

Outline

- 1) Causality
- 2) Potential Outcomes
- 3) Randomized Experiments

Causality

Correlation \neq Causation

So far, course has focused on drawing inferences about probability dist.s

Ex $D_i = 1 \{ \text{Student } i \text{ gets scholarship} \}$

$Y_i = \text{College GPA of student } i$

Given iid sample of (D_i, Y_i) pairs, could

- estimate $\mathbb{E}[Y_i | D_i]$
- get interval for $\theta = \mathbb{E}[Y_i | D_i = 1] - \mathbb{E}[Y_i | D_i = 0]$
- test $H_0: \text{dist}(Y_i | D_i = 1) = \text{dist}(Y_i | D_i = 0)$

But Suppose CI for θ is $[0.7, 0.8]$

Still can't conclude scholarship caused GPA boost

(Why not?)

People often interpret regression coeff.s causally, "all other things equal" - usually w/o justification

We can never draw causal conclusions from
joint dist. alone

But causal questions are often very
interesting, drive policy, etc.

When can we draw causal conclusions?

Two questions about causality:

Hard question: what caused an outcome?

"Why couldn't he find a job?"

- very few things are monocausal
- causes may interact w/ each other

Easier: what effect does an intervention have?

"Does a specific job training program
help people find jobs?"

- more tractable — could imagine experiment
- only meaningful relative to specific ^(control) counterfactual
(same person, no training program)

"Effect of a cause" not "Cause of an effect"

Potential Outcomes

Idea: Define causal effect in terms of outcomes that would happen in alt. universe

Ex (Cont'd):

$Y_i(1)$ = GPA of student i with schol.

$Y_i(0)$ = GPA " w/o schol.

Observed $Y_i = Y_i(D_i)$ (SUTVA)

(Implicit: Y_i only depends on D_i (not D_j for $j \neq i$))

Problem: Only one potential outcome observed

Define average treatment effect (ATE)

$$\text{as } \delta = \mathbb{E}[Y_i(1) - Y_i(0)]$$

Different from θ :

$$\theta = \mathbb{E}[Y_i(1) | D_i = 1] - \mathbb{E}[Y_i(0) | D_i = 0]$$

$$= \underbrace{\mathbb{E}[Y_i(1) - Y_i(0) | D_i = 1]}_{\text{ATE}} + \underbrace{\mathbb{E}[Y_i(0) | D_i = 1] - \mathbb{E}[Y_i(0) | D_i = 0]}$$

Randomized Experiments

selection bias

Suppose D_i assigned at random,
w/o regard to any attributes

$$\text{Then, } E[Y_i(1) | D_i=1] = E[Y_i(1)]$$

$$\Rightarrow \theta = \delta$$

Can use usual two-sample methods,
interpret causally.

For example, can get unbiased
estimate of $\delta = \theta$:

$$\hat{\delta} = \frac{1}{n_1} \sum_{i: D_i=1} Y_i - \frac{1}{n_0} \sum_{i: D_i=0} Y_i$$

Usually can't observe both potential outcomes
in social science applications

Randomization Tests

What if we want to test whether treatment has any effect (on any unit)?

Fisher's sharp null:

$$H_0: Y_i(1) = Y_i(0) \quad \forall i = 1, \dots, n$$

Note: statement about (unobserved) aspect of these n units (no sampling model)

Idea: Use randomness of treatment assignment to our benefit

Under H_0 , $Y_i(D_i^*) = Y_i(D_i)$ under any treatment assignment

We know dist. of $D \Rightarrow$ know dist of (null) any test stat $T(D, Y(D))$

For $b = 1, \dots, B$:

$$D^{*b} = \text{permute}(D)$$

$$\hat{\delta}^{*b} = \frac{1}{n_1} \sum_{i: D_i^{*b} = 1} Y_i - \frac{1}{n_0} \sum_{i: D_i^{*b} = 0} Y_i$$

