

**SUPPLEMENTARY MATERIALS FOR
PROPENSITY SCORE WEIGHTING FOR CAUSAL
INFERENCE WITH MULTIPLE TREATMENTS**

BY FAN LI AND FAN LI

Yale University and Duke University

SUPPLEMENT A: ON TRANSITIVITY

With multiple treatments, a desirable property of a given class of estimands is *transitivity*. For pairwise comparisons, lack of transitivity often implies that comparisons of treatments are based on different populations. As a result, non-transitivity may lead to incompatible pairwise contrasts; for example, it is possible that treatment A is favored over treatment B, treatment B is favored over treatment C, but treatment C is found to better than treatment A at the same time. Below we provide a formal definition of transitivity and offer two related remarks.

DEFINITION 1. *The class of causal estimands $\mathbb{T}(h, \mathbb{A}) = \{\tau^h(\mathbf{a}) : \mathbf{a} \in \mathbb{A} \subset \mathbb{R}^J\}$ is transitive if the following equivariance relationship holds: $\tau^h(\mathbf{a}) + \tau^h(\mathbf{a}') = \tau^h(\mathbf{a}'')$ whenever $\mathbf{a}, \mathbf{a}', \mathbf{a}'' \in \mathbb{A}$ and $\mathbf{a} + \mathbf{a}' = \mathbf{a}''$.*

REMARK 1. *Fixing a tilting function h , the class of estimands specifying all pairwise contrasts, namely, $\mathbb{T}(h, \mathbb{S})$ is transitive. For example, with $h(\mathbf{X}) = 1$, the class of pairwise ATE estimands is transitive; with $h(\mathbf{X}) = e_j(\mathbf{X}) \prod_{k=1}^J E_k(\mathbf{X})$, the class of ATT estimands in [Lopez and Gutman \(2017\)](#) is also transitive.*

REMARK 2. *The union of $\mathbb{T}(h_1, \mathbb{S})$ and $\mathbb{T}(h_2, \mathbb{S})$ or that of their subsets is generally non-transitive for $h_1 \neq h_2$. This explains why several existing classes of estimands are non-transitive, including the class of ATT estimands of [Lechner \(2001\)](#), $\{\mathbb{E}[Y_i(j) - Y_i(j') | Z_i = j \text{ or } Z_i = j'] : j < j'\}$. The reason is that each individual estimand corresponds to a distinct tilting function $h_{j,j'}(\mathbf{X}) = (e_j(\mathbf{X}) + e_{j'}(\mathbf{X})) / e_1(\mathbf{X})$, and therefore this class of estimands is the union of $\binom{J}{2}$ elements, each of which is contained in $\mathbb{T}(h_{j,j'}, \mathbb{S})$ for some $j < j'$.*

SUPPLEMENT B: PROOF OF PROPOSITIONS

For proving the Propositions, we assume regularity conditions on $m_j(\mathbf{X}) = \mathbb{E}[Y(j) | \mathbf{X}]$ and $v_j(\mathbf{X}) = \mathbb{V}(Y(j) | \mathbf{X})$ necessary to ensure that the integrals

are well defined.

Proof of Proposition 1. By definition of the generalized propensity score, we must have $\mathbb{E}[\mathbf{1}\{Z = j\}/e_j(\mathbf{X})|\mathbf{X}] = 1$ for all $j \in \mathbb{Z}$. Then the average of the potential outcomes in target population h

$$\begin{aligned}
 m_j^h &= \frac{\int_{\mathbb{X}} m_j(\mathbf{X})f(\mathbf{X})h(\mathbf{X})\mu(d\mathbf{X})}{\int_{\mathbb{X}} f(\mathbf{X})h(\mathbf{X})\mu(d\mathbf{X})} \\
 &= \frac{\int_{\mathbb{X}} \mathbb{E}[\mathbf{1}\{Z = j\}Y(j)(h(\mathbf{X})/e_j(\mathbf{X}))|\mathbf{X}]f(\mathbf{X})\mu(d\mathbf{X})}{\int_{\mathbb{X}} \mathbb{E}[\mathbf{1}\{Z = j\}(h(\mathbf{X})/e_j(\mathbf{X}))|\mathbf{X}]f(\mathbf{X})\mu(d\mathbf{X})} \\
 \text{(A.1)} \quad &= \frac{\int_{\mathbb{X}} \mathbb{E}[\mathbf{1}\{Z = j\}Y(j)w_j(\mathbf{X})|\mathbf{X}]f(\mathbf{X})\mu(d\mathbf{X})}{\int_{\mathbb{X}} \mathbb{E}[\mathbf{1}\{Z = j\}w_j(\mathbf{X})|\mathbf{X}]f(\mathbf{X})\mu(d\mathbf{X})}
 \end{aligned}$$

where the second equation holds due to the weak unconfoundedness assumption, $Y(j) \perp \mathbf{1}\{Z = j\}|\mathbf{X}$ (Imbens, 2000). Because $D_{ij} = \mathbf{1}\{Z_i = j\}$, it follows that the estimators, $n^{-1} \sum_{i=1}^n D_{ij}Y_iw_j(\mathbf{X}_i)$ and $n^{-1} \sum_{i=1}^n D_{ij}w_j(\mathbf{X}_i)$, consistently estimate the numerator and denominator of (A.1). Therefore, $\hat{m}_j^h = \sum_{i=1}^n D_{ij}Y_iw_j(\mathbf{X}_i) / \sum_{i=1}^n D_{ij}w_j(\mathbf{X}_i)$ is consistent for m_j^h , and $\hat{\tau}^h(\mathbf{a}) = \sum_{j=1}^J a_j \hat{m}_j^h$ must be consistent for $\tau^h(\mathbf{a}) = \sum_{j=1}^J a_j m_j^h$.

Proof of Proposition 2. By SUTVA (Imbens and Rubin, 2015), we write

$$\hat{\tau}_{\mathbf{a}}^h = \sum_{j=1}^J a_j \frac{\sum_{i=1}^n D_{ij}Y_iw_j(\mathbf{X}_i)}{\sum_{i=1}^n D_{ij}w_j(\mathbf{X}_i)} = \sum_{j=1}^J a_j \frac{\sum_{i=1}^n D_{ij}Y_i(j)w_j(\mathbf{X}_i)}{\sum_{i=1}^n D_{ij}w_j(\mathbf{X}_i)}.$$

Conditional on the assignment $\underline{\mathbf{Z}}$ and sample design $\underline{\mathbf{X}}$, only the potential outcomes are random. Therefore the residual variance of $\hat{\tau}^h(\mathbf{a})$ is

$$\begin{aligned}
 \mathbb{V}[\hat{\tau}^h(\mathbf{a})|\underline{\mathbf{Z}}, \underline{\mathbf{X}}] &= \sum_{j=1}^J a_j^2 \frac{\sum_{i=1}^n v_j(\mathbf{X}_i)D_{ij}w_j^2(\mathbf{X}_i)}{[\sum_{i=1}^n D_{ij}w_j(\mathbf{X}_i)]^2} \\
 &= \sum_{j=1}^J a_j^2 \frac{\sum_{i=1}^n \{v_j(\mathbf{X}_i)/e_j(\mathbf{X}_i)\}\{D_{ij}/e_j(\mathbf{X}_i)\}h^2(\mathbf{X}_i)}{[\sum_{i=1}^n \{D_{ij}/e_j(\mathbf{X}_i)\}h(\mathbf{X}_i)]^2}.
 \end{aligned}$$

