

# MS&E 228: Applied Causal Inference Powered by ML and AI

## Lecture 4: Estimation, Confidence Intervals, and Doubly Robust Learning

Vasilis Syrgkanis

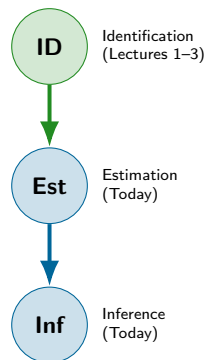
Stanford University

Winter 2026

**Readings:** *Applied Causal Inference Powered by ML and AI*, §9.

# Goals for Today

1. Move from **identification** (infinite-data formulas) to **estimation & inference** (finite- $n$  procedures).
2. Formalize what we want from estimators: **convergence rates** and **confidence intervals**.
3. Understand why the **plug-in g-estimator** can give **misleading uncertainty** when  $\hat{g}$  is learned.
4. Build the **doubly robust / debiased estimator** + **cross-fitting** and get valid CIs using modern ML.



# Warm-up: The Plug-in g-Estimator

Outcome regression from Lecture 2

# The Plug-in g-Estimator (Outcome Regression)

## Assumptions

Conditional ignorability + overlap:

$$(Y(0), Y(1)) \perp\!\!\!\perp D \mid X, \quad 0 < \Pr(D = 1 \mid X) < 1.$$

## Identification (g-formula):

$$\theta_0 \equiv \text{ATE} = \mathbb{E}[g_0(1, X) - g_0(0, X)], \quad g_0(d, x) = \mathbb{E}[Y \mid D = d, X = x].$$

## Operationalization:

$$\hat{\theta}_{\text{plug-in}} = \frac{1}{n} \sum_{i=1}^n \left( \hat{g}(1, X_i) - \hat{g}(0, X_i) \right)$$

# What Questions Do We Want to Answer?

Suppose we repeat the analysis on fresh i.i.d. samples of size  $n$ .

- ▶ **Accuracy:** How fast does  $\hat{\theta}$  approach  $\theta_0$  as  $n$  grows?
- ▶ **Variability:** How much does  $\hat{\theta}$  change across repeated samples?
- ▶ **Confidence:** Can we build an interval around  $\hat{\theta}$  that contains  $\theta_0$ , in approximately 95% of our repetitions?

## Why CIs are Central in Causal Inference

We cannot directly validate causal estimates on held-out data.  
CIs are the main operational analogue of “generalization error” in prediction.

# Formalizing “Convergence Rate”

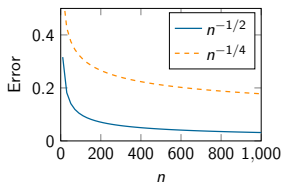
Let  $\theta_0$  be the target scalar parameter and  $\hat{\theta}_n$  an estimator.

## (Mean-square) Convergence Rate

We say  $\hat{\theta}_n$  has rate  $r_n$  if

$$\mathbb{E}[(\hat{\theta}_n - \theta_0)^2] \lesssim r_n^2.$$

**Typical benchmark:**  $r_n = n^{-1/2}$  (“parametric” / CLT rate).



- ▶ If  $\hat{\theta}_n - \theta_0 = O_p(n^{-1/2})$ , we can hope for tight CIs.
- ▶ If rate is slower, uncertainty may dominate even for large samples.

# Formalizing Confidence Intervals

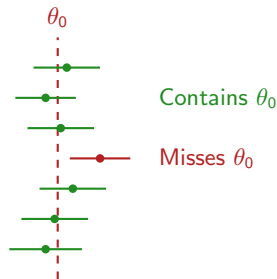
A (random) interval  $CI_n = [L_n, U_n]$  is a **95% confidence interval** if

$$\Pr(\theta_0 \in CI_n) \approx 0.95.$$

## Interpretation (Frequentist Coverage)

If we re-run the whole analysis many times on fresh samples, about 95% of the time the interval will contain the true parameter.

- ▶ CIs inform decisions: “treat only if we are confident it won’t harm.”
- ▶ A narrow CI means high precision; a wide CI means the data are not very informative.



# How We Usually Get CIs: Asymptotic Normality

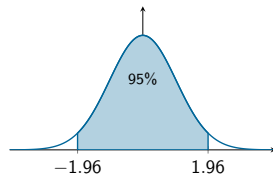
## Common path:

1. Show **asymptotic normality**:

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \Rightarrow \mathcal{N}(0, \sigma^2).$$

2. Estimate  $\sigma^2$  with a **standard error**  $\widehat{\text{se}}$ .
3. Form an approximate 95% CI:

$$\boxed{\hat{\theta}_n \pm 1.96 \widehat{\text{se}}}$$



$$\begin{aligned}\Pr(\theta_0 \in \hat{\theta} \pm 1.96 \widehat{\text{se}}) \\ &\approx \Pr(\sqrt{n}(\hat{\theta} - \theta_0) \in \pm 1.96\sigma) \\ &\approx \Pr(N(0, \sigma^2) \in \pm 1.96\sigma) = 0.95\end{aligned}$$

## Key Question

When is it valid to treat  $\hat{\theta}_n$  as “approximately normal”?



