From experiments to observational data

STSCI / INFO / ILRST 3900

16 Sep 2025

Learning goals for today

At the end of class, you will be able to.

- ► Tie analysis of observational data to an idealized experiment
- ► Ask good questions which
 - ▶ involve treatments that exist (positivity assumption)
 - ► involve precise treatments (consistency assumption)

After class:

► Optional: Hernán, M. 2016.

"Does water kill? A call for less casual causal inferences." Annals of Epidemiology 26(10):674–680.

Logistics

- ► PSET 1 peer review due tonight
- ► First in-class quiz 9/18
- ► PSET 2 due 9/19

- ► Marginal exchangeability: $Y_i^{a=1}$, $Y_i^{a=0} \perp A_i$ holds in conditionally experiments
- ▶ Is almost never true in observational data

- ► Marginal exchangeability: $Y_i^{a=1}$, $Y_i^{a=0} \perp A_i$ holds in conditionally experiments
- ▶ Is almost never true in observational data

- ► Conditional exchangeability: $Y_i^{a=1}$, $Y_i^{a=0} \perp A_i \mid L$ holds in conditionally randomized experiments
- ► We've typically discussed *L* being a single variable, but it could also be a set of variables
- ▶ But does it ever hold in observational data?

What is the effect of college degree on income at age 35

- ▶ $A_i = 1$ if four year college degree; $A_i = 0$ if no college degree
- ► Suppose we have information on parental income
 - $ightharpoonup L_i = 0$: parents have high income
 - ► $L_i = 1$: parents have low income
- ▶ Does conditional exchangeability hold given parental income?

What is the effect of college degree on income at age 35

- ▶ $A_i = 1$ if four year college degree; $A_i = 0$ if no college degree
- Suppose we have information on parental income
 - $ightharpoonup L_i = 0$: parents have high income
 - ▶ $L_i = 1$: parents have low income
- ▶ Does conditional exchangeability hold given parental income?
- ► What additional information would you gather to make conditional exchangeability plausible?

► Even if gathering data was possible for every covariate we want, when do we stop?

- ► Even if gathering data was possible for every covariate we want, when do we stop?
- ► Never 100% sure that conditional exchangeability holds
- ► Is it reasonable?

- ► Even if gathering data was possible for every covariate we want, when do we stop?
- ▶ Never 100% sure that conditional exchangeability holds
- ▶ Is it reasonable?
- ► In observational data, conditional exchangeability is an assumption we make (but can't typically verify)
- Requires expert knowledge

- ► Even if gathering data was possible for every covariate we want, when do we stop?
- ▶ Never 100% sure that conditional exchangeability holds
- ▶ Is it reasonable?
- ► In observational data, conditional exchangeability is an assumption we make (but can't typically verify)
- ► Requires expert knowledge
- ► Causal claims are data + outside knowledge

Formulating causal questions

Asking "good" causal questions involve

- ▶ Positivity condition: Treatments that exist
- ► Consistency: Treatments that are precise
- ► Accounts for interference

Good causal questions involve **treatments that exist**

Employer 1	Employer 2
100 employees	200 employees
Face-to-face interaction	Work in individual offices
75% randomized to vaccine 25% randomized to no vaccine	50% randomized to vaccine 50% randomized to no vaccine

How do you estimate the average effect over all 300 employees?

Employer 1	Employer 2
100 employees	200 employees
Face-to-face interaction	Work in individual offices
100% randomized to vaccine 0% randomized to no vaccine	50% randomized to vaccine 50% randomized to no vaccine

How do you estimate the average effect over all 300 employees?

If units are exchangeable given a confounder L, then to estimate $E(Y^a)$ we need **positivity** to hold

$$\mathsf{P}(A=a\mid \vec{L}=\vec{\ell})>0$$



Source: Wikimedia A, B, C



Source: Wikimedia A, B, C

Would the bulbs in Ithaca bloom if it did not freeze all winter?



Source: Wikimedia A, B, C

Would the bulbs in Ithaca bloom if it did not freeze all winter?

Confounder L Ithaca

Treatment a Did not freeze

Outcome Y^a Blooms?



Source: Wikimedia A, B, C



Source: Wikimedia

Would the bulbs in Ithaca bloom if it did not freeze all winter?

Confounder L Ithaca

Treatment a Did not freeze Outcome Y^a Blooms?

Sarah has no MD training. Would Sarah earn more money if she were a surgeon?



Source: Wikimedia A, B, C



Source: Wikimedia

Would the bulbs in Ithaca bloom if it did not freeze all winter?

Confounder *L* Ithaca

Treatment a Did not freeze Outcome Y^a Blooms?

Sarah has no MD training. Would Sarah earn more money if she were a surgeon?

Confounder L No MD training

Treatment a Surgeon Outcome Y^a Earnings

We can choose causal questions so that positivity holds.

$$P(A=a\mid \vec{L}=\vec{\ell})>0$$

- lacktriangle in each population subgroup $ec{L}=ec{\ell}$
- ▶ only study treatment values *a* that can actually happen

Good causal questions involve **precise treatments**

Consistency.

$$Y = Y^A$$

- 1. holds for precise treatments
- 2. holds with clarity about interference among units

Imagine you are a high school counselor.