Averaging over the joint distribution of $\underline{\mathbf{Z}}$ and $\underline{\mathbf{X}}$, we observe by the Weak Law of Large Numbers that

$$\frac{1}{n} \sum_{i=1}^n \{D_{ij}/e_j(\mathbf{X}_i)\}h(\mathbf{X}_i) \xrightarrow{P} \int_{\mathbb{X}} \mathbb{E}[\mathbf{1}\{Z = j\}/e_j(\mathbf{X})|\mathbf{X}]h(\mathbf{X})f(\mathbf{X})\mu(d\mathbf{X}) = C_h,$$

and

$$\begin{aligned}
 & \frac{1}{n} \sum_{i=1}^n \{v_j(\mathbf{X}_i)/e_j(\mathbf{X}_i)\} \{D_{ij}/e_j(\mathbf{X}_i)\} h^2(\mathbf{X}_i) \\
 & \xrightarrow{p} \int_{\mathbb{X}} v_j(\mathbf{X})/e_j(\mathbf{X}) \mathbb{E}[\mathbb{1}\{Z=j\}/e_j(\mathbf{X})|\mathbf{X}] h^2(\mathbf{X}) f(\mathbf{X}) \mu(d\mathbf{X}) \\
 & = \int_{\mathbb{X}} \{v_j(\mathbf{X})/e_j(\mathbf{X})\} h^2(\mathbf{X}) f(\mathbf{X}) \mu(d\mathbf{X})
 \end{aligned}$$

An application of the Slutsky's Theorem shows $n \cdot \mathbb{V}[\hat{\tau}^h(\mathbf{a})|\underline{\mathbf{Z}}, \underline{\mathbf{X}}] \xrightarrow{p} Q(\mathbf{a}, h)$, where $Q(\mathbf{a}, h)$ is a constant defined in Proposition 2. The uniform integrability assumption for the family of random variables $\{\mathbb{V}[\hat{\tau}^h(\mathbf{a})|\underline{\mathbf{Z}}, \underline{\mathbf{X}}], n \geq 1\}$ then gives the desired L_1 convergence result.

Proof of Proposition 3. For notational simplicity, we use the $\mathbb{E}[\cdot]$ operator to represent $\int_{\mathbb{X}} \cdot f(\mathbf{X}) \mu(d\mathbf{X})$. Under homoscedasticity, $v_j(\mathbf{X}) = v$,

$$\begin{aligned}
 Q(\mathbf{a}, h) &= (v/C_h^2) \int_{\mathbb{X}} \left(\sum_{j=1}^J a_j^2/e_j(\mathbf{X}) \right) h^2(\mathbf{X}) f(\mathbf{X}) \mu(d\mathbf{X}) \\
 &= (v/C_h^2) \mathbb{E} \left\{ h^2(\mathbf{X}) \left(\sum_{j=1}^J a_j^2/e_j(\mathbf{X}) \right) \right\}.
 \end{aligned}$$

Applying the Cauchy-Schwarz inequality, we have

$$\begin{aligned}
 C_h^2 = [\mathbb{E}\{h(\mathbf{X})\}]^2 &= \left[\mathbb{E} \left\{ h(\mathbf{X}) \left(\sum_{j=1}^J a_j^2/e_j(\mathbf{X}) \right)^{1/2} \left(\sum_{j=1}^J a_j^2/e_j(\mathbf{X}) \right)^{-1/2} \right\} \right]^2 \\
 &\leq \mathbb{E} \left\{ h^2(\mathbf{X}) \left(\sum_{j=1}^J a_j^2/e_j(\mathbf{X}) \right) \right\} \mathbb{E} \left\{ \left(\sum_{j=1}^J a_j^2/e_j(\mathbf{X}) \right)^{-1} \right\},
 \end{aligned}$$

and the equality is attained when $h = \tilde{h}(\mathbf{X}) \propto \left(\sum_{j=1}^J a_j^2/e_j(\mathbf{X}) \right)^{-1}$. This implies that

$$\mathbb{E} \left\{ h^2(\mathbf{X}) \left(\sum_{j=1}^J a_j^2/e_j(\mathbf{X}) \right) \right\} / C_h^2 \geq \left[\mathbb{E} \left\{ \left(\sum_{j=1}^J a_j^2/e_j(\mathbf{X}) \right)^{-1} \right\} \right]^{-1} = C_{\tilde{h}}^{-1},$$

which gives $Q(\mathbf{a}, \tilde{h}) = v/C_{\tilde{h}}$.

SUPPLEMENT C: PROOF OF THEOREM 1

From the multinomial logistic model, we have for $i = 1, \dots, n$,

$$\begin{aligned} e_1(\mathbf{X}_i) &= Pr(Z_i = 1|\mathbf{X}_i) = \frac{1}{1 + \sum_{k=2}^J \exp(\alpha_k + \mathbf{X}_i^T \boldsymbol{\beta}_k)}. \\ e_j(\mathbf{X}_i) &= Pr(Z_i = j|\mathbf{X}_i) = \frac{\exp(\alpha_j + \mathbf{X}_i^T \boldsymbol{\beta}_j)}{1 + \sum_{k=2}^J \exp(\alpha_k + \mathbf{X}_i^T \boldsymbol{\beta}_k)}, \quad j = 2, \dots, J. \end{aligned}$$

Since $D_{ij} = \mathbb{1}\{Z_i = j\}$, it is straightforward to show that the log likelihood function

$$l(\boldsymbol{\theta}) = \sum_{i=1}^n l_i(\boldsymbol{\theta}) = \sum_{i=1}^n \left[\sum_{j=2}^J \left\{ D_{ij}(\alpha_j + \mathbf{X}_i^T \boldsymbol{\beta}_j) \right\} - \log \left\{ 1 + \sum_{k=2}^J \exp(\alpha_k + \mathbf{X}_i^T \boldsymbol{\beta}_k) \right\} \right]$$

When the estimation of model parameters is carried out by maximum likelihood, the first-order condition is obtained by differentiating the log likelihood with respect to $\boldsymbol{\theta}$,

$$\begin{aligned} \mathbf{0} &= \mathbf{S}_{\boldsymbol{\theta}} = \sum_{i=1}^n \mathbf{S}_{\boldsymbol{\theta},i} = \sum_{i=1}^n \frac{\partial}{\partial \boldsymbol{\theta}} l_i(\boldsymbol{\theta}) \\ \text{(A.2)} &= \sum_{i=1}^n \left(\frac{\partial}{\partial \alpha_2} l_i(\boldsymbol{\theta}), \dots, \frac{\partial}{\partial \alpha_J} l_i(\boldsymbol{\theta}), \frac{\partial}{\partial \boldsymbol{\beta}_2^T} l_i(\boldsymbol{\theta}), \dots, \frac{\partial}{\partial \boldsymbol{\beta}_J^T} l_i(\boldsymbol{\theta}) \right)^T, \end{aligned}$$

where for $l = 2, \dots, J$,

$$\frac{\partial}{\partial \beta_l} l_i(\boldsymbol{\theta}) = \mathbf{X}_i \frac{\partial}{\partial \alpha_l} l_i(\boldsymbol{\theta}) = \mathbf{X}_i \{D_{il} - e_l(\mathbf{X}_i)\}.$$

