# Lec 9: Estimation of Average Treatment Effect

Yanjun Han

Nov 14, 2023

# Potential outcome model.

For a binary treatment $W \in \{0,1\}$, an individual $i$ has two potential outcomes $Y_i(1)$ and $Y_i(0)$ ← the outcome individual $i$ would have experienced had he/she received the treatment or not, respectively

---

**Average Treatment Effect (ATE):**
$$\tau = \mathbb{E}[Y_i(1) - Y_i(0)]$$

---

A typical dataset: $\{(X_i, W_i, Y_i)\}_{i=1}^{n}$:
- $W_i \in \{0,1\}$: indicator of treatment / control
- $Y_i \in \mathbb{R}$: observed outcome $Y_i = Y_i(W_i)$
- $X_i \in \mathbb{R}^p$ (optional): feature of individual $i$

(Optional material: SUTVA — stable unit treatment value assumption
"the potential outcomes for any unit do not vary with the treatments assigned to each other unit, and, for each unit, there are no different forms or versions of each treatment level, which lead to different potential outcomes", e.g.
   — you taking the aspirin cannot have an affect on my headache
   — different aspirins should have the same strength                )

# Randomised control trials (RCT)  (no $X_i$)

Assumption: $\begin{cases} W_i \perp\!\!\!\perp (Y_i(0), Y_i(1)) & \text{(random treatment assignment)} \\ \text{each } i \text{ has the same marginal prob. of getting treated} \end{cases}$

Difference-in-mean estimation.

$$\hat{\tau}_{DM} = \frac{1}{n_1} \sum_{W_i=1} Y_i - \frac{1}{n_0} \sum_{W_i=0} Y_i \;,\; \text{where} \; n_j = \#\{i : W_i = j\}$$

Unbiasedness of $\hat{\tau}_{DM}$:

$$\mathbb{E}\left[\frac{1}{n_1}\sum_{W_i=1} Y_i\right] = \mathbb{E}\left[\frac{1}{n_1}\sum_{i=1}^{\hat{n}} W_i Y_i\right]$$

$$= \mathbb{E}\left[\frac{1}{n_1}\sum_{i=1}^{\hat{n}} W_i Y_i(1)\right] \quad (\text{SUTVA})$$

$$= \mathbb{E}\left[\mathbb{E}\left[\frac{1}{n_1}\sum_{i=1}^{n} W_i Y_i(1) \;\Big|\; \{Y_i(0), Y_i(1)\}_{i=1}^{\hat{n}}, n\right]\right]$$

$$= \mathbb{E}\left[\frac{1}{n_1}\sum_{i=1}^{\hat{n}} Y_i(1)\cdot \mathbb{E}[W_i \mid n_i]\right] \quad \begin{array}{l}(\text{random}\\ \text{treatment}\\ \text{assignment})\end{array}$$

$$= \mathbb{E}\left[\frac{1}{n_1}\sum_{i=1}^{n} Y_i(1)\cdot \frac{\hat{n}}{n}\right] \quad \begin{array}{l}(\text{same}\\ \text{marginal prob.})\end{array}$$

$$= \mathbb{E}[Y_i(1)]$$

$$\Rightarrow \mathbb{E}[\hat{\tau}_{DM}] = \mathbb{E}[Y_i(1)] - \mathbb{E}[Y_i(0)] = \tau.$$

## Propensity score.

Question. What happens if we combine two RCTs, but with different treatment probabilities?

Failure of $\hat{\tau}_{DM}$: Simpson's Paradox

Example: discourage teenagers from smoking
in Palo Alto, CA (5%) & NYC (20%)

| Palo Alto | Non-S. | Smoker |
|---|---|---|
| Treat. | 95 | 5 |
| Control | 1700 | 200 |

$+$

| NYC | Non-S. | Smoker |
|---|---|---|
| Treat. | 255 | 145 |
| Control | 800 | 800 |

$=$

| All | Non-S. | Smoker |
|---|---|---|
| Treat. | 350 | 150 |
| Control | 2500 | 1000 |

19:1  vs.  8.5:1

treatment effect: +

1.76:1 vs. 1:1

treatment effect: +

2.33:1 vs. 2.5:1

treatment effect: − (!!)