# **Some Important Preliminaries**

Basics of Probability and Statistics

# Empirical Averages

When given  $n$  i.i.d. samples  $\{X_1, \dots, X_n\}$  of a random vector  $X$ , denote:

$$\mathbb{E}_n[X] = \frac{1}{n} \sum_{i=1}^n X_i \quad (\text{empirical average})$$

# Convergence in Probability and Distribution

We will use the following intuitive notations to avoid complicated math expressions:

## Convergence in Probability

**Intuitively:** for large  $n$  the quantity  $A_n$  is a good approximation to the quantity  $\mu$

$$A_n \approx \mu$$

**Formally:** the random sequence  $A_n$  converges in probability to the quantity  $\mu$ .

## Convergence in Distribution

**Intuitively:** for large  $n$  the random vector  $A_n \in \mathbb{R}^d$  is approximately distributed like a Gaussian with mean 0 and covariance  $V$ .

$$A_n \stackrel{a}{\sim} \mathcal{N}(0, V)$$

**Formally:** If  $\mathcal{R}$  is the set of all rectangles in  $d$  dimensions, then

$$\sup_{R \in \mathcal{R}} |\Pr(A_n \in R) - \Pr(\mathcal{N}(0, V) \in R)| \approx 0.$$

# LLN and CLT

When given  $n$  i.i.d. samples  $\{X_1, \dots, X_n\}$  of a random vector  $X$ , denote:

$$\mathbb{E}_n[X] = \frac{1}{n} \sum_{i=1}^n X_i \quad (\text{empirical average})$$

## (Weak) Law of Large Numbers

$$\mathbb{E}_n[X] \approx \mathbb{E}[X]$$

## Central Limit Theorem

$$\sqrt{n}(\mathbb{E}_n[X] - \mathbb{E}[X]) \overset{d}{\sim} N(0, \text{Var}(X))$$

## Concrete Example: RCTs

Two-means estimator + CLT

## RCT: Two-means Estimator

In an RCT, treatment is randomized:  $(Y(0), Y(1)) \perp\!\!\!\perp D$ .

$$\hat{\theta}_{\text{RCT}} = \bar{Y}_1 - \bar{Y}_0, \quad \bar{Y}_1 = \frac{1}{n_1} \sum_{i:D_i=1} Y_i, \quad \bar{Y}_0 = \frac{1}{n_0} \sum_{i:D_i=0} Y_i.$$

# Central Limit Theorem (CLT) Intuition

Each group mean  $\bar{Y}_d$  is an average of  $n_d$  i.i.d. outcomes with mean  $\mathbb{E}[Y \mid D = d]$ .

1. By the Central Limit Theorem, each  $\bar{Y}_d$  is approximately distributed like a Gaussian  $\mathcal{N}(\mu_d, \sigma_d^2)$  with mean and variance

$$\mu_d = \mathbb{E}[Y \mid D = d], \quad \sigma_d^2 = \frac{\text{Var}(Y \mid D = d)}{n_d}.$$

2.  $\hat{\theta}_{\text{RCT}} = \bar{Y}_1 - \bar{Y}_0$  is approximately the difference of two independent Gaussians  $\mathcal{N}(\mu_1, \sigma_1^2)$ ,  $\mathcal{N}(\mu_0, \sigma_0^2)$  (since samples in the two empirical averages are disjoint)

Hence,  $\hat{\theta}_{\text{RCT}}$  behaves like a Gaussian with mean  $\mu_1 - \mu_0$  and variance  $\sigma_0^2 + \sigma_1^2$ .

# Asymptotic Variance & Standard Error in an RCT

**A more formal statement.** Under mild conditions,

$$\sqrt{n}(\hat{\theta}_{\text{RCT}} - \theta_0) \overset{a}{\sim} \mathcal{N}(0, \sigma_{\text{RCT}}^2),$$

where

$$\sigma_{\text{RCT}}^2 = \frac{\text{Var}(Y \mid D = 1)}{\pi} + \frac{\text{Var}(Y \mid D = 0)}{1 - \pi}, \quad \pi = \Pr(D = 1).$$

**Plug-in standard error estimate:**

$$\widehat{\text{se}} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_0^2}{n_0}}, \quad s_d^2 = \frac{1}{n_d} \sum_{i:D_i=d} (Y_i - \bar{Y}_d)^2.$$

**95% CI for RCT**

$$\hat{\theta}_{\text{RCT}} \pm 1.96 \cdot \sqrt{\frac{s_1^2}{n_1} + \frac{s_0^2}{n_0}}$$



