

Chapter 1

Basics of Causal Inference

Zhan Gao, February 11, 2026

Adapted from [Zhentao Shi's lecture notes](#) and heavily dependent on [Wager \(2025\)](#).

1.1 Potential Outcome Framework

This a thought experiment. An individual i has two potential outcomes $Y_i(1)$ and $Y_i(0)$, and the individual treatment effect is the difference $\Delta_i = Y_i(1) - Y_i(0)$. However, no one can step into the same river twice. One and only one of the potential outcomes will be realized, and therefore Δ_i is unobservable. Given the treatment status $D_i \in \{0, 1\}$, in reality what we can be observed is the realized outcome

$$Y_i = Y_i(D_i) = D_i Y_i(1) + (1 - D_i) Y_i(0). \quad (1.1)$$

Stable Unit Treatment Values Assumption. (1.1) is often referred in the literature as the *Stable Unit Treatment Values Assumption* (SUTVA) assumption, that assumes that the treatment on one individual does not interfere the outcome of others. SUTVA may be questionable in the presence of interference through social interactions, for example see [Leung \(2022\)](#). In this note, we maintain the SUTVA assumption, and for now suppress subscript i for simplicity.

We want to use the observable Y to learn the *average treatment effect* (ATE)

$$\text{ATE} := \mathbb{E}[\Delta] = \mathbb{E}[Y(1)] - \mathbb{E}[Y(0)].$$

Randomized controlled trials. If the treatment is randomly assigned by flipping a coin

(the coin does not need to be even), then

$$\begin{pmatrix} Y(1) \\ Y(0) \end{pmatrix} \perp D. \quad (1.2)$$

This is the simplest framework for identifying causal effects, as well as the most reliable level of proof in terms of *interval validity*. [Abadie and Cattaneo \(2018\)](#) refer to RCTs as the “gold standard”.

Characterization of ATE. (1.2) implies

$$\mathbb{E}[Y|D=1] = \mathbb{E}[DY(1) + (1-D)Y(0)|D=1] = \mathbb{E}[Y(1)|D=1] = \mathbb{E}[Y(1)].$$

where the last equality holds by randomized treatment assignment. Since (Y, D) are observable, the LHS is operational. Independence between $Y(1)$ and D ensures that the conditional expectation $\mathbb{E}[Y(1)|D=1]$ equals the unconditional expectation $\mathbb{E}[Y(1)]$. Similarly,

$$\mathbb{E}[Y|D=0] = \mathbb{E}[Y(0)|D=0] = \mathbb{E}[Y(0)].$$

Under RCT, we have an operational formula for ATE:

$$\text{ATE} := \tau = \mathbb{E}[Y|D=1] - \mathbb{E}[Y|D=0].$$

Given the data, we mimic the population average to compute

$$\widehat{\text{ATE}} := \widehat{\tau} = \frac{1}{n_1} \sum_{\{i:D_i=1\}} Y_i - \frac{1}{n_0} \sum_{\{i:D_i=0\}} Y_i$$

where $n_1 = \sum_{i=1}^n \mathbf{1}\{D_i = 1\}$ and $n_0 = \sum_{i=1}^n \mathbf{1}\{D_i = 0\}$. This is often referred as the *difference-in-means* estimator.

The above approach uses conditioning, which is intuitive. There is an alternative, yet equivalent way to ATE, using only unconditional quantities. Notice that

$$DY = D^2Y(1) + D(1-D)Y(0) = DY(1)$$

since for either $D \in \{0, 1\}$, we have $D^2 = D$ and $D(1-D) = 0$. Again, Y depends on D but $Y(1)$ is independent of D , and therefore

$$\mathbb{E}[DY] = \mathbb{E}[DY(1)] = \mathbb{E}[D]\mathbb{E}[Y(1)] = \Pr[D=1]\mathbb{E}[Y(1)]$$

and therefore

$$\mathbb{E}[Y(1)] = \frac{\mathbb{E}[DY]}{\Pr[D=1]} = \frac{\mathbb{E}[\mathbb{I}(D=1)Y]}{\Pr[D=1]}$$

if $\Pr[D=1] \neq 0$. The denominator $\Pr[D=1]$ is called the *propensity score*: the probability that a person is assigned into the treatment group.

Similarly, if $\Pr[D=0] \neq 0$ we have

$$\mathbb{E}[Y(0)] = \frac{\mathbb{E}[(1-D)Y]}{\Pr[D=0]} = \frac{\mathbb{E}[\mathbb{I}(D=0)Y]}{\Pr[D=0]}.$$

Therefore, if $\Pr[D=1] \in (0, 1)$ the ATE is

$$\text{ATE} = \frac{\mathbb{E}[DY]}{\Pr[D=1]} - \frac{\mathbb{E}[(1-D)Y]}{\Pr[D=0]}. \quad (1.3)$$

ATE is the difference of the two ratios.

Given data, we can compute $\mathbb{E}[DY] / \Pr[D=1]$ as

$$\frac{(\sum_{i=1}^n D_i Y_i)/n}{n_1/n} = \frac{\sum_{i=1}^n D_i Y_i}{n_1} = \frac{1}{n_1} \sum_{\{i: D_i=1\}} Y_i.$$

It is easy to see that, the sample version is the same either we compute via the conditioning or the propensity score. Their equivalence is analogous to the fact that the conditional density can be written as the ratio of the joint density and the marginal density.

Theorem 1 (Wager (2025) Theorem 1.2). *Suppose $(Y_i(1), Y_i(0))$ are i.i.d. draws from some super-population \mathbb{P} . Treatments are assigned based on a Bernoulli trial:*

$$D_i \stackrel{\text{i.i.d.}}{\sim} \text{Bern}(\pi), \quad 0 < \pi < 1.$$

Then

$$\sqrt{n}(\hat{\tau} - \tau) \xrightarrow{d} \mathcal{N}\left(0, \frac{\sigma_1^2}{\pi} + \frac{\sigma_0^2}{1-\pi}\right), \quad (1.4)$$

where $\sigma_d^2 := \text{Var}(Y_i(d))$ for $d \in \{0, 1\}$.

1.1.1 Stratification

At a granular level, the researcher also observe some confounding factor (causal inference term) (alternatively, in plain stat term it is called covariate) X .

The conditional ATE (CATE)

$$\tau(x) := \text{ATE}(x) = \mathbb{E}[\Delta|X = x] = \mathbb{E}[Y(1)|x] - \mathbb{E}[Y(0)|x]$$

can vary across different realizations of X . For example, an vaccine is more effective to children than adults. In RCT, if the random treatment assignment is regardless of the age group, we have $(Y(1), Y(0))$ depend X , while $D \perp X$, then we maintains (1.2).

Stratification. For different age groups, we may want to use different treatment assignment probability. For example, we put 70% of the children into the treatment, while we put 40% of the adults into the treatment. In this case, D also depends on X . RCT can be implemented in a stratified approach: inside each age group, the researcher runs a RCT by flipping a coin to assign treatment.

To formalize the discussion, suppose the covariate $X \in \mathcal{X}$ takes discrete values, where $|\mathcal{X}| = m < \infty$. Assume further that

$$\begin{pmatrix} Y(1) \\ Y(0) \end{pmatrix} \perp D | X = x \text{ for all } x \in \mathcal{X}. \quad (1.5)$$

Caveat on aggregation. Notice that when D depends on x ,

$$\mathbb{E}\{\mathbb{E}[Y|D = d, x]\} \neq \mathbb{E}[Y|D = d]$$

and thus

$$\text{ATE} \neq \mathbb{E}[Y|D = 1] - \mathbb{E}[Y|D = 0].$$

With that being said, the difference-in-means estimator

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

no longer provides a proper estimator for ATE.

