

# Causality

<https://www.cambridge.org/core/books/causal-inference-for-statistics-social-and-biomedical-sciences/71126BE90C58F1A431FE9B2DD07938AB>

[https://scholar.harvard.edu/imbens/files/efficient\\_estimation\\_of\\_average\\_treatment\\_effects\\_using\\_the\\_estimated\\_propensity\\_score.pdf](https://scholar.harvard.edu/imbens/files/efficient_estimation_of_average_treatment_effects_using_the_estimated_propensity_score.pdf)

<https://onlinelibrary.wiley.com/doi/abs/10.1111/ectj.12097>

<https://www.pnas.org/content/116/10/4156>

<https://arxiv.org/pdf/1712.04912.pdf>

# Prediction and causality

- A central goal of ML is to predict an outcome given variables describing a situation
  - Given patient characteristics, will their outcome improve?
- Most decision-making problems revolve around a decision / intervention / treatment
  - What would happen if we changed the system?
  - Given patient characteristics, will their outcome improve if **they follow a new diet**?
- We want to develop a scientific understanding of a decision
  - If you predict housing demand based on price, then a prediction model will say high price means high demand

# Prediction and causality

- Causal inference is a multi-disciplinary field spanning across economics, epidemiology, and statistics
- Focus is on questions about **counterfactuals**
  - What structure of data do we need to answer this question?
  - How do we interpret the key estimands?
- ML models can predict outcomes; when can it predict counterfactuals?
  - How can we leverage flexible ML models to infer causality?

# Potential outcomes

- Framework for explicitly modeling counterfactuals
- $A$ : binary treatment assignment (1: treated, 0: control)
- $Y(1)$  and  $Y(0)$  are potential outcomes
- $X$  is observed covariates

**First goal:** Estimate average treatment effect

$$\tau := \mathbb{E}[Y(1) - Y(0)]$$

**Problem:** We only observe  $Y := Y(A)$

# Observational studies

- Randomization is sometimes infeasible or prohibitively expensive
  - e.g. post-market drug surveillance, effect of air pollution on long-term health outcomes
- Experimentation can be risky in high-stakes scenarios
  - operational scenarios: new inventory policy for Amazon, new pricing algorithm for Uber
- May want to use existing large-scale data collected under some data-generating policy (e.g. legacy system)

# Observational studies

- Historically, many important findings from observational data
  - “citrus fruit curing scurvy described in the 1700s or insulin as a treatment for diabetes in the 1920s long preceded the advent of the modern randomized clinical trial.”
  - “these methods had in common a reliable method of diagnosis, a predictable clinical course, and a large and obvious effect of the treatment.” [Corrigan-Curay et al. 2018]
- These results need to be contextualized and viewed with more skepticism than RCTs

# SUTVA

- Throughout we implicitly assumed there is only a single version of the treatment that gets applied to all treated units
  - This may not be true if drugs go stale in storage, or dosages differ
- We also assumed there is *no interference between units*
  - Whether or not individual  $i$  is treated has no impact on the treatment effect of another individual  $j$
  - This can also fail in many real-world scenarios
- Together these assumptions are called stable unit treatment value assumption (SUTVA)

# Interference

- Any two-sided platform faces interference between units
- Consider the following scenario:
  - Lyft A/B tests a new promotion strategy for drivers
  - Each driver is randomized into treatment or control
  - It is observed that drivers finish a lot more rides with the promotion
  - So they decide this promotion is worth spending resources on
- But the estimate turned out to be an **overestimate**, not worth the cost of the promotion. Why?



# Interference

- Both treated and control drivers see the same set of demand
- If promotion incentivizes treated drivers to work more for less nominal fares, this cannibalizes demand that would usually go to control drivers
- Interference occurs in a number of different settings
  - Two-sided platforms: Airbnb, ridesharing, ad auctions
  - Network effects: e.g. adoption of new education technology
- When this happens, the potential outcomes now depend on all possible  $2^n$  treatment assignments
  - Very active area of research

# No unobserved confounding

- Previous regression-based direct method still works if there are no unobserved confounders (also called ignorability)

**Assumption.**  $Y(1), Y(0) \perp A \mid X$

- Observed treatment assignments are based on covariate information alone (+ random noise)
  - Treatment assignment does not use information about counterfactuals
- Strong assumption. Often violated in practice.
  - e.g. doctors often use unrecorded info to prescribe treatments

# No unobserved confounding

**Assumption.**  $Y(1), Y(0) \perp A \mid X$

- Under **no unobserved confounding**,

$$\mathbb{E}_P[Y(1)] = \mathbb{E}_P[\mathbb{E}[Y(1) \mid X]] = \mathbb{E}[\mathbb{E}_P[Y(1) \mid X, Z = 1]]$$

- Directly regress  $Y$  on  $X$  for treated units ( $Z=1$ ) to get  $\mathbb{E}_{\hat{P}}[Y \mid X, Z = 1]$

# Overlap

- We need enough samples for both control and treatment throughout the covariate space
  - This governs the effective sample size
- Propensity score  $e^\star(X) := \mathbb{P}(A = 1 \mid X)$
- Assume that there exists  $\epsilon > 0$  such that  $\epsilon \leq e^\star(X) \leq 1 - \epsilon$  almost surely
- This means I have at least  $\epsilon n$  number of samples for fitting the two outcome models

# Overlap

- This breaks if data is generated by a deterministic policy
  - e.g. always assign the drug (treatment) when age  $> 50$
- We need sufficient amount of randomness in treatment assignment in all covariate regions
- Governs difficulty of estimation. Often violated in practice.

# Assessing overlap

- “If the covariate distributions are similar, as they would be, in expectation, in the setting of a completely randomized experiment, there is less reason to be concerned about the sensitivity of estimates to the specific method chosen than if these distributions are substantially different.”
- “On the other hand, even if unconfoundedness holds, it may be that there are regions of the covariate space with relatively few treated units or relatively few control units, and, as a result, inferences for such regions rely largely on extrapolation and are therefore less credible than inferences for regions with substantial overlap in covariate distributions.”
- Imbens and Rubin

# Assessing overlap

- Overlap governs effective sample size
  - Even approaches that don't require propensity weighting is affected under this fundamental restriction
- Causal inference literature has developed various “supplementary analysis” tools for assessing credibility of empirical claims
- One of the most common conventions is to plot the propensity scores of treated and control groups

# Assessing overlap

- Difference in covariate distributions between treatment and control group is summarized by the propensity score
- Let  $f_1(X)$  be the density of  $X$  in the treatment group (similarly  $f_0(X)$ )
- Let  $p := \mathbb{P}(A = 1)$

$$\text{Var}(e^\star(X)) = p(1 - p)(\mathbb{E}[e^\star(X) | A = 1] - \mathbb{E}[e^\star(X) | A = 0])$$

$$= p^2(1 - p)^2 \cdot \mathbb{E} \left[ \left( \frac{f_1(X) - f_0(X)}{pf_1(X) + (1 - p)f_0(X)} \right)^2 \right]$$