## Appendix: A bit more formal

Deal with the fact that number of treated and control people is also random and formally argue that treated and control means are independent random variables.

$$\hat{\theta}_d := \bar{Y}_d = \frac{\sum_i Y_i \cdot \mathbb{I}\{D_i = d\}}{\sum_i \mathbb{I}\{D_i = d\}} = \frac{\mathbb{E}_n[Y \cdot \mathbb{I}\{D = d\}]}{\mathbb{E}_n[D]}, \quad \theta_{0,d} := \mathbb{E}[Y \mid D = d]$$

We have:

$$\begin{aligned} \sqrt{n}(\hat{\theta}_d - \theta_{0,d}) &= \sqrt{n} \left( \frac{\mathbb{E}_n[Y \mathbb{I}\{D = d\}]}{\mathbb{E}_n[\mathbb{I}\{D = d\}]} - \theta_{0,d} \frac{\mathbb{E}_n[\mathbb{I}\{D = d\}]}{\mathbb{E}_n[\mathbb{I}\{D = d\}]} \right) \\ &= \frac{\sqrt{n} \mathbb{E}_n[(Y - \theta_{0,d}) \cdot \mathbb{I}\{D = d\}]}{\mathbb{E}_n[\mathbb{I}\{D = d\}]} \\ &\approx \frac{\sqrt{n} \mathbb{E}_n[(Y - \theta_{0,d}) \cdot \mathbb{I}\{D = d\}]}{\mathbb{E}[\mathbb{I}\{D = d\}]} \end{aligned} \quad (\text{LLN})$$

## Appendix: A bit more formal

Define:  $X := (X_0, X_1)$   $X_d := \frac{(Y - \theta_{0,d}) \cdot \mathbb{I}\{D = d\}}{\mathbb{E}[\mathbb{I}\{D = d\}]}$

Then, we have proven that:

$$\sqrt{n}\{\hat{\theta}_d - \theta_{0,d}\}_{d \in \{0,1\}} \approx \sqrt{n}\mathbb{E}_n[X]$$

Note that  $\mathbb{E}[X] = 0$ . Thus, by the central limit theorem:

$$\sqrt{n}\{\hat{\theta}_d - \theta_{0,d}\}_{d \in \{0,1\}} \stackrel{a}{\sim} \mathcal{N}(0, \text{Var}(X)) \quad \text{Var}(X) = \begin{pmatrix} \mathbb{E}[X_0^2] & \mathbb{E}[X_0 \cdot X_1] \\ \mathbb{E}[X_0 \cdot X_1] & \mathbb{E}[X_1^2] \end{pmatrix}$$

Diagonal terms are of the form: for  $d \in \{0,1\}$

$$\frac{\mathbb{E}[(Y - \theta_{0,d})^2 \mathbb{I}\{D = d\}]}{\mathbb{E}[\mathbb{I}\{D = d\}]^2} = \frac{\mathbb{E}[(Y - \theta_{0,d})^2 \mid D = d] \Pr(D = d)}{\Pr(D = d)^2} = \frac{\text{Var}(Y \mid D = d)}{\Pr(D = d)}$$

and cross terms are zero because they contain  $\mathbb{I}\{D = 0\} \cdot \mathbb{I}\{D = 1\} = 0$ .

## Poll Everywhere

If we *double* the sample size in each arm of an RCT (so  $n_1, n_0$  both double), what happens to the standard error of  $\bar{Y}_1 - \bar{Y}_0$  (roughly)?

- A. It halves
- B. It shrinks by a factor of  $\sqrt{2}$
- C. It stays the same



# Back to Observational Data

Why plug-in outcome regression can mislead

## Naive CLT for the Plug-in g-Estimator

The plug-in estimator is an empirical average, so one might try:

$$\sqrt{n}(\hat{\theta}_{\text{plug-in}} - \theta_0) \stackrel{?}{\Rightarrow} \mathcal{N}\left(0, \text{Var}(\hat{g}(1, X) - \hat{g}(0, X))\right).$$

and calculate confidence intervals with  $\hat{se} = \sqrt{\text{Var}_n(\hat{g}(1, X) - \hat{g}(0, X)) / n}$ .

### What's Wrong with this approach?

This treats  $\hat{g}$  as **fixed** and **very accurate**. But  $\hat{g}$  is **learned from the same data** and **varies substantially across resamples**.