Under H_0 , $\hat{\delta}, \hat{\delta}^{*1}, \dots, \hat{\delta}^{*B}$ are

exchangeable \Rightarrow use $p = \frac{1}{B+1} (1 + \#\{\hat{\delta}^{*b} \geq \hat{\delta}\})$

Can extend to CI for constant treatment effect:

$$H_0^{\delta}: Y_i(1) - Y_i(0) = \delta \quad \forall i$$

$$\text{Under } H_0, Y_i(D_i^{*b}) - Y_i(D_i) = \delta (D_i^{*b} - D_i)$$

$$\begin{aligned} T^{*b} &= \frac{1}{n_1} \sum_{D_i^{*b} = 1} (Y_i + \delta(1 - D_i)) \\ &\quad - \frac{1}{n_0} \sum_{D_i^{*b} = 0} (Y_i - \delta D_i) \end{aligned}$$

Experimental Design

Often can improve precision of estimator if we randomize more carefully

e.g. Suppose pre-treatment covariate X_i is highly predictive of outcome

say, $X_i = \text{age}$

$Y_i = 1 \{ \text{recovers from illness} \}$

If treatment assigned unif. randomly, treatment group will get more or less older people by chance
Could be main driver of variance!

Better: Match 2 oldest, next 2, ..., 2 youngest

Randomize w/in pairs

→ Different randomization test.

Unconfoundedness

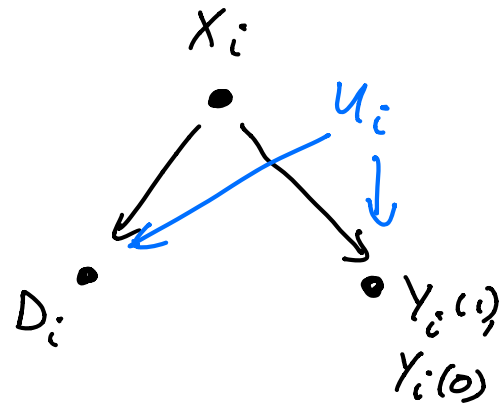
In observational studies, we don't randomise,
just observe world as it is.