# Assessing overlap

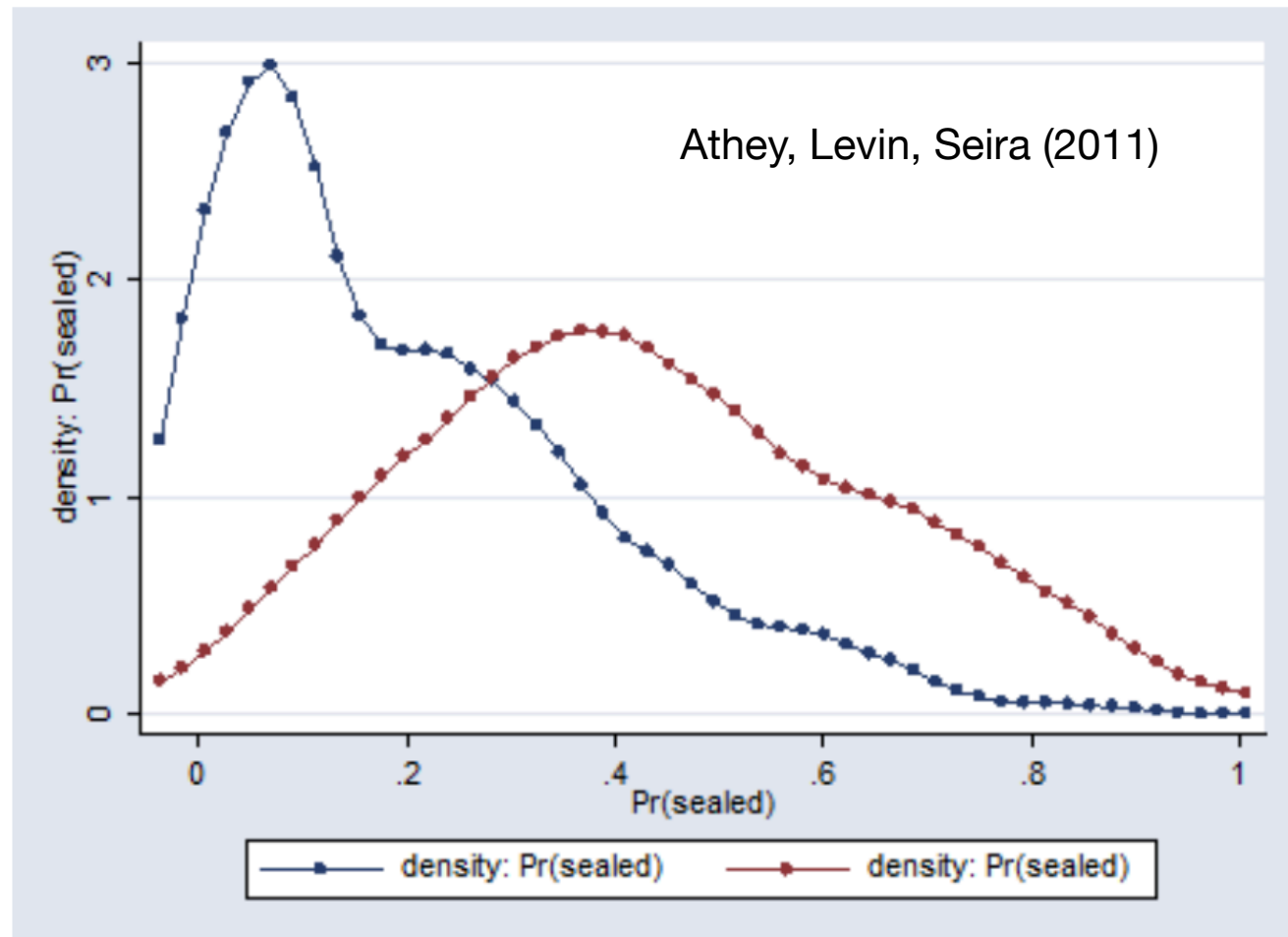
- A common visualization is to look at the pdf of the propensity score across treatment groups
- Plot approximates pdfs of the distribution  $\mathbb{P}(e^\star(X) \in \cdot \mid A = a)$
- For each  $q \in (0,1)$ , plot fraction of observations in the treatment group with  $e^\star(x) = q$  (and similarly for control)

# Assessing overlap

- Athey, Levin, Seira (2011) studied timber auctions
  - Award timber harvest contracts via first price sealed auction or open ascending auction
- Idaho: randomized with different probabilities across different regions
- California: determined by small vs. large sales volume; cutoff varies by region

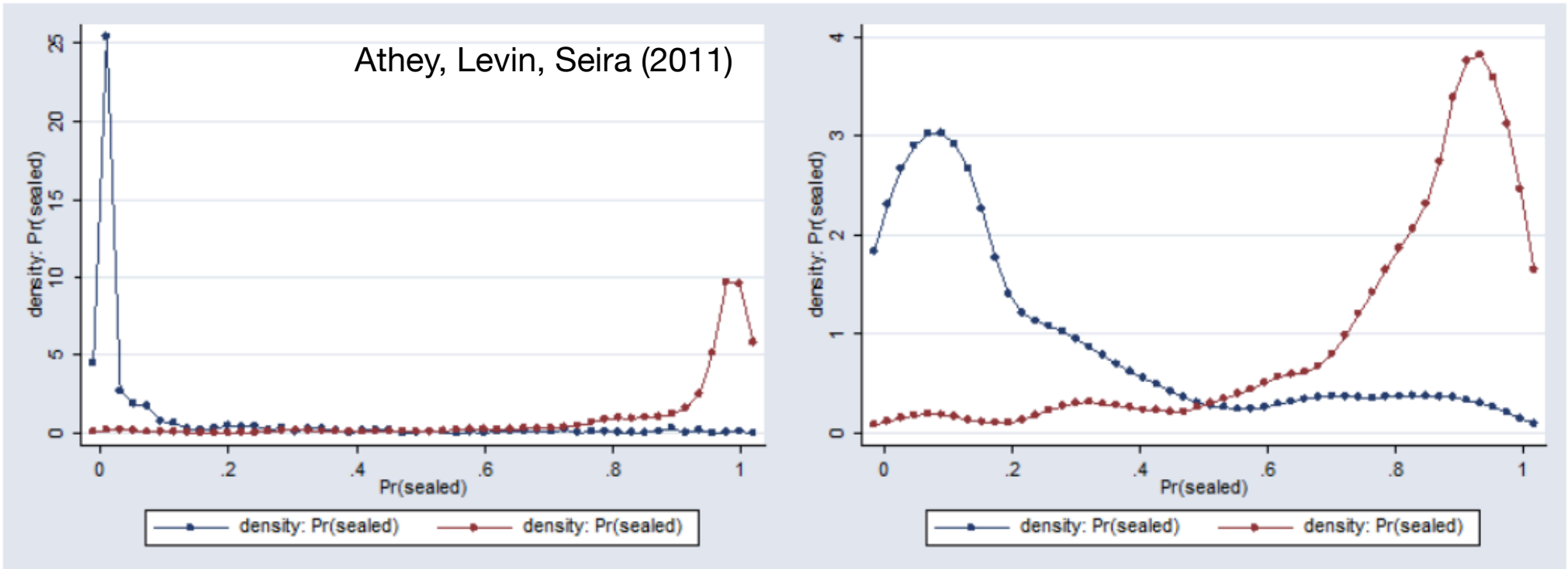
# Idaho

Very few observations with extreme propensity scores



# California

Untrimmed v. trimmed so that  $e(x) \in [.025, .975]$



# Estimators

# Direct method

- By no unobserved confounding,

$$\begin{aligned}\mu_a^\star(X) &:= \mathbb{E}[Y(a) \mid X] = \mathbb{E}[Y(a) \mid X, A = a] \\ &= \mathbb{E}[Y \mid X, A = a] \leftarrow \text{observable}\end{aligned}$$

