

SUPPLEMENT TO “WHO SHOULD BE TREATED? EMPIRICAL WELFARE
MAXIMIZATION METHODS FOR TREATMENT CHOICE”
(*Econometrica*, Vol. 86, No. 2, March 2018, 591–616)

TORU KITAGAWA
Cemmap and Department of Economics, University College London

ALEKSEY TETENOV
Department of Economics, University of Bristol

APPENDIX A: LEMMAS AND PROOFS

A.1. Notations and Basic Lemmas

LET $Z_i = (Y_i, D_i, X_i) \in \mathcal{Z}$. The *subgraph* of a real-valued function $f : \mathcal{Z} \mapsto \mathbb{R}$ is the set

$$\text{SG}(f) \equiv \{(z, t) \in \mathcal{Z} \times \mathbb{R} : 0 \leq t \leq f(z) \text{ or } f(z) \leq t \leq 0\}.$$

The following lemma establishes a link between the VC-dimension of a class of subsets in the covariate space \mathcal{X} and the VC-dimension of a class of subgraphs of functions on $\mathcal{Z} = \mathbb{R} \times \{0, 1\} \times \mathcal{X}$ (their subgraphs will be in $\mathcal{Z} \times \mathbb{R}$).

LEMMA A.1: *Let \mathcal{G} be a VC-class of subsets of \mathcal{X} with VC-dimension $v < \infty$. Let g and h be two given functions from \mathcal{Z} to \mathbb{R} . Then the set of functions from \mathcal{Z} to \mathbb{R}*

$$\mathcal{F} = \{f_G(z) = g(z) \cdot 1\{x \in G\} + h(z)1\{x \notin G\} : G \in \mathcal{G}\}$$

is a VC-subgraph class of functions with VC-dimension less than or equal to v .

PROOF: Let $z_i = (y_i, d_i, x_i)$. By the assumption, no set of $(v + 1)$ points in \mathcal{X} could be shattered by \mathcal{G} . Take an arbitrary set of $(v + 1)$ points in $\mathcal{Z} \times \mathbb{R}$, $A = \{(z_1, t_1), \dots, (z_{v+1}, t_{v+1})\}$. Denote the collection of subgraphs of \mathcal{F} by $\text{SG}(\mathcal{F}) \equiv \{\text{SG}(f_G), G \in \mathcal{G}\}$. We want to show that $\text{SG}(\mathcal{F})$ does not shatter A .

If, for some $i \in \{1, \dots, (v + 1)\}$, $(z_i, t_i) \in \text{SG}(g) \cap \text{SG}(h)$, then $\text{SG}(\mathcal{F})$ cannot pick out all of the subsets of A because the i th point is included in any $S \in \text{SG}(\mathcal{F})$. Similarly, if, for some $i \in \{1, \dots, (v + 1)\}$, $(z_i, t_i) \in \text{SG}(g)^c \cap \text{SG}(h)^c$, then point i cannot be included in any $S \in \text{SG}(\mathcal{F})$.

The remaining case is that, for each i , either $(z_i, t_i) \in \text{SG}(g) \cap \text{SG}(h)^c$ or $(z_i, t_i) \in \text{SG}(g)^c \cap \text{SG}(h)$ holds. Indicate the former case by $\delta_i = 0$ and the latter case by $\delta_i = 1$. The points with $\delta_i = 0$ could be picked by $\text{SG}(f_G)$ if and only if $x_i \notin G$. The points with $\delta_i = 1$ could be picked if and only if $x_i \in G$. Given that \mathcal{G} is a VC-class with VC-dimension v , there exists a subset X_0 of $\{x_1, \dots, x_{v+1}\}$ such that $X_0 \neq (\{x_1, \dots, x_{v+1}\} \cap G)$ for any $G \in \mathcal{G}$. Then there could be no set $S \in \text{SG}(\mathcal{F})$ that picks out the set (possibly empty)

$$\{(z_i, t_i) : (x_i \in X_0 \text{ and } \delta_i = 1) \text{ or } (x_i \notin X_0 \text{ and } \delta_i = 0)\}, \quad (\text{A.1})$$

Toru Kitagawa: t.kitagawa@ucl.ac.uk
Aleksey Tetenov: a.tetenov@bristol.ac.uk

because this set of points could only be picked out by $\text{SG}(f_G)$ if $(\{x_1, \dots, x_{v+1}\} \cap G) = X_0$. Hence, \mathcal{F} is a VC-subgraph class of functions with VC-dimension less than or equal to v . Q.E.D.

In addition to the notations introduced in the main text, the following notations are used throughout the supplementary material. The empirical probability distribution based on an i.i.d. size n sample of $Z_i = (Y_i, D_i, X_i)$ is denoted by P^n . $L_2(P)$ metric for f is denoted by $\|f\|_{L_2(P)} = [\int_{\mathcal{Z}} f^2 dP]^{1/2}$, and the sup-metric of f is denoted by $\|f\|_\infty$. Positive constants that only depend on the class of data generating processes, not on the sample size nor the VC-dimension, are denoted by c_1, c_2, c_3, \dots . The universal constants are denoted by the capital letter C_1, C_2, C_3, \dots .

In what follows, we present lemmas that will be used in the proofs of Theorems 2.1 and 2.3. Lemmas A.2 and A.3 are classical inequalities whose proofs can be found, for instance, in Lugosi (2002).

LEMMA A.2—Hoeffding’s Lemma: *Let X be a random variable with $EX = 0$, $a \leq X \leq b$. Then, for $s > 0$,*

$$E(e^{sX}) \leq e^{s^2(b-a)^2/8}.$$

LEMMA A.3: *Let $\lambda > 0$, $n \geq 2$, and let Y_1, \dots, Y_n be real-valued random variables such that, for all $s > 0$ and $1 \leq i \leq n$, $E(e^{sY_i}) \leq e^{s^2\lambda^2/2}$ holds. Then,*

- (i) $E\left(\max_{i \leq n} Y_i\right) \leq \lambda\sqrt{2\ln n}$,
- (ii) $E\left(\max_{i \leq n} |Y_i|\right) \leq \lambda\sqrt{2\ln(2n)}$.

The next two lemmas give maximal inequalities that bound the mean of a supremum of centered empirical processes indexed by a VC-subgraph class of functions. The first maximal inequality (Lemma A.4) is standard in the empirical process literature, and it yields our Theorem 2.1 as a corollary. Though its proof can be found elsewhere (e.g., Dudley (1999), van der Vaart and Wellner (1996)), we present it here for the sake of completeness and for later reference in the proof of Lemma A.5. The second maximal inequality (Lemma A.5) concerns the class of functions whose diameter is constrained by the $L_2(P)$ -norm. Lemma A.5 will be used in the proof of Theorem 2.3. A lemma similar to our Lemma A.5 appears in Massart and Nédélec (2006, Lemma A.3).

LEMMA A.4: *Let \mathcal{F} be a class of uniformly bounded functions, that is, there exists $\bar{F} < \infty$ such that $\|f\|_\infty \leq \bar{F}$ for all $f \in \mathcal{F}$. Assume that \mathcal{F} is a VC-subgraph class with VC-dimension $v < \infty$. Then, there is a universal constant C_1 such that*

$$E_{P^n} \left[\sup_{f \in \mathcal{F}} |E_n(f) - E_P(f)| \right] \leq C_1 \bar{F} \sqrt{\frac{v}{n}}$$

holds for all $n \geq 1$.