Example 1. (Wager, 2025, Chapter 2, p.17) We are interested in evaluating whether providing teenagers with cash incentives can discourage smoking. To this end, an experiment was carried out in two locations: Palo Alto, CA, and Geneva, Switzerland.

Palo Alto	Non-Smoker	Smoker	Ratio (Smoker / Total)
Treatment	152	5	5 / (152 + 5) = 0.032
Control	2362	122	122 / (2362 + 122) = 0.049

Geneva	Non-Smoker	Smoker	Ratio (Smoker / Total)
Treatment	581	350	$350 / (581 + 350) = 0.376$
Control	2278	1979	$1979 / (2278 + 1979) = 0.465$

The results indicate that the treatment reduces the smoking rate among teenagers in both locations. However, looking at aggregate data can be misleading:

Palo Alto + Geneva	Non-Smoker	Smoker	Ratio (Smoker / Total)
Treatment	733	401	$401 / (733 + 401) = 0.354$
Control	4640	2101	$2101 / (4640 + 2101) = 0.312$

This echos the *Simpson's paradox*. Intuitively, we need to re-weight the two sample taking into consideration that Genevans are both more likely to be treated and more likely to smoke. This results in the following estimation:

$$\begin{aligned}\hat{\tau}_{PA} &= \frac{5}{152+5} - \frac{122}{2362+122} \approx -1.7\% \\ \hat{\tau}_{GVA} &= \frac{350}{350+581} - \frac{1979}{2278+1979} \approx -8.9\% \\ \hat{\tau} &= \frac{2641}{2641+5188} \hat{\tau}_{PA} + \frac{5188}{2641+5188} \hat{\tau}_{GVA} \approx -6.5\%\end{aligned}$$

Stratified Estimation. Example 1 suggests we ought to estimate the ATE by aggregating the estimates of CATE by

$$\hat{\tau}^S = \sum_{x \in \mathcal{X}} \frac{n_x}{n} \underbrace{\left(\frac{1}{n_{x1}} \sum_{\{X_i=x, D_i=1\}} Y_i - \frac{1}{n_{x0}} \sum_{\{X_i=x, D_i=0\}} Y_i \right)}_{\hat{\tau}(x)} = \sum_{x \in \mathcal{X}} \frac{n_x}{n} \hat{\tau}(x), \quad (1.6)$$

where $n_x = |\{i : X_i = x\}|$ and $n_{nd} = |\{i : D_i = d, X_i = x\}|$.

Theorem 2 (Wager (2025) Theorem 2.1). *Suppose $(Y_i(1), Y_i(0), D_i, X_i) \sim_{\text{i.i.d.}} \mathbb{P}$ for some super-population \mathbb{P} . $X_i \in \mathcal{X}$ where $|\mathcal{X}| = m < \infty$. $\mathbb{E}(Y_i(d)^2 | X_i) < \infty$. Assume stratified design (1.5) and SUTVA hold, and there is nontrivial treatment variation in \mathcal{X} so that $p(x) = \Pr(D_i = 1 | X_i = x) \in (0, 1)$ for all $x \in \mathcal{X}$. Then*

$$\sqrt{n} (\hat{\tau}^S - \tau) = \mathcal{N}(0, V^S),$$

where

$$\text{Var}(\tau(X_i)) + \mathbb{E} \left[\frac{\sigma_1^2(X_i)}{p(X_i)} + \frac{\sigma_0^2(X_i)}{1-p(X_i)} \right]. \quad (1.7)$$

For a straightforward proof, see Wager (2025, pp. 19–20).

1.1.2 Inverse Propensity Weighting

The intuition from the stratified design can be in general formalized as the *unconfoundedness condition* (or *conditional independence assumption*):

$$\begin{pmatrix} Y(1) \\ Y(0) \end{pmatrix} \perp D \mid X. \quad (1.8)$$

Qualitatively, (1.8) requires that we have measured enough covariates to capture the dependence between treatment assignments and potential outcomes so that D cannot “peek” at $(Y(1), Y(0))$. $x \in \mathcal{X}$ is now allowed to be either discrete or continuous. When \mathcal{X} is continuous, we cannot estimate $\hat{\tau}(x)$ for each $x \in \mathcal{X}$ as in the stratified design, since there will not be enough observations at each value of x to compute the ATE reliably.

Characterization of ATE via Inverse Propensity Weighting. We take another route to do the computation via the propensity score as for (1.3). Denote the (conditional) propensity score as $p(x) = \Pr[D = 1 \mid X = x]$. Then

$$\tau(x) := \text{ATE}(x) = \frac{\mathbb{E}[DY|x]}{p(x)} - \frac{\mathbb{E}[(1-D)Y|x]}{1-p(x)}, \quad (1.9)$$

given that $0 < p(x) < 1$. Compared to (1.18), the (1.9) makes the role of $p(x)$ explicit. If we use (1.9) instead of (1.18), then

$$\text{ATE} = \mathbb{E}[\text{ATE}(x)] = \mathbb{E}\left[\frac{\mathbb{E}[DY|x]}{p(x)} - \frac{\mathbb{E}[(1-D)Y|x]}{1-p(x)}\right], \quad (1.10)$$

where the different $p(x)$ is explicitly accounted. Although (1.19) and (1.10) are mathematically equivalent, when we use the sample average to mimic the population average, (1.10) is easier to work with as it conditions only on the random variable X , but not D .

The unconfoundedness condition (1.8) seems impractical at the first glance. however, as shown in [Rosenbaum and Rubin \(1983\)](#), the propensity score can work as a dimension reduction tools which can make (1.8) more tractable: If (1.8) holds, then

$$\begin{pmatrix} Y(1) \\ Y(0) \end{pmatrix} \perp D \mid p(X). \quad (1.11)$$

It suffices to control for $p(X)$ rather than X to account for the non-random treatment assignment. It is easy to verify that ([Wager, 2025](#), p.21)

$$\Pr[D = d \mid Y(0), Y(1), p(X)] = \Pr[D = d \mid X]. \quad (1.12)$$

Suppose the functional form of $p(x)$ is known, and we can use the *inverse propensity weighting* (IPW) estimator

$$\tilde{\tau}^{\text{IPW}} = \frac{1}{n} \sum_{i=1}^n \left[\frac{D_i Y_i}{p(X_i)} - \frac{(1 - D_i) Y_i}{1 - p(X_i)} \right]. \quad (1.13)$$

Notice

$$\begin{aligned} \mathbb{E} \left[\frac{D_i Y_i}{p(X_i)} \right] &= \mathbb{E} \left\{ \mathbb{E} \left[\frac{D Y}{p(X)} \middle| X \right] \right\} = \mathbb{E} \left\{ \frac{1}{p(X)} \mathbb{E} [D Y | X] \right\} = \mathbb{E} \left\{ \frac{1}{p(X)} \mathbb{E} [D Y(1) | X] \right\} \\ &= \mathbb{E} \left\{ \frac{\mathbb{E} [D | x]}{p(X)} \mathbb{E} [Y(1) | X] \right\} = \mathbb{E} \{ \mathbb{E} [Y(1) | X] \} = \mathbb{E} [Y(1)], \end{aligned}$$

where we apply the unconfoundedness condition (1.8) from the first to the second line, and similarly

$$\mathbb{E} \left[\frac{(1 - D_i) Y_i}{1 - p(X_i)} \right] = \mathbb{E} [Y(0)].$$