# A Concrete Failure Case: Random Forest Plug-in + Naive SE

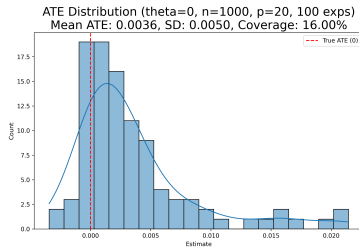
Synthetic DGP with true ATE  $\theta_0 = 0$ . Using a flexible outcome model (RF) inside the plug-in estimator, then forming a *naive* standard error as if predictions were fixed.

$$X \sim N(0, I_p),$$

$$D \sim \text{Bernoulli}(p = 0.5 + \text{clip}(X[0], -0.4, 0.4)),$$

$$Y = \theta \cdot D + X[0] + X[1] + N(0, 1)$$

```
1 # One-shot plug-in + naive SE (INCORRECT!)
2 def est(X, D, Y):
3     g = RandomForestRegressor(min_samples_leaf=20)
4     g.fit(np.c_[D, X], Y)
5     mu = g.predict(np.c_[np.ones(len(X)), X]) \
6         - g.predict(np.c_[np.zeros(len(X)), X])
7     ate_hat = mu.mean()
8     # WRONG: treats g as fixed
9     se_naive = mu.std(ddof=1) / np.sqrt(len(X))
10    return ate_hat, se_naive
```



Biased! Naive CIs under-cover.

# Two Problems to Solve

## Problem 1: Data Reuse

$\hat{g}$  is fit on the same data  
used in the average



### Fix: Cross-fitting

Out-of-fold predictions

## Problem 2: First-Order Sensitivity

Error in  $\hat{g}$  leads to error in  $\hat{\theta}$   
of the *same order*



### Fix: Orthogonal/Debiased

formula (insensitive to nuisance errors)

# In Class Activity

## Poll Everywhere

*How can we fix the problem? Try out some ideas and share back.*

Hints: change the random forest to something else, change how you use your data to train your random forest to avoid data reuse, add some extra factor to the standard error

Report Results Here



Notebook QR Code  
Or Click here!



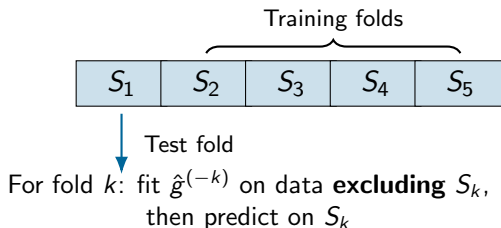


# **Solving the Data Reuse Problem: Cross-fitting**

Don't predict on samples you used to train.

## Cross-Fitting (Out-of-Fold Predictions)

Split indices into  $K$  folds  $S_1, \dots, S_K$ .



- Each  $i$  gets prediction  $\hat{g}^{(-k(i))}(D_i, X_i)$  from a model that did *not* train on  $i$ ; where  $k(i)$  is the index of the fold that  $i$ -th sample belongs.

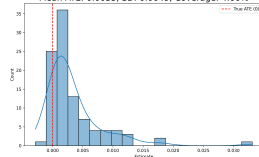
### Key Benefit

Restores (approximate) i.i.d. behavior needed for LLN/CLT arguments with flexible ML.

# Python Pseudocode: Cross-Fitted Plug-in g-Estimator

```
1 # K-fold cross-fitting for plug-in g-estimator
2 cv = KFold(n_splits=K, shuffle=True, random_state=123)
3 mu = np.zeros(n)
4 for train, test in cv.split(X):
5     g = clone(model)
6     g.fit(np.c_[D[train], X[train]], Y[train])
7     mu[test] = g.predict(np.c_[np.ones(len(test)), X[test]]) \
8                 - g.predict(np.c_[np.zeros(len(test)), X[test]])
9
10 ate_hat = mu.mean()
11 # WRONG: treats g as fixed
12 se_naive = mu.std(ddof=1) / np.sqrt(len(X))
```

Cross-Fitted ATE Distribution (theta=0, n=1000, p=20, 100 exps)  
Mean ATE: 0.0035, SD: 0.0049, Coverage: 4.00%



Notebook with Cfit



Or Click here!

**Cross-fitting fixes data reuse. But we still need to fix first-order sensitivity.**

# **Solving the First Order Sensitivity Problem: Debiasing**

Make the formula insensitive to nuisance errors

# Regression Residuals as a “Bias Detector”

## Key Idea

When  $\hat{g}$  is wrong, residuals  $Y - \hat{g}(D, X)$  can reveal systematic prediction mistakes.

**We will add a correction term that uses these residuals.**

If  $g_0(d, x) = \mathbb{E}[Y \mid D = d, X = x]$  then

$$\mathbb{E}[Y - g_0(D, X) \mid D, X] = 0.$$

So for any function  $a(D, X)$ , by the tower rule:

$$\mathbb{E}[a(D, X) \cdot (Y - g_0(D, X))] = \mathbb{E}[a(D, X) \cdot \mathbb{E}[Y - g_0(D, X) \mid D, X]] = 0$$

$$\boxed{\mathbb{E}[Y - \hat{g}(D, X) \mid D, X]} \neq 0 \Rightarrow \text{bias signal!}$$

# A Family of Candidate “Debiased” Formulas

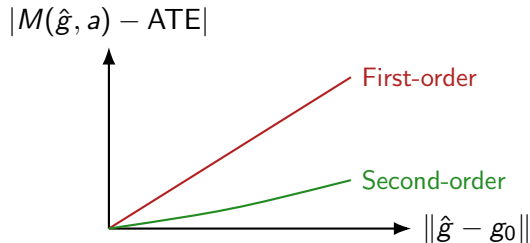
Consider

$$M(g, a) = \mathbb{E}[g(1, X) - g(0, X)] + \mathbb{E}[a(D, X)(Y - g(D, X))].$$

If  $g = g_0$ , the second term is zero, so  $M(g_0, a) = \text{ATE}$  for *any*  $a$ .

