

Using Double Machine Learning to Rank Treatments

Apoorva Lal, Winston Chou, and Jordan Schafer

Netflix



Problem Setup

We observe an outcome and set of treatments $(Y_i, \mathbf{W}_i)_{i=1}^N$, and seek to rank these treatments by their treatment effects τ . For example, \mathbf{W}_i may consist of different actions taken by a user on a digital platform, and the platform seeks to promote those actions that have the largest treatment effects on a north-star metric Y_i .

A common approach to this problem is to write a Partially Linear Model of the outcome as a function of each treatment and pre-treatment covariates $Y_i = \tau W_i + g(\mathbf{X}_i) + \varepsilon_i$ and to run the residuals-on-residuals regression:

$$Y_i - \mathbb{E}[Y_i | \mathbf{X}_i] = \tilde{\tau}(W_i - \mathbb{E}[W_i | \mathbf{X}_i]) + \eta_i$$

The treatments are then ranked according to the partially ordered set $(\tilde{\tau}_{\text{plr}}, \leq)$.

We show via a simple example and numerical experiments that, when treatment effects are heterogeneous, $\tilde{\tau}_{\text{plr}}$ does not equal the Average Treatment Effect (ATE) in general, and consequently:

$$(\tilde{\tau}_{\text{plr}}, \leq) \neq (\tilde{\tau}_{\text{ATE}}, \leq).$$

Other estimators like Augmented IPW do not fall prey to this problem and should be preferred when researchers seek to rank treatments according to their ATEs, as opposed to treatment effects on implicit & potentially incomparable populations.

PLM Estimates a Weighted Average Treatment Effect

A well-known result from Angrist (1998) states that, under the linearity of the propensity score and arbitrary treatment effect heterogeneity, linear regression recovers:

$$\text{plim } \hat{\tau} = \mathbb{E}[\gamma(X)\tau(X)]$$

where $\gamma(X) = \frac{\mathbb{V}[W | X]}{\mathbb{E}[\mathbb{V}[W | X]]}$.

When the treatment W is binary, the weights simplify to

$$\gamma(X) = \frac{p(X)(1-p(X))}{\mathbb{E}[p(X)(1-p(X))]}$$

These weights are strictly positive, and largest for units with propensity scores close to 0.5.

An immediate implication is that, unless the treatments in \mathbf{W}_i are randomized according to the same mechanism, the regression weights $\gamma(X)$ will be non-uniform across treatments, and therefore treatment effects are not comparable across treatments.

See also Goldsmith-Pinkham et al. (2024) on *contamination bias*.

A Simple Example

Consider a single binary covariate x distributed uniformly in the population ($P(X = 1) = 0.5$) and binary treatments W_1, W_2 . We seek to rank W_1 and W_2 according to their ATEs. The propensity scores are as follows