We have $\tilde{\tau}^{\text{IPW}}$ is an unbiased estimator of ATE:

$$\mathbb{E} [\tilde{\tau}^{\text{IPW}}] = \mathbb{E} [Y(1)] - \mathbb{E} [Y(0)] = \text{ATE}.$$

Theorem 3. Suppose $(Y_i(1), Y_i(0), D_i, X_i) \sim_{\text{i.i.d.}} \mathbb{P}$. Assume the unconfoundedness condition (1.8) and SUTVA hold. $\mathbb{E} (Y_i(d)^2 | X_i) < \infty$ and there $\eta > 0$ s.t.

$$\epsilon < p(x) < 1 - \epsilon \text{ for all } x \in \mathcal{X}. \quad (1.14)$$

Then,

$$\sqrt{n} (\tilde{\tau}^{\text{IPW}} - \tau) = \mathcal{N} (0, V^{\text{IPW}}),$$

where

$$V^{\text{IPW}} = \text{Var} (\tau (X_i)) + \mathbb{E} \left[\frac{\sigma_1^2 (X_i)}{p(X_i)} + \frac{\sigma_0^2 (X_i)}{1 - p(X_i)} \right] + \mathbb{E} \left[\frac{(\mu_0 (X_i) + (1 - p(X_i)) \tau (X_i))^2}{p(X_i) (1 - p(X_i))} \right]. \quad (1.15)$$

Refer to [Wager \(2025\)](#), p.23 - 24) for a proof.

Overlap Assumption. (1.14), together with the finite second moments of the outcomes, ensures all moments involved in (1.15) are finite. This is known as the *overlap* assumption. In general, we need to guarantee there is enough randomness in treatment assignment to justify the treatment effect estimation, i.e. D_i cannot be perfectly predicted by X_i . Technically,

if we strengthen the finite second moment assumption to assume outcomes are uniformly bounded, then (1.14) can be relaxed to

$$\mathbb{E} \left(\frac{1}{p(X_i)(1-p(X_i))} \right) < \infty,$$

which often referred as weak overlap condition.

Experimental v.s. observational studies. There are two conceptually distinct ways in which potential outcomes can satisfy (1.5). The foregoing discussion primarily addressed the case where data arise from an experiment employing stratified treatment assignment. In this experimental setting, nature first draws $\{Y_i(0), Y_i(1), X_i\} \sim \mathbb{P}$. The experimenter then assigns treatment according to $D_i \sim \text{Bern}(p(X_i))$, where the propensity score $p(\cdot)$ is deliberately chosen and thus known. Alternatively, we may encounter data $(Y_i(1), Y_i(0), D_i, X_i) \sim \mathbb{P}$ that do not originate from a controlled experiment; in such cases, the unconfoundedness assumption is invoked to enable causal identification. These latter scenarios are commonly referred to as *observational studies* or *natural experiments*. In such settings, justifying the unconfoundedness assumption is often challenging. Furthermore, the propensity score $p(x)$ is unknown to the researcher and must be estimated from the data.

Feasible IPW. In observational studies, we often don't know the exact propensity score $p(x)$. To construct the *feasible* IPW estimator, we need to plug in the estimates of the propensity score, $\hat{p}(x)$, in (1.13):

$$\hat{\tau}^{\text{IPW}} = \frac{1}{n} \sum_{i=1}^n \left(\frac{D_i Y_i}{\hat{p}(X_i)} - \frac{(1 - D_i) Y_i}{1 - \hat{p}(X_i)} \right). \quad (1.16)$$

In practice, estimating the propensity score accurately can be difficult. Theoretical guarantees as in Theorem 3 can break down if we plug in generic estimates of the propensity without further underlying assumptions. Robustness of methods to errors in the propensity scores is important.

Inefficiency of $\hat{\tau}^{\text{IPW}}$. Let's revisit the stratified design, and compare the asymptotic variance of the IPW estimator $\hat{\tau}^{\text{IPW}}$ and the stratification estimator $\hat{\tau}^S$, i.e. V^S in (1.7) and V^{IPW} in (1.15):

$$V^{\text{IPW}} = V^S + \mathbb{E} \underbrace{\left[\frac{(\mu_0(X_i) + (1 - p(X_i))\tau(X_i))^2}{p(X_i)(1 - p(X_i))} \right]}_{\geq 0} \geq V^S,$$

i.e. $\hat{\tau}^{\text{IPW}}$ is strictly less efficient than $\hat{\tau}^S$ unless $\mu_0(X_i) + (1 - p(X_i))\tau(X_i) = 0$ almost surely.

Zooming into $\hat{\tau}^S$, it turns on $\hat{\tau}^S$ is actually the *feasible* IPW estimator

$$\hat{\tau}^S = \hat{\tau}^{\text{IPW}},$$

where $\hat{p}(x) = \frac{n_{x1}}{n_x}$. Surprisingly, this yields the seemingly paradoxical result that the *oracle* IPW estimator is actually less efficient than the *feasible* IPW estimator when \mathcal{X} is discrete. [Hirano et al. \(2003\)](#) formalize this insight for the case of continuous covariates, demonstrating that the IPW estimator, based on a nonparametrically estimated propensity score obtained from a sufficiently smooth, growing sieve approximation, attains the semiparametric efficiency bound.

Example 2 (Missing-at-random with binary X : why \hat{p} can reduce variance). We consider the missing-at-random setup used in Section 3 of [Hirano et al. \(2003\)](#). Let $X \in \{0, 1\}$ be a binary covariate, $D \in \{0, 1\}$ indicate whether Y is observed, and suppose the (observation) propensity score is

$$p(x) = \Pr(D = 1 \mid X = x) = \frac{1}{2},$$

and $D \perp Y \mid X$. The outcome satisfies the structural equation

$$Y = \mu(X) + \varepsilon, \quad \mathbb{E}[\varepsilon \mid X] = 0, \quad \text{Var}(\varepsilon \mid X) = \text{Var}(Y \mid X).$$

The target is the population mean $\theta = \mathbb{E}[Y] = \mathbb{E}[\mu(X)]$.