## Goal

Choose  $a = a_0$  so that  $M(g, a_0)$  is **insensitive** to small perturbations in  $g$ .



With the right  $a_0$ , error becomes  
second-order (product of RMSEs)

## Orthogonality Condition for $a_0$

Let  $\Delta(D, X)$  be any “error function” and consider  $g_t = g_0 + t\Delta$ . Then

$$\left. \frac{\partial}{\partial t} M(g_t, a) \right|_{t=0} = \mathbb{E}[\Delta(1, X) - \Delta(0, X)] - \mathbb{E}[a(D, X) \Delta(D, X)].$$

We want this derivative to be **zero for all**  $\Delta$ , i.e.

$$\boxed{\mathbb{E}[a_0(D, X) \Delta(D, X)] = \mathbb{E}[\Delta(1, X) - \Delta(0, X)] \quad \forall \Delta}$$

**Intuition:**  $a_0(D, X)$  is such that weighted average of “errors in prediction of  $Y$ ” translates them into errors in the ATE!

# The “Magic” Choice: Horvitz–Thompson Weights

Let  $p_0(x) = \Pr(D = 1 \mid X = x)$  and define for any propensity function  $p$

$$H_p(D, X) = \frac{D}{p(X)} - \frac{1 - D}{1 - p(X)}$$

## Key Identity

For any function  $\Delta(D, X)$ ,

$$\mathbb{E}[H_{p_0}(D, X) \Delta(D, X)] = \mathbb{E}[\Delta(1, X) - \Delta(0, X)].$$

So  $a_0(D, X) = H_{p_0}(D, X)$  satisfies the orthogonality condition!

**Remember from Lecture 3:** These are the IPW weights!



# Proof Sketch of the Key Identity

Fix any measurable  $\nu$ .

$$\mathbb{E}\left[\frac{D}{p_0(X)}\Delta(D, X) \mid X\right] = \frac{\mathbb{E}[D \mid X]}{p_0(X)} \Delta(1, X) = \Delta(1, X),$$

$$\mathbb{E}\left[\frac{1-D}{1-p_0(X)}\Delta(D, X) \mid X\right] = \frac{\mathbb{E}[1-D \mid X]}{1-p_0(X)} \Delta(0, X) = \nu(0, X).$$

Subtract and then take expectation over  $X$ .

**Key insight:** This identity holds for *any* function  $\Delta$ , not just the true regression  $g_0$ . It's a property of the HT weights themselves.

# Doubly Robust (DR) Formula for ATE

Plug  $a_0 = H_0$  into  $M(g, a)$ :

$$\text{ATE} = \mathbb{E}\left[g_0(1, X) - g_0(0, X) + H_{p_0}(D, X)(Y - g_0(D, X))\right].$$

Define the per-sample **DR score**:

**g-formula term + IPW correction term**

$$\psi_0(Z) = \underbrace{g_0(1, X) - g_0(0, X)}_{\text{outcome regression}} + \underbrace{H_{p_0}(D, X) \cdot (Y - g_0(D, X))}_{\text{bias correction}}$$

Then  $\text{ATE} = \mathbb{E}[\psi_0(Z)]$ .

# Why “Doubly Robust”?

Consider the population functional

$$\Theta(g, p) = \mathbb{E} \left[ g(1, X) - g(0, X) + H_p(D, X)(Y - g(D, X)) \right].$$

**If propensity is correct, i.e.,  $p = p_0$**

Even if regression  $g$  is wrong:

$$\begin{aligned} \Theta(g, p_0) &= \mathbb{E}[Y H_{p_0}(D, X)] \\ &= \text{IPW estimand} = \text{ATE} \checkmark \end{aligned}$$

**If regression is correct, i.e.,  $g = g_0$**

Even if propensity  $p$  is wrong:

$$\begin{aligned} &\text{Residual term has mean zero} \\ \Theta(g_0, p) &= \mathbb{E}[g_0(1, X) - g_0(0, X)] \\ &= \text{g-formula} = \text{ATE} \checkmark \end{aligned}$$