Implication: propensity score plays a central role !

Propensity score:   $e(x) = \mathbb{P}(W_i = 1 \mid X_i = x)$

Assumptions: 1. unconfoundedness:  $(Y_i(0), Y_i(1)) \perp\!\!\!\perp W_i \mid X_i$

(no unexplained feature affects both $W_i$ & $(Y_i(0), Y_i(1))$)

2. overlap:   $\eta \le e(x) \le 1-\eta$  for all $x$.

Inverse-propensity weighting (IPW).

Theorem.  $\mathbb{E}\left[ \underbrace{\dfrac{WY}{e(x)} - \dfrac{(1-W)Y}{1-e(x)} - \tau}_{f_{\tau, e}(W, X, Y): \text{ estimating function}} \right] = 0$

Pf.   $\mathbb{E}\left[\dfrac{WY}{e(x)}\right] = \mathbb{E}\left[\dfrac{WY(1)}{e(x)}\right]$   (SUTVA)

$= \mathbb{E}\left\{ \mathbb{E}\left[\dfrac{WY(1)}{e(x)} \mid x\right] \right\}$

$= \mathbb{E}\left\{ \dfrac{1}{e(x)} \mathbb{E}[W \mid x]\, \mathbb{E}[Y(1) \mid x] \right\}$  (unconfoundedness)

$= \mathbb{E}\left\{ \mathbb{E}[Y(1) \mid x] \right\}$   ($e(x) = \mathbb{P}(W=1 \mid x)$)

$= \mathbb{E}[Y(1)]$ .  ☐

IPW estimator: given an estimate $\hat{e}(x)$ for $e(x)$, then

$$\frac{1}{n} \sum_{i=1}^{n} \left( \frac{W_i Y_i}{\hat{e}(x_i)} - \frac{(1-W_i) Y_i}{1-\hat{e}(x_i)} - \hat{\tau}_{IPW} \right) = 0$$

$$\implies \hat{\tau}_{IPW} = \frac{1}{n} \sum_{i=1}^{n} \left( \frac{W_i Y_i}{\hat{e}(x_i)} - \frac{(1-W_i) Y_i}{1-\hat{e}(x_i)} \right)$$

Pros:      consistent ($\hat{\tau}_{IPW} \to \tau$ as $n \to \infty$)

Cons:      potentially large variance ;
              not robust to nuisance estimation error $\hat{e}(x) - e(x)$

## Double robust estimation: Augmented IPW (AIPW).

__Model__. $\begin{cases} Y = \mu_W(x) + \varepsilon_W & , \quad \mathbb{E}[\varepsilon_0 | W, x] = 0, \; \mathbb{E}[\varepsilon_1 | W, x] = 0. \\ W \sim \text{Bern}(e(x)) \end{cases}$

Target parameter: $\tau = \mathbb{E}[\mu_1(x) - \mu_0(x)]$

Nuisance parameter: mean outcomes $\mu_0(x), \mu_1(x)$
                               propensity score $e(x)$

__AIPW estimator__. Given nuisance estimates $(\hat{\mu}_1(x), \hat{\mu}_0(x), \hat{e}(x))$:

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

Interpretation: 1. from IPW, subtract the mean outcomes
                       $(\hat{\mu}_0(x_i), \hat{\mu}_1(x_i))$ from $Y_i$ ;
                2. from $\frac{1}{n} \sum_{i=1}^{n} (\hat{\mu}_1(x_i) - \hat{\mu}_0(x_i))$, debias using
                    IPW applied to the regression residuals.