Denote

$$n_x = \sum_{i=1}^n \mathbf{1}\{X_i = x\}, \quad n_{x1} = \sum_{i=1}^n D_i \mathbf{1}\{X_i = x\}.$$

Using the true propensity score yields

$$\tilde{\theta} = \frac{1}{n} \sum_{i=1}^n \frac{D_i Y_i}{p(X_i)},$$

which is the Horvitz–Thompson estimator ([Horvitz and Thompson, 1952](#)). To see where inefficiency comes from, decompose the estimator using $Y_i = \mu(X_i) + \varepsilon_i$:

$$\begin{aligned} \tilde{\theta} &= \frac{1}{n} \sum_{i=1}^n \frac{D_i \mu(X_i)}{p(X_i)} + \frac{1}{n} \sum_{i=1}^n \frac{D_i \varepsilon_i}{p(X_i)} \\ &= \frac{1}{n} \sum_{x \in \{0,1\}} \mu(x) \frac{\sum_{i=1}^n D_i \mathbf{1}\{\mathbf{X}_i = \mathbf{x}\}}{p(x)} + \frac{1}{n} \sum_{i=1}^n \frac{D_i \varepsilon_i}{p(X_i)} \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{n} \sum_{x \in \{0,1\}} \mu(x) \underbrace{\frac{n_x}{p(x)}}_{:= \tilde{n}_x} + \frac{1}{n} \sum_{i=1}^n \frac{D_i \varepsilon_i}{p(X_i)} \\
&= \sum_{x \in \{0,1\}} \frac{n_x}{n} \mu(x) + \frac{1}{n} \sum_{x \in \{0,1\}} (\tilde{n}_x - n_x) \mu(x) + \frac{1}{n} \sum_{i=1}^n \frac{D_i \varepsilon_i}{p(X_i)} \tag{1.17}
\end{aligned}$$

- The first term in (1.17) represents the “full-sample” mean of $\mu(X)$, which can be further expanded as

$$\begin{aligned}
\sum_{x \in \{0,1\}} \frac{n_x}{n} \mu(x) &= \frac{1}{n} \sum_{i=1}^n \mu(X_i) \\
&= \sum_{x \in \{0,1\}} \Pr(X = x) \mu(x) + \sum_{x \in \{0,1\}} \left(\frac{n_x}{n} - \Pr(X = x) \right) \mu(x) \\
&= \theta + \sum_{x \in \{0,1\}} \left(\frac{n_x}{n} - \Pr(X = x) \right) \mu(x).
\end{aligned}$$

This term appropriately re-weights each group defined by $x \in \{0, 1\}$ in the finite sample according to its sample fraction $\frac{n_x}{n}$. It is centered at θ and its contribution to the variance arises solely from sampling variability in n_x , which in finite samples follows a binomial distribution with parameters n and $\Pr(X = x)$.

- The second term in (1.17) arises because the oracle propensity score does not guarantee any *finite-sample* balance. Due to the randomness of missingness within each group, the observed subsample may over- or under-represent a group, as captured by the discrepancy $\tilde{n}_x - n_x$. This imbalance introduces additional variation into the estimator.
- The third term in (1.17) contains the idiosyncratic noises.

we can compute the variance of the estimator,

$$\begin{aligned}
n \text{Var}(\tilde{\theta}) &= n \left[\mathbb{E} \left(\text{Var} \left(\tilde{\theta} \mid X \right) \right) + \text{Var} \left(\mathbb{E} \left(\tilde{\theta} \mid X \right) \right) \right] \\
&= \text{Var}(\mu(X)) + \mathbb{E} \left(\frac{\text{Var}(Y \mid X)}{p(X)} \right) + \mathbb{E} \left(\frac{1 - p(X)}{p(X)} \mu(X)^2 \right)
\end{aligned}$$

- The term $\text{Var}(\mu(X))$ is the variance contribution of the mean of the outcome.
- The term $\mathbb{E} \left(\frac{\text{Var}(Y \mid X)}{p(X)} \right)$ represents the irreducible idiosyncratic noise ε .
- The term $\mathbb{E} \left(\frac{1 - p(X)}{p(X)} \mu(X)^2 \right)$ reflects the variance contribution of the imbalance term.

Note that in (1.17), replacing $p(x)$ by the estimated propensity score $\hat{p}(x) = \frac{n_{x1}}{n_x}$ leads to

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n \frac{D_i Y_i}{\hat{p}(X_i)}.$$

Now $\hat{n}_x = \frac{n_{x1}}{\hat{p}(x)} = n_x$, the imputed weights exactly balance the realized subsamples, so we can follow the calculation in (1.17) and get Therefore,

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n \mu(X_i) + \frac{1}{n} \sum_{i=1}^n \frac{D_i \varepsilon_i}{\hat{p}(X_i)},$$

in which the imbalance term cancels out. By the standard CLT argument, we can show¹

$$\sqrt{n} \hat{\theta} \xrightarrow{d} \mathcal{N} \left(0, \text{Var}(\mu(X)) + \mathbb{E} \left(\frac{\text{Var}(Y|X)}{p(X)} \right) \right).$$

Thus the efficiency gain is exactly the removal of the $\mathbb{E} \left(\frac{1-p(X)}{p(X)} \mu(X)^2 \right)$ component.

1.1.3 Doubly Robust Estimator

Characterization of ATE via Outcome Regression. The CATE is made operational via

$$\begin{aligned} \text{ATE}(x) &= \mathbb{E}[Y(1)|X=x] - \mathbb{E}[Y(0)|X=x] \\ &= \mathbb{E}[Y(1)|D(x)=1, x] - \mathbb{E}[Y(0)|D(x)=0, X_i=x] \\ &= \mathbb{E}[Y|D=1, x] - \mathbb{E}[Y|D=0, x]. \end{aligned} \tag{1.18}$$

Inside each age group, we just need to compute the difference between the averages of those treated and those untreated, respectively. Now, suppose the researcher has CATE at hand, and she wants to aggregate the CATE across subgroups into an overall ATE. Then by the law of iterated expectations:

$$\text{ATE} = \mathbb{E}[\text{ATE}(X)] = \mathbb{E}\{\mathbb{E}[Y|D=1, X] - \mathbb{E}[Y|D=0, X]\} = \mathbb{E}[\mu_1(X) - \mu_0(X)], \tag{1.19}$$

¹Observe that $\hat{\theta} = \sum_{x \in \{0,1\}} \frac{n_x}{n_{x1}} \left(\frac{1}{n} \sum_{i=1}^n D_i \mathbf{1}\{X_i=x\} \varepsilon_i \right)$. To analyze the asymptotic distribution, apply the Central Limit Theorem (CLT) separately to each term within the parentheses for $x \in \{0,1\}$, and then invoke Slutsky's theorem to conclude.

where $\mu_d(x) = \mathbb{E}[Y|D=d, X=x]$. This leads to a simple and consistent (*but not necessarily optimal*) nonparametric regression estimator for ATE:

$$\hat{\tau}^{\text{REG}} = \frac{1}{n} \sum_{i=1}^n (\hat{\mu}_1(X_i) - \hat{\mu}_0(X_i)), \quad (1.20)$$

where $\hat{\mu}_j(\cdot)$ is a consistent nonparametric estimators for $\mu_j(\cdot)$.

Augmented Inverse Propensity Weighting (AIPW). Given the two characterizations, ATE can be estimated by either the IPW estimator (1.16) or the outcome regression estimator (1.20). Both estimators require the estimation of nuisance components ($p(\cdot)$ in $\hat{\tau}^{\text{IPW}}$ and $\mu_d(\cdot)$ in $\hat{\tau}^{\text{REG}}$). Consequently, the validity of these estimators relies heavily on the accuracy with which these nuisance components are estimated. Though the original motivation was mainly theoretical relying the *semiparametric efficiency theory*², it is natural to ask whether it is possible to combine both strategies to mitigate the bias and improve efficiency; see brief discussion in [Ding \(2024, Section 12.2\)](#) and references therein. This leads to the augmented inverse propensity weighting (AIPW) estimator ([Robins et al., 1994](#)),

$$\hat{\tau}^{\text{AIPW}} = \frac{1}{n} \sum_{i=1}^n \left[\hat{\mu}_1(X_i) - \hat{\mu}_0(X_i) + \frac{D_i(Y_i - \hat{\mu}_1(X_i))}{\hat{p}(X_i)} - \frac{(1 - D_i)(Y_i - \hat{\mu}_0(X_i))}{1 - \hat{p}(X_i)} \right] \quad (1.21)$$

The AIPW estimator can be interpreted as first making a best effort at estimating τ by modeling the conditional means $\mu_0(x)$ and $\mu_1(x)$. Then, it adjusts for any remaining bias in the estimated conditional means $\hat{\mu}_d(x)$ by applying IPW to the regression residuals. Statistically, AIPW not only inherits the robustness properties of both the regression and IPW estimators, but it can also improve upon both.