**Correct if either  $\hat{g}$  or  $\hat{p}$  is correct!**

## ATE Error is Second-Order (Product of RMSEs)

For any estimates  $\hat{g}, \hat{p}$  of the regression and the propensity:

$$\Theta(\hat{g}, \hat{p}) - \Theta(g_0, p_0) = \mathbb{E} \left[ (\hat{g}(D, X) - g_0(D, X)) \cdot (H_{p_0}(D, X) - H_{\hat{p}}(D, X)) \right],$$

By Cauchy–Schwarz inequality and assuming *strict overlap* ( $p(X) \in [\epsilon, 1 - \epsilon]$  for  $\epsilon > 0$ ):

$$|\Theta(\hat{g}, \hat{p}) - \Theta(g_0, p_0)| \lesssim \text{RMSE}(\hat{g}) \cdot \text{RMSE}(\hat{p})$$

### Consequence

If both nuisance learners achieve rate  $n^{-1/4}$  in RMSE, their product is  $n^{-1/2}$ , compatible with **root- $n$  inference** for  $\theta$ .

## Poll Everywhere

Suppose  $\text{RMSE}(\hat{g}) = O_p(n^{-1/4})$  and  $\text{RMSE}(\hat{p}) = O_p(n^{-1/4})$ . Roughly, what order is the DR estimation error contributed by nuisance estimation?

- A.  $O_p(n^{-1/4})$
- B.  $O_p(n^{-1/2})$
- C.  $O_p(n^{-1})$



# Estimator and Confidence Intervals

Asymptotic properties of the Doubly Robust Estimator

# The Doubly Robust Estimator (with Cross-Fitting)

Split your samples in  $K$  folds for a small constant  $K$

For each fold  $k$ ,

1. estimate  $\hat{g}^{(-k)}$  and  $\hat{p}^{(-k)}$  out-of-fold
2. for every sample  $i$  in the  $k$ -th fold compute the “plug-in” DR score estimate:

$$\hat{\psi}_i = \hat{g}^{(-k)}(1, X_i) - \hat{g}^{(-k)}(0, X_i) + H_{\hat{p}^{(-k)}}(D_i, X_i) (Y_i - \hat{g}^{(-k)}(D_i, X_i))$$

Using all the data, calculate the ATE estimate:

$$\hat{\theta}_{\text{DR}} = \frac{1}{n} \sum_{i=1}^n \hat{\psi}_i$$

# Asymptotic Normality Theorem (ATE, DR + Cross-Fitting)

## Main Theorem (DR ATE asymptotic normality)

Under mild regularity conditions, conditional ignorability and **strict** overlap, if  $\text{RMSE}(\hat{g}), \text{RMSE}(\hat{p}) \approx 0$  and

$$\sqrt{n} \text{RMSE}(\hat{g}) \text{RMSE}(\hat{p}) \approx 0, \quad (\text{Product Rate Condition})$$

then the Doubly Robust Estimator with cross-fitting satisfies

$$\sqrt{n}(\hat{\theta}_{\text{DR}} - \theta_0) \Rightarrow \mathcal{N}(0, \text{Var}(\psi_0(Z))).$$

where  $\psi_0(Z) = g_0(1, X) - g_0(0, X) + H_{p_0}(D, X)(Y - g_0(D, X))$ .

## Variance and standard error estimation:

$$\hat{V} = \text{Var}_n(\hat{\psi}_i) = \frac{1}{n} \sum_{i=1}^n (\hat{\psi}_i - \hat{\theta}_{\text{DR}})^2, \quad \hat{\text{se}} = \sqrt{\hat{V}/n}, \quad 95\% \text{ CI} = \hat{\theta}_{\text{DR}} \pm 1.96\hat{\text{se}}.$$



## Pseudocode: DR ATE with Cross-Fitting + CI (S-Learner Variant)

```
1 cv = KFold(n_splits=K, shuffle=True, random_state=123)
2 g0hat, g1hat, ghat, phat = np.zeros(n), np.zeros(n), np.zeros(n), np.zeros(n)
3
4 for train, test in cv.split(X):
5     # Outcome regression  $g(d,x) = E[Y|D=d, X=x]$ 
6     g = clone(model)
7     g.fit(np.c_[D[train], X[train]], Y[train])
8     g0hat[test] = g.predict(np.c_[np.ones(len(test)), X[test]])
9     g1hat[test] = g.predict(np.c_[np.zeros(len(test)), X[test]])
10    ghat[test] = g.predict(np.c_[D[test], X[test]])
11    # Propensity  $p(x) = P(D=1|X=x)$ 
12    phat[test] = clone(model_p).fit(X[train], D[train]).predict_proba(X[test])[:,1]
13
14 phat = np.clip(phat, 0.01, 0.99) # clipping for stability
15 Hhat = D / phat - (1 - D) / (1 - phat)
16 psi = (g1hat - g0hat) + Hhat * (Y - ghat) # DR scores
17 ate, se = psi.mean(), psi.std() / np.sqrt(n)
18 ci95 = (ate - 1.96 * se, ate + 1.96 * se)
```