We further let $\mathbf{I}_{\boldsymbol{\theta}\boldsymbol{\theta}} = -\mathbb{E}[\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} l_i(\boldsymbol{\theta})]$ be the information matrix, whose exact form can be expressed in a similar fashion but is omitted here for brevity. We denote a consistent estimator for this information by $\hat{\mathbf{I}}_{\boldsymbol{\theta}\boldsymbol{\theta}}$. Under standard regularity conditions (Lehmann, 1983), the stochastic expansion for the maximum likelihood estimator is

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) = \mathbf{I}_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{S}_{\boldsymbol{\theta},i} + o_p(1),$$

where $o_p(1)$ is asymptotically negligible as $n \rightarrow \infty$.

With the multinomial logistic model, the generalized overlap weights are expressed as functions of $\boldsymbol{\theta}$:

$$\begin{aligned} w_1(\mathbf{X}_i) &= w_1(\mathbf{X}_i; \boldsymbol{\theta}) = \frac{1}{1 + \sum_{k=2}^J \exp(-\alpha_k - \mathbf{X}_i^T \boldsymbol{\beta}_k)} \\ w_j(\mathbf{X}_i) &= w_j(\mathbf{X}_i; \boldsymbol{\theta}) = \frac{\exp(-\alpha_j - \mathbf{X}_i^T \boldsymbol{\beta}_j)}{1 + \sum_{k=2}^J \exp(-\alpha_k - \mathbf{X}_i^T \boldsymbol{\beta}_k)}, \quad j = 2, \dots, J, \end{aligned}$$

and the derivative of the weights takes the form

$$\dot{w}_j(\mathbf{X}_i) \equiv \frac{\partial}{\partial \boldsymbol{\theta}} w_j(\mathbf{X}_i) = \left(\frac{\partial}{\partial \alpha_2} w_j(\mathbf{X}_i), \dots, \frac{\partial}{\partial \alpha_J} w_j(\mathbf{X}_i), \frac{\partial}{\partial \boldsymbol{\beta}_2^T} w_j(\mathbf{X}_i), \dots, \frac{\partial}{\partial \boldsymbol{\beta}_J^T} w_j(\mathbf{X}_i) \right)^T,$$

where for $j = 1, \dots, J$ and $l = 2, \dots, J$,

$$\frac{\partial}{\partial \boldsymbol{\beta}_l} w_j(\mathbf{X}_i) = \mathbf{X}_i \frac{\partial}{\partial \alpha_l} w_j(\mathbf{X}_i) = \mathbf{X}_i \{w_j(\mathbf{X}_i) w_l(\mathbf{X}_i) - \delta_{jl} w_l(\mathbf{X}_i)\},$$

and $\delta_{jl} = \mathbb{1}\{j = l\}$.

For $j = 1, \dots, J$, the plug-in weighting estimator \hat{m}_j^h can be regarded as the solution of the following estimating equation

$$\sum_{i=1}^n \mathbf{U}(\hat{m}_j^h, \hat{\boldsymbol{\theta}}) = \sum_{i=1}^n D_{ij} (Y_i - \hat{m}_j^h) w_j(\mathbf{X}_i; \hat{\boldsymbol{\theta}}) = \mathbf{0}.$$

Under standard regularity conditions (van der Vaart, 1998), a first-order Taylor expansion of the unbiased estimating equations around the truth leads to

$$\begin{aligned} \sqrt{n}(\hat{m}_j^h - m_j^h) &= \varpi^{-1} \left\{ \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{U}(m_j^h, \boldsymbol{\theta}) + \mathbf{H}_j^T \sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \right\} + o_p(1) \\ \text{(A.3)} \quad &= \varpi^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n \left\{ D_{ij} (Y_i - m_j^h) w_j(\mathbf{X}_i) + \mathbf{H}_j^T \mathbf{I}_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1} \mathbf{S}_{\boldsymbol{\theta},i} \right\} + o_p(1), \end{aligned}$$

where $\varpi = \mathbb{E}[D_{ij} w_j(\mathbf{X}_i)] = \mathbb{E}[h(\mathbf{X}_i)]$, and $\mathbf{H}_j = \mathbb{E}[D_{ij} (Y_i - m_j^h) \dot{w}_j(\mathbf{X}_i)] = \mathbb{E}[(Y_i - m_j^h) e_j(\mathbf{X}_i) \dot{w}_j(\mathbf{X}_i)]$. Therefore, given any fixed coefficient $\mathbf{a} = (a_1, \dots, a_J)'$, we have

$$\sqrt{n}\{\hat{\tau}^h(\mathbf{a}) - \tau^h(\mathbf{a})\} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \left\{ \frac{1}{\varpi} \sum_{j=1}^J a_j \psi_{ij} \right\} + o_p(1),$$

where we define $\psi_{ij} = D_{ij} (Y_i - m_j^h) w_j(\mathbf{X}_i) + \mathbf{H}_j^T \mathbf{I}_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1} \mathbf{S}_{\boldsymbol{\theta},i}$. Since the triplets $\{Y_i, \mathbf{X}_i, Z_i\}'s$ are assumed i.i.d., an application of the standard Central Limit Theorem gives,

$$\sqrt{n}\{\hat{\tau}^h(\mathbf{a}) - \tau^h(\mathbf{a})\} \xrightarrow{d} \mathcal{N} \left(0, \varpi^{-2} \mathbb{E} \left\{ \sum_{j=1}^J a_j \psi_{ij} \right\}^2 \right).$$

In practice, we use the empirical sandwich estimator to consistently estimate the large-sample variance (Stefanski and Boos, 2002); the variance of $\hat{\tau}^h(\mathbf{a})$ is estimated by

$$\frac{1}{(n\hat{\varpi})^2} \sum_{i=1}^n \left\{ \sum_{j=1}^J a_j \hat{\psi}_{ij} \right\}^2,$$

where

$$\begin{aligned} \hat{\psi}_{ij} &= D_{ij}(Y_i - \hat{m}_j^h) w_j(\mathbf{X}_i; \hat{\boldsymbol{\theta}}) + \hat{\mathbf{H}}_j^T \hat{\mathbf{I}}_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1} \hat{\mathbf{S}}_{\boldsymbol{\theta},i}, \\ \hat{\varpi} &= \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{k=1}^J 1/\hat{e}_k(\mathbf{X}_i) \right\}^{-1}, \\ \hat{\mathbf{H}}_j &= \frac{1}{n} \sum_{i=1}^n D_{ij}(Y_i - \hat{m}_j^h) \dot{w}_j(\mathbf{X}_i; \hat{\boldsymbol{\theta}}), \end{aligned}$$

and $\hat{\mathbf{S}}_{\boldsymbol{\theta},i}$ is the estimated individual score function (A.2) from the propensity model. For pairwise comparisons, we substitute \mathbf{a} with $\boldsymbol{\lambda}_{j,j'}$ to obtain the results in Theorem 1.