Assume $(X_i, D_i, Y_i(0), Y_i(1)) \stackrel{iid}{\sim} P$

pre-treatment
covariate

treatment

P.O.s

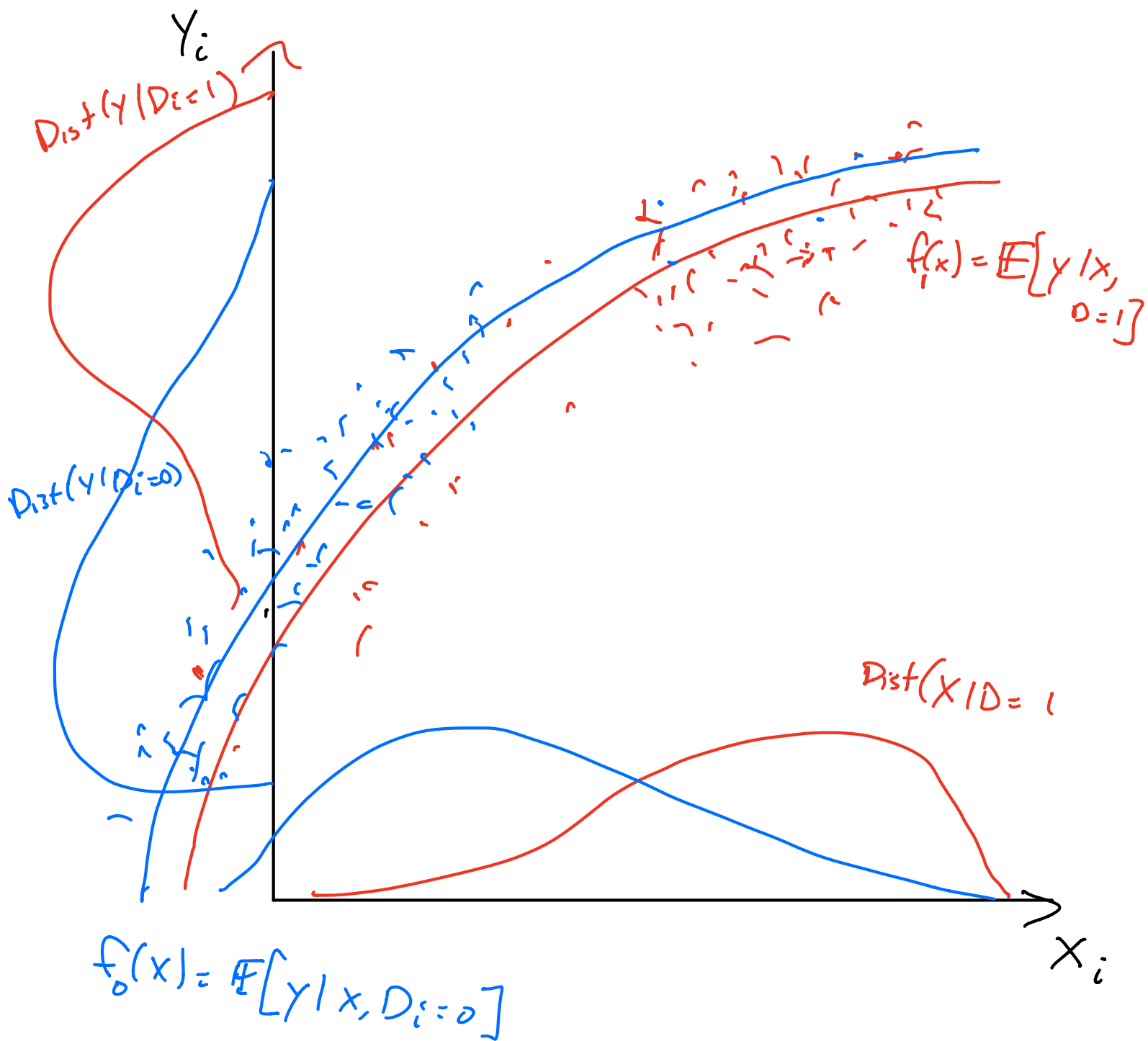


we want $\delta = \mathbb{E}_P[Y_i(1) - Y_i(0)]$

Suppose $(Y_i(1), Y_i(0)) \perp\!\!\!\perp D_i \mid X_i$

then δ becomes identifiable from obs. data.

$$\begin{aligned}\mathbb{E}[Y_i(1) - Y_i(0)] &= \mathbb{E}[\mathbb{E}[Y_i(1) \mid X_i] - \mathbb{E}[Y_i(0) \mid X_i]] \\ &= \mathbb{E}\left[\underbrace{\mathbb{E}[Y_i \mid X_i, D_i=1]}_{\text{observed } Y_i} - \mathbb{E}[Y_i \mid X_i, D_i=0]\right] \\ &= \mathbb{E}[f_1(X_i) - f_0(X_i)] \\ f_d(X_i) &= \mathbb{E}[Y_i \mid X_i, D_i=d] \text{ identifiable.}\end{aligned}$$



Example of Simpson's Paradox: apparent
 'effect' is reversed after we condition
 on something.

Propensity Scores

Sometimes can't randomize ourselves, but
can estimate how "nature" randomized
treatment assignment.

Suppose we observe covariates X_i , know

$$e(x) = P(D_i = 1 \mid X_i = x) \quad \begin{array}{l} \text{Propensity} \\ \text{score function} \end{array}$$

And assume unconfoundedness: $(e(x) \in (0, 1))$

$$(Y_i(0), Y_i(1)) \perp\!\!\!\perp D_i \mid X_i$$

Then also true that $(Y_i(0), Y_i(1)) \perp\!\!\!\perp D_i \mid e(X_i)$

Why?

$$\begin{aligned} P(Y_i \in A, D_i = 1 | e(x_i)) &= E \left[P(Y_i \in A, D_i = 1 | x_i) | e(x_i) \right] \\ &= E \left[P(Y_i \in A | x_i) \underbrace{P(D_i = 1 | x_i)}_{= e(x_i)} | e(x_i) \right] \\ &= e(x_i) E \left[P(Y_i \in A | x_i) | e(x_i) \right] \\ &= e(x_i) P(Y_i \in A | e(x_i)) \end{aligned}$$

We could use regression adjustment

$$\mu_d = E \left[E \left\{ Y_i | e(x_i), D_i = d \right\} \right]$$

But can actually make exact adj.
instead.