PROOF: Introduce (Z'_1, \dots, Z'_n) , an independent copy of $(Z_1, \dots, Z_n) \sim P^n$. We denote the probability law of (Z'_1, \dots, Z'_n) by $P^{n'}$, its expectation by $E_{P^{n'}}(\cdot)$, and the sample average with respect to (Z'_1, \dots, Z'_n) by $E'_n(\cdot)$. Define i.i.d. Rademacher variables

$\sigma^n \equiv (\sigma_1, \dots, \sigma_n)$ such that $\Pr(\sigma_1 = -1) = \Pr(\sigma_1 = 1) = 1/2$ and they are independent of $Z_1, Z'_1, \dots, Z_n, Z'_n$. Then,

$$\begin{aligned}
E_{p^n} \left[\sup_{f \in \mathcal{F}} |E_n(f) - E_P(f)| \right] &= E_{p^n} \left[\sup_{f \in \mathcal{F}} |E_{p^{n'}}[E_n(f) - E'_n(f) | Z_1, \dots, Z_n]| \right] \\
&\leq E_{p^n} \left[\sup_{f \in \mathcal{F}} E_{p^{n'}} [|E_n(f) - E'_n(f)| | Z_1, \dots, Z_n] \right] \\
&(\because \text{Jensen's inequality}) \\
&\leq E_{p^n, p^{n'}} \left[\sup_{f \in \mathcal{F}} |E_n(f) - E'_n(f)| \right] \\
&= \frac{1}{n} E_{p^n, p^{n'}} \left\{ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n (f(Z_i) - f(Z'_i)) \right| \right\} \\
&= \frac{1}{n} E_{p^n, p^{n'}, \sigma^n} \left\{ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i (f(Z_i) - f(Z'_i)) \right| \right\} \tag{A.2} \\
&(\because f(Z_i) - f(Z'_i) \sim \sigma_i (f(Z_i) - f(Z'_i)) \text{ for all } i) \\
&\leq \frac{1}{n} E_{p^n, p^{n'}, \sigma^n} \left\{ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i f(Z_i) \right| + \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i f(Z'_i) \right| \right\} \\
&= \frac{2}{n} E_{p^n, \sigma^n} \left[\sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i f(Z_i) \right| \right] \\
&= \frac{2}{n} E_{p^n} \left\{ E_{\sigma^n} \left[\sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i f(Z_i) \right| \middle| Z_1, \dots, Z_n \right] \right\}.
\end{aligned}$$

Fix Z_1, \dots, Z_n , and define $\mathbf{f} \equiv (f(Z_1), \dots, f(Z_n)) = (f_1, \dots, f_n)$, which is a vector of length n collecting the value of $f \in \mathcal{F}$ evaluated at each of (Z_1, \dots, Z_n) . Let $\mathbf{F} \equiv \{\mathbf{f} : f \in \mathcal{F}\} \subset \mathbb{R}^n$, which is a bounded set in \mathbb{R}^n with radius \bar{F} , since \mathcal{F} is the set of uniformly bounded functions with $|f(\cdot)| \leq \bar{F}$. Introduce the Euclidean norm to \mathbf{F} ,

$$\rho(\mathbf{f}, \mathbf{f}') = \left(\frac{1}{n} \sum_{i=1}^n (f_i - f'_i)^2 \right)^{1/2}.$$

Let $\mathbf{f}^{(0)} = (0, \dots, 0)$, and $\mathbf{f}^* = (f_1^*, \dots, f_n^*)$ be a random element in \mathbf{F} maximizing $|\sum_{i=1}^n \sigma_i f_i|$. Let $B_0 = \{\mathbf{f}^{(0)}\}$ and construct $\{B_k : k = 1, \dots, \bar{K}\}$ a sequence of covers of \mathbf{F} , such that $B_k \subset \mathbf{F}$ is a minimal cover with radius $2^{-k} \bar{F}$ and $B_{\bar{K}} = \mathbf{F}$. Note that such $\bar{K} < \infty$ exists at given n and (Z_1, \dots, Z_n) . Define also $\{\mathbf{f}^{(k)} \in B_k : k = 1, \dots, \bar{K}\}$ be a random sequence such that $\mathbf{f}^{(k)} \in \arg \min_{\mathbf{f} \in B_k} \rho(\mathbf{f}, \mathbf{f}^*)$. Since B_k is a cover with radius $2^{-k} \bar{F}$, $\rho(\mathbf{f}^{(k)}, \mathbf{f}^*) \leq 2^{-k} \bar{F}$ holds. In addition, we have

$$\rho(\mathbf{f}^{(k-1)}, \mathbf{f}^{(k)}) \leq \rho(\mathbf{f}^{(k)}, \mathbf{f}^*) + \rho(\mathbf{f}^{(k-1)}, \mathbf{f}^*) \leq 3 \cdot 2^{-k} \bar{F}.$$

By a telescope sum,

$$\sum_{i=1}^n \sigma_i f_i^* = \sum_{i=1}^n \sigma_i f_i^{(0)} + \sum_{k=1}^{\bar{K}} \sum_{i=1}^n \sigma_i (f_i^{(k)} - f_i^{(k-1)}) = \sum_{k=1}^{\bar{K}} \sum_{i=1}^n \sigma_i (f_i^{(k)} - f_i^{(k-1)}).$$

We hence obtain

$$\begin{aligned} E_{\sigma^n} \left| \sum_{i=1}^n \sigma_i f_i^* \right| &\leq \sum_{k=1}^{\bar{K}} E_{\sigma^n} \left| \sum_{i=1}^n \sigma_i (f_i^{(k)} - f_i^{(k-1)}) \right| \\ &\leq \sum_{k=1}^{\bar{K}} E_{\sigma^n} \max_{\mathbf{f} \in B_k, \mathbf{g} \in B_{k-1}: \rho(\mathbf{f}, \mathbf{g}) \leq 3 \cdot 2^{-k} \bar{F}} \left| \sum_{i=1}^n \sigma_i (f_i - g_i) \right|. \end{aligned} \quad (\text{A.3})$$

We apply Lemma A.2 to obtain

$$\begin{aligned} E_{\sigma^n} (e^{s \sum_{i=1}^n \sigma_i (f_i - g_i)}) &= \prod_{i=1}^n E_{\sigma_i} [e^{s \sigma_i (f_i - g_i)}] \leq \prod_{i=1}^n e^{s^2 (f_i - g_i)^2 / 2} \\ &= \exp(s^2 n \rho^2(\mathbf{f}, \mathbf{g}) / 2) \\ &\leq \exp(s^2 n (3 \cdot 2^{-k} \bar{F})^2 / 2). \end{aligned}$$

An application of Lemma A.3(ii) with $\lambda = 3\sqrt{n} \cdot 2^{-k} \bar{F}$ and $n = |B_k| |B_{k-1}| \leq |B_k|^2$ then yields

$$\begin{aligned} E_{\sigma^n} \max_{\mathbf{f} \in B_k, \mathbf{g} \in B_{k-1}: \rho(\mathbf{f}, \mathbf{g}) \leq 3 \cdot 2^{-k} \bar{F}} \left| \sum_{i=1}^n \sigma_i (f_i - g_i) \right| &\leq 3\sqrt{n} \cdot 2^{-k} \bar{F} \sqrt{2 \ln 2 |B_k|^2} \\ &= 3\sqrt{n} \cdot 2^{-k} \bar{F} \sqrt{2 \ln 2 N(2^{-k} \bar{F}, \mathbf{F}, \rho)^2} \\ &= 6\sqrt{n} \cdot 2^{-k} \bar{F} \sqrt{\ln 2^{1/2} N(2^{-k} \bar{F}, \mathbf{F}, \rho)}, \end{aligned}$$

where $N(r, \mathbf{F}, \rho)$ is the covering number of \mathbf{F} with radius r in terms of norm ρ . Accordingly,

$$\begin{aligned} E_{\sigma^n} \left| \sum_{i=1}^n \sigma_i f_i^* \right| &\leq \sum_{k=1}^{\bar{K}} 6\sqrt{n} \cdot 2^{-k} \bar{F} \sqrt{\ln 2^{1/2} N(2^{-k} \bar{F}, \mathbf{F}, \rho)} \\ &\leq 12\sqrt{n} \sum_{k=1}^{\infty} 2^{-(k+1)} \bar{F} \sqrt{\ln 2^{1/2} N(2^{-k} \bar{F}, \mathbf{F}, \rho)} \\ &\leq 12\sqrt{n} \int_0^1 \bar{F} \sqrt{\ln 2^{1/2} N(\varepsilon \bar{F}, \mathbf{F}, \rho)} d\varepsilon, \end{aligned} \quad (\text{A.4})$$

where the last line follows from the fact that $N(\varepsilon \bar{F}, \mathbf{F}, \rho)$ is decreasing in ε .