	$W_1 = 0$	$W_2 = 1$
$X = 0$	0.01	0.5
$X = 1$	0.5	0.01

The treatment effects for the two treatments for the two strata are

	τ_1	τ_2
$X = 0$	-3	-2
$X = 1$	3	3
ATE	0	0.5

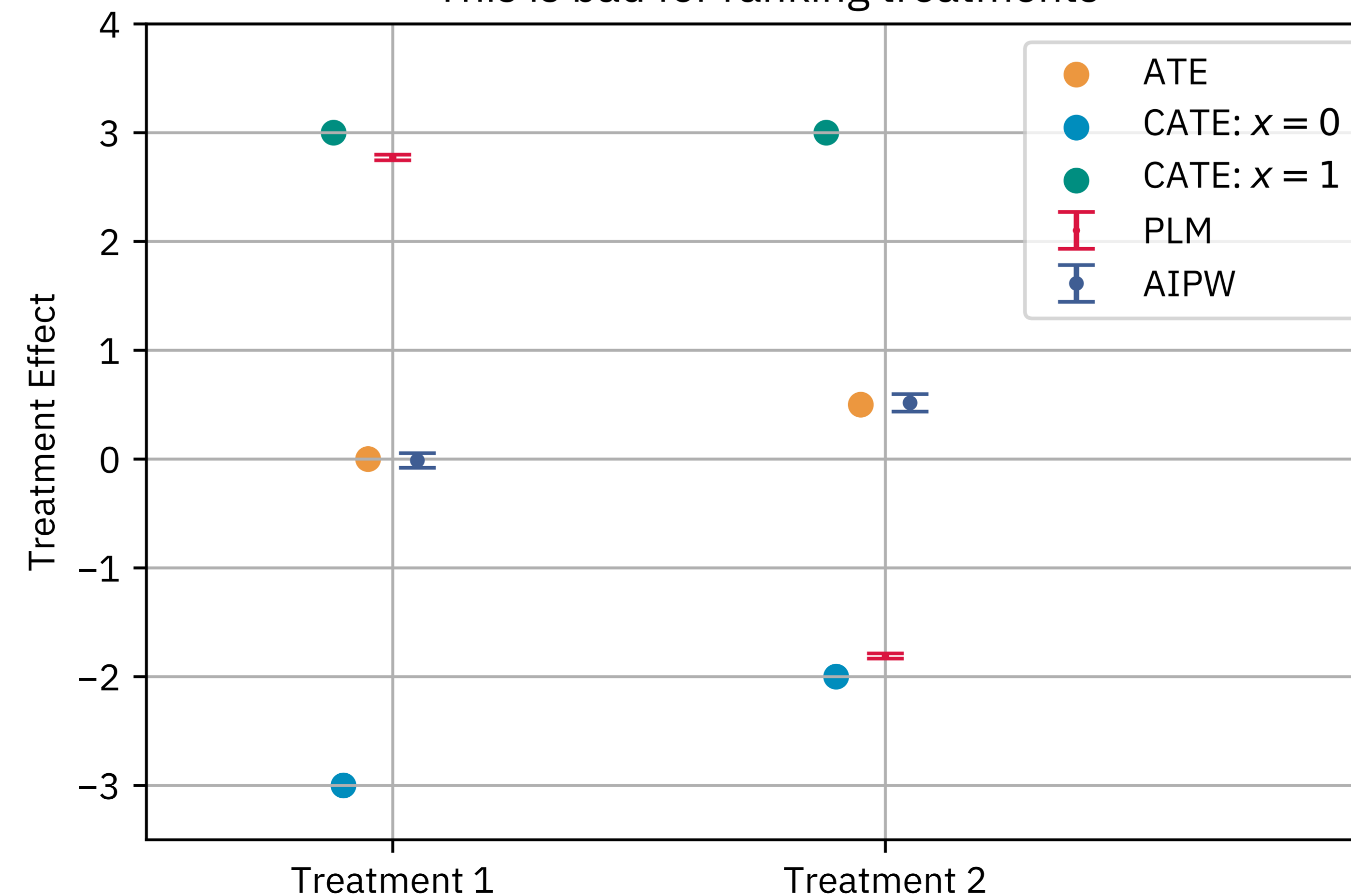
Note that W_2 has the larger ATE, but the worse treatment effect for units with propensity score $P(W_2|X = 0) = 0.5$. In contrast, W_1 has the smaller ATE, but the better treatment effect for units with propensity score $P(W_1|X = 1) = 0.5$.

The Partially Linear Model (PLM) does not recover the ATEs or the correct ranking:

$$\tilde{\tau}_1 = \frac{-3 \cdot 0.01 \cdot 0.99 + 3 \cdot 0.5 \cdot 0.5}{0.01 \cdot 0.99 + 0.5 \cdot 0.5} = 2.7714$$

$$\tilde{\tau}_2 = \frac{-2 \cdot 0.5 \cdot 0.5 + 3 \cdot 0.01 \cdot 0.99}{0.01 \cdot 0.99 + 0.5 \cdot 0.5} = -1.8095$$