For completeness, we next offer three remarks regarding variance estimation.

REMARK 3. *One could similarly characterize the asymptotic distribution of a collection of estimators specified by different contrast coefficients. Briefly, let the coefficient matrix $\mathbf{A}_{J \times R} = (\mathbf{a}_1, \dots, \mathbf{a}_R)$, where the vector \mathbf{a} 's are distinct from one another. For pairwise comparisons, each vector \mathbf{a} is a distinct element in the set \mathbb{S} . Write $\boldsymbol{\tau} = (\tau^h(\mathbf{a}_1), \dots, \tau^h(\mathbf{a}_R))'$, and $\hat{\boldsymbol{\tau}} = (\hat{\tau}^h(\mathbf{a}_1), \dots, \hat{\tau}^h(\mathbf{a}_R))'$ as the corresponding weighting estimators. Further denote $\boldsymbol{\psi}_i = (\psi_{i1}, \dots, \psi_{iJ})'$, and it can be shown using similar arguments that*

$$\sqrt{n}(\hat{\boldsymbol{\tau}}^h - \boldsymbol{\tau}^h) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \varpi^{-1} \mathbf{A}^T \boldsymbol{\psi}_i + o_p(1) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \varpi^{-2} \mathbf{A}^T \mathbb{E}\{\boldsymbol{\psi}_i \boldsymbol{\psi}_i^T\} \mathbf{A}).$$

The covariance for $\hat{\boldsymbol{\tau}}^h$ can then be estimated by the empirical sandwich estimator

$$\hat{\mathbf{V}}(\hat{\boldsymbol{\tau}}^h) = (n\hat{\varpi})^{-2} \mathbf{A}^T \left\{ \sum_{i=1}^n \hat{\boldsymbol{\psi}}_i \hat{\boldsymbol{\psi}}_i^T \right\} \mathbf{A}.$$

REMARK 4. *Although the above derivation focuses on the generalized overlap weights, a more general presentation for other members of the balancing weights is possible, provided that the balancing weights is a differentiable function in the generalized propensity scores. This differentiability condition rules out the generalized matching weights, which is smooth but not everywhere differentiable and so closed-form variance requires parametric smooth approximation (Li and Greene, 2013) (such approximations could be challenging with multiple treatments since the weight function have infinite-many non-differentiable points). In particular, if we choose the balancing weights as the inverse probability weights, in which case $h(\mathbf{X}) = 1$ and the target population is the combined population from all groups, the above derivation can be repeated by substituting the correct forms of $w_j(\mathbf{X}_i)$ and $\dot{w}_j(\mathbf{X}_i)$. For example, the inverse probability weights are*

$$\begin{aligned} w_1(\mathbf{X}_i) &= 1/e_1(\mathbf{X}_i) = 1 + \sum_{k=2}^J \exp(\alpha_k + \mathbf{X}_i^T \boldsymbol{\beta}_k) \\ w_j(\mathbf{X}_i) &= 1/e_j(\mathbf{X}_i) = \frac{1 + \sum_{k=2}^J \exp(\alpha_k + \mathbf{X}_i^T \boldsymbol{\beta}_k)}{\exp(\alpha_j + \mathbf{X}_i^T \boldsymbol{\beta}_j)}, \quad j = 2, \dots, J, \end{aligned}$$

and the derivative of the weights takes the form

$$\frac{\partial}{\partial \boldsymbol{\beta}_l} w_j(\mathbf{X}_i) = \mathbf{X}_i \frac{\partial}{\partial \alpha_l} w_j(\mathbf{X}_i) = \mathbf{X}_i \{w_j(\mathbf{X}_i)/w_l(\mathbf{X}_i) - \delta_{jl} w_l(\mathbf{X}_i)\},$$

for $j = 1, \dots, J$ and $l = 2, \dots, J$. Of note, this empirical sandwich variance for $h(\mathbf{X}) = 1$ extends the one proposed by Lunceford and Davidian (2004) for binary treatments, and is used to obtain the interval estimates for IPW in the main manuscript.

REMARK 5. *We have focused on the case with a multinomial logistic propensity score model, but in fact the derivation can be made more general to accommodate other propensity score models that admit a regular and asymptotically linear estimator for the model parameters (Tsiatis, 2006). This condition permits a stochastic expansion for $\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})$, which can then be substituted into (A.3) to obtain the corresponding sandwich variance estimator. In particular, one could replace the multinomial logistic model with a multinomial Probit model, which is another commonly used regression model to accommodate categorical responses.*

SUPPLEMENT D: ADDITIONAL SIMULATION RESULTS

In the second simulation with $J = 6$ groups, we specify the parameters to simulate both adequate and lack of overlap. Specifically, we spec-

ify $\beta_2^T = \kappa_2 \times (1, 1, 2, 1, 1, 1)$, $\beta_3^T = \kappa_3 \times (1, 1, 1, 1, 1, -5)$, $\beta_4^T = \kappa_4 \times (1, 1, 1, 1, 1, 5)$, $\beta_5^T = \kappa_5 \times (1, 1, 1, -2, 1, 1)$ and $\beta_6^T = \kappa_6 \times (1, 1, 1, -2, -1, 1)$. We use $(\kappa_2, \kappa_3, \kappa_4, \kappa_5, \kappa_6) = (0.1, 0.15, 0.2, 0.25, 0.3)$ to simulate a scenario with adequate overlap and $(\kappa_2, \kappa_3, \kappa_4, \kappa_5, \kappa_6) = (0.4, 0.6, 0.8, 1, 1.2)$ to represent a challenging scenario with strong propensity tails. The intercepts are chosen so that the marginal treatment proportions are fixed around $(0.12, 0.16, 0.12, 0.25, 0.2, 0.15)$. Finally, the coefficients for the outcome model is specified as $\gamma_1^T = (-1.5, 1, 1, 1, 1, 1)$, $\gamma_2^T = (-4, 2, 3, 1, 2, 2)$, $\gamma_3^T = (4, 3, 1, 2, -1, -1)$, $\gamma_4^T = (1, 4, 1, 2, -1, -1)$, $\gamma_5^T = (3.5, 5, 1, 2, -1, -1)$ and $\gamma_6^T = (3.5, 6, 1, 2, -1, -1)$. The total sample size is fixed at $n = 6000$ for $J = 6$. Visual inspections of the overlap in each simulation scenario are provided in Supplementary Figures 1-4. Simulation results for $J = 6$ are presented in Supplementary Figures 5 and 6.