- Fit  $\mu_a^\star(X)$  via the loss minimization problem

$$\text{minimize}_{\mu_a \in \mathfrak{M}_a} \mathbb{E}[(Y - \mu_a(X))^2 \mid A = a]$$

- ATE estimator  $\hat{\tau}_{\text{DM}} := \frac{1}{n} \sum_{i=1}^n \hat{\mu}_1(X_i) - \hat{\mu}_0(X_i)$

- Good if the outcome models are easy to learn

# Inverse propensity weighting

- What if the outcome models are very complex and difficult to estimate?
- A natural approach is to reweight samples to correct for confounding bias
  - Essentially importance sampling
  - Reweight treated units to look like “everyone”
- First, estimate the propensity score  $e^{\star}(X) := \mathbb{P}(A = 1 | X)$ 
  - e.g. run logistic regression to predict A given X

# Inverse propensity weighting

$$\hat{\tau}_{\text{IPW}} := \frac{1}{n} \sum_{i=1}^n \left( \frac{A_i}{\hat{e}(X_i)} Y_i - \frac{1 - A_i}{1 - \hat{e}(X_i)} Y_i \right)$$

- Can work well if propensity score is simple to estimate
- But estimating this well over the entire covariate space can be difficult
  - Calibration is hard, especially in high-dimensions
- When overlap doesn't hold, importance weights blow up



# Unbiasedness

Propensity score  $e^*(x) := \mathbb{E}[A|x] = \mathbb{P}(A=1|x)$

$$\hat{\tau}_{IPW} := \frac{1}{n} \sum_i \left( \frac{A_i Y_i}{e^*(x_i)} - \frac{(1-A_i) Y_i}{1-e^*(x_i)} \right)$$

$$\begin{aligned} \mathbb{E} \left[ \frac{AY}{e^*(x)} \right] &= \mathbb{E} \left[ \mathbb{E} \left[ \frac{AY}{e^*(x)} \mid x \right] \right] = \mathbb{E} \left[ \mathbb{E} \left[ \frac{AY}{e^*(x)} \mid x, A=1 \right] e^*(x) \right] \\ &= \mathbb{E} \left[ \mathbb{E} [AY \mid x, A=1] \right] = \mathbb{E} \left[ \mathbb{E} [Y(1) \mid x, A=1] \right] \\ &= \mathbb{E} \left[ \mathbb{E} [Y(1) \mid x] \right] \quad \text{by no unobserved confounding} \end{aligned}$$

Similarly  $\mathbb{E} \left[ \frac{(1-A)Y}{1-e^*(x)} \right] = \mathbb{E} \left[ \mathbb{E} [Y(0) \mid x] \right] = \mathbb{E} [Y(0)].$

# Estimating propensity score

But  $e^*$  is often unknown. How do we estimate  $e^*$ ?

Logistic regression

$$\log \frac{e(x)}{1-e(x)} = f(x; \theta) \quad \text{e.g. } f(x; \theta) = \theta^T x$$

Then solve

$$\min_{\theta \in \Theta} \mathbb{E} \left[ \log \left( 1 + \exp \left( \underbrace{\{-1, 1\}}_{\psi} \cdot f(x; \theta) \right) \right) \right] \quad \downarrow \text{MLE}$$

If model is well-specified, then  $\hat{e} \approx e^*$ .

# Inverse propensity weighting

- Can work well if propensity score is simple to estimate
- But estimating this well over the entire covariate space can be difficult
  - Calibration is hard, especially in high-dimensions
- When overlap doesn't hold, importance weights blow up

# Debiasing

- **Idea: Correct plug-in estimator using the first-order error**

$$\psi(\hat{P}) - \psi(P) = \nabla \psi(\hat{P})^\top (\hat{P} - P) + \text{Rem}_2$$

- Debaised estimator

$$\psi(\hat{P}) - \nabla \psi(\hat{P})^\top (\hat{P} - P)$$

- Automatically achieves second-order error  $\text{Rem}_2$

# Debiasing

- Outcome model  $\mu_1^\star(X) := \mathbb{E}[Y | X, A = 1]$ ,  
propensity score  $e^\star(X) := \mathbb{P}(A = 1 | X)$
- Debiasing gives **doubly robust** estimator

$$\mathbb{E}[Y(1)] = \mathbb{E} \left[ \mu_1^\star(X) + \frac{A}{\mathbb{P}(A = 1 | X)} (Y - \mu_1^\star(X)) \right]$$

- Propensity weight residuals to debias the direct method
- Accurate if you can do either well; **insensitive** to errors in nuisance estimates

# Debiasing

- Outcome model  $\mu_1^\star(X) := \mathbb{E}[Y | X, A = 1]$ ,  
propensity score  $e^\star(X) := \mathbb{P}(A = 1 | X)$

$$\hat{\tau}_{\text{AIPW}} := \frac{1}{n} \sum_{i=1}^n (\hat{\mu}_1(X_i) - \hat{\mu}_0(X_i)) + \frac{1}{n} \sum_{i=1}^n \left( \frac{A_i}{\hat{e}(X_i)} (Y_i - \hat{\mu}_1(X_i)) - \frac{1 - A_i}{1 - \hat{e}(X_i)} (Y_i - \hat{\mu}_0(X_i)) \right)$$

# Control variate

- Control variate is a centered RV, such that if you want to estimate  $E[Y]$ , you estimate  $E[Y + V]$  instead
- You get a variance reduction whenever  $\text{Var}(V) - 2\text{Cov}(Y, V) < 0$

$$\begin{aligned} & \text{Var}\left(\frac{AY}{e^{\eta(x)}} + \left(1 - \frac{A}{e^{\eta(x)}}\right) H_1^{\eta}(x)\right) \\ &= \text{Var}\left(\frac{AY}{e^{\eta(x)}}\right) + \text{Var}\left(\left(1 - \frac{A}{e^{\eta(x)}}\right) H_1^{\eta}(x)\right) + 2\text{Cov}\left(\frac{AY}{e^{\eta(x)}}, \left(1 - \frac{A}{e^{\eta(x)}}\right) H_1^{\eta}(x)\right) \end{aligned}$$