Weak Double Robustness. $\hat{\tau}^{\text{AIPW}}$ is consistent if *either* $\hat{\mu}_d(x)$ are consistent *or* $\hat{p}(x)$ is consistent. To see this, first consider the case where $\hat{\mu}_d(x)$ is consistent, i.e., $\hat{\mu}_d(x) \approx \mu_d(x)$. Then,

$$\begin{aligned} \hat{\tau}^{\text{AIPW}} &= \underbrace{\frac{1}{n} \sum_{i=1}^n (\hat{\mu}_1(X_i) - \hat{\mu}_0(X_i))}_{=\hat{\tau}^{\text{REG}}} \\ &\quad + \underbrace{\frac{1}{n} \sum_{i=1}^n \left(\frac{D_i}{\hat{p}(X_i)} (Y_i - \hat{\mu}_1(X_i)) - \frac{1 - D_i}{1 - \hat{p}(X_i)} (Y_i - \hat{\mu}_0(X_i)) \right)}_{\approx \text{mean-zero noise}}, \end{aligned}$$

because $\mathbb{E}(Y - \hat{\mu}_d(X) | X, D) \approx 0$ under unconfoundedness, which zeros out the estimation

²See [Bodhisattva Sen's note](#) for a accessible introduction to semiparametric efficiency theory, and refer to [Bickel et al. \(1993\)](#) for the rigorous treatment.

error in the propensity score. Conversely, suppose $\hat{p}(x)$ is consistent, i.e. $\hat{p}(x) \approx p(x)$, then

$$\begin{aligned}\hat{\tau}^{\text{AIPW}} &= \underbrace{\frac{1}{n} \sum_{i=1}^n \left(\frac{D_i Y_i}{\hat{p}(X_i)} - \frac{(1 - D_i) Y_i}{1 - \hat{p}(X_i)} \right)}_{= \hat{\tau}^{\text{IPW}}} \\ &\quad + \underbrace{\frac{1}{n} \sum_{i=1}^n \left(\hat{\mu}_1(X_i) \left(1 - \frac{D_i}{\hat{p}(X_i)} \right) - \hat{\mu}_0(X_i) \left(1 - \frac{1 - D_i}{1 - \hat{p}(X_i)} \right) \right)}_{\approx \text{mean-zero noise}},\end{aligned}$$

since $\mathbb{E}(1 - D/\hat{p}(X) | x) \approx 0$. Again, this can zero out the estimation error in the outcome regressions even if $\hat{\mu}_d(x)$ is inconsistent.

Remark 1. Weak robustness ensures only the consistency of $\hat{\tau}^{\text{AIPW}}$. If both $\hat{\mu}_d(x)$ and $\hat{p}(x)$ are consistently estimated using appropriate nonparametric or machine learning methods, then weak double robustness does not, by itself, provide further advantages. In practice, however, we are often interested not just in consistency, but also in the rate of convergence and in uncertainty quantification. Modern machine learning algorithms used to estimate $\mu_d(x)$ and $p(x)$ often cannot achieve the parametric \sqrt{n} rate but rather slower rate like $n^{1/4}$ unless strong assumptions are imposed. This natural question is how robust is the performance of the AIPW estimator to the quality of the nuisance estimates $\mu_d(x)$ and $p(x)$, compared to the oracle AIPW estimator in which both nuisance components are known.

Oracle AIPW. The *oracle* AIPW estimator is constructed with true outcome means $\mu_1(x)$ and $\mu_0(x)$ and propensity score $p(x)$:

$$\tilde{\tau}^{\text{AIPW}} = \frac{1}{n} \sum_{i=1}^n \left[\mu_1(X_i) - \mu_0(X_i) + \frac{D_i (Y_i - \mu_1(X_i))}{p(X_i)} - \frac{1 - D_i (Y_i - \mu_0(X_i))}{1 - p(X_i)} \right]. \quad (1.22)$$

Theorem 4 (Wager (2025) Proposition 3.1). *Under the assumptions of Theorem 3,*

$$\sqrt{n} (\tilde{\tau}^{\text{AIPW}} - \tau) = \mathcal{N}(0, V^*),$$

where

$$V^* = \text{Var}(\tau(X_i)) + \mathbb{E} \left[\frac{\sigma_1^2(X_i)}{p(X_i)} + \frac{\sigma_0^2(X_i)}{1 - p(X_i)} \right]. \quad (1.23)$$

It is easy to show the result by applying central limit theorem; see Wager (2025, Proposition 3.1) for a proof.

Remark 2. Note that $V^* = V^S$, but under general conditions beyond stratified design with discrete covariates. In fact, V^* is the *efficiency bound* for nonparametric ATE estimation

under unconfoundedness. [Wager \(2025, Theorem 3.4\)](#) provides a proof sketch using the argument of [Chamberlain \(1992\)](#). Given $\tilde{\tau}^{\text{AIPW}}$ achieves the *optimal* asymptotic behavior, for any other operational estimator, such as $\hat{\tau}^{\text{AIPW}}$, we aim to establish the asymptotic equivalence that $\sqrt{n} (\hat{\tau} - \tilde{\tau}^{\text{AIPW}}) = o_p(1)$.

Population bias decomposition of $\hat{\tau}^{\text{AIPW}}$. Let $\eta := (\mu_0, \mu_1, p)$ collect nuisance functions and let $\eta_0 := (\mu_{0,0}, \mu_{1,0}, p_0)$ denote the truth, where

$$p_0(x) := \Pr(D = 1 \mid X = x), \mu_{d,0}(x) := \mathbb{E}(Y \mid D = d, X = x), d \in \{0, 1\}.$$

Let $W := (Y, D, X)$. Define the AIPW score

$$\psi(W; \eta) := \mu_1(X) - \mu_0(X) + \frac{D}{p(X)}(Y - \mu_1(X)) - \frac{1 - D}{1 - p(X)}(Y - \mu_0(X)), \quad (1.24)$$

and the associated moment function

$$m(W; \tau, \eta) := \psi(W; \eta) - \tau. \quad (1.25)$$

Under unconfoundedness and overlap, the target $\tau_0 := \mathbb{E}(Y(1) - Y(0))$ satisfies the *population moment condition*

$$\mathbb{E}[m(W; \tau_0, \eta_0)] = 0. \quad (1.26)$$

Remark 3. Note that $\hat{\tau}^{\text{AIPW}} = \frac{1}{n} \sum_{i=1}^n \psi(W_i; \hat{\eta})$ and $\tilde{\tau}^{\text{AIPW}} = \frac{1}{n} \sum_i^n \psi(W_i; \eta_0)$.