A statistician tells you

The probability of receiving a BA in 6 years would be higher if a student initially enrolled in the State University of New York instead of a community college

$$\mathsf{P}\bigg(\mathsf{B}\mathsf{A}^{\mathsf{Enroll} \; \mathsf{in} \; \mathsf{SUNY}}\bigg) > \mathsf{P}\bigg(\mathsf{B}\mathsf{A}^{\mathsf{Enroll} \; \mathsf{in} \; \mathsf{Community} \; \mathsf{College}}\bigg)$$

How would you advise students?







6-year graduation rate

BINGHAMTON UNIVERSITY STATE UNIVERSITY OF NEW YORK

78%





The treatment value Enroll in SUNY is not sufficiently precise

6-year graduation rate

BINGHAMTON
UNIVERSITY
STATE UNIVERSITY OF NEW YORK

83%

78%



74%



The treatment value Enroll in SUNY is not sufficiently precise

 $\mathsf{BA}^{\mathsf{Binghamton}}
eq \mathsf{BA}^{\mathsf{Stony}\ \mathsf{Brook}}$ $eq \mathsf{BA}^{\mathsf{Buffalo}}$ $eq \mathsf{BA}^{\mathsf{Albany}}$

6-year graduation rate

BINGHAMTON
UNIVERSITY
STATE UNIVERSITY OF NEW YORK

83%

78%



74%



The treatment value Enroll in SUNY is not sufficiently precise

 $\mathsf{BA}^{\mathsf{Binghamton}}
eq \mathsf{BA}^{\mathsf{Stony}\ \mathsf{Brook}}$ $eq \mathsf{BA}^{\mathsf{Buffalo}}$ $eq \mathsf{BA}^{\mathsf{Albany}}$

To advise the student, a precise treatment is more helpful

6-year graduation rate

BINGHAMTON UNIVERSITY STATE UNIVERSITY OF NEW YORK

78%



74%



Consistency assumption: $Y = Y^A$

More credible when A is very precise

- it is clear how to run a hypothetical experiment
- ▶ is is clear how to inform policy

Example:

- ▶ if a = SUNY, then Y^a is vague. To which SUNY should you send the student?
- ▶ if a = Binghamton, then Y^a is clearer

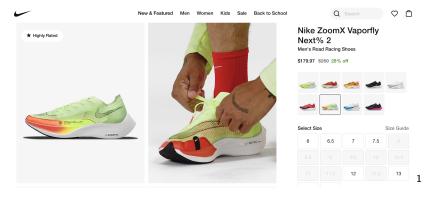
A good read:

Hernán, M. 2016.

"Does water kill? A call for less casual causal inferences."

Annals of Epidemiology 26(10):674-680.

Good causal questions involve clarity about interference



¹Image source: Nike

You and a friend race in your normal shoes.

You and a friend race in your normal shoes. It is extremely close.

You and a friend race in your normal shoes.

It is extremely close.

You barely lose.

You and a friend race in your normal shoes.

It is extremely close.

You barely lose.

$$Y_{You} = Lose$$

You and a friend race in your normal shoes.

It is extremely close.

You barely lose.

$$Y_{You} = Lose$$

What if you had the springy shoes?

You and a friend race in your normal shoes.

It is extremely close.

You barely lose.

$$Y_{You} = Lose$$

What if you had the springy shoes?

$$Y_{\mathsf{You}}^{\mathsf{You} \ \mathsf{wear} \ \mathsf{springy} \ \mathsf{shoes}} = \mathsf{Win}$$

You and a friend race in your normal shoes.

It is extremely close.

You barely lose.

$$Y_{You} = Lose$$

What if you had the springy shoes?

$$Y_{You}^{You \ wear \ springy \ shoes} = Win$$

But what if your friend also wears them?

You and a friend race in your normal shoes.

It is extremely close.

You barely lose.

$$Y_{You} = Lose$$

What if you had the springy shoes?

$$Y_{\mathsf{You}}^{\mathsf{You} \; \mathsf{wear} \; \mathsf{springy} \; \mathsf{shoes}} = \mathsf{Win}$$

But what if your friend also wears them?

$$Y_{\mathsf{You}}^{\mathsf{You}}$$
 wear springy shoes, Your friend wears springy shoes $= \mathsf{Lose}$

$$Y_{\mathsf{You}}^{\mathsf{You}}$$
 wears springy shoes, Your friend wear normal shoes $= \mathsf{Win}$

Good causal questions: In math

We should study treatments that exist

(positivity)

$$\mathsf{P}(A=a\mid \vec{L}=\vec{\ell})>0$$

with potential outcomes that are well-defined

(consistency)

$$Y = Y^A$$

Well-defined potential outcomes involve precise treatments

BA^{Binghamton} instead of BA^{SUNY}

and incorporate interference when it exists

 $Y^{a_{you},a_{your friend}}$ instead of $Y^{a_{you}}$

Learning goals for today

At the end of class, you will be able to.

- ► Tie analysis of observational data to an idealized experiment
- ► Ask good questions which
 - ▶ involve treatments that exist (positivity assumption)
 - ► involve precise treatments (consistency assumption)

After class:

► Optional: Hernán, M. 2016.

"Does water kill? A call for less casual causal inferences." Annals of Epidemiology 26(10):674–680.