## Double machine learning in practice.

1. Split the dataset into $K$ folds;
2. For $k = 1, \dots, K$, use all data but the $k$-th fold to estimate $(\hat{\mu}_1^{(-k)}(x), \hat{\mu}_0^{(-k)}(x), \hat{e}^{(-k)}(x))$, possibly via machine learning;
3. Estimate ATE by   ← $i$ belongs to $k_i$-th fold

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

## Theoretical properties.

$$f_{(\mu_1, \mu_0, e, \tau)}(W, X, Y) = \mu_1(x) - \mu_0(x) + W \frac{Y - \mu_1(x)}{e(x)} - (1-W) \frac{Y - \mu_0(x)}{1 - e(x)} - \tau$$

Claim 1: $f$ is an estimating function, i.e.
$$\mathbb{E}[f_{(\mu_1, \mu_0, e, \tau)}(W, X, Y)] = 0.$$

Pf. $\mathbb{E}\left[ W \frac{Y - \mu_1(x)}{e(x)} \right]$

$$= \mathbb{E}\left[ W \frac{Y(1) - \mu_1(x)}{e(x)} \right] \quad \text{(SUTVA)}$$

$$= \mathbb{E}\left[ \frac{W \varepsilon_1}{e(x)} \right]$$

$$= \mathbb{E}\left\{ \mathbb{E}\left[ \frac{W \varepsilon_1}{e(x)} \mid W, x \right] \right\}$$

$$= 0 \quad \left( \mathbb{E}[\varepsilon_1 \mid W, X] = 0, \text{ or unconfoundedness} \right)$$

$$\implies \mathbb{E}[f] = \mathbb{E}[\mu_1(x) - \mu_0(x)] - \tau$$
$$= \mathbb{E}[\varepsilon_0 - \varepsilon_1] = 0. \qquad \square$$

**Claim 2:** $f$ is Neyman orthogonal, i.e.

$$\mathbb{E}\left[\nabla_g f_{(\mu_1, \mu_0, e, \tau)}(W, X, Y)\right] = 0, \quad \forall\, g \in \{\mu_0, \mu_1, e\}.$$

(Implication: nuisance estimation errors only have second-order effects on the estimation of $\tau$)

Pf. (1) $g = \mu_1$:
$$\mathbb{E}[\nabla_{\mu_1} f] = \mathbb{E}\left[1 - \frac{W}{e(X)} \mid X\right]$$
$$= 1 - \mathbb{E}\left[\frac{W}{e(X)} \mid X\right]$$
$$= 0 \quad (\mathbb{P}(W = 1 \mid X) = e(X))$$

(2) $g = \mu_0$:
$$\mathbb{E}[\nabla_{\mu_0} f] = \mathbb{E}\left[-1 + \frac{1-W}{1-e(X)} \mid X\right]$$
$$= -1 + \mathbb{E}\left[\frac{1-W}{1-e(X)} \mid X\right]$$
$$= 0 \quad (\mathbb{P}(W = 0 \mid X) = 1 - e(X))$$

(3) $g = e$:
$$\mathbb{E}[\nabla_e f] = \mathbb{E}\left[-\frac{W(Y - \mu_1(X))}{e(X)^2} + \frac{(1-W)(Y - \mu_0(X))}{(1-e(X))^2} \mid X\right]$$
$$= \mathbb{E}\left[-\frac{W\varepsilon_1}{e(X)^2} + \frac{(1-W)\varepsilon_0}{(1-e(X))^2} \mid X\right] \text{ (SUTVA)}$$
$$= \mathbb{E}\left\{\mathbb{E}\left[-\frac{W\varepsilon_1}{e(X)^2} + \frac{(1-W)\varepsilon_0}{(1-e(X))^2} \mid W, X\right] \mid X\right\}$$
$$= 0 \quad (\mathbb{E}[\varepsilon_1, \varepsilon_0 \mid W, X] = 0, \text{ or}$$
$$\text{unconfoundedness}) \quad \square$$

**Claim 3:** $f$ is (weakly) double robust, i.e.

$$\mathbb{E}\left[f_{(\hat{\mu}_1, \hat{\mu}_0, \hat{e}, \tau)}(W, X, Y)\right] = 0 \quad \text{if} \quad \begin{array}{c}(\hat{\mu}_1, \hat{\mu}_0) = (\mu_1, \mu_0)\\ \underline{OR} \quad \hat{e} = e.\end{array}$$

(Implication: AIPW is consistent if either $(\hat{\mu}_1(X), \hat{\mu}_0(X))$ are consistent, or $\hat{e}(X)$ is consistent)