Remark 4. In the following discussion, we first establish identities for *fixed* (non-random) functions $\hat{\eta}$. When we later substitute an estimated nuisance $\hat{\eta} = \hat{\eta}(\mathcal{D}_n)$ constructed from a sample $\mathcal{D}_n := \{W_i\}_{i=1}^n$, the same identities are used in the following precise sense: let $W^* = (Y^*, D^*, X^*)$ be an independent draw from the population, $W^* \perp \mathcal{D}_n$; then all “population” expectations involving $\hat{\eta}$ should be read as conditional expectations over W^* given \mathcal{D}_n , e.g.

$$\mathbb{E}[m(W; \tau_0, \hat{\eta})] \quad \text{means} \quad \mathbb{E}[m(W^*; \tau_0, \hat{\eta}(\mathcal{D}_n)) \mid \mathcal{D}_n].$$

This distinction matters because if $\hat{\eta}$ is trained on the full sample, then $\hat{\eta}(X_i)$ may depend on (Y_i, D_i) through the fitting procedure; hence one generally cannot treat $\hat{\eta}(X_i)$ as measurable w.r.t. $\sigma(X_i)$ in *in-sample* conditional expectations. Cross-fitting introduced later restores the needed conditional independence fold-by-fold.

Fix $\tilde{\eta} := (\tilde{\mu}_0, \tilde{\mu}_1, \tilde{p})$ with $\varepsilon \leq \tilde{p}(X) \leq 1 - \varepsilon$ a.s. Define $\delta_d(X) := \tilde{\mu}_d(X) - \mu_{d,0}(X)$. Because $\tilde{\mu}_d(X)$ and $\tilde{p}(X)$ are $\sigma(X)$ -measurable (they are *fixed* functions of X), we may condition on

X and pull them out of conditional expectations.

$$\begin{aligned}\mathbb{E} \left[\frac{D}{\tilde{p}(X)} (Y - \tilde{\mu}_1(X)) \mid X \right] &= \frac{1}{\tilde{p}(X)} \mathbb{E} [D \mathbb{E} (Y - \tilde{\mu}_1(X) \mid D, X) \mid X] \\ &= \frac{1}{\tilde{p}(X)} \mathbb{E} [D (\mu_1(X) - \tilde{\mu}_1(X)) \mid X] \\ &= \frac{p_0(X)}{\tilde{p}(X)} (\mu_1(X) - \tilde{\mu}_1(X)) = -\frac{p_0(X)}{\tilde{p}(X)} \delta_1(X),\end{aligned}$$

and similarly

$$\mathbb{E} \left[\frac{1-D}{1-\tilde{p}(X)} (Y - \tilde{\mu}_0(X)) \mid X \right] = -\frac{1-p_0(X)}{1-\tilde{p}(X)} \delta_0(X).$$

Substituting these into (1.24) gives

$$\begin{aligned}\mathbb{E} [\psi(W; \tilde{\eta}) \mid X] &= \tilde{\mu}_1(X) - \tilde{\mu}_0(X) - \frac{p_0(X)}{\tilde{p}(X)} \delta_1(X) + \frac{1-p_0(X)}{1-\tilde{p}(X)} \delta_0(X) \\ &= \mu_1(X) - \mu_0(X) + \left(1 - \frac{p_0(X)}{\tilde{p}(X)}\right) \delta_1(X) + \left(\frac{1-p_0(X)}{1-\tilde{p}(X)} - 1\right) \delta_0(X).\end{aligned}$$

Simplying the above, we obtain the exact conditional representation

$$\mathbb{E} [\psi(W; \tilde{\eta}) \mid X] = \mu_{1,0}(X) - \mu_{0,0}(X) + (\tilde{p}(X) - p_0(X)) \left\{ \frac{\tilde{\mu}_1(X) - \mu_{1,0}(X)}{\tilde{p}(X)} + \frac{\tilde{\mu}_0(X) - \mu_{0,0}(X)}{1-\tilde{p}(X)} \right\}.$$

Taking expectations over X yields the population bias identity

$$\mathbb{E} [\psi(W; \tilde{\eta})] - \tau_0 = \mathbb{E} \left[(\tilde{p}(X) - p_0(X)) \left\{ \frac{\tilde{\mu}_1(X) - \mu_{1,0}(X)}{\tilde{p}(X)} + \frac{\tilde{\mu}_0(X) - \mu_{0,0}(X)}{1-\tilde{p}(X)} \right\} \right]. \quad (1.27)$$

Then,

$$\mathbb{E} [m(W; \tau_0, \tilde{\eta})] = \mathbb{E} \left[(\tilde{p}(X) - p_0(X)) \left\{ \frac{\delta_1(X)}{\tilde{p}(X)} + \frac{\delta_0(X)}{1-\tilde{p}(X)} \right\} \right]. \quad (1.28)$$

Equation (1.28) immediately yields:

$$\mathbb{E} [m(W; \tau_0, \tilde{\eta})] = 0 \quad \text{if either } \tilde{\mu}_d = \mu_{d,0} \text{ a.s. for } d = 0, 1, \quad \text{or} \quad \tilde{p} = p_0 \text{ a.s.}$$

This is the population version of *double robustness*. Moreover, (1.28) shows that, in general, the population bias is a *product of nuisance errors*.

Remark 5. Even if $\hat{\eta} = \hat{\eta}(\mathcal{D}_n)$ is estimated using the full sample (no sample splitting), the algebra behind (1.28) remains valid *conditionally* for an independent draw W^* . Specifically,

with $W^* \perp \mathcal{D}_n$,

$$\mathbb{E}[m(W^*; \tau_0, \hat{\eta}) \mid \mathcal{D}_n] = \mathbb{E}\left[\left(\hat{p}(X^*) - p_0(X^*)\right) \left\{ \frac{\hat{\mu}_1(X^*) - \mu_{1,0}(X^*)}{\hat{p}(X^*)} + \frac{\hat{\mu}_0(X^*) - \mu_{0,0}(X^*)}{1 - \hat{p}(X^*)} \right\} \mid \mathcal{D}_n\right]. \quad (1.29)$$

Equation (1.29) is the precise sense in which we apply the “exact bias decomposition” when $\tilde{\eta} = \hat{\eta}$. However, (1.29) does *not* justify in-sample manipulations such as $\mathbb{E}[D_i/\hat{p}(X_i) \mid X_i] = p_0(X_i)/\hat{p}(X_i)$, because $\hat{p}(X_i)$ need not be $\sigma(X_i)$ -measurable when trained on the same sample.

Neyman orthogonality. For a perturbation direction $h := (h_0, h_1, h_p)$, define a path $\eta_t := \eta_0 + th$. We say the moment $m(W; \tau, \eta)$ is *Neyman orthogonal at (τ_0, η_0)* if

$$\frac{d}{dt}\mathbb{E}[m(W; \tau_0, \eta_t)]\Big|_{t=0} = 0 \quad \text{for all admissible directions } h. \quad (1.30)$$

Equivalently, the (Gateaux) derivative of $\eta \mapsto \mathbb{E}[m(W; \tau_0, \eta)]$ vanishes at η_0 ³.

Step 3: compute derivatives for AIPW. It is convenient to rewrite ψ as

$$\psi(W; \eta) = \frac{DY}{p(X)} - \frac{(1-D)Y}{1-p(X)} + \mu_1(X) \left(1 - \frac{D}{p(X)}\right) - \mu_0(X) \left(1 - \frac{1-D}{1-p(X)}\right). \quad (1.31)$$

From (1.25) and (1.31), note that $\partial_\eta \mathbb{E}[m] = \partial_\eta \mathbb{E}[\psi]$ since τ_0 is fixed.