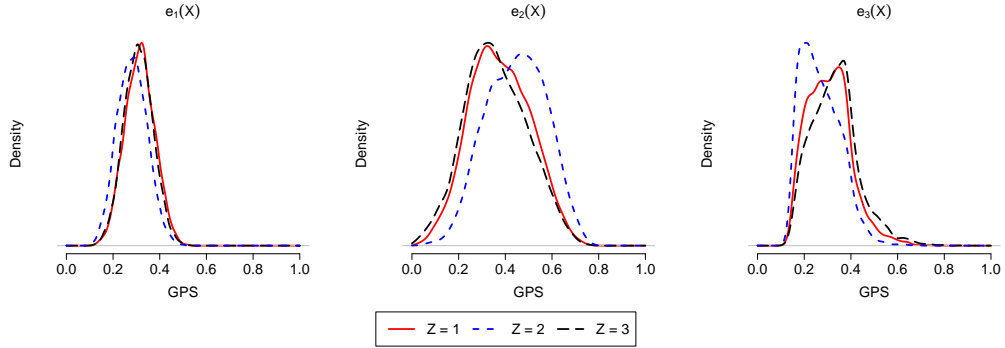


FIG 1. *Distribution of the generalized propensity scores in the simulation with $J = 3$ groups and adequate overlap.*

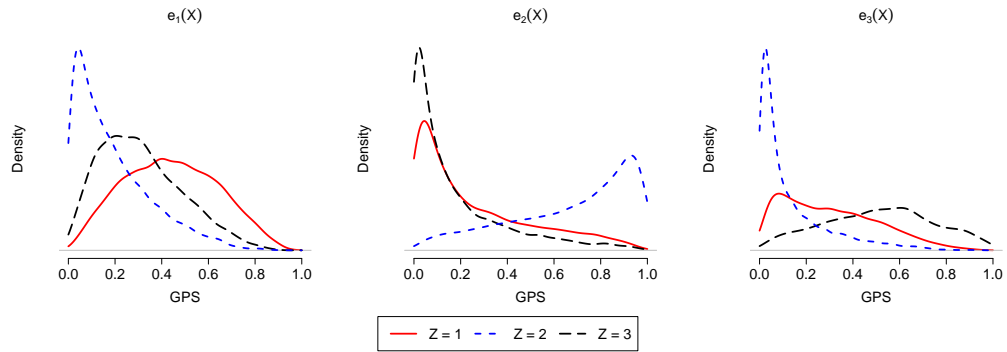


FIG 2. Distribution of the generalized propensity scores in the simulation with $J = 3$ groups and lack of overlap.

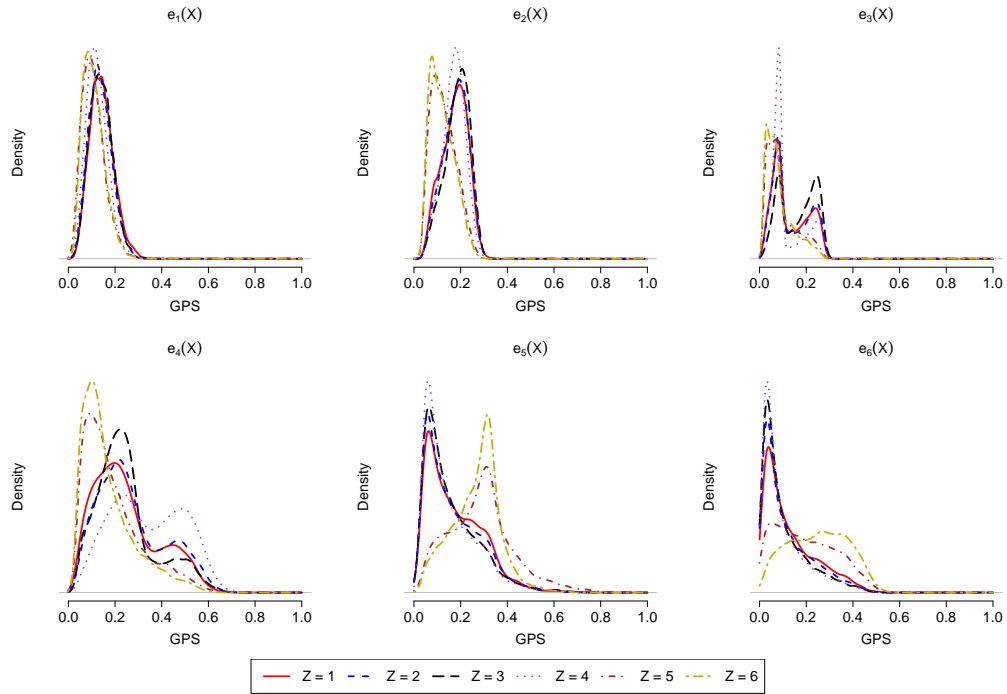


FIG 3. Distribution of the generalized propensity scores in the simulation with $J = 6$ groups and adequate overlap.

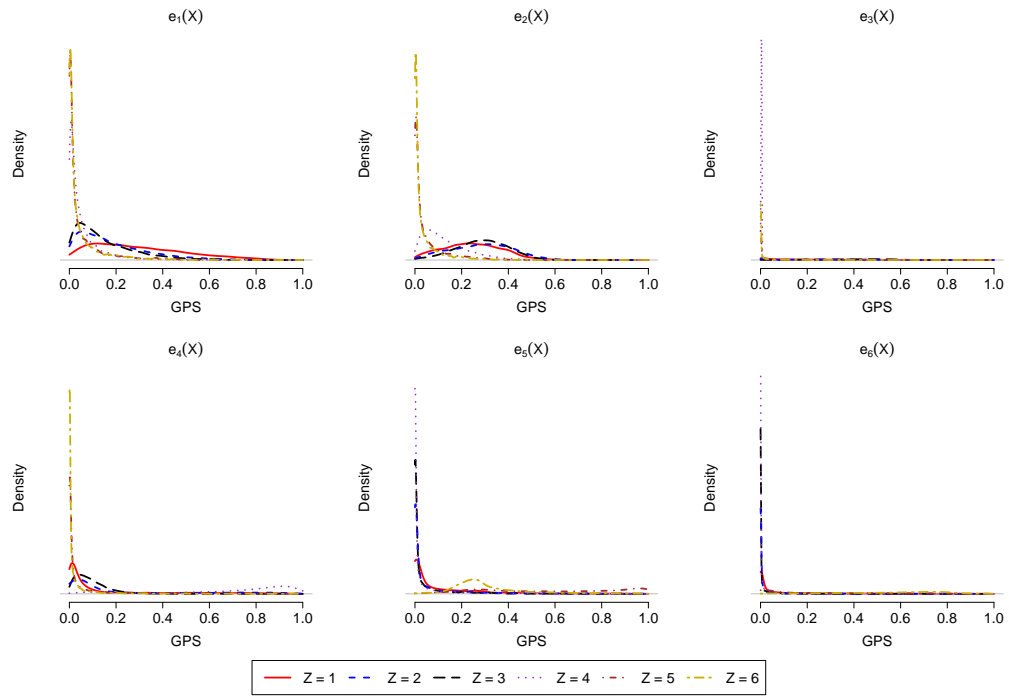


FIG 4. *Distribution of the generalized propensity scores in the simulation with $J = 6$ groups and lack of overlap.*