# Control variate

$$\begin{aligned}
 \text{Var} \left( \left(1 - \frac{A}{e^*(x)}\right) M_i^*(x) \right) &= \mathbb{E} \left[ M_i^*(x)^2 \cdot \left(1 - \frac{A}{e^*(x)}\right)^2 \right] \\
 &= \mathbb{E} \left[ M_i^*(x)^2 \cdot \left\{ e^*(x) \cdot \mathbb{E} \left[ \left(1 - \frac{A}{e^*(x)}\right)^2 \mid X, A=1 \right] \right. \right. \\
 &\quad \left. \left. + (1 - e^*(x)) \mathbb{E} \left[ \left(1 - \frac{A}{e^*(x)}\right)^2 \mid X, A=0 \right] \right\} \right] \\
 &= \mathbb{E} \left[ M_i^*(x)^2 \cdot \left\{ \frac{(e^*(x)-1)^2}{e^*(x)} + 1 - e^*(x) \right\} \right] \\
 &= \mathbb{E} \left[ M_i^*(x)^2 \cdot \frac{1}{e^*(x)} \cdot \left\{ \cancel{e^*(x)^2} - 2e^*(x) + 1 + e^*(x) - \cancel{e^*(x)^2} \right\} \right] \\
 &= \mathbb{E} \left[ M_i^*(x)^2 \left( -1 + \frac{1}{e^*(x)} \right) \right].
 \end{aligned}$$



# Control variate

$$\begin{aligned} \text{Cov} \left( \frac{AY}{e^{\eta(x)}}, \left(1 - \frac{A}{e^{\eta(x)}}\right) M_i^*(x) \right) &= \mathbb{E} \left[ \frac{AY}{e^{\eta(x)}} \left(1 - \frac{A}{e^{\eta(x)}}\right) M_i^*(x) \right] \\ &= \mathbb{E} \left[ e^{\eta(x)} \cdot \mathbb{E} \left[ \frac{AY}{e^{\eta(x)}} \left(1 - \frac{A}{e^{\eta(x)}}\right) M_i^*(x) \mid X, A=1 \right] \right] \\ &= \mathbb{E} \left[ \mathbb{E} \left[ \left(1 - \frac{1}{e^{\eta(x)}}\right) \cdot Y(1) \cdot M_i^*(x) \mid X, A=1 \right] \right] \\ &= \mathbb{E} \left[ \left(1 - \frac{1}{e^{\eta(x)}}\right) \cdot M_i^*(x) \cdot \mathbb{E}[Y(1) \mid X] \right] \quad \text{by no unobs. conf.} \\ &= \mathbb{E} \left[ \left(1 - \frac{1}{e^{\eta(x)}}\right) M_i^*(x)^2 \right] \end{aligned}$$

# Control variate

$$\begin{aligned}
 \text{Cov} \left( \frac{AY}{e^{\theta(x)}}, \left(1 - \frac{A}{e^{\theta(x)}}\right) M_i^{\theta}(x) \right) &= \mathbb{E} \left[ \frac{AY}{e^{\theta(x)}} \left(1 - \frac{A}{e^{\theta(x)}}\right) M_i^{\theta}(x) \right] \\
 &= \mathbb{E} \left[ e^{\theta(x)} \cdot \mathbb{E} \left[ \frac{AY}{e^{\theta(x)}} \left(1 - \frac{A}{e^{\theta(x)}}\right) M_i^{\theta}(x) \mid X, A=1 \right] \right] \\
 &= \mathbb{E} \left[ \mathbb{E} \left[ \left(1 - \frac{1}{e^{\theta(x)}}\right) \cdot Y(1) \cdot M_i^{\theta}(x) \mid X, A=1 \right] \right] \\
 &= \mathbb{E} \left[ \left(1 - \frac{1}{e^{\theta(x)}}\right) \cdot M_i^{\theta}(x) \cdot \mathbb{E}[Y(1) \mid X] \right] \quad \text{by no covs. conf.} \\
 &= \mathbb{E} \left[ \left(1 - \frac{1}{e^{\theta(x)}}\right) M_i^{\theta}(x)^2 \right]
 \end{aligned}$$

$$\begin{aligned}
 \text{So } \text{Var} \left( \frac{AY}{e^{\theta(x)}} + \left(1 - \frac{A}{e^{\theta(x)}}\right) M_i^{\theta}(x) \right) & \quad \text{variance reduction} \\
 = \text{Var} \left( \frac{AY}{e^{\theta(x)}} \right) - \mathbb{E} \left[ \left( \frac{1}{e^{\theta(x)}} - 1 \right) M_i^{\theta}(x)^2 \right] & \quad \checkmark
 \end{aligned}$$

# Nuisance parameters

- Outcome model  $\mu_a^\star(X) := \mathbb{E}[Y | X, A = a]$ ,  
propensity score  $e^\star(X) := \mathbb{P}(A = 1 | X)$
- If a good parametric model exists, then can estimate at the usual  $1/\sqrt{n}$  rates
- In general, these are infinite dimensional objects. Can be difficult to estimate.

# Semiparametrics

- We only care about estimating the ATE
  - One-dimensional estimand, infinite dimensional nuisance parameters
- Estimation accuracy of nuisance parameters is good only insofar as it helps with estimating the ATE
- Due to its high-dimensional nature, often difficult to estimate nuisances at parametric rates
- Goal: semiparametric estimators that are insensitive to errors in nuisance estimates

# Doubly robust

- One main advantage of AIPW is that even if one of the nuisance parameter models are **misspecified**, you can still get correct asymptotic behavior
- Consistent estimator of the ATE so long as ***either*** outcome models or propensity score can be estimated consistently

# Heuristic derivation

$$\begin{aligned} & \mathbb{E} \left[ Y(1) - \frac{1}{n} \sum_{i=1}^n \hat{\mu}_1(X_i) + \frac{A_i}{\hat{e}(X_i)} (Y_i - \hat{\mu}_1(X_i)) \right] \\ &= \mathbb{E} \left[ \left( 1 - \frac{e^\star(X)}{\hat{e}(X)} \right) (\mu_1^\star(X) - \hat{\mu}(X)) \right] \end{aligned}$$

# Orthogonality

- When is a semiparametric estimator insensitive to errors in nuisance estimates?
- Directional derivative of functional wrt nuisance parameters at true value is near-zero
- Ensures that a little perturbation in nuisance parameters near the truth values does not affect functional

# Orthogonality

- Directional (Gateaux) derivative of functional w.r.t. nuisance parameters around the true values  $\gamma^* = (\mu_1^*, \mu_0^*, e^*, h^*)$  is zero

Let  $\eta = (\mu_0, \mu_1, e)$  be the tuple of nuisance param.