To bound (A.4) from above, we apply a uniform entropy bound for the covering number. In Theorem 2.6.7 of van der Vaart and Wellner (1996), by setting $r = 2$ and Q at the empirical probability measure of (Z_1, \dots, Z_n) , we have,

$$N(\varepsilon \bar{F}, \mathbf{F}, \rho) \leq K(v+1)(16e)^{(v+1)} \left(\frac{1}{\varepsilon}\right)^{2v}, \quad (\text{A.5})$$

where $K > 0$ is a universal constant. Plugging this into (A.4) leads to

$$\begin{aligned} E_\sigma \left| \sum_{i=1}^n \sigma_i f_i^* \right| &\leq 12\bar{F}\sqrt{n} \int_0^1 \sqrt{\ln(2^{1/2}K) + \ln(v+1) + (v+1)\ln(16e) - 2v\ln \varepsilon} d\varepsilon \\ &\leq 12\bar{F}\sqrt{nv} \int_0^1 \sqrt{\ln(2^{1/2}K) + \ln 2 + 2\ln(16e) - 2\ln \varepsilon} d\varepsilon \\ &= C'\bar{F}\sqrt{nv}, \end{aligned} \quad (\text{A.6})$$

where $C' = 12 \int_0^1 \sqrt{\ln(2^{1/2}K) + \ln 2 + 2\ln(16e) - 2\ln \varepsilon} d\varepsilon < \infty$. Combining (A.6) with (A.2) and setting $C_1 = 2C'$ lead to the conclusion. Q.E.D.

LEMMA A.5: *Let \mathcal{F} be a class of uniformly bounded functions with $\|f\|_\infty \leq \bar{F} < \infty$ for all $f \in \mathcal{F}$. Assume that \mathcal{F} is a VC-subgraph class with VC-dimension $v < \infty$. Assume further that $\sup_{f \in \mathcal{F}} \|f\|_{L_2(P)} \leq \delta$. Then, there exists a positive universal constant C_2 such that*

$$E_{P^n} \left[\sup_{f \in \mathcal{F}} (E_n(f) - E_P(f)) \right] \leq C_2 \delta \sqrt{\frac{v}{n}}$$

holds for all $n \geq C_1^2 \bar{F}^2 v / \delta^2$, where C_1 is the universal constant defined in Lemma A.4.

PROOF: By the same symmetrization argument and the same use of Rademacher variables as in the proof of Lemma A.4, we have

$$E_{P^n} \left[\sup_{f \in \mathcal{F}} (E_n(f) - E_P(f)) \right] \leq \frac{2}{n} E_{P^n} \left\{ E_{\sigma^n} \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^n \sigma_i f(Z_i) \mid Z_1, \dots, Z_n \right] \right\}. \quad (\text{A.7})$$

Fix the values of Z_1, \dots, Z_n , and define $\mathbf{f}, \mathbf{f}^{(0)}, \mathbf{F}$, and norm $\rho(\mathbf{f}, \mathbf{f}')$ as in the proof of Lemma A.4. Let \mathbf{f}^* be a maximizer of $\sum_{i=1}^n \sigma_i f(Z_i)$ in \mathbf{F} and let $\delta_n = \sup_{\mathbf{f} \in \mathbf{F}} \rho(\mathbf{f}^{(0)}, \mathbf{f}) \leq \bar{F}$. Let $B_0 = \{\mathbf{f}^{(0)}\}$ and construct $\{B_k : k = 1, \dots, \bar{K}\}$ a sequence of covers of \mathbf{F} , such that $B_k \subset \mathbf{F}$ is a minimal cover with radius $2^{-k} \delta_n$ and $B_{\bar{K}} = \mathbf{F}$. We define $\{\mathbf{f}^{(k)} \in B_k : k = 1, \dots, \bar{K}\}$ to be a random sequence such that $\mathbf{f}^{(k)} \in \arg \min_{\mathbf{f} \in B_k} \rho(\mathbf{f}, \mathbf{f}^*)$. By applying the chaining argument in the proof of Lemma A.4, Lemma A.3(i), and the uniform bound of the covering number (A.5), we obtain

$$E_\sigma \sum_{i=1}^n \sigma_i f_i^* \leq 12\sqrt{n} \int_0^1 \delta_n \sqrt{\log N(\varepsilon \delta_n, \mathbf{F}, \rho)} d\varepsilon \leq 2^{-1} C_1 \delta_n \sqrt{nv}$$

for the universal constant C_1 defined in the proof of Lemma A.4. Hence, from (A.7), we have

$$\begin{aligned} E_{P^n} \left[\sup_{f \in \mathcal{F}} (E_n(f) - E_P(f)) \right] &\leq C_1 \sqrt{\frac{v}{n}} E_{P^n}(\delta_n) = C_1 \sqrt{\frac{v}{n}} E_{P^n} \left(\left[\sup_{f \in \mathcal{F}} E_n(f^2) \right]^{1/2} \right) \\ &\leq C_1 \sqrt{\frac{v}{n}} \left[E_{P^n} \left(\sup_{f \in \mathcal{F}} E_n(f^2) \right) \right]^{1/2}. \end{aligned} \quad (\text{A.8})$$

Note that $E_n(f^2)$ is bounded by

$$\begin{aligned} E_n(f^2) &= E_n(f^2 - E_P(f^2)) + E_P(f^2) \\ &= E_n[(f - \|f\|_{L_2(P)})(f + \|f\|_{L_2(P)})] + \|f\|_{L_2(P)}^2 \\ &\leq 2\bar{F} E_n[f - \|f\|_{L_2(P)}] + \|f\|_{L_2(P)}^2 \\ &\leq 2\bar{F} E_n[f - E_P(f)] + \|f\|_{L_2(P)}^2 \\ &\quad (\because \|f\|_{L_2(P)} \geq E_P(f) \text{ by the Cauchy-Schwarz inequality}). \end{aligned}$$

Combining this inequality with (A.8) yields

$$E_{P^n} \left[\sup_{f \in \mathcal{F}} (E_n(f) - E_P(f)) \right] \leq C_1 \sqrt{\frac{v}{n}} \sqrt{2\bar{F} E_{P^n} \left[\sup_{f \in \mathcal{F}} (E_n(f) - E_P(f)) \right]} + \delta^2.$$

Solving this inequality for $E_{P^n}[\sup_{f \in \mathcal{F}} (E_n(f) - E_P(f))]$ leads to

$$E_{P^n} \left[\sup_{f \in \mathcal{F}} (E_n(f) - E_P(f)) \right] \leq \bar{F} C_1^2 \sqrt{\frac{v}{n}} \left(\sqrt{\frac{v}{n}} + \sqrt{\frac{v}{n} + \frac{\delta^2}{\bar{F}^2 C_1^2}} \right).$$

For $\frac{v}{n} \leq \frac{\delta^2}{\bar{F}^2 C_1^2}$, that is, $n \geq \frac{C_1^2 \bar{F}^2 v}{\delta^2}$, the upper bound can be further bounded by $(1 + \sqrt{2}) \times C_1 \delta \sqrt{\frac{v}{n}}$, so the conclusion of the lemma follows with $C_2 = (1 + \sqrt{2}) C_1$. Q.E.D.

A.2. Proofs of Theorems 2.1 and 2.2

PROOF OF THEOREM 2.1: Define

$$f(Z_i; G) = \left[\frac{Y_i D_i}{e(X_i)} \cdot 1\{X_i \in G\} + \frac{Y_i(1 - D_i)}{1 - e(X_i)} \cdot 1\{X_i \notin G\} \right],$$

and the class of functions on \mathcal{Z}

$$\mathcal{F} = \{f(\cdot; G) : G \in \mathcal{G}\}.$$

With these notations, we can express inequality (2.3) in the main text as

$$W_G^* - W(\hat{G}_{\text{EWM}}) \leq 2 \sup_{f \in \mathcal{F}} |E_n(f) - E_P(f)|. \quad (\text{A.9})$$

Note that Assumption 2.1 (BO) and (SO) imply that \mathcal{F} has uniform envelope $\bar{F} = M/(2\kappa)$. Also, by Assumption 2.1 (VC) and Lemma A.1, \mathcal{F} is a VC-subgraph class of functions with VC-dimension at most v . We apply Lemma A.4 to (A.9) to obtain

$$E_{P^n}[W_{\mathcal{G}}^* - W(\hat{G}_{\text{EWM}})] \leq C_1 \frac{M}{\kappa} \sqrt{\frac{v}{n}}.$$

Since this upper bound does not depend on $P \in \mathcal{P}(M, \kappa)$, the upper bound is uniform over $\mathcal{P}(M, \kappa)$. *Q.E.D.*

PROOF OF THEOREM 2.2: In obtaining the rate lower bound, we normalize the support of outcomes to $Y_{1,i}, Y_{0,i} \in [-\frac{1}{2}, \frac{1}{2}]$. That is, we focus on bounding $\sup_{P \in \mathcal{P}(1, \kappa)} E_{P^n}[W_{\mathcal{G}}^* - W(G_n)]$. The lower bound of the original welfare loss $\sup_{P \in \mathcal{P}(M, \kappa)} E_{P^n}[W_{\mathcal{G}}^* - W(G_n)]$ is obtained by multiplying by M the lower bound of $\sup_{P \in \mathcal{P}(1, \kappa)} E_{P^n}[W_{\mathcal{G}}^* - W(G_n)]$.