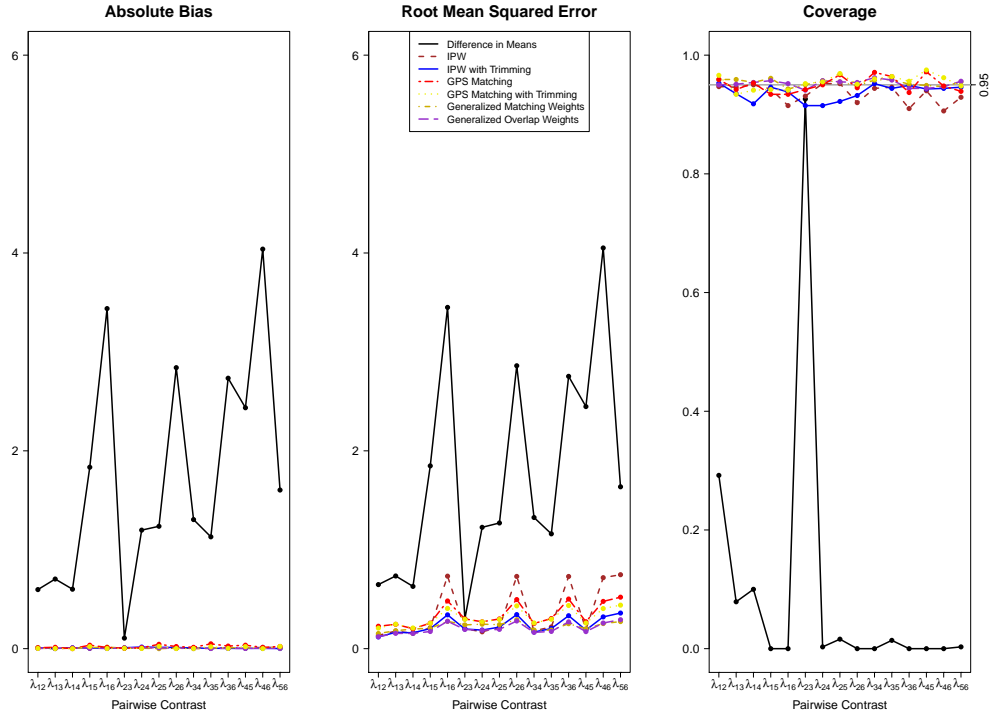


FIG 5. Simulation results with $J = 6$ treatment groups and $(\kappa_2, \kappa_3, \kappa_4, \kappa_5, \kappa_6) = (0.1, 0.15, 0.2, 0.25, 0.3)$, i.e., with adequate overlap. Optimal trimming excludes 3% ~ 7% of the total sample. For a given approach, each one of the 15 causal comparisons is represented by the contrast $\lambda_{j,j'}$ for notational simplicity.

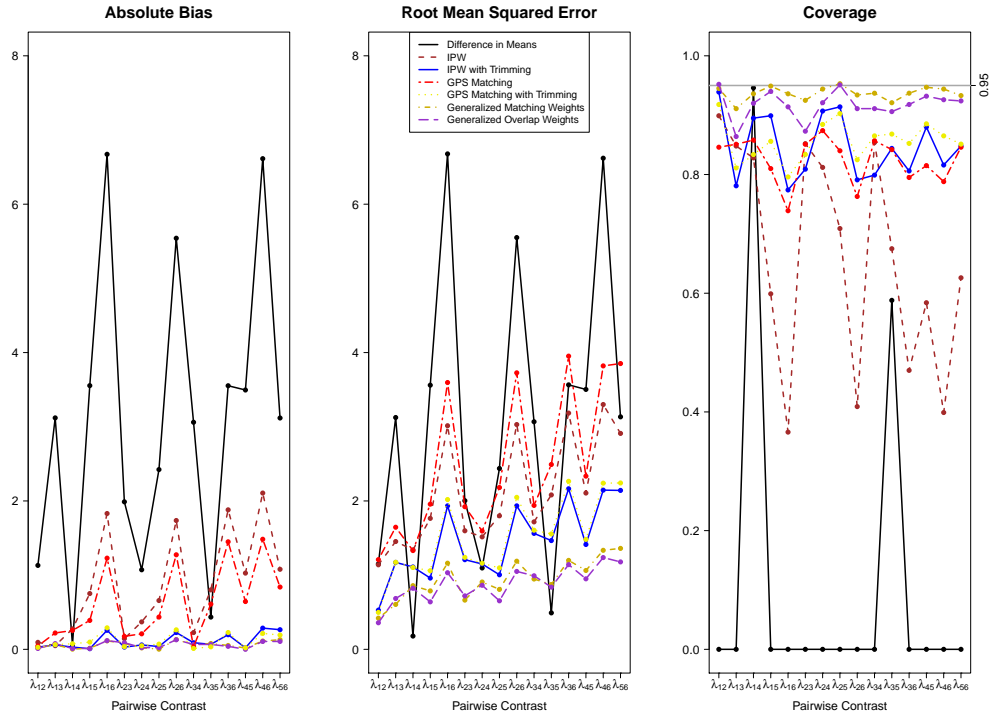


FIG 6. Simulation results with $J = 6$ treatment groups and $(\kappa_2, \kappa_3, \kappa_4, \kappa_5, \kappa_6) = (0.4, 0.6, 0.8, 1, 1.2)$, i.e., with strong propensity tails. Optimal trimming excludes 52% \sim 74% of the total sample. For a given approach, each one of the 15 causal comparisons is represented by the contrast $\lambda_{j,j'}$ for notational simplicity.

SUPPLEMENT E: SAMPLE R CODE FOR IMPLEMENTING
GENERALIZED OVERLAP WEIGHTS

We provide sample R code to illustrate the application of the generalized overlap weighting scheme along with the implementation of the empirical sandwich variance, based on a simulated observational data set with $J = 3$ treatments. We use the following code to simulate a data set for an illustrative analysis. The data generating process is described in Section 5 of the main manuscript.

```

1 # Consider three treatments and sample size 1500
2 set.seed(2019)
3 require(MASS)
4 J = 3
5 n = 1500
6
7 # Simulate pre-treatment covariates
8 # The simulation setup is similar to Yang et al. (2016)
9 # X1-X3 are multivariate normal covariates
10 # X4 is uniform variable
11 # X5 is a Chi-squared variable
12 # X6 is binary
13 vars = c(2,1,1)
14 covars = c(1,-1,-.5)
15 mu = c(0,0,0)
16 tau = 1
17 Sigma = diag(vars)
18 Sigma[2,1] = Sigma[1,2] = covars[1]
19 Sigma[3,1] = Sigma[1,3] = covars[2]
20 Sigma[3,2] = Sigma[2,3] = covars[3]
21
22 X13 = mvrnorm(n, mu=mu, Sigma=Sigma, empirical = FALSE)
23 X4 = runif(n,-3,3)
24 X5 = rchisq(n, df=1)
25 X6 = rbinom(n, size=1, prob=.5)
26 X16 = cbind(X13, X4, X5, X6)
27 X06 = cbind(1, X13, X4, X5, X6)
28
29 # Assignment mechanism
30 beta1 = c(0,0,0,0,0,0,0)
31 beta2 = c(0.344, 0.2*c(1,1,1,-1,-1,1))
32 beta3 = c(-0.178, 0.1*c(1,1,1,1,1,1))
33 xb2 = c(X06 %*% beta2)
34 xb3 = c(X06 %*% beta3)
35 exb2 = exp(xb2)
36 exb3 = exp(xb3)
37 e1 = 1 / (1+exp(xb2)+exp(xb3))
38 e2 = exp(xb2)/(1+exp(xb2)+exp(xb3))