A statistical functional  $\mathbb{E}_P \psi(D; \eta)$  is Neyman orthogonal on  $\Lambda$  if

$$\left. \frac{d}{dt} \mathbb{E}_P \psi(D; \eta^* + t(\eta - \eta^*)) \right|_{t=0} = 0 \quad \forall \eta \in \Lambda$$

↑ Directional derivative at true parameter  $\eta^*$

$D = (X, Y, A)$ ,  $\eta = (\mu_0, \mu_1, e)$ . Let

$$\psi_{\text{AIPW}}(D; \eta) := \mu_1(x) - \mu_0(x) + \frac{A}{e(x)} (Y - \mu_1(x)) - \frac{1-A}{1-e(x)} (Y - \mu_0(x))$$



# Orthogonality

Assume  $\frac{d}{dr}$  &  $\mathbb{E}_P$  are interchangeable throughout.

$$\frac{d}{dr} \mathbb{E}_P[(M_i + r(M_i - M_i^*)) (x)] \Big|_{r=0} = \mathbb{E}[(M_i - M_i^*) (x)]$$

$$\begin{aligned} & \frac{d}{dr} \mathbb{E}_P \left[ \frac{A}{e^r + r(e - e^r)} (\gamma - (M_i + r(M_i - M_i^*))) (x) \right] \Big|_{r=0} \\ &= \mathbb{E} \left[ - \frac{A(e - e^r)}{(e^r + r(e - e^r))^2} (\gamma - M_i^*(x)) \right] \Big|_{r=0} - \mathbb{E} \left[ \frac{A}{e^r(x)} (M_i - M_i^*)(x) \right] \\ &= \mathbb{E} \left[ - \frac{A(e - e^r)}{e^r(x)^2} (\gamma - M_i^*(x)) \right] - \mathbb{E} \left[ \frac{A}{e^r(x)} (M_i - M_i^*)(x) \right] \\ &= \mathbb{E} \left[ e^r(x) \mathbb{E} \left[ - \frac{A(e - e^r)}{e^r(x)^2} (\gamma - M_i^*(x)) \mid X, A=1 \right] \right] - \mathbb{E}[(M_i - M_i^*)(x)] \end{aligned}$$

# Orthogonality

$$= -\mathbb{E}\left[\frac{(e-e^*)(x)}{e^*(x)} \mathbb{E}[Y(i) - M_i^*(x) | X, A=i]\right] - \mathbb{E}[(M_i - M_i^*)(x)]$$

$$= -\mathbb{E}[(M_i - M_i^*)(x)] \quad \text{by no unobserved confounding}$$

$$\Rightarrow \frac{d}{dr} \mathbb{E}_P[(M_i^* + r(M_i - M_i^*))(x)] \Big|_{r=0} + \frac{d}{dr} \mathbb{E}_P\left[\frac{A}{e^* + r(e - e^*)} (Y - (M_i^* + r(M_i - M_i^*))(x))\right] \Big|_{r=0} = 0$$

Similarly for  $(1-A)$ -related terms.

We conclude  $\mathbb{E}_P \Psi_{\text{HAW}}$  is Neyman orthogonal at  $\eta^*$ ,  
for all  $\eta$  s.t.  $\frac{d}{dr}$  &  $\mathbb{E}$  are interchangeable.

# Why orthogonality?

- Allows getting central limit rates on ATE estimation even when we can only estimate nuisance parameters at slower rates
- In addition to no unobserved confounding,  $e^\star(X), \hat{e}(X) \in [\epsilon, 1 - \epsilon]$ , we assume the following rate condition

$$\|\hat{e} - e^\star\|_{P,2} (\|\hat{\mu}_1 - \mu_1^\star\|_{P,2} + \|\hat{\mu}_0 - \mu_0^\star\|_{P,2}) = o_p(n^{-1/2})$$

- This allows us to trade-off errors between nuisance parameters. Only their product needs to go down at this rate!

# Central limit result

- CLT for the semiparametric AIPW, even when nuisance estimates converge at slower-than-parametric rates

$$\sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n \psi_{\text{AIPW}}(X_i, Y_i, A_i; \hat{\mu}_0, \hat{\mu}_1, \hat{e}) - \tau \right) \Rightarrow N(0, \sigma_{\text{AIPW}}^2)$$

where  $\sigma_{\text{AIPW}}^2 := \text{Var} \left( \psi_{\text{AIPW}}(X, Y, A; \mu_0^*, \mu_1^*, e^*) \right)$

- This is the oracle asymptotic variance; when the true nuisance parameters are known
- AIPW achieves optimal asymptotic efficiency

# Sketch of asymptotics

We use the following decomposition.

$$\frac{1}{n} \sum \psi(D_i; \hat{\eta}) - \mathbb{E} \psi(D; \eta^*) = \underbrace{\frac{1}{n} \sum \psi(D_i; \hat{\eta}) - \mathbb{E}_{\mathbb{P}_{\mathbb{D}}} \psi(D; \hat{\eta})}_{\textcircled{1}} + \underbrace{\mathbb{E}_{\mathbb{P}_{\mathbb{D}}} \psi(D; \hat{\eta}) - \mathbb{E} \psi(D; \eta^*)}_{\textcircled{2}}$$

By the triangular CLT,  $\sqrt{n} \cdot \textcircled{1} \Rightarrow N(0, \sigma_{\text{ML}}^2)$  if  $\hat{\eta} \rightarrow \eta^*$

So it suffices to show that  $\sqrt{n} \cdot \textcircled{2} \xrightarrow{p} 0$ .

# Sketch of asymptotics

Define  $Q(r) := \mathbb{E}_{D_n, p} [\psi(D, \eta^* + r(\hat{\eta} - \eta^*))]$ .

Then  $\textcircled{2} = Q(1) - Q(0)$ . If  $r \mapsto Q(r)$  is cont. diff.

then  $= Q'(r)$  for some  $r \in [0, 1]$ .

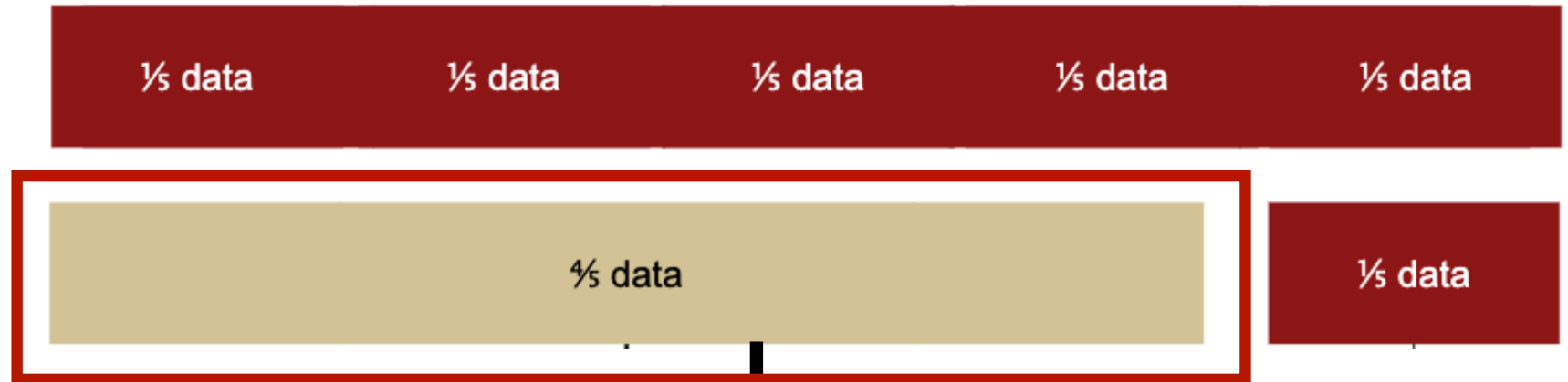
So if we can show  $\sup_{r \in [0, 1]} Q'(r) = o_p(n^{-\frac{1}{2}}) \dots (*)$