We consider a suitable subclass $\mathcal{P}^* \subset \mathcal{P}(1, \kappa)$, for which the worst-case welfare loss can be bounded from below by a distribution-free term that converges at rate $n^{-1/2}$. The construction of \mathcal{P}^* proceeds as follows. First, let $x_1, \dots, x_v \in \mathcal{X}$ be v points that are shattered by \mathcal{G} . We constrain P_X (the marginal distribution of X) to being supported only on (x_1, \dots, x_v) . We put the equal mass $1/v$ at x_i , $i \leq v$. Thus-constructed marginal distribution of X is common in \mathcal{P}^* . Let the distribution of treatment indicator D be independent of (Y_1, Y_0, X) , and D follows the Bernoulli distribution with $\Pr(D = 1) = 1/2$. Let $\mathbf{b} = (b_1, \dots, b_v) \in \{0, 1\}^v$ be a bit vector used to index a member of \mathcal{P}^* , that is, \mathcal{P}^* consists of a finite number of DGPs. For each $j = 1, \dots, v$, and depending on \mathbf{b} , construct the following conditional distribution of Y_1 given $X = x_j$: if $b_j = 1$,

$$Y_1 = \begin{cases} \frac{1}{2} & \text{with prob. } \frac{1}{2} + \gamma, \\ -\frac{1}{2} & \text{with prob. } \frac{1}{2} - \gamma, \end{cases} \quad (\text{A.10})$$

and, if $b_j = 0$,

$$Y_1 = \begin{cases} \frac{1}{2} & \text{with prob. } \frac{1}{2} - \gamma, \\ -\frac{1}{2} & \text{with prob. } \frac{1}{2} + \gamma, \end{cases} \quad (\text{A.11})$$

where $\gamma \in [0, \frac{1}{2}]$ is chosen properly in a later step of the proof. As for Y_0 's conditional distribution, we consider the degenerate distribution at $Y_0 = 0$ at every $X = x_j$, $j = 1, \dots, v$. That is, when $b_j = 1$, $\tau(x_j) = \gamma$, and when $b_j = 0$, $\tau(x_j) = -\gamma$. For each $\mathbf{b} \in \{0, 1\}^v$, $P_{\mathbf{b}} \in \mathcal{P}(1, \kappa)$ clearly holds. We accordingly define a subclass of $\mathcal{P}(1, \kappa)$ by $\mathcal{P}^* = \{P_{\mathbf{b}} : \mathbf{b} \in \{0, 1\}^v\}$.

With knowledge of $P_{\mathbf{b}} \in \mathcal{P}^*$, the optimal treatment assignment rule is

$$G_{\mathbf{b}}^* = \{x_j : b_j = 1, j \leq v\},$$

which is feasible $G_{\mathbf{b}}^* \in \mathcal{G}$ by the construction of the support points of X . The maximized social welfare is

$$W(G_{\mathbf{b}}^*) = v^{-1} \gamma \left(\sum_{j=1}^v b_j \right).$$

Let \hat{G} be an arbitrary treatment choice rule depending on sample (Z_1, \dots, Z_n) , and $\hat{\mathbf{b}} \in \{0, 1\}^v$ be a binary vector whose j th element is $\hat{b}_j = 1\{x_j \in \hat{G}\}$. Consider $\pi(\mathbf{b})$ a prior distribution for \mathbf{b} such that b_1, \dots, b_v are i.i.d. and $b_1 \sim \text{Ber}(1/2)$. The welfare loss satisfies the following inequalities:

$$\begin{aligned}
\sup_{P \in \mathcal{P}(1, \kappa)} E_{P^n} [W_G^* - W(\hat{G})] &\geq \sup_{P_{\mathbf{b}} \in \mathcal{P}^*} E_{P_{\mathbf{b}}^n} [W(G_{\mathbf{b}}^*) - W(\hat{G})] \\
&\geq \int_{\mathbf{b}} E_{P_{\mathbf{b}}^n} [W(G_{\mathbf{b}}^*) - W(\hat{G})] d\pi(\mathbf{b}) \\
&= \gamma \int_{\mathbf{b}} E_{P_{\mathbf{b}}^n} [P_X(G_{\mathbf{b}}^* \Delta \hat{G})] d\pi(\mathbf{b}) \\
&= \gamma \int_{\mathbf{b}} \int_{Z_1, \dots, Z_n} P_X(\{b(X) \neq \hat{b}(X)\}) dP_{\mathbf{b}}^n(Z_1, \dots, Z_n) d\pi(\mathbf{b}) \\
&\geq \inf_{\hat{G}} \gamma \int_{\mathbf{b}} \int_{Z_1, \dots, Z_n} P_X(\{b(X) \neq \hat{b}(X)\}) dP_{\mathbf{b}}^n(Z_1, \dots, Z_n) d\pi(\mathbf{b}),
\end{aligned}$$

where $b(X)$ and $\hat{b}(X)$ are elements of \mathbf{b} and $\hat{\mathbf{b}}$, respectively, such that $b(x_j) = b_j$ and $\hat{b}(x_j) = \hat{b}_j$. Note that the infimum over assignment rules \hat{G} can be seen as the minimization problem of the Bayes risk with the loss function corresponding to the classification error for predicting binary random variable $b(X)$. Hence, a minimizer of the Bayes risk is attained by the Bayes classifier,

$$\hat{G}^* = \left\{ x_j : \pi(b_j = 1 | Z_1, \dots, Z_n) \geq \frac{1}{2}, j \leq v \right\},$$

where $\pi(b_j = 1 | Z_1, \dots, Z_n)$ is the posterior probability for $b_j = 1$. The minimized Bayes risk is given by

$$\begin{aligned}
&\gamma \int_{Z_1, \dots, Z_n} E_X [\min\{\pi(b(X) = 1 | Z_1, \dots, Z_n), 1 - \pi(b(X) = 1 | Z_1, \dots, Z_n)\}] d\tilde{P}^n \\
&= v^{-1} \gamma \int_{Z_1, \dots, Z_n} \sum_{j=1}^v [\min\{\pi(b_j = 1 | Z_1, \dots, Z_n), 1 - \pi(b_j = 1 | Z_1, \dots, Z_n)\}] d\tilde{P}^n,
\end{aligned} \tag{A.12}$$

where \tilde{P}^n is the marginal likelihood of $\{(Y_{1,i}, Y_{0,i}, D_i, X_i) : i = 1, \dots, n\}$ with prior $\pi(\mathbf{b})$. For each $j = 1, \dots, (v)$, let

$$\begin{aligned}
k_j^+ &= \# \left\{ i : X_i = x_j, Y_i D_i = \frac{1}{2} \right\}, \\
k_j^- &= \# \left\{ i : X_i = x_j, Y_i D_i = -\frac{1}{2} \right\}.
\end{aligned}$$

The posterior for $b_j = 1$ can be written as

$$\pi(b_j = 1 | Z_1, \dots, Z_n) = \begin{cases} \frac{1}{2} & \text{if } \#\{i : X_i = x_j, D_i = 1\} = 0, \\ \frac{\left(\frac{1}{2} + \gamma\right)^{k_j^+} \left(\frac{1}{2} - \gamma\right)^{k_j^-}}{\left(\frac{1}{2} + \gamma\right)^{k_j^+} \left(\frac{1}{2} - \gamma\right)^{k_j^-} + \left(\frac{1}{2} + \gamma\right)^{k_j^-} \left(\frac{1}{2} - \gamma\right)^{k_j^+}} & \text{otherwise.} \end{cases}$$