Derivative w.r.t. μ_1 . Let $\mu_{1,t} := \mu_{1,0} + th_1$ and hold $(\mu_0, p) = (\mu_{0,0}, p_0)$ fixed. Then by (1.31),

$$\psi(W; \mu_{0,0}, \mu_{1,t}, p_0) - \psi(W; \eta_0) = t h_1(X) \left(1 - \frac{D}{p_0(X)}\right).$$

Hence

$$\begin{aligned} \frac{d}{dt}\mathbb{E}[\psi(W; \mu_{0,0}, \mu_{1,t}, p_0)]\Big|_{t=0} &= \mathbb{E}\left[h_1(X) \left(1 - \frac{D}{p_0(X)}\right)\right] \\ &= \mathbb{E}\left[h_1(X) \mathbb{E}\left(1 - \frac{D}{p_0(X)} \mid X\right)\right] \\ &= \mathbb{E}\left[h_1(X) \left(1 - \frac{\mathbb{E}(D \mid X)}{p_0(X)}\right)\right] = 0. \end{aligned} \quad (1.32)$$

³Refer to [Luenberger \(1997, Section 7.2\)](#) for formal definitions, and other chapters for textbook treatment of foundations of functional spaces.

Derivative w.r.t. μ_0 . Similarly, letting $\mu_{0,t} := \mu_{0,0} + th_0$ and holding $(\mu_1, p) = (\mu_{1,0}, p_0)$ fixed,

$$\psi(W; \mu_{0,t}, \mu_{1,0}, p_0) - \psi(W; \eta_0) = -t h_0(X) \left(1 - \frac{1-D}{1-p_0(X)} \right),$$

so

$$\begin{aligned} \frac{d}{dt} \mathbb{E}[\psi(W; \mu_{0,t}, \mu_{1,0}, p_0)] \Big|_{t=0} &= -\mathbb{E} \left[h_0(X) \left(1 - \frac{1-D}{1-p_0(X)} \right) \right] \\ &= -\mathbb{E} \left[h_0(X) \mathbb{E} \left(1 - \frac{1-D}{1-p_0(X)} \middle| X \right) \right] = 0. \end{aligned} \quad (1.33)$$

Derivative w.r.t. p . Now let $p_t := p_0 + th_p$ and hold $(\mu_0, \mu_1) = (\mu_{0,0}, \mu_{1,0})$ fixed. Using the original form of ψ ,

$$\psi(W; \mu_{0,0}, \mu_{1,0}, p_t) = \mu_{1,0}(X) - \mu_{0,0}(X) + \frac{D}{p_t(X)} (Y - \mu_{1,0}(X)) - \frac{1-D}{1-p_t(X)} (Y - \mu_{0,0}(X)),$$

and differentiating at $t = 0$ yields

$$\frac{d}{dt} \psi(W; \mu_{0,0}, \mu_{1,0}, p_t) \Big|_{t=0} = -\frac{D h_p(X)}{p_0(X)^2} (Y - \mu_{1,0}(X)) - \frac{(1-D) h_p(X)}{(1-p_0(X))^2} (Y - \mu_{0,0}(X)).$$

Taking expectations and conditioning on X gives

$$\begin{aligned} \frac{d}{dt} \mathbb{E}[\psi(W; \mu_{0,0}, \mu_{1,0}, p_t)] \Big|_{t=0} &= -\mathbb{E} \left[\frac{h_p(X)}{p_0(X)^2} \mathbb{E}(D(Y - \mu_{1,0}(X)) \mid X) \right] \\ &\quad - \mathbb{E} \left[\frac{h_p(X)}{(1-p_0(X))^2} \mathbb{E}((1-D)(Y - \mu_{0,0}(X)) \mid X) \right]. \end{aligned} \quad (1.34)$$

Now, by the law of iterated expectations,

$$\mathbb{E}(D(Y - \mu_{1,0}(X)) \mid X) = \mathbb{E}(D \mathbb{E}(Y - \mu_{1,0}(X) \mid D, X) \mid X) = \mathbb{E}(D \cdot 0 \mid X) = 0,$$

and similarly $\mathbb{E}((1-D)(Y - \mu_{0,0}(X)) \mid X) = 0$. Substituting into (1.34) yields

$$\frac{d}{dt} \mathbb{E}[\psi(W; \mu_{0,0}, \mu_{1,0}, p_t)] \Big|_{t=0} = 0. \quad (1.35)$$

Combining (1.32), (1.33), and (1.35), we obtain

$$\frac{d}{dt} \mathbb{E}[m(W; \tau_0, \eta_0 + th)] \Big|_{t=0} = \frac{d}{dt} \mathbb{E}[\psi(W; \eta_0 + th)] \Big|_{t=0} = 0 \quad \text{for all } h,$$

which verifies (1.30). Hence, the AIPW moment is Neyman orthogonal at (τ_0, η_0) .

The Taylor expansion around η_0 has the form

$$\mathbb{E}[m(W; \tau_0, \hat{\eta})] = \mathbb{E}[m(W; \tau_0, \eta_0)] + \partial_\eta \mathbb{E}[m(W; \tau_0, \eta)] \Big|_{\eta=\eta_0} [\hat{\eta} - \eta_0] + R_2(\hat{\eta}, \eta_0), \quad (1.36)$$

where $R_2(\hat{\eta}, \eta_0)$ is a second-order remainder. By (1.26), the first term is zero. By Neyman orthogonality (1.30), the linear term is also zero. Thus,

$$\mathbb{E}[m(W; \tau_0, \hat{\eta})] = R_2(\hat{\eta}, \eta_0), \quad (1.37)$$

so any population bias from plugging in $\hat{\eta}$ is *second order*.⁴ Indeed, the exact bias identity applied conditionally (cf. (1.29)) yields

$$\mathbb{E}[m(W; \tau_0, \hat{\eta})] = \mathbb{E}\left[\left(\hat{p}(X) - p_0(X)\right) \left\{ \frac{\hat{\mu}_1(X) - \mu_{1,0}(X)}{\hat{p}(X)} + \frac{\hat{\mu}_0(X) - \mu_{0,0}(X)}{1 - \hat{p}(X)} \right\}\right], \quad (1.38)$$

again with the convention that this is a conditional population expectation when $\hat{\eta}$ is random. In particular, if $\hat{p}(X) \in [\varepsilon, 1 - \varepsilon]$ a.s., then

$$|\mathbb{E}[m(W; \tau_0, \hat{\eta})]| \leq \frac{1}{\varepsilon} \mathbb{E}[|\hat{p}(X) - p_0(X)| |\hat{\mu}_1(X) - \mu_{1,0}(X) + \hat{\mu}_0(X) - \mu_{0,0}(X)|], \quad (1.39)$$

and by Cauchy–Schwarz,

$$|\mathbb{E}[m(W; \tau_0, \hat{\eta})]| \leq \frac{1}{\varepsilon} \|\hat{p} - p_0\|_{L^2(P_X)} (\|\hat{\mu}_1 - \mu_{1,0}\|_{L^2(P_X)} + \|\hat{\mu}_0 - \mu_{0,0}\|_{L^2(P_X)}), \quad (1.40)$$

where $\|f\|_{L^2(P_X)} := (\mathbb{E}[f(X)^2])^{1/2}$ denotes the $L^2(P_X)$ norm.