By orthogonality,  $Q'(0) = 0$ . Since  $Q'(\cdot)$  is sufficiently smooth, then we can argue that  $(*)$  holds under the rate condition

# Cross-fitting

- Instead of sample-splitting, we can alternate the role of main and auxiliary samples over multiple splits

**Cross-fitting**  
[Chernozhukov '18]



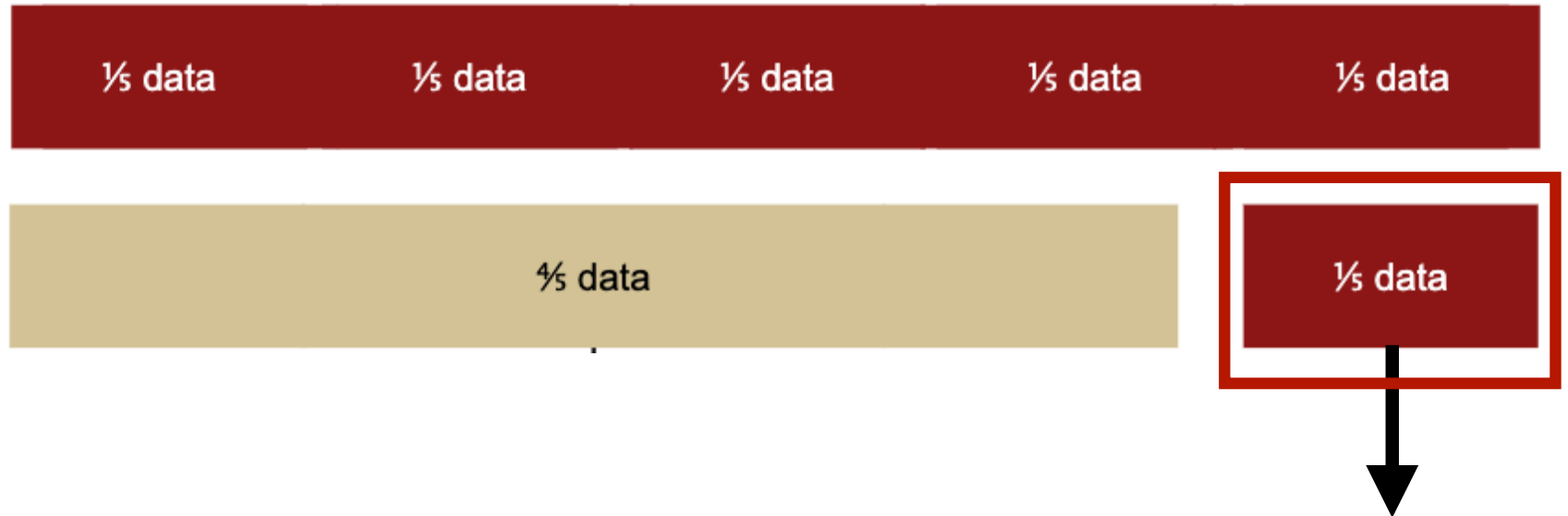
$$\hat{\mu}_a(X) \approx \mathbb{E}[Y(a) \mid X = x], \quad a \in \{0, 1\}$$

$$\hat{e}(X) \approx \mathbb{P}(A = 1 \mid X)$$

- Estimate nuisance parameters on the auxiliary sample

# Cross-fitting

Cross-fitting  
[Chernozhukov '18]



$$\hat{\tau}_1 := \frac{1}{n} \sum_{i=1}^n \hat{\mu}_1(X_i) - \hat{\mu}_0(X_i) + \frac{A_i}{\hat{e}(X_i)} (Y - \mu_1(X_i)) - \frac{1 - A_i}{1 - \hat{e}(X_i)} (Y - \mu_0(X_i))$$

- Estimate ATE by plugging in nuisance estimates



# Cross-fitting

Cross-fitting  
[Chernozhukov '18]



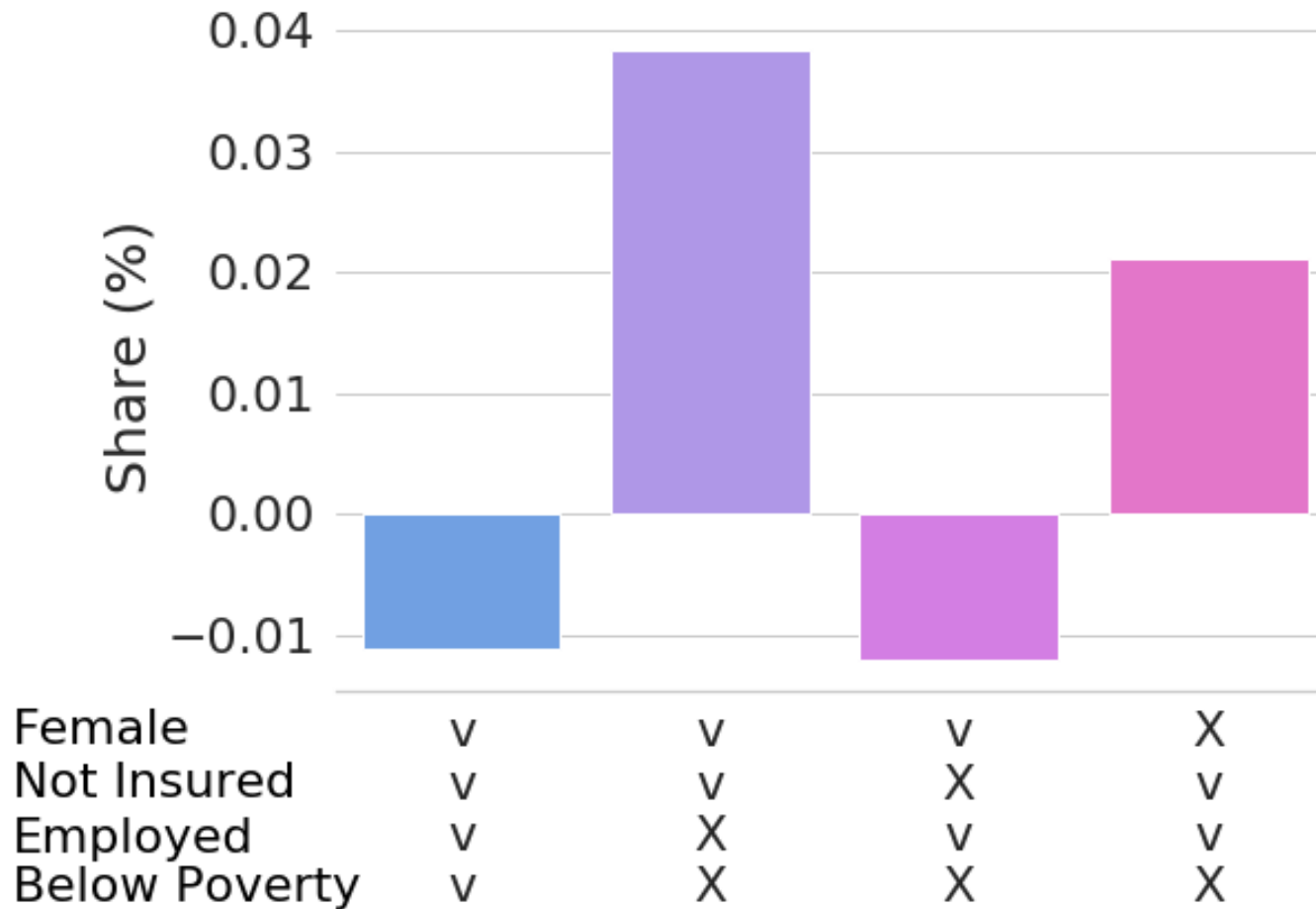
$$\hat{\tau} = \frac{1}{5} \left( \hat{\tau}_1 + \hat{\tau}_2 + \hat{\tau}_3 + \hat{\tau}_4 + \hat{\tau}_5 \right)$$