Hence,

$$\begin{aligned} & \min\{\pi(b_j = 1 | Z_1, \dots, Z_n), 1 - \pi(b_j = 1 | Z_1, \dots, Z_n)\} \\ &= \frac{\min\left\{\left(\frac{1}{2} + \gamma\right)^{k_j^+} \left(\frac{1}{2} - \gamma\right)^{k_j^-}, \left(\frac{1}{2} + \gamma\right)^{k_j^-} \left(\frac{1}{2} - \gamma\right)^{k_j^+}\right\}}{\left(\frac{1}{2} + \gamma\right)^{k_j^+} \left(\frac{1}{2} - \gamma\right)^{k_j^-} + \left(\frac{1}{2} + \gamma\right)^{k_j^-} \left(\frac{1}{2} - \gamma\right)^{k_j^+}} \\ &= \frac{\min\left\{1, \left(\frac{\frac{1}{2} + \gamma}{\frac{1}{2} - \gamma}\right)^{k_j^+ - k_j^-}\right\}}{1 + \left(\frac{\frac{1}{2} + \gamma}{\frac{1}{2} - \gamma}\right)^{k_j^+ - k_j^-}} \\ &= \frac{1}{1 + a^{|k_j^+ - k_j^-|}}, \quad \text{where } a = \frac{1 + 2\gamma}{1 - 2\gamma} > 1. \end{aligned} \tag{A.13}$$

Since $k_j^+ - k_j^- = \sum_{i: X_i = x_j} 2Y_i D_i$, plugging (A.13) into (A.12) yields

$$\begin{aligned} v^{-1} \gamma \sum_{j=1}^v E_{\tilde{p}^n} \left[\frac{1}{1 + a^{|\sum_{i: X_i = x_j} 2Y_i D_i|}} \right] &\geq \frac{\gamma}{2v} \sum_{j=1}^v E_{\tilde{p}^n} \left[\frac{1}{a^{|\sum_{i: X_i = x_j} 2Y_i D_i|}} \right] \\ &\geq \frac{\gamma}{2v} \sum_{j=1}^v a^{-E_{\tilde{p}^n} |\sum_{i: X_i = x_j} 2Y_i D_i|}, \end{aligned}$$

where $E_{\tilde{p}^n}(\cdot)$ is the expectation with respect to the marginal likelihood of $\{(Y_{1,i}, Y_{0,i}, D_i, X_i), i = 1, \dots, n\}$. The second line follows by $a > 1$, and the third line follows by Jensen's inequality. Given our prior specification for \mathbf{b} , the marginal distribution of $Y_{1,i}$ is $\Pr(Y_{1,i} = 1/2) = \Pr(Y_{1,i} = -1/2) = 1/2$, so

$$E_{\tilde{p}^n} \left| \sum_{i: X_i = x_j} 2Y_i D_i \right| = E_{\tilde{p}^n} \left| \sum_{i=1: X_i = x_j, D_i = 1} 2Y_{1,i} \right| = \sum_{k=0}^n \binom{n}{k} \left(\frac{1}{2v}\right)^k \left(1 - \frac{1}{2v}\right)^{n-k} E \left| B\left(k, \frac{1}{2}\right) - \frac{k}{2} \right|$$

holds, where $B(k, \frac{1}{2})$ is the binomial random variable with parameters k and $\frac{1}{2}$. By noting

$$\begin{aligned} E \left| B\left(k, \frac{1}{2}\right) - \frac{k}{2} \right| &\leq \sqrt{E\left(B\left(k, \frac{1}{2}\right) - \frac{k}{2}\right)^2} \quad (\because \text{Cauchy-Schwarz inequality}) \\ &= \sqrt{\frac{k}{4}}, \end{aligned}$$

we obtain

$$\begin{aligned} E_{\tilde{p}^n} \left| \sum_{i: X_i = x_j} 2Y_i D_i \right| &\leq \sum_{k=0}^n \binom{n}{k} \left(\frac{1}{2v}\right)^k \left(1 - \frac{1}{2v}\right)^{n-k} \sqrt{\frac{k}{4}} = E \sqrt{\frac{B\left(n, \frac{1}{2v}\right)}{4}} \\ &\leq \sqrt{\frac{n}{8v}} \quad (\because \text{Jensen's inequality}). \end{aligned}$$

Hence, the Bayes risk is bounded from below by

$$\begin{aligned} \frac{\gamma}{2} a^{-\sqrt{\frac{n}{8v}}} &\geq \frac{\gamma}{2} \exp\left\{-\left(a-1\right)\sqrt{\frac{n}{8v}}\right\} \quad (\because 1+x \leq e^x \forall x) \\ &= \frac{\gamma}{2} \exp\left\{-\frac{4\gamma}{1-2\gamma}\sqrt{\frac{n}{8v}}\right\}. \end{aligned} \tag{A.14}$$

This lower bound of the Bayes risk has the slowest convergence rate when γ is set to be proportional to $n^{-1/2}$. Specifically, let $\gamma = \sqrt{\frac{v}{n}}$. Then, we have

$$\frac{\gamma}{2} \exp\left\{-\frac{4\gamma}{1-2\gamma}\sqrt{\frac{n}{8v}}\right\} = \frac{1}{2}\sqrt{\frac{v}{n}} \exp\left\{-\frac{\sqrt{2}}{1-2\gamma}\right\} \geq \frac{1}{2}\sqrt{\frac{v}{n}} \exp\{-2\sqrt{2}\} \quad \text{if } 1-2\gamma \geq \frac{1}{2}.$$

The condition $1-2\gamma \geq \frac{1}{2}$ is equivalent to $n \geq 16v$. Multiplying M to this lower bound completes the proof. *Q.E.D.*

A.3. Proofs of Theorems 2.3 and 2.4

The next lemma is the concentration inequality of [Bousquet \(2002\)](#).

LEMMA A.6: *Let \mathcal{F} be a countable family of measurable functions, such that $\sup_{f \in \mathcal{F}} E_P(f^2) \leq \delta^2$ and $\sup_{f \in \mathcal{F}} \|f\|_\infty \leq \bar{F}$ for some constants δ and \bar{F} . Let $S = \sup_{f \in \mathcal{F}} (E_n(f) - E_P(f))$. Then, for every positive t ,*

$$P^n \left(S - E_{P^n}(S) \geq \sqrt{\frac{2[\delta^2 + 4\bar{F}E_{P^n}(S)]t}{n}} + \frac{2\bar{F}t}{3n} \right) \leq \exp(-t).$$

In proving Theorem 2.3, it is convenient to work with the *normalized welfare difference*,

$$d(G, G') \equiv \frac{\kappa}{M} [W(G) - W(G')],$$

and its sample analogue

$$d_n(G, G') \equiv \frac{\kappa}{M} [W_n(G) - W_n(G')]. \quad (\text{A.15})$$

By Assumption 2.1 (BO) and (SO), both $d(G, G')$ and $d_n(G, G')$ are bounded in $[-1, 1]$, and the normalized welfare difference relates to the original welfare loss of decision set G as

$$d(G_{\text{FB}}^*, G) = \frac{\kappa}{M} [W(G_{\text{FB}}^*) - W(G)] \in [0, 1]. \quad (\text{A.16})$$

Hence, the welfare loss upper bound of \hat{G}_{EWM} can be obtained by multiplying M/κ by the upper bound of $d(G_{\text{FB}}^*, \hat{G}_{\text{EWM}})$.

Note that $d(G_{\text{FB}}^*, G)$ can be bounded from above by $P_X(G_{\text{FB}}^* \Delta G)$, since

$$d(G_{\text{FB}}^*, G) = \frac{\kappa}{M} \int_{G_{\text{FB}}^* \Delta G} |\tau(X)| dP_X \leq \kappa P_X(G_{\text{FB}}^* \Delta G) \leq P_X(G_{\text{FB}}^* \Delta G). \quad (\text{A.17})$$

On the other hand, with Assumption 2.2 (MA) imposed, $P_X(G_{\text{FB}}^* \Delta G)$ can be bounded from above by a function of $d(G_{\text{FB}}^*, G)$, as the next lemma shows. We borrow this lemma from [Tsybakov \(2004\)](#).