Inverse Propensity Weighting

Suppose we actually know

$$e(x) = \mathbb{P}(D_i = 1 \mid X_i = x)$$

Could use it to consistently est. δ :

$$\hat{\delta} = \frac{1}{n} \sum_{D_i=1} \frac{Y_i}{e(x_i)} - \frac{1}{n} \sum_{D_i=0} \frac{Y_i}{1-e(x_i)}$$

not n_1

Informally: for each $(X_i, Y_i(0), Y_i(1))$ we have
an $e(x_i)$ chance of observing $Y_i(1)$,
(otherwise $D_i=0$ and $Y_i(1)$ missing).
 \Rightarrow upweight observed ones accordingly

(works as long as $e(x) \in (0,1) \quad \forall x$)

More formal:

$$\frac{1}{n} \sum_{D_i=1} \frac{Y_i}{e(X_i)} = \frac{1}{n} \sum_{i=1}^n \frac{Y_i(1) D_i}{e(X_i)}$$

$$\mathbb{E} \left[\frac{Y_i(1) D_i}{e(X_i)} \mid X_i \right] = \mathbb{E} [Y_i(1) \mid X_i] \cdot \underbrace{\mathbb{E} \left[\frac{D_i}{e(X_i)} \mid X_i \right]}_{=1}$$

$$\Rightarrow \mathbb{E} \frac{Y_i D_i}{e(X_i)} = \mathbb{E} Y_i(1)$$

Similarly, $\frac{1}{n} \sum_{D_i=0} \frac{Y_i}{1-e(X_i)} = \frac{1}{n} \sum_{i=1}^n \frac{Y_i(0)(1-D_i)}{1-e(X_i)}$

$$\Rightarrow \mathbb{E} \frac{Y_i (1-D_i)}{1-e(X_i)} = \mathbb{E} Y_i(0)$$

$$\text{So, } \mathbb{E} \hat{\delta} = \mathbb{E} (Y_i(1) - Y_i(0)) = \delta$$

Practical issues: if $e(X_i) \approx 0$ or ≈ 1

for some X_i , IPW estimator can have extremely high variance. Called "poor overlap"

Estimated Prop. Scores

Problem: we don't know $e(X_i)$

Could estimate, e.g. using logistic regression

Different e 😞

$$e(x) = \frac{e^{\beta'x}}{1 + e^{\beta'x}}$$

Use
$$\hat{e}(x) = \frac{e^{\hat{\beta}'x}}{1 + e^{\hat{\beta}'x}}$$

Bias $\rightarrow 0$, $\hat{\delta}$ still consistent

If propensity score model is misspecified,

Bias $\not\rightarrow 0$, $\hat{\delta}$ not consistent

Unsurprising since our model was wrong

More surprising: we can (sometimes) correct it!

What is the bias? Assume to

Write $\hat{\mu}_1 = \frac{1}{n} \sum \frac{Y_i D_i}{e(X_i)}$, $\hat{\mu}_0 = \frac{1}{n} \sum \frac{Y_i (1-D_i)}{1-e(X_i)}$

then $\hat{\delta} = \hat{\mu}_1 - \hat{\mu}_0$

$$\text{Bias}(\hat{\mu}_1) = \mathbb{E} \left[\frac{Y_i(1) D_i}{\hat{e}(X_i)} - Y_i(1) \right]$$