- Same procedure for direct method, IPW
- Similar central limit result follows as before

# Heterogeneous treatment effects

- Treatment effect often varies with user / patient / agent characteristics (covariates)
- Example: Oregon Health Insurance Experiment
  - Evaluate effect of Medicaid on low-income adults on emergency department (ED) visits in 2008
  - Precursory study to federal Medicaid expansion in 2014, which cost \$553 billion/year
  - Insurance allows visits ED, but access to preventive care may also reduce need of ED visits

# Oregon Health Insurance Experiment



# CATE

- To estimate personalized treatment effects, we want to estimate the **conditional average treatment effect (CATE)**

$$\tau(X) := \mathbb{E}[Y(1) - Y(0) \mid X]$$

- Few different ways to estimate this using black-box ML models
- Again, key challenging is missing data
  - We never observed counterfactuals

# S-Learner

Fit a single model for all treatment options

$$Y \leftarrow (X, A)$$

$$\min_{\theta} \mathbb{E}[(Y - f_{\theta}(A, X))^2]$$

Then  $\hat{\tau}(x) := f_{\theta}(1, x) - f_{\theta}(0, x)$

- Shared feature representation, assuming similar model class for both treatment and control

# T-Learner

Fit two models, for treatment & control

$$\min_{\theta} \mathbb{E}[(Y - \mu_{\theta,1}(x))^2 | A=1]$$

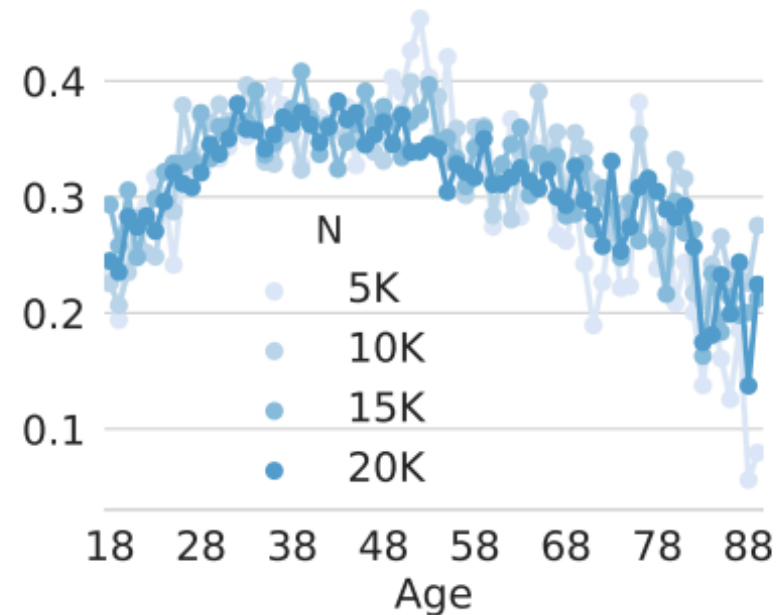
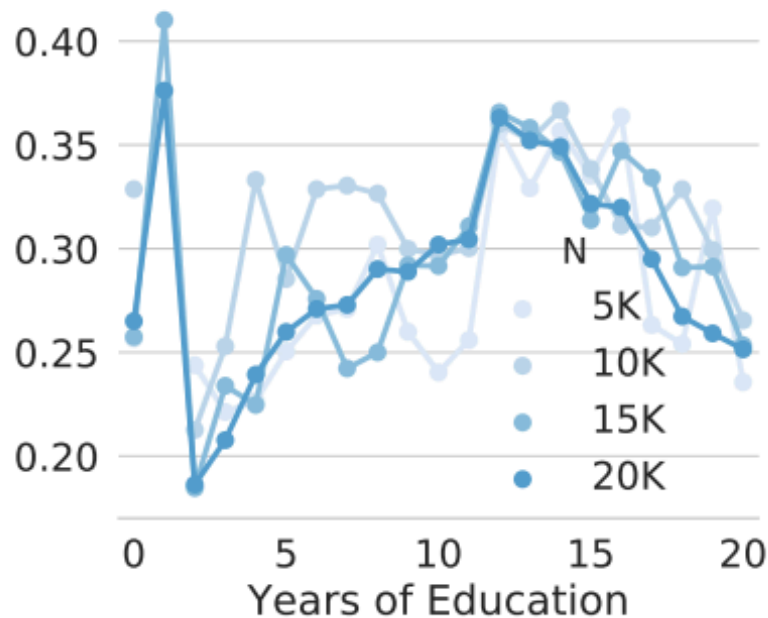
$$\min_{\theta} \mathbb{E}[(Y - \mu_{\theta,0}(x))^2 | A=0]$$

$$\hat{\tau}(x) := \hat{\mu}_{\theta,1}(x) - \hat{\mu}_{\theta,0}(x)$$

- Can fit different models over treatment options

# Welfare attitudes experiment

- Evaluate effect of wording on survey results (“welfare” vs “assistance to the poor”)
- Resoundingly positive treatment effects, but significant heterogeneity across covariates



# X-Learner

Kunzel et al. (2018)