LEMMA A.7: *Suppose Assumption 2.2 (MA) holds with margin coefficient $\alpha \in (0, \infty)$. Then*

$$P_X(G_{\text{FB}}^* \Delta G) \leq c_1(M, \kappa, \eta, \alpha) d(G_{\text{FB}}^*, G)^{\frac{\alpha}{1+\alpha}}$$

holds for all $G \in \mathcal{G}$, where $c_1(M, \kappa, \eta, \alpha) = \left(\frac{M}{\kappa\eta\alpha}\right)^{\frac{\alpha}{1+\alpha}} (1 + \alpha)$.

PROOF: Let $A = \{x : |\tau(x)| > t\}$ and consider the following inequalities:

$$\begin{aligned} W(G_{\text{FB}}^*) - W(G) &= \int_{G_{\text{FB}}^* \Delta G} |\tau(X)| dP_X \geq \int_{G_{\text{FB}}^* \Delta G} |\tau(X)| 1\{X \in A\} dP_X \\ &\geq t P_X((G_{\text{FB}}^* \Delta G) \cap A) \geq t [P_X(G_{\text{FB}}^* \Delta G) - P_X(A^c)] \\ &\geq t \left[P_X(G_{\text{FB}}^* \Delta G) - \left(\frac{t}{\eta}\right)^\alpha \right], \end{aligned}$$

where the final line uses the margin condition. The right-hand side is maximized at $t = \eta(1 + \alpha)^{-\frac{1}{\alpha}} [P_X(G_{\text{FB}}^* \Delta G)]^{\frac{1}{\alpha}} \leq \eta$, so it holds that

$$W(G_{\text{FB}}^*) - W(G) \geq \eta \alpha \left(\frac{1}{1 + \alpha}\right)^{\frac{1+\alpha}{\alpha}} [P_X(G_{\text{FB}}^* \Delta G)]^{\frac{1+\alpha}{\alpha}}.$$

This, in turn, implies

$$P_X(G_{\text{FB}}^* \Delta G) \leq \left(\frac{M}{\kappa\eta\alpha}\right)^{\frac{\alpha}{1+\alpha}} (1 + \alpha) d(G_{\text{FB}}^*, G)^{\frac{\alpha}{1+\alpha}}. \quad \text{Q.E.D.}$$

PROOF OF THEOREM 2.3: Let $a = \sqrt{kt}\varepsilon_n$ with $k \geq 1$, $t \geq 1$, and $\varepsilon_n > 0$, where $t \geq 1$ is arbitrary, k is a constant that we choose later, and ε_n is a sequence indexed by sample size

n whose proper choice will be discussed in a later step. The normalized welfare loss can be bounded by

$$d(G_{\text{FB}}^*, \hat{G}_{\text{EWM}}) \leq d(G_{\text{FB}}^*, \hat{G}_{\text{EWM}}) - d_n(G_{\text{FB}}^*, \hat{G}_{\text{EWM}}),$$

as $d_n(G_{\text{FB}}^*, \hat{G}_{\text{EWM}}) \leq 0$ by Assumption 2.2 (FB). Define a class of functions induced by $G \in \mathcal{G}$:

$$\begin{aligned} \mathcal{H} &\equiv \{h(Z_i; G) : G \in \mathcal{G}\}, \\ h(Z_i; G) &\equiv \frac{\kappa}{M} \left(\frac{Y_i D_i}{e(X_i)} - \frac{Y_i(1-D_i)}{1-e(X_i)} \right) [1\{X_i \in G\} - 1\{X_i \in G_{\text{FB}}^*\}]. \end{aligned}$$

By Assumption 2.1 (VC) and Lemma A.1, \mathcal{H} is a VC-subgraph class with VC-dimension at most $v < \infty$ with envelope $\bar{H} = 1$. Using $h(Z_i; G)$, we can write $d(G_{\text{FB}}^*, G) = -E_P(h(Z_i; G))$. Since $d(G_{\text{FB}}^*, G) \geq 0$ for all $G \in \mathcal{G}$, it holds that $-E_P(h) \geq 0$ for all $h \in \mathcal{H}$. Since we have

$$d(G_{\text{FB}}^*, \hat{G}_{\text{EWM}}) - d_n(G_{\text{FB}}^*, \hat{G}_{\text{EWM}}) = E_n(h(Z_i; \hat{G}_{\text{EWM}})) - E_P(h(Z_i; \hat{G}_{\text{EWM}}))$$

and $d_n(G_{\text{FB}}^*, \hat{G}_{\text{EWM}}) \leq 0$, the normalized welfare loss can be bounded by

$$\begin{aligned} d(G_{\text{FB}}^*, \hat{G}_{\text{EWM}}) &\leq E_n(h(Z_i; \hat{G}_{\text{EWM}})) - E_P(h(Z_i; \hat{G}_{\text{EWM}})) \\ &\leq V_a [d(G_{\text{FB}}^*, \hat{G}_{\text{EWM}}) + a^2], \end{aligned}$$

where

$$V_a = \sup_{h \in \mathcal{H}} \left\{ \frac{E_n(h) - E_P(h)}{-E_P(h) + a^2} \right\} = \sup_{h \in \mathcal{H}} \left\{ E_n \left(\frac{h}{-E_P(h) + a^2} \right) - E_P \left(\frac{h}{-E_P(h) + a^2} \right) \right\}.$$

On event $V_a < \frac{1}{2}$, $d(G_{\text{FB}}^*, \hat{G}_{\text{EWM}}) \leq a^2$ holds, so this implies

$$P^n(d(G_{\text{FB}}^*, \hat{G}_{\text{EWM}}) \geq a^2) \leq P^n \left(V_a \geq \frac{1}{2} \right). \quad (\text{A.18})$$

In what follows, our aim is to construct an exponential inequality for $P^n(V_a \geq \frac{1}{2})$ involving only t , and we make use of such exponential tail bound to bound $E_{P^n}(d(G_{\text{FB}}^*, \hat{G}_{\text{EWM}}))$.

To apply Bousquet's inequality (Lemma A.6) to V_a , note first that

$$\begin{aligned} E_P \left(\left(\frac{h}{-E_P(h) + a^2} \right)^2 \right) &\leq \frac{P_X(G_{\text{FB}}^* \Delta G)}{(-E_P(h) + a^2)^2} \leq c_1 \frac{[-E_P(h)]^{\frac{\alpha}{1+\alpha}}}{(-E_P(h) + a^2)^2} \\ &(\because \text{by Lemma A.7 and } d(G_{\text{FB}}^*, G) = -E_P(h(Z_i; G))) \\ &\leq c_1 \sup_{\varepsilon \geq 0} \frac{\varepsilon^{\frac{2\alpha}{1+\alpha}}}{(\varepsilon^2 + a^2)^2} \leq c_1 \frac{1}{a^2} \sup_{\varepsilon \geq 0} \frac{\varepsilon^{\frac{2\alpha}{1+\alpha}}}{\varepsilon^2 + a^2} \leq c_1 \frac{1}{a^2} \sup_{\varepsilon \geq 0} \left(\frac{\varepsilon^{\frac{\alpha}{1+\alpha}}}{\varepsilon \vee a} \right)^2 \\ &\leq c_1 \frac{1}{a^4} a^{\frac{2\alpha}{1+\alpha}}, \end{aligned}$$

where c_1 is a constant that depends only on $(M, \kappa, \eta, \alpha)$ as defined in Lemma A.7. We, on the other hand, have

$$\sup_{h \in \mathcal{H}} \left| \sup_Z \frac{h}{-E_P(h) + a^2} \right| \leq \frac{1}{a^2}.$$

Hence, Lemma A.6 gives, with probability larger than $1 - \exp(-t)$,

$$V_a \leq E_{P^n}(V_a) + \sqrt{\frac{2[c_1 a^{\frac{2\alpha}{1+\alpha}-2} + 4E_{P^n}(V_a)]t}{na^2}} + \frac{2t}{3na^2}. \quad (\text{A.19})$$