```

```

39 e3 = exp(xb3)/(1+exp(xb2)+exp(xb3))
40 e = cbind(e1,e2,e3) # true propensity scores
41
42 # Simulate observed treatment
43 D = matrix(NA, n, J)
44 colnames(D) = c("D1", "D2", "D3")
45 for(k in 1:n){
46   D[k,] = rmultinom(1, 1, prob = e[k,])
47 }
48 Z = D[, "D1"] + 2*D[, "D2"] + 3*D[, "D3"]
49
50 # True potential outcome models
51 u = rnorm(n)
52 gamma1 = c(-1.5, 1, 1, 1, 1, 1, 1)
53 gamma2 = c(-4, 2, 3, 1, 2, 2, 2)
54 gamma3 = c(3, 3, 1, 2, -1, -1, -1)
55 EY1 = c(X06 %%% gamma1)
56 EY2 = c(X06 %%% gamma2)
57 EY3 = c(X06 %%% gamma3)
58 EY = cbind(EY1, EY2, EY3)
59 Y = rowSums(EY*D) + u

```

We now create a data set with only observed outcomes, treatments and pre-treatment covariates, and examine the first few rows to get a sense of the data structure.

```

1 # Create the analysis datasets
2 analdata = data.frame(Y=Y, Z=Z, X=X16)
3 colnames(analdata) = c("Y", "Z", "X1", "X2", "X3", "X4", "X5",
4   "X6")
5 # Peek at the data
6 round(head(analdata), 3)
7
8 #
9 # 1 7.673 2 0.869 1.696 0.261 2.310 0.156 0
10 # 2 -3.701 2 -1.142 0.370 0.467 -1.628 0.041 1
11 # 3 -11.345 2 -1.975 -1.866 1.247 -0.468 0.819 1
12 # 4 3.534 2 1.857 0.463 -0.003 2.195 0.015 0
13 # 5 0.524 1 -2.555 -0.790 -0.120 2.398 2.391 0
14 # 6 7.198 3 1.472 0.767 0.349 1.229 0.819 0

```

We could repeat the above data generating process for a large number of times, and numerically approximate the true values of the pairwise ATO (by averaging out the Monte Carlo errors in repeated simulations). We will omit the details here but indicate that the true pairwise ATO quantities are $\tau^h(\lambda_{1,2}) = 1.03$, $\tau^h(\lambda_{1,3}) = -1.36$ and $\tau^h(\lambda_{2,3}) = -2.39$. From now on, we

will remain agnostic to the true data generating process, and estimate the pairwise ATO using only the observed data object `analdata`. We estimate the generalized propensity scores using a correctly specified multinomial logistic regression.

```

1 require(nnet)
2 require(numDeriv)
3
4 # Create a n x J matrix of treatment indicators
5 # Each row corresponds to an observation
6 # Each column represents a treatment category
7 D = cbind(as.numeric(analdata$Z == 1),
8           as.numeric(analdata$Z == 2),
9           as.numeric(analdata$Z == 3))
10
11 # Estimate the GPS model
12 n = length(analdata$Y)
13 gps.fit = multinom(Z ~ X1 + X2 + X3 + X4 + X5 + X6, data =
14               analdata, maxit = 500, Hess = TRUE, trace = FALSE)
15 Thetah = t(coef(gps.fit)) # matrix coefficient
16 thetah = c(Thetah) # vec operator
17 IthetaInv = n*vcov(gps.fit) # extract the covariance
18           matrix
19 e = gps.fit$fitted.values # estimated propensity scores
20 eInv = 1/e # inverse probability weights
21 h = as.numeric(1/rowSums(eInv)) # obtain the optimal tilting
22           function
23 w = eInv * h # create the generalized
24           overlap weights

```

In practice, we need to check whether the estimated generalized overlap weights adequately balance the weighted covariates across groups (therefore assess the adequacy of the generalized propensity score model). If the weighted balance is unsatisfactory, we may need to revise the propensity score model by including additional higher-order terms or interactions to make improvements. We defined the following two functions to calculate the *population standardized difference* (PSD) and the pairwise *absolute standardized difference* (ASD), defined in Section 3.2.

```

1 # Population standardized difference
2 PSDfun = function(Z, D, covM, h, w){
3   # Z: treatment value
4   # D: matrix of treatment indicators
5   # covM: covariate matrix, n by p
6   # h: value of the h function, n by 1
7   # w: matrix of weights, n by J
8

```

```

9   p = ncol(covM)
10  J = ncol(D)
11  psdM = matrix(NA, J, p)
12  for(k in 1:p){
13    x = as.numeric(covM[,k])
14    xbar.p = sum(rowSums(D*x*h))/sum(rowSums(D*h))
15    xbar.j = colSums(D*x*w)/colSums(D*w)
16    s2x = mean(tapply(x, Z, var))
17    psdM[,k] = abs(xbar.j - xbar.p)/sqrt(s2x)
18  }
19  return(t(psdM))
20 }
21
22 # create contrast matrix
23 J = 3
24 Jc = choose(J, 2) # number of contrasts
25 loc = t(combn(J,2))
26 C = matrix(0, Jc, J)
27 for(i in 1:Jc){
28   C[i,loc[i,1]] = 1
29   C[i,loc[i,2]] = -1
30 }
31
32 # Absolute standardized difference
33 ASDfun = function(Z, D, C, covM, h, w){
34   # Z: treatment value
35   # D: matrix of treatment indicators
36   # C: contrast matrix
37   # covM: covariate matrix, n by p
38   # h: value of the h function
39   # w: matrix of weights, n by J
40
41   p = ncol(covM)
42   Jc = nrow(C)
43   asdM = matrix(NA, Jc, p)
44
45   # loop through each pair
46   for(k in 1:p){
47     x = as.numeric(covM[,k])
48     xbar.j = colSums(D*x*w)/colSums(D*w)
49     s2x = mean(tapply(x, Z, var))
50     asdM[,k] = abs(C %*% xbar.j)/sqrt(s2x)
51   }
52   return(t(asdM))
53 }