Pf. (1) If $(\hat\mu_1, \hat\mu_0) = (\mu_1, \mu_0)$ : same argument in Claim 1

(2) If $\hat{e} = e$, rewrite

$\overbrace{\phantom{xxxxxxxxxxxxx}}$ $\mathbb{E}[\cdot] = 0$ by IPW analysis

$$f_{(\hat\mu_1, \hat\mu_0, e, \tau)}(W, X, Y) = \frac{WY}{e(x)} - \frac{(1-W)Y}{1-e(x)} - \tau$$

$$- (W - e(x))\left(\frac{\hat\mu_1(x)}{e(x)} - \frac{\hat\mu_0(x)}{1-e(x)}\right)$$

$\underbrace{\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxxxxx}}$

$$\mathbb{E}[\cdot] = \mathbb{E}\{\mathbb{E}[\cdot | X]\} = 0$$

since $\mathbb{P}(W = 1 | X) = e(x)$.

$\square$

## Derivation of AIPW (Optional)

<u>First derivation</u> : use efficient influence

(see J. Hahn, "On the role of propensity score in efficient semiparametric estimation of average treatment effects", Econometrica, 1998)

<u>Second derivation</u> : find the projection of IPW

$$f_{\tau e}(W, X, Y) = \frac{WY}{e(x)} - \frac{(1-W)Y}{1-e(x)} - \tau$$

to the orthogonal complement of $L$, where

$$L = \{ g(W, X, Y) : \mathbb{E}[g | X, Y(0), Y(1)] = 0 \}.$$

<u>Lemma 1</u>   $L = \{ (W - e(x)) h(x)$  for general $h \}$

Pf. Obviously $\mathbb{E}[(W - e(x)) h(x) | X, Y(0), Y(1)] = h(x) \mathbb{E}[W - e(x) | X] = 0.$

Now we show that any $g(W, X, Y) \in L$ must take this form.

$$\mathbb{E}[g | X, Y(0), Y(1)] = e(x) \underbrace{g(1, X, Y(1))}_{\equiv g_1(x)} + (1-e(x)) \underbrace{g(0, X, Y(0))}_{\equiv g_0(x)} \equiv 0$$

$$\Rightarrow \frac{g_1(x)}{1-e(x)} = -\frac{g_0(x)}{e(x)} =: h(x)$$

$$\Rightarrow g(W,X,Y) = \begin{cases} g_1(x) & \text{if } W=1 \\ g_0(x) & \text{if } W=0 \end{cases} = (W-e(x))h(x) \qquad \square.$$

<u>Lemma 2.</u> $\text{Proj}_{\perp}(f_{\tau,e}(W,X,Y)) = f_{(\mu_0,\mu_1,e,\tau)}(W,X,Y),$
the estimating function of AIPW.

Pf. Aim to find $h_o(x)$ s.t,

$$\mathbb{E}\left[\left(\frac{WY}{e(x)} - \frac{(1-W)Y}{1-e(x)} - \tau - (W-e(x))h_o(x)\right) \times (W-e(x))h(x)\right] = 0$$
$$\forall h$$

$$\Rightarrow 0 = \mathbb{E}\left[\left(\frac{WY}{e(x)} - \frac{(1-W)Y}{1-e(x)} - \tau - (W-e(x))h_o(x)\right)(W-e(x)) \mid x\right]$$

$$= \mu_1(x)(1-e(x)) - \mu_0(x)(-e(x)) - e(x)(1-e(x))h_o(x)$$

$$\Rightarrow h_o(x) = \frac{\mu_1(x)}{e(x)} + \frac{\mu_0(x)}{1-e(x)}$$

Therefore, $\text{Proj}_{\perp}(f_{\tau,e}(W,X,Y))$

$$= \frac{WY}{e(x)} - \frac{(1-W)Y}{1-e(x)} - \tau - (W-e(x))\left(\frac{\mu_1(x)}{e(x)} + \frac{\mu_0(x)}{1-e(x)}\right)$$

$$= \mu_1(x) - \mu_0(x) - \tau + W\frac{Y-\mu_1(x)}{e(x)} - (1-W)\frac{Y-\mu_0(x)}{1-e(x)}$$

$$= f_{(\mu_0,\mu_1,e,\tau)}(W,X,Y) \qquad \square$$