Next, we derive an upper bound of $E_{P^n}(V_a)$ by applying the maximal inequality of Lemma A.5. Let $r > 1$ be arbitrary and consider partitioning \mathcal{H} by $\mathcal{H}_0, \mathcal{H}_1, \dots$, where $\mathcal{H}_0 = \{h \in \mathcal{H} : -E_P(h) \leq a^2\}$ and $\mathcal{H}_j = \{h \in \mathcal{H} : r^{2(j-1)}a^2 < -E_P(h) \leq r^{2j}a^2\}$, $j = 1, 2, \dots$. Then,

$$\begin{aligned} V_a &\leq \sup_{h \in \mathcal{H}_0} \left\{ \frac{E_n(h) - E_P(h)}{-E_P(h) + a^2} \right\} + \sum_{j \geq 1} \sup_{h \in \mathcal{H}_j} \left\{ \frac{E_n(h) - E_P(h)}{-E_P(h) + a^2} \right\} \\ &\leq \frac{1}{a^2} \left[\sup_{h \in \mathcal{H}_0} (E_n(h) - E_P(h)) + \sum_{j \geq 1} (1 + r^{2(j-1)})^{-1} \sup_{h \in \mathcal{H}_j} (E_n(h) - E_P(h)) \right] \\ &\leq \frac{1}{a^2} \left[\sup_{-E_P(h) \leq a^2} (E_n(h) - E_P(h)) + \sum_{j \geq 1} (1 + r^{2(j-1)})^{-1} \sup_{-E_P(h) \leq r^{2j}a^2} (E_n(h) - E_P(h)) \right]. \end{aligned} \quad (\text{A.20})$$

Since it holds that $\|h\|_{L_2(P)}^2 \leq P_X(G_{\text{FB}}^* \Delta G) \leq c_1(M, \kappa, \eta, \alpha)[-E_P(h)]^{\frac{\alpha}{1+\alpha}}$, where the latter inequality follows from Lemma A.7, $-E_P(h) \leq r^{2j}a^2$ implies $\|h\|_{L_2(P)} \leq c_1^{1/2} r^{\frac{\alpha}{1+\alpha}j} a^{\frac{\alpha}{1+\alpha}}$. Hence, (A.20) can be further bounded by

$$\begin{aligned} V_a &\leq \frac{1}{a^2} \left[\sup_{\|h\|_{L_2(P)} \leq c_1^{1/2} a^{\frac{\alpha}{1+\alpha}}} (E_n(h) - E_P(h)) \right. \\ &\quad \left. + \sum_{j \geq 1} (1 + r^{2(j-1)})^{-1} \sup_{\|h\|_{L_2(P)} \leq c_1^{1/2} r^{\frac{\alpha}{1+\alpha}j} a^{\frac{\alpha}{1+\alpha}}} (E_n(h) - E_P(h)) \right]. \end{aligned}$$

We apply Lemma A.5 to each supremum term, and obtain

$$E_{P^n}(V_a) \leq C_2 \frac{c_1^{\frac{1}{2}}}{a^2} \sqrt{\frac{v}{n}} a^{\frac{\alpha}{1+\alpha}} \sum_{j \geq 0} \frac{r^{\frac{\alpha}{1+\alpha}j}}{1 + r^{2(j-1)}} \leq C_2 c_1^{\frac{1}{2}} \sqrt{\frac{v}{n}} a^{\frac{\alpha}{1+\alpha}-2} \left(\frac{r^2}{1 - r^{-\frac{2+\alpha}{1+\alpha}}} \right) \leq c_2 \sqrt{\frac{v}{n}} a^{\frac{\alpha}{1+\alpha}-2}$$

for

$$n \geq \frac{C_1 v}{c_1 a^{\frac{2\alpha}{1+\alpha}}} \iff a \geq \left(\frac{C_1}{c_1} \right)^{\frac{1+\alpha}{2\alpha}} \left(\frac{v}{n} \right)^{\frac{1+\alpha}{2\alpha}}, \quad (\text{A.21})$$

where C_1 and C_2 are universal constants defined in Lemmas A.4 and A.5, and $c_2 = C_2 c_1^{\frac{1}{2}} \left(\frac{r^2}{1 - r^{-\frac{2+\alpha}{1+\alpha}}} \right) \vee 1$ is a constant greater than or equal to 1 and depends only on

$(M, \kappa, \eta, \alpha)$, as $r > 1$ is fixed. We plug this upper bound into (A.19) to obtain

$$V_a \leq c_2 \sqrt{\frac{v}{n}} a^{\frac{\alpha}{1+\alpha}-2} + \sqrt{\frac{2 \left[c_1 a^{\frac{2\alpha}{1+\alpha}-2} + 4c_2 \sqrt{\frac{v}{n}} a^{\frac{\alpha}{1+\alpha}-2} \right] t}{na^2}} + \frac{2t}{3na^2}. \quad (\text{A.22})$$

Choose ε_n as the root of $c_2 \sqrt{\frac{v}{n}} a^{\frac{\alpha}{1+\alpha}-2} = 1$, that is,

$$\varepsilon_n = \left(c_2 \sqrt{\frac{v}{n}} \right)^{\frac{1+\alpha}{2+\alpha}}. \quad (\text{A.23})$$

Note that the right-hand side of (A.22) is decreasing in a , and $a \geq \varepsilon_n$ by the construction. Hence, if ε_n satisfies inequality (A.21), that is,

$$n \geq c_2^{-\alpha} \left(\frac{C_1}{c_1} \right)^{1+\frac{\alpha}{2}} v,$$

which can be reduced to an innocuous restriction $n \geq 1$ by inflating, if necessary, c_1 large enough, we can substitute ε_n for a to bound the right-hand side of (A.22). In particular, by noting

$$\begin{aligned} c_2 \sqrt{\frac{v}{n}} a^{\frac{\alpha}{1+\alpha}-2} &\leq \frac{\varepsilon_n}{a} = \frac{1}{\sqrt{kt}} \leq \frac{1}{\sqrt{k}} \quad \text{and} \\ a^{\frac{2\alpha}{1+\alpha}-2} &= a^{2(\frac{\alpha}{1+\alpha}-2)} a^2 \leq [\varepsilon_n^{\frac{\alpha}{1+\alpha}-2}]^2 \varepsilon_n^2 = c_2^{-2} v^{-1} n \varepsilon_n^2, \end{aligned}$$

the right-hand side of (A.22) can be bounded by

$$\begin{aligned} V_a &\leq \frac{1}{\sqrt{k}} + \sqrt{2 \frac{c_1 c_2^{-2} v^{-1} n \varepsilon_n^2 + 8}{nk \varepsilon_n^2}} + \frac{2}{3nk \varepsilon_n^2} \\ &= \frac{1}{\sqrt{k}} + \sqrt{\frac{2c_1 c_2^{-2} v^{-1}}{k} + \frac{8}{nk \varepsilon_n^2}} + \frac{2}{3nk \varepsilon_n^2} \\ &\leq \frac{1}{\sqrt{k}} + \sqrt{\frac{2c_1 c_2^{-2} v^{-1}}{k} + \frac{8}{k}} + \frac{2}{3k} \quad \text{for } n \varepsilon_n^2 \geq 1. \end{aligned} \quad (\text{A.24})$$

Note that condition $n \varepsilon_n^2 \geq 1$ used to derive the last line is valid for all n , since it is equivalent to $n \geq c_2^{-2(1+\alpha)} v^{-(1+\alpha)}$, which holds for all $n \geq 1$ since $c_2 \geq 1$ and $v \geq 1$. By choosing k large enough so that the right-hand side of (A.24) is less than $\frac{1}{2}$, we can conclude

$$\Pr \left(V_a < \frac{1}{2} \right) \geq 1 - \exp(-t). \quad (\text{A.25})$$

Hence, (A.18) yields

$$P^n(d(G_{\text{FB}}^*, \hat{G}_{\text{EWM}}) \geq kt \varepsilon_n^2) \leq \exp(-t)$$

for all $t \geq 1$. From this exponential bound, we obtain

$$\begin{aligned}
E_{P^n}(d(G_{\text{FB}}^*, \hat{G}_{\text{EWM}})) &= \int_0^\infty P^n(d(G_{\text{FB}}^*, \hat{G}_{\text{EWM}}) > t') dt' \\
&\leq \int_0^{k\varepsilon_n^2} P^n(d(G_{\text{FB}}^*, \hat{G}_{\text{EWM}}) \geq t') dt' + \int_{k\varepsilon_n^2}^\infty P^n(d(G_{\text{FB}}^*, \hat{G}_{\text{EWM}}) \geq t') dt' \\
&\leq k\varepsilon_n^2 + k\varepsilon_n^2 e^{-1} \\
&= (1 + e^{-1})kc_2^{\frac{2(1+\alpha)}{2+\alpha}} \left(\frac{v}{n}\right)^{\frac{1+\alpha}{2+\alpha}}.
\end{aligned}$$

So, setting $c = \frac{M}{\kappa}(1 + e^{-1})kc_2^{\frac{2(1+\alpha)}{2+\alpha}}$ leads to the conclusion. *Q.E.D.*

PROOF OF THEOREM 2.4: As in the proof of Theorem 2.2, we work with the normalized outcome support, $Y_{1,i}, Y_{0,i} \in [-\frac{1}{2}, \frac{1}{2}]$. With the normalized outcome, we can assume without loss of generality that constant η of the margin assumption satisfies $\eta \leq 1$.

