beta})^2}{\Delta_A^2}$$



General Power Calculations

$$n = 2\sigma^2 \frac{(\xi_{1-\alpha/2} + \xi_{1-\beta})^2}{\Delta_A^2}$$



- $\xi_{1-\alpha/2}$ is the critical value obtained from a normal distribution for Type I error equal to α .

General Power Calculations

$$n = 2\sigma^2 \frac{(\xi_{1-\alpha/2} + \xi_{1-\beta})^2}{\Delta_A^2}$$



- $\xi_{1-\alpha/2}$ is the critical value obtained from a normal distribution for Type I error equal to α .
- $\xi_{1-\beta}$ is the critical value from a normal distribution for Type II error rate β .

General Power Calculations

$$n = 2\sigma^2 \frac{(\xi_{1-\alpha/2} + \xi_{1-\beta})^2}{\Delta_A^2}$$



- $\xi_{1-\alpha/2}$ is the critical value obtained from a normal distribution for Type I error equal to α .
- $\xi_{1-\beta}$ is the critical value from a normal distribution for Type II error rate β .
- Δ_A is the effect size, typically 20% of the standard deviation of the outcome in the population

General Power Calculations

$$n = 2\sigma^2 \frac{(\xi_{1-\alpha/2} + \xi_{1-\beta})^2}{\Delta_A^2}$$



- $\xi_{1-\alpha/2}$ is the critical value obtained from a normal distribution for Type I error equal to α .
- $\xi_{1-\beta}$ is the critical value from a normal distribution for Type II error rate β .
- Δ_A is the effect size, typically 20% of the standard deviation of the outcome in the population
- σ^2 is the true variance of the outcome in the population, typically replaced with an estimate from a sample

- We replace σ^2 with an estimate from each subgroup



Our Modification



- We replace σ^2 with an estimate from each subgroup
- Shiny app gives two estimates, one if you were to use auxiliary data in analysis, and one without

Our Modification



- We replace σ^2 with an estimate from each subgroup
- Shiny app gives two estimates, one if you were to use auxiliary data in analysis, and one without
- "Without auxiliary data" estimate is variance of outcome for that subgroup



- We replace σ^2 with an estimate from each subgroup
- Shiny app gives two estimates, one if you were to use auxiliary data in analysis, and one without
- "Without auxiliary data" estimate is variance of outcome for that subgroup
- "With auxiliary data" estimate is variance of the residuals, $(y_i - \hat{y}_i)$, where \hat{y}_i are out-of-bag predictions from model

Defining Subgroups

Three options:

1. Categorical Variable

- Divide based on levels of categorical variable
- Can create your own categorical variables



Defining Subgroups

Three options:

1. Categorical Variable

- Divide based on levels of categorical variable
- Can create your own categorical variables

2. Numerical Variable

- Divide into 10 (adjustable) equally sized groups
- May need to divide into fewer if there isn't enough variation



Defining Subgroups

Three options:

1. Categorical Variable

- Divide based on levels of categorical variable
- Can create your own categorical variables

2. Numerical Variable

- Divide into 10 (adjustable) equally sized groups
- May need to divide into fewer if there isn't enough variation

3. Best-Worst Case Scenario

- Divide based on how predictive we expect the auxiliary model to be for that group
- Good starting point





Shiny App Demo



References

Gagnon-Bartsch, Johann A., Adam C. Sales, Edward Wu, Anthony F. Botelho, John A. Erickson, Luke W. Miratrix and Neil T. Heffernan. 2023. "Precise unbiased estimation in randomized experiments using auxiliary observational data." *Journal of Causal Inference* 11(1):20220011.

URL: <https://www.degruyter.com/document/doi/10.1515/jci-2022-0011/html>

Künzel, Sören R, Jasjeet S Sekhon, Peter J Bickel and Bin Yu. 2019. "Metalearners for estimating heterogeneous treatment effects using machine learning." *Proceedings of the National Academy of Sciences* 116(10):4156–4165.

Lowe, Jaylin, Charlotte Mann, Jiaying Wang, Adam Sales and Johann Gagnon-Bartsch. Forthcoming. "Power Calculations for Randomized Controlled Trials with Auxiliary Observational Data." *EDM* 2024 .



Mann, Charlotte, Jiaying Wang, Adam Sales and Johann Gagnon-Bartsch.

Forthcoming. "Using Publicly Available Auxiliary Data to Improve Precision of Treatment Effect Estimation in a Randomized Efficacy Trial." *EDM 2024* .

Pane, John F., Beth Ann Griffin, Daniel F. McCaffrey and Rita Karam. 2014.

"Effectiveness of Cognitive Tutor Algebra I at Scale." *Educational Evaluation and Policy Analysis* 36(2):127–144.

URL: <https://doi.org/10.3102/0162373713507480>



Pham, Duy, Kirk Vanacore, Adam Sales and Johann Gagnon-Bartsch.

Forthcoming. "LOOL: Towards Personalization with Flexible Robust Estimation of Heterogeneous Treatment Effects." *EDM 2024* .

Sales, Adam C, Ethan B Prihar, Johann A Gagnon-Bartsch and Neil T Heffernan.

2023. "Using Auxiliary Data to Boost Precision in the Analysis of A/B Tests on an Online Educational Platform: New Data and New Results." *arXiv preprint arXiv:2306.06273* .

Wager, Stefan and Susan Athey. 2018. "Estimation and Inference of Heterogeneous Treatment Effects using Random Forests." *Journal of the American Statistical Association* 113(523):1228–1242.

Wu, Edward and Johann A. Gagnon-Bartsch. 2018. "The LOOP Estimator: Adjusting for Covariates in Randomized Experiments." *Evaluation Review* 42(4):458–488. Publisher: SAGE Publications Inc.

URL: <https://doi.org/10.1177/0193841X18808003>

Wu, Edward and Johann A. Gagnon-Bartsch. 2021. "Design-Based Covariate Adjustments in Paired Experiments." *Journal of Educational and Behavioral Statistics* 46(1):109–132. Publisher: American Educational Research Association.

URL: <https://doi.org/10.3102/1076998620941469>