```

We compared the covariate balance before and after weighting using the following code. The covariate balance is satisfactory after weighting, which assures the adequacy of the propensity score model.


```

1 # Before weighting: PSD and ASD may exceed 10%
2 round(PSDfun(Z=Z, D=D, covM=X, h=rep(1,n), w=rep(1,n)), 3)
3 #      [,1] [,2] [,3]
4 # [1,] 0.182 0.122 0.006
5 # [2,] 0.240 0.144 0.033
6 # [3,] 0.085 0.042 0.026
7 # [4,] 0.083 0.254 0.294
8 # [5,] 0.003 0.141 0.207
9 # [6,] 0.058 0.034 0.008
10 round(ASDfun(Z=Z, D=D, C=C, covM=X, h=rep(1,n), w=rep(1,n)), 3)
11 #      [,1] [,2] [,3]
12 # [1,] 0.304 0.188 0.116
13 # [2,] 0.384 0.273 0.111
14 # [3,] 0.127 0.111 0.016
15 # [4,] 0.338 0.210 0.548
16 # [5,] 0.144 0.204 0.348
17 # [6,] 0.092 0.066 0.026
18
19 # After weighting: PSD and ASD all below 10%
20 round(PSDfun(Z=Z, D=D, covM=X, h=h, w=w), 3)
21 #      [,1] [,2] [,3]
22 # [1,] 0.018 0.016 0.024
23 # [2,] 0.019 0.002 0.005
24 # [3,] 0.007 0.004 0.013
25 # [4,] 0.006 0.004 0.007
26 # [5,] 0.002 0.034 0.011
27 # [6,] 0.013 0.007 0.008
28
29 round(ASDfun(Z=Z, D=D, C=C, covM=X, h=h, w=w), 3)
30 #      [,1] [,2] [,3]
31 # [1,] 0.002 0.006 0.008
32 # [2,] 0.021 0.014 0.006
33 # [3,] 0.004 0.006 0.010
34 # [4,] 0.003 0.001 0.004
35 # [5,] 0.035 0.009 0.044
36 # [6,] 0.006 0.005 0.001

```

The following code is used to estimate the pairwise ATO using the generalized overlap weights. The variance is obtained using the empirical sandwich estimator, which is also used to construct the 95% confidence interval. Recall that the true pairwise ATO quantities are $\tau^h(\boldsymbol{\lambda}_{1,2}) = 1.03$, $\tau^h(\boldsymbol{\lambda}_{1,3}) = -1.36$ and $\tau^h(\boldsymbol{\lambda}_{2,3}) = -2.39$; the estimated pairwise ATO, $\hat{\tau}^h(\boldsymbol{\lambda}_{1,2}) = 1.08$, $\hat{\tau}^h(\boldsymbol{\lambda}_{1,3}) = -1.19$ and $\hat{\tau}^h(\boldsymbol{\lambda}_{2,3}) = -2.27$, are close to the true quantities in this example.

```

1 # Point estimation
2 mhat = as.numeric(colSums(D*Y*w)/colSums(D*w))

```

```

3 tau = as.matrix(C**mhat)
4
5 # Variance and interval estimation
6 omega = mean(h)
7
8 # Calculate gradient of weights
9 Hmat = NULL
10 for(j in 1:J){
11   wj = function(theta){
12     Theta = matrix(theta, ncol(X)+1, ncol(D)-1)
13     Eta = cbind(1,X) ** cbind(0,Theta)
14     return(as.numeric(exp(-Eta[,j])) / as.numeric(rowSums(exp(-
15       Eta))))
16   }
17   wdotj = jacobian(wj,thetah)
18   Hmat = rbind(Hmat, c(colMeans(D[,j] * (Y - mhat[j]) * wdotj)
19     ))
20 }
21 # Multinomial logistic score function
22 loglik = function(theta){
23   Theta = matrix(theta, ncol(X)+1, ncol(D)-1)
24   Eta = cbind(1,X) ** Theta
25   ltheta = as.numeric(rowSums(D[,,-1]*Eta)-log(1 + rowSums(exp(
26     Eta))))
27   return(ltheta)
28 }
29 Sthetah = jacobian(loglik, thetah)
30
31 # Covariance matrix of pairwise ATO estimates
32 # This a multivariate version of Theorem 1
33 # and relevant details are provided in Remark 3 (Supplement C)
34 YMat = matrix(rep(Y, J), n, J)
35 mhatMat = matrix(rep(mhat, each=n), n, J)
36 Psi = t(D*(YMat - mhatMat)*w) + Hmat ** IthetaInv ** t(
37   Sthetah)
38 Sigmah = diag(C ** tcrossprod(Psi) ** t(C) / (n * omega)^2)
39 se = sqrt(Sigmah)
40 lcl = tau - qnorm(0.975) * se
41 ucl = tau + qnorm(0.975) * se
42 results = cbind(tau, lcl, ucl)
43 colnames(results) = c("Point estimates", "95% Lower limit", "
44   95% Upper limit")
45 rownames(results) = c("1-2", "1-3", "2-3")
46 round(results, 3)
47 #      Point estimates 95% Lower limit 95% Upper limit
48 # 1-2          1.080          0.758          1.402
49 # 1-3         -1.188         -1.515         -0.861
50 # 2-3         -2.268         -2.756         -1.779

```

REFERENCES

- IMBENS, G. W. (2000). The role of the propensity score in estimating dose-response functions. *Biometrika* **87** 706–710.
- IMBENS, G. W. and RUBIN, D. B. (2015). *Causal Inference for Statistics, Social, and Biomedical Sciences*. Cambridge University Press, New York, NY.
- LECHNER, M. (2001). Identification and estimation of causal effects of multiple treatments under the conditional independence assumption. In *Econometric Evaluations of Active Labor Market Policies in Europe* (M. Lechner and F. Pfeiffer, eds.) 43–58. Heidelberg: Physica.
- LEHMANN, E. (1983). *Theory of Point Estimation*. Springer, New York.
- LI, L. and GREENE, T. (2013). A weighting analogue to pair matching in propensity score analysis. *International Journal of Biostatistics* **9** 1-20.
- LOPEZ, M. J. and GUTMAN, R. (2017). Estimation of causal effects with multiple treatments: A review and new ideas. *Statistical Science* **32** 432–454.
- LUNCEFORD, J. K. and DAVIDIAN, M. (2004). Stratification and weighting via the propensity score in estimation of causal treatment effects: a comparative study. *Statistics in Medicine* **23** 2937–2960.
- STEFANSKI, L. A. and BOOS, D. D. (2002). The calculus of M-estimation. *American Statistician* **56** 29–38.
- TSIATIS, A. A. (2006). *Semiparametric Theory and Missing Data*. Springer, New York, NY.
- VAN DER VAART, A. W. (1998). *Asymptotic Statistics. Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press.

F. LI
 DEPARTMENT OF BIOSTATISTICS
 YALE UNIVERSITY
 135 COLLEGE ST
 NEW HAVEN, CONNECTICUT 06510
 USA
 E-MAIL: fan.f.li@yale.edu

F. LI
 DEPARTMENT OF STATISTICAL SCIENCE
 DUKE UNIVERSITY
 122 OLD CHEMISTRY BUILDING
 DURHAM, NORTH CAROLINA 27708
 USA
 E-MAIL: ffi@stat.duke.edu