PLMs do not estimate ATEs
This is bad for ranking treatments

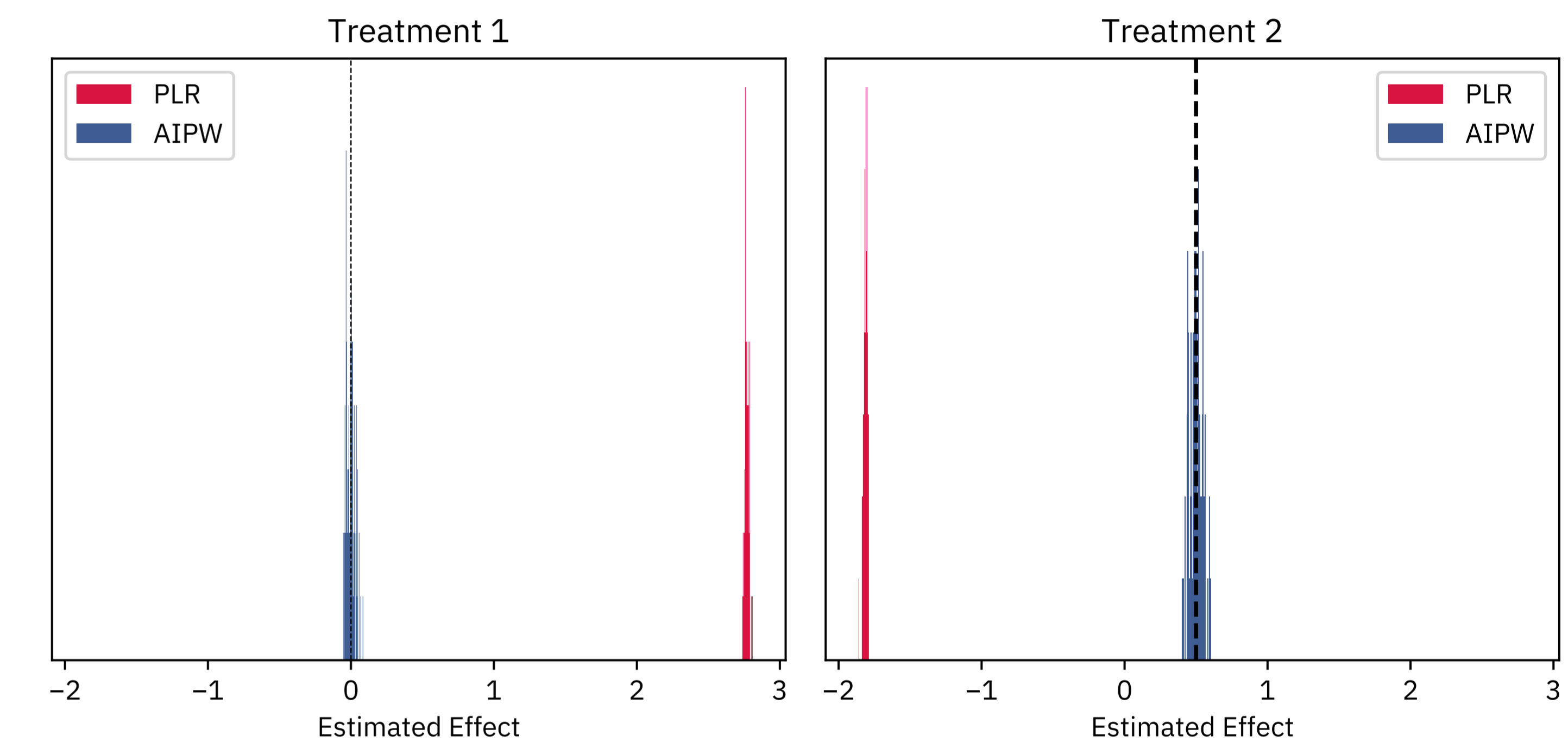


Solution: Weight Explicitly Rather Than Implicitly

Rather than implicitly weight by the propensity score, we propose to use estimators that explicitly target the ATE under heterogeneity in assignment mechanisms and treatment effects, such as the Augmented IPW estimator:

$$\hat{\tau}^{\text{AIPW}} = \left(\hat{\mu}^1(x) + \frac{W_i}{\hat{\pi}(x)}(y_i - \hat{\mu}^1(x)) \right) - \left(\hat{\mu}^0(x) + \frac{1 - W_i}{1 - \hat{\pi}(x)}(y_i - \hat{\mu}^0(x)) \right)$$

Applying AIPW to our example recovers the ATEs, and therefore the correct ranking of treatments.



References

- [1] Angrist, J. D. Estimating the labor market impact of voluntary military service using social security data on military applicants. 249–288.
- [2] Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W., and Robins, J. Double/debiased machine learning for treatment and structural parameters. C1–C68.
- [3] Goldsmith-Pinkham, P. S., Hull, P., and Kolesár, M. Contamination bias in linear regressions.
- [4] Robinson, P. M. Root-N-consistent semiparametric regression. 931–954.



Netflix Tech Blog



Netflix Jobs