## Pseudocode: DR ATE with Cross-Fitting + CI (T-Learner Variant)

```
1 cv = KFold(n_splits=K, shuffle=True, random_state=123)
2 g0hat, g1hat, phat = np.zeros(n), np.zeros(n), np.zeros(n)
3
4 for train, test in cv.split(X):
5     # Outcome regression  $g(d,x) = E[Y|D=d, X=x]$ 
6     g0hat[test] = model_y.fit(X[train][D[train]==0],
7                               Y[train][D[train]==0]).predict(X[test])
8     g1hat[test] = model_y.fit(X[train][D[train]==1],
9                               Y[train][D[train]==1]).predict(X[test])
10    # Propensity  $p(x) = P(D=1|X=x)$ 
11    phat[test] = clone(model_p).fit(X[train], D[train]).predict_proba(X[test])[:,1]
12
13 phat = np.clip(phat, 0.01, 0.99) # clipping for stability
14 ghat = g1hat * D + g0hat * (1-D)
15 Hhat = D / phat - (1 - D) / (1 - phat)
16 psi = (g1hat - g0hat) + Hhat * (Y - ghat) # DR scores
17 ate, se = psi.mean(), psi.std() / np.sqrt(n)
18 ci95 = (ate - 1.96 * se, ate + 1.96 * se)
```

## Big Picture

We get *root-n* inference for  $\theta$  without requiring asymptotic normality of ML predictions—only RMSE rates (that are slower than parametric rates).

## Poll Everywhere

Suppose we know the propensity score and we use  $\hat{p} = p_0$ . Then what rate of convergence do we need for the outcome regression so that the doubly robust estimator is asymptotically normal?

- A.  $O_p(n^{-1/4})$
- B.  $O_p(n^{-1/2})$
- C.  $O_p(n^{-1})$
- C.  $O_p(1)$
- C. No assumption.



## **Doubly Robust Estimator for the ATT**

# ATT Variant (Doubly Robust + Cross-Fitting)

Target:

$$\text{ATT} = \mathbb{E}[Y(1) - Y(0) \mid D = 1].$$

ATT g-formula:

$$\text{ATT} = \mathbb{E}[Y - g_0(X) \mid D = 1] \qquad g_0(X) = \mathbb{E}[Y \mid D = 0, X]$$

ATT g-formula estimator:

$$\widehat{\text{ATT}} = \mathbb{E}_n[Y - \hat{g}(X) \mid D = 1] = \frac{1}{n_1} \sum_{i:D_i=1} (Y_i - \hat{g}(X_i))$$

A naive standard error of  $\sqrt{\text{Var}_n(Y - \hat{g}(X) \mid D = 1)/n_1}$  would ignore errors in  $\hat{g}$ .

## Debiasing for ATT

Same debiasing idea: we can add weighted regression residuals:

$$ATT = \mathbb{E}[Y - g_0(X) \mid D = 1] + \mathbb{E}[a(X) \cdot (Y - g_0(X)) \mid D = 0]$$

To cancel the first order effect of  $g$  on ATT, we need for any “error” function  $\Delta(X)$ :

$$-\mathbb{E}[\Delta(X) \mid D = 1] = \mathbb{E}[a(X)\Delta(X) \mid D = 0] \quad (\text{“derivative” of ATT formula})$$

which leads to the ATT inverse propensity weights we saw in the last lecture:

$$a(X) = w_0(X) := \frac{p_0(X)}{1 - p_0(X)} \frac{1 - \pi}{\pi} \quad \pi = \mathbb{E}[D] = \Pr(D = 1)$$

# The Doubly Robust Formula for the ATT

Pluggin in  $w_0$  in the debiased formula, we get:

$$ATT = \mathbb{E}[Y - g_0(X) \mid D = 1] + \mathbb{E}[w_0(X) \cdot (Y - g_0(X)) \mid D = 0]$$

which can also be re-written as unconditional expectations:

$$ATT = \frac{\mathbb{E}[D \cdot (Y - g_0(X))]}{\mathbb{E}[D]} + \frac{\mathbb{E}[(1 - D) \cdot w_0(X) \cdot (Y - g_0(X))]}{\mathbb{E}[1 - D]}$$

which simplifies to:

## Population DR Formula (simplified)

$$ATT = \frac{\mathbb{E}\left[ D(Y - g_0(X)) - (1 - D) \frac{p_0(X)}{1 - p_0(X)} (Y - g_0(X)) \right]}{\mathbb{E}[D]}.$$



# The Doubly Robust Formula and its Empirical Plug-in Analogue

## Population DR Formula (simplified)

$$ATT = \frac{\mathbb{E} \left[ D (Y - g_0(X)) - (1 - D) \frac{p_0(X)}{1 - p_0(X)} (Y - g_0(X)) \right]}{\mathbb{E}[D]}.$$