The AIPW estimator can be written as $\hat{\tau}^{\text{AIPW}} = \mathbb{P}_n[\psi(W; \hat{\eta})]$. A basic decomposition is

$$\hat{\tau}^{\text{AIPW}} - \tau_0 = (\mathbb{P}_n - \mathbb{P})\psi(W; \eta_0) + \underbrace{\mathbb{P}[\psi(W; \hat{\eta}) - \psi(W; \eta_0)]}_{\text{population bias}} + \underbrace{(\mathbb{P}_n - \mathbb{P})(\psi(W; \hat{\eta}) - \psi(W; \eta_0))}_{\text{overfitting}}. \quad (1.41)$$

The first term is the usual sampling fluctuation at the truth. By (1.38), orthogonality (and the exact bias identity) control the *second term*: it is second order and vanishes if either nuisance is correct (weak double robustness in terms of consistency), and it is $o_p(n^{-1/2})$ under a product-rate condition (for root- n inference).

The subtlety is the *third term* in (1.41). Even if the population Neyman condition holds,

⁴If $\hat{\eta}$ is data-dependent, read $\mathbb{E}[m(W; \tau_0, \hat{\eta})]$ as $\mathbb{E}[m(W^*; \tau_0, \hat{\eta}) \mid \mathcal{D}_n]$

this third term can fail to be small if $\hat{\eta}$ is estimated on the same sample used in \mathbb{P}_n . The reason is that $\psi(\cdot; \hat{\eta})$ is then a *data-dependent* function evaluated on the same data, so controlling $(\mathbb{P}_n - \mathbb{P})f_n$ with $f_n = \psi(\cdot; \hat{\eta}) - \psi(\cdot; \eta_0)$ requires strong uniform complexity conditions (e.g. Donsker/entropy conditions) that typically fail for modern ML.

Sample Splitting and Cross-fitting. To deal with the third term in (1.41), we rely on sample splitting and cross-fitting to restore conditional i.i.d. structure and make the third term negligible under mild conditions. Partition $\{1, \dots, n\}$ into K folds I_1, \dots, I_K . For each fold k , estimate nuisances $\hat{\eta}^{(-k)}$ using only the training sample I_k^c , and evaluate on the held-out fold I_k . Define the cross-fitted AIPW estimator

$$\hat{\tau}^{\text{cf}} := \frac{1}{n} \sum_{k=1}^K \sum_{i \in I_k} \psi(W_i; \hat{\eta}^{(-k)}). \quad (1.42)$$

Let $\mathbb{P}_{n,k} := |I_k|^{-1} \sum_{i \in I_k}$ be the fold-specific empirical measure and let $\Delta_k(W) := \psi(W; \hat{\eta}^{(-k)}) - \psi(W; \eta_0)$. Then

$$\hat{\tau}^{\text{cf}} - \tau_0 = \sum_{k=1}^K \frac{|I_k|}{n} (\mathbb{P}_{n,k} - \mathbb{P}) \psi(W; \eta_0) + \sum_{k=1}^K \frac{|I_k|}{n} \mathbb{P}[\Delta_k(W)] + \sum_{k=1}^K \frac{|I_k|}{n} (\mathbb{P}_{n,k} - \mathbb{P}) \Delta_k(W). \quad (1.43)$$

The middle term is the population bias term and is handled by the same orthogonality/product-of-errors argument as before. The crucial point is the last term: conditional on the training data I_k^c , the function $\Delta_k(\cdot)$ is fixed, and the held-out observations $\{W_i : i \in I_k\}$ are i.i.d. and independent of $\hat{\eta}^{(-k)}$. Therefore, conditional on I_k^c ,

$$\begin{aligned} \mathbb{E}[(\mathbb{P}_{n,k} - \mathbb{P}) \Delta_k(W) \mid I_k^c] &= 0, \\ \text{Var}((\mathbb{P}_{n,k} - \mathbb{P}) \Delta_k(W) \mid I_k^c) &= \frac{\text{Var}(\Delta_k(W) \mid I_k^c)}{|I_k|}. \end{aligned}$$

In particular, using $\text{Var}(Z) \leq \mathbb{E}[Z^2]$,

$$\mathbb{E}[(\mathbb{P}_{n,k} - \mathbb{P}) \Delta_k(W)^2 \mid I_k^c] \leq \frac{1}{|I_k|} \mathbb{E}[\Delta_k(W)^2 \mid I_k^c]. \quad (1.44)$$

Thus,

$$(\mathbb{P}_{n,k} - \mathbb{P}) \Delta_k(W) = O_p\left(\frac{\|\Delta_k\|_{L^2(P)}}{\sqrt{|I_k|}}\right). \quad (1.45)$$

Consequently, if $\|\Delta_k\|_{L^2(P)} = o_p(1)$ (a mild stability/consistency requirement), then the last

term in (1.43) is $o_p(n^{-1/2})$ when $|I_k| \asymp n$. This is the key technical role of cross-fitting: it turns the difficult empirical-process term into an ordinary conditional LLN/CLT problem.

Quick summary of the key ingredients.

- *Neyman orthogonality.* The Neyman orthogonality condition (1.30) ensures that the population moment has no first-order sensitivity to nuisance errors, so the population bias is second order. For AIPW, this second-order term is exactly the product-of-errors expression (1.38), interpreted conditionally when $\hat{\eta}$ is random (cf. Step 0 and (1.29)).
- *Cross-fitting.* Orthogonality alone does *not* control the empirical-process/overfitting term in (1.41) when $\hat{\eta}$ is estimated on the same sample. Cross-fitting enforces foldwise independence between the fitted nuisances and the evaluation observations, yielding the variance bound (1.44) and the rate control (1.45) without imposing strong uniform complexity (Donsker) conditions on the first-stage estimators.

Bibliography

Abadie, A. and M. D. Cattaneo (2018). Econometric methods for program evaluation. *Annual Review of Economics* 10(1), 465–503.

Bickel, P. J., C. A. Klaassen, P. J. Bickel, Y. Ritov, J. Klaassen, J. A. Wellner, and Y. Ritov (1993). *Efficient and adaptive estimation for semiparametric models*, Volume 4. Springer.

Chamberlain, G. (1992). Efficiency bounds for semiparametric regression. *Econometrica: Journal of the Econometric Society*, 567–596.

Ding, P. (2024). *A first course in causal inference*. Chapman and Hall/CRC.

Hirano, K., G. W. Imbens, and G. Ridder (2003). Efficient estimation of average treatment effects using the estimated propensity score. *Econometrica* 71(4), 1161–1189.

Horvitz, D. G. and D. J. Thompson (1952). A generalization of sampling without replacement from a finite universe. *Journal of the American statistical Association* 47(260), 663–685.

Leung, M. P. (2022). Causal inference under approximate neighborhood interference. *Econometrica* 90(1), 267–293.

Luenberger, D. G. (1997). *Optimization by vector space methods*. John Wiley & Sons.

Robins, J. M., A. Rotnitzky, and L. P. Zhao (1994). Estimation of regression coefficients when some regressors are not always observed. *Journal of the American statistical Association* 89(427), 846–866.

Rosenbaum, P. R. and D. B. Rubin (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika* 70(1), 41–55.

Wager, S. (2025). *Causal inference: A statistical learning approach*.