- Regress on the imputed treatment effect  $Y(1) - Y(0)$
- Fit T-learner models and compute imputed treatment effects

$$Y_i - \hat{\mu}_{\theta,0}(X_i) \text{ if } A_i = 1, \hat{\mu}_{\theta,1}(X_i) - Y_i \text{ if } A_i = 0$$

- Fit another set of models  $\hat{\tau}_1, \hat{\tau}_0$  on the two category of imputed values, take

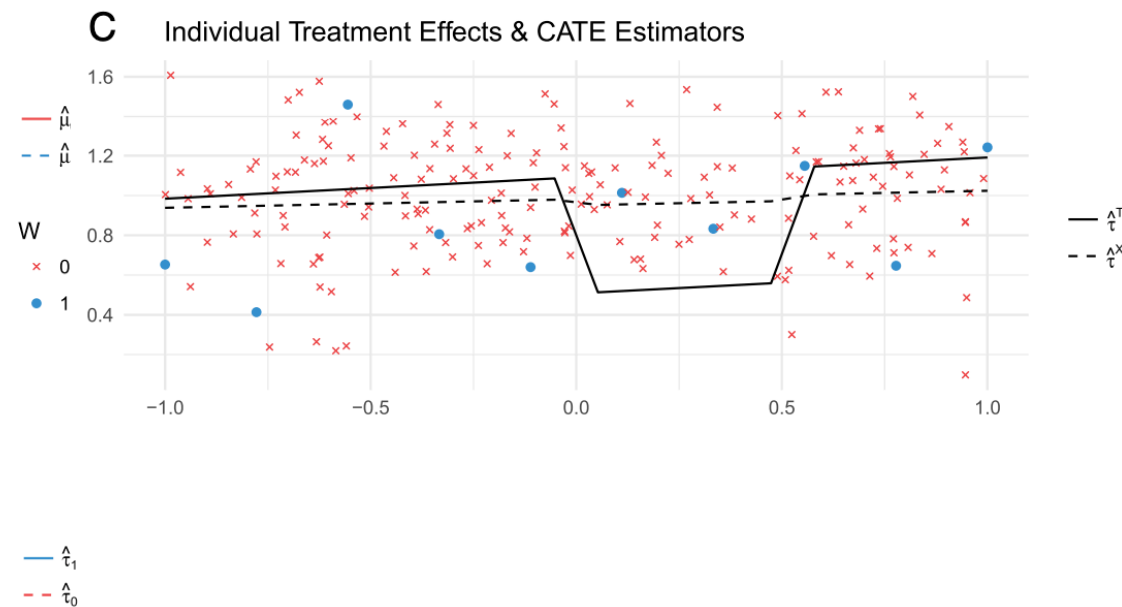
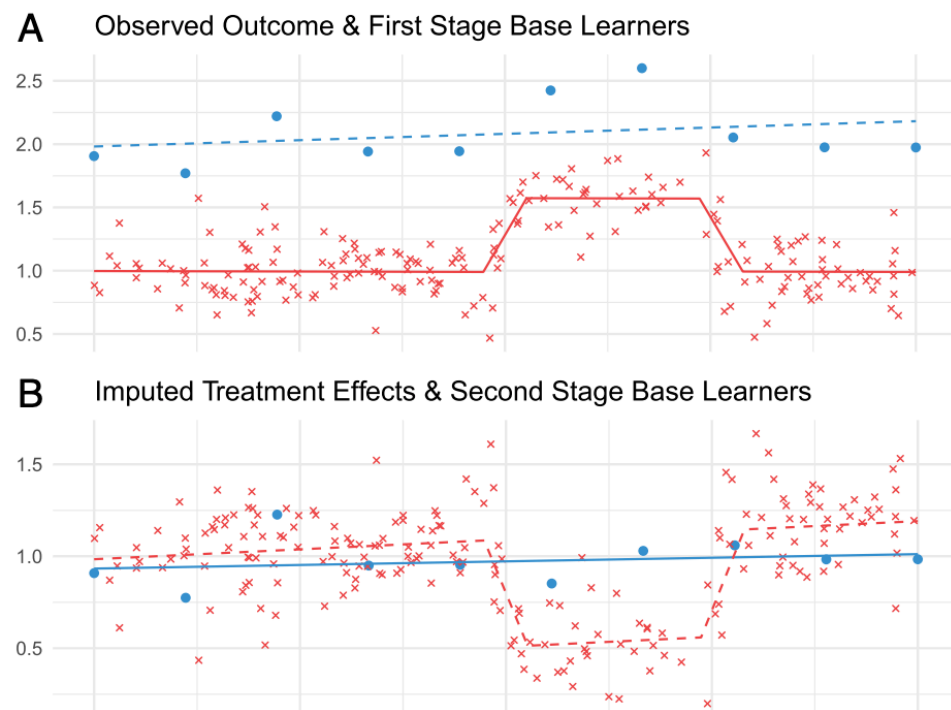
$$\hat{\tau}(X) := \hat{e}(X)\hat{\tau}_0(X) + (1 - \hat{e}(X))\hat{\tau}_1(X)$$



# X-Learner

Kunzel et al. (2018)

- Usually, number of samples in treatment >> those in control
- Advantageous if CATE is much smoother than individual outcome functions



# R-Learner

Nie and Wager (2020)

Robinson's decomposition Define  $\varepsilon(a) := Y(a) - \overbrace{(M_0^*(x) + a\tau(x))}^{\text{red bracket}}$ ,  $\varepsilon := \varepsilon(A)$ .

From no unobserved confounding,

$$\begin{aligned} \mathbb{E}[\varepsilon | X] &= e^*(x) \mathbb{E}[\varepsilon(1) | X, A=1] + (1 - e^*(x)) \mathbb{E}[\varepsilon(0) | X, A=0] \\ &= e^*(x) \mathbb{E}[Y(1) - M_1^*(x) | X, A=1] + (1 - e^*(x)) \mathbb{E}[Y(0) - M_0^*(x) | X, A=0] \\ &= e^*(x) \mathbb{E}[Y(1) - M_1^*(x) | X] + (1 - e^*(x)) \mathbb{E}[Y(0) - M_0^*(x) | X] \\ &= 0 \end{aligned}$$

By def of  $\varepsilon$ ,  $Y = M_0(x) + A\tau(x) + \varepsilon \dots (*)$

Define  $m^*(x) := \mathbb{E}[Y | X]$ . Taking  $\mathbb{E}[\cdot | X]$  on both sides

$$m^*(x) = M_0(x) + e^*(x)\tau(x).$$

# R-Learner

Nie and Wager (2020)

Subtracting this from (\*), we arrive at

$$Y - \hat{m}(x) = (A - \hat{e}(x))z(x) + \epsilon.$$

So if we fit  $\hat{m}, \hat{e}$  on heldout data,

we can now solve

$$\min_{\theta} \mathbb{E}[(Y - \hat{m}(x) - (A - \hat{e}(x))z_{\theta}(x))^2]$$

to get  $\hat{z}_{\theta}(x)$ .