## Empirical plug-in analogue

$$\widehat{ATT} = \frac{\mathbb{E}_n \left[ D (Y - \hat{g}(X)) - (1 - D) \frac{\hat{p}(X)}{1 - \hat{p}(X)} (Y - \hat{g}(X)) \right]}{\mathbb{E}_n[D]}.$$

# The Doubly Robust ATT Estimator (with Cross-Fitting)

Split your samples in  $K$  folds for a small constant  $K$

For each fold  $k$ ,

1. estimate  $\hat{g}^{(-k)}$  by regression  $Y$  on  $X$  using only out-of-fold un-treated samples
2. estimate  $\hat{p}^{(-k)}$  out-of-fold by running a classification of  $D$  on  $X$
3. for every sample  $i$  in the  $k$ -th fold compute the “plug-in” DR score estimate:

$$\hat{m}_i = D_i \cdot (Y_i - \hat{g}^{(-k)}(X_i)) + (1 - D_i) \cdot \frac{\hat{p}^{(-k)}(D_i, X_i)}{1 - \hat{p}^{(-k)}(D_i, X_i)} (Y_i - \hat{g}^{(-k)}(X_i))$$

Using all the data, calculate the ATT estimate as:

$$\hat{\theta}_{\text{ATT, DR}} = \frac{\frac{1}{n} \sum_{i=1}^n \hat{m}_i}{\frac{1}{n} \sum_{i=1}^n D_i}$$

# ATT: Doubly Robust Asymptotic Normality

$$ATT = \frac{\mathbb{E}[m_0(Z)]}{\mathbb{E}[D]} \quad m_0(Z) := D(Y - g_0(X)) - (1 - D) \frac{p_0(X)}{1 - p_0(X)} (Y - g_0(X)),$$

## Theorem (DR ATT asymptotic normality)

Under mild regularity conditions, conditional ignorability and **strict** one-sided overlap, if  $RMSE(\hat{g}), RMSE(\hat{p}) \approx 0$  and

$$\sqrt{n} RMSE(\hat{g}) RMSE(\hat{p}) \approx 0 \quad (\text{product rate condition})$$

then the Doubly Robust ATT Estimator with cross-fitting satisfies

$$\sqrt{n}(\widehat{ATT} - ATT) \Rightarrow \mathcal{N}(0, V_{ATT}), \quad V_{ATT} = \text{Var}(\phi_0(Z)).$$

where  $\phi_0$  is defined as:

$$\phi_0(Z) = \frac{m_0(Z) - ATT \cdot D}{\mathbb{E}[D]}$$

## Pseudocode: DR ATT with Cross-Fitting + CI

```
1 cv = KFold(n_splits=K, shuffle=True, random_state=123)
2 g0hat, phat = np.zeros(n), np.zeros(n)
3 for train_idx, test_idx in cv.split(X):
4     # Fit outcome model  $g(X) = E[Y|D=0, X]$  on UNTREATED training samples only
5     untreated_train = train_idx[D[train_idx] == 0]
6     g = clone(model_y)
7     g.fit(X[untreated_train], Y[untreated_train])
8     g0hat[test_idx] = g.predict(X[test_idx])
9     # Propensity  $p(x) = P(D=1|X=x)$ 
10    phat[test] = clone(model_p).fit(X[train], D[train]).predict_proba(X[test])[:,1]
11 phat = np.clip(phat, 0, 0.99) # one-sided clipping
12 # Compute the DR scores for ATT
13 ipw_weight = phat / (1 - phat)
14 m = D * (Y - g0hat) - (1 - D) * ipw_weight * (Y - g0hat)
15 # Compute ATT Estimate, Standard Error and 95% CI
16 att = m.mean() / D.mean()
17 phi = (m - att * D) / D.mean()
18 se = phi.std() / np.sqrt(n)
19 ci95 = (att - 1.96 * se, att + 1.96 * se)
```

## Summary: The 4 Key Takeaways

1. Identification gives an infinite-data formula; **estimation** turns it into a finite- $n$  procedure.
2. CIs come from **asymptotic normality** + **variance estimation** (standard errors).
3. Plug-in g-estimation can fail with ML due to **data reuse** and **first-order sensitivity**.
4. **Cross-fitting** + **doubly robust scores** deliver valid root- $n$  CIs under only product RMSE conditions.

## Next Time

- ▶ How do modern ML methods achieve small RMSE for outcome regression or propensity estimation?
- ▶ Bias-variance tradeoffs, regularization, model selection. How to select among models and what diagnostic metrics to report?
- ▶ Why  $n^{-1/4}$  RMSE is often attainable in high dimensions with structure.

Lasso

Random Forest

Gradient Boosting

Neural Networks

Ensembling/Stacking

AutoML