Let $\alpha \in (0, \infty)$ and $\eta \in (0, 1]$ be given. Similarly to the proof of Theorem 2.2, we consider constructing a suitable subclass $\mathcal{P}^* \subset \mathcal{P}(1, \kappa, \eta, \alpha)$. Let $x_1, \dots, x_v \in \mathcal{X}$ be v points that are shattered by \mathcal{G} , and let γ be a positive number satisfying $\gamma \leq \min\{\eta, \frac{1}{2}\}$, whose proper choice will be given later. We fix the marginal distribution of X at the one supported only on (x_1, \dots, x_v) and having the probability mass function

$$\begin{aligned}
P_X(X_i = x_j) &= \frac{1}{v-1} \left(\frac{\gamma}{\eta}\right)^\alpha \quad \text{for } j = 1, \dots, (v-1), \quad \text{and} \\
P_X(X_i = x_v) &= 1 - \left(\frac{\gamma}{\eta}\right)^\alpha.
\end{aligned}$$

Thus-constructed marginal distribution of X is common in \mathcal{P}^* . As in the proof of Theorem 2.2, we specify D to be independent of (Y_1, Y_0, X) and follow the Bernoulli distribution with $\Pr(D = 1) = 1/2$. Let $\mathbf{b} = (b_1, \dots, b_{v-1}) \in \{0, 1\}^{v-1}$ be a binary vector that uniquely indexes a member of \mathcal{P}^* , and, accordingly, write $\mathcal{P}^* = \{P_{\mathbf{b}} : \mathbf{b} \in \{0, 1\}^{v-1}\}$. For each $j = 1, \dots, (v-1)$, we specify the conditional distribution of Y_1 given $X = x_j$ to be (A.10) if $b_j = 1$ and (A.11) if $b_j = 0$. For $j = v$, the conditional distribution of Y_1 given $X = x_v$ is degenerate at $Y_1 = \frac{1}{2}$. As for the conditional distribution of Y_0 given $X = x_j$, we consider the degenerate distribution at $Y_0 = 0$ for $j = 1, \dots, (v-1)$, and the degenerate distribution at $Y_0 = -\frac{1}{2}$ for $X = x_v$. In this specification of \mathcal{P}^* , it holds that

$$P_X(|\tau(X)| \leq t) = \begin{cases} 0 & \text{for } t \in [0, \gamma), \\ \left(\frac{\gamma}{\eta}\right)^\alpha & \text{for } t \in [\gamma, 1), \\ 1 & \text{for } t \geq 1 \end{cases}$$

for every $P_{\mathbf{b}} \in \mathcal{P}^*$. Since $\gamma \leq \eta$, $P_X(|\tau(X)| \leq t) \leq (t/\eta)^\alpha$ holds for all $t \in [0, \eta]$. Furthermore, by the construction of the support points, for every $P_{\mathbf{b}} \in \mathcal{P}^*$, the first-best decision rule $G_{\mathbf{b}}^* = \{x_j : j < v, b_j = 1\} \cup \{x_v\}$ is contained in \mathcal{G} . Hence, $\mathcal{P}^* \subset \mathcal{P}_{\text{FB}}(1, \kappa, \eta, \alpha)$ holds.

Let $\pi(\mathbf{b})$ be a prior distribution for \mathbf{b} such that b_1, \dots, b_{v-1} are i.i.d. and $b_1 \sim \text{Ber}(1/2)$. The maximized social welfare is

$$W(G_{\mathbf{b}}^*) = \frac{\gamma}{v-1} \left(\frac{\gamma}{\eta} \right)^\alpha \left(\sum_{j=1}^{v-1} b_j \right) + \left[1 - \left(\frac{\gamma}{\eta} \right)^\alpha \right].$$

Let \hat{G} be an arbitrary treatment choice rule as a function of (Z_1, \dots, Z_n) , and $\hat{\mathbf{b}} \in \{0, 1\}^v$ be a binary vector whose j th element is $\hat{b}_j = 1\{x_j \in \hat{G}\}$.

The welfare loss can be bounded from below as follows:

$$\begin{aligned} & \sup_{P \in \mathcal{P}(1, \kappa, \eta, \alpha)} E_{P^n} [W_G^* - W(\hat{G})] \\ & \geq \sup_{P_{\mathbf{b}} \in \mathcal{P}^*} E_{P_{\mathbf{b}}^n} [W(G_{\mathbf{b}}^*) - W(\hat{G})] \\ & \geq \int_{\mathbf{b}} E_{P_{\mathbf{b}}^n} [W(G_{\mathbf{b}}^*) - W(\hat{G})] d\pi(\mathbf{b}) \geq \int_{\mathbf{b}} E_{P_{\mathbf{b}}^n} [W(G_{\mathbf{b}}^*) - W(\hat{G} \cup \{x_v\})] d\pi(\mathbf{b}) \\ & = \gamma \int_{\mathbf{b}} E_{P_{\mathbf{b}}^n} [P_X(\{(G_{\mathbf{b}}^* \Delta \hat{G}) \cap \{x_1, \dots, x_{v-1}\}\})] d\pi(\mathbf{b}) \\ & = \gamma \int_{\mathbf{b}} \int_{Z_1, \dots, Z_n} P_X(\{b(X) \neq \hat{b}(X)\} \cap \{x_1, \dots, x_{v-1}\}) dP_{\mathbf{b}}^n(Z_1, \dots, Z_n) d\pi(\mathbf{b}) \\ & \geq \inf_{G_n} \gamma \int_{\mathbf{b}} \int_{Z_1, \dots, Z_n} P_X(\{b(X) \neq b_n(X)\}) dP_{\mathbf{b}}^n(Z_1, \dots, Z_n) d\pi(\mathbf{b}), \end{aligned}$$

where the second line follows since $W(G_{\mathbf{b}}^*) - W(\hat{G}) \geq W(G_{\mathbf{b}}^*) - W(\hat{G} \cup \{x_v\})$ holds for every \mathbf{b} and \hat{G} . The infimum in the last line is taken over decision sets $G_n = \{x_j : b_n(x_j) = 1\}$ that are constrained to contain $\{x_v\}$, that is, $b_n(x_v) = 1$.

By the same reasonings as in obtaining (A.12), the lower bound of the welfare loss as viewed as the Bayes risk can be expressed as

$$\begin{aligned} & \sup_{P \in \mathcal{P}(1, \kappa, \eta, \alpha)} E_{P^n} [W(G^*) - W(\hat{G})] \\ & \geq \frac{\gamma}{v-1} \left(\frac{\gamma}{\eta} \right)^\alpha \int_{Z_1, \dots, Z_n} \sum_{j=1}^{v-1} [\min\{\pi(b_j = 1 | Z_1, \dots, Z_n), 1 - \pi(b_j = 1 | Z_1, \dots, Z_n)\}] d\tilde{P}^n. \end{aligned}$$

Repeating the same bounding arguments as in the proof of Theorem 2.2, a lower bound of the Bayes risk analogous to (A.14) is obtained by

$$\sup_{P \in \mathcal{P}(1, \kappa, \eta, \alpha)} E_{P^n} [W(G^*) - W(\hat{G})] \geq \frac{\gamma}{2} \left(\frac{\gamma}{\eta} \right)^\alpha \exp \left\{ -\frac{4\gamma}{1-2\gamma} \sqrt{\frac{n}{8(v-1)}} \left(\frac{\gamma}{