Assume for simplicity $\hat{e}(\cdot) \perp (X_i, D_i, Y_i)_{i=1}^n$
(e.g. data splitting)

$$\begin{aligned} \mathbb{E} \left[\frac{Y_i(1) D_i}{\hat{e}(X_i)} - Y_i(1) \mid X_i \right] &= \mathbb{E} \left[Y_i(1) \cdot \frac{D_i - \hat{e}(X_i)}{\hat{e}(X_i)} \mid X_i \right] \\ &= \underbrace{\mathbb{E}[Y_i(1) \mid X_i]}_{f_1(X_i)} \mathbb{E} \left[\frac{D_i - \hat{e}(X_i)}{\hat{e}(X_i)} \mid X_i \right] \end{aligned}$$

$$\Rightarrow \text{Bias}(\hat{\mu}_1) = \mathbb{E} \left[f_1(X_i) \cdot \frac{D_i - \hat{e}(X_i)}{\hat{e}(X_i)} \right]$$

Similarly, $\text{Bias}(\hat{\mu}_0) = \mathbb{E} \left[f_0(X_i) \cdot \frac{\hat{e}(X_i) - D_i}{1 - e(X_i)} \right]$

We can estimate/correct Bias using regression!

Double Robustness

(Assume for simplicity $\hat{f}_0, \hat{f}_1, \hat{e} \perp\!\!\!\perp \text{data}$)

Let $\hat{\delta}^{DR} = \hat{M}_1^{DR} - \hat{M}_0^{DR}$, where

$$\hat{M}_1^{DR} = \frac{1}{n} \sum_{i=1}^n \frac{Y_i D_i}{\hat{e}(x_i)} - \underbrace{\frac{1}{n} \sum_{i=1}^n \frac{D_i - \hat{e}(x_i)}{\hat{e}(x_i)} \hat{f}_1(x_i)}_{\widehat{\text{Bias}}(\hat{M}_1)}$$

$$\hat{M}_0^{DR} = \frac{1}{n} \sum_{i=1}^n \frac{Y_i (1-D_i)}{1 - \hat{e}(x_i)} - \underbrace{\frac{1}{n} \sum_{i=1}^n \frac{\hat{e}(x_i) - D_i}{1 - \hat{e}(x_i)} \hat{f}_0(x_i)}_{\widehat{\text{Bias}}(\hat{M}_0)}$$

Claim If either

- $\hat{e}(x) = e(x) (+ o_p(1))$ or
- $\hat{f}_d(x) = f_d(x) (+ o_p(1))$, for $d=0,1$

Then $\mathbb{E} \hat{\delta}^{DR} = \delta (+ o_p(1))$

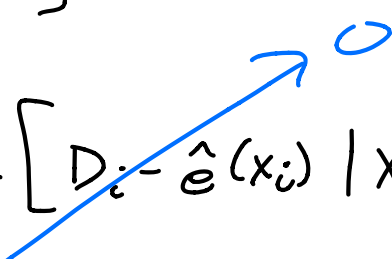
Proof

1. Suppose $\hat{f}_d(x) = f_d(x)$ for $d = 0, 1$

By previous derivation,

$$\mathbb{E} \hat{\mu}_d^{\text{DR}} = \mathbb{E} \hat{\mu}_d - \text{Bias}(\hat{\mu}_d) = \mu_d$$

2. Suppose $\hat{e}(x) = e(x)$, \hat{f}_0, \hat{f}_1 possibly wrong

$$\begin{aligned} \mathbb{E} \left[\hat{f}_1(x_i) \frac{D_i - \hat{e}(x_i)}{\hat{e}(x_i)} \mid x_i \right] \\ = \frac{\hat{f}_1(x_i)}{\hat{e}(x_i)} \cdot \mathbb{E} \left[\cancel{D_i - \hat{e}(x_i)} \mid x_i \right] \end{aligned}$$


$$\text{So } \mathbb{E} \hat{\mu}_1^{\text{DR}} = \mathbb{E} \hat{\mu}_1 = \mu_1$$

Similar for $\hat{\mu}_0^{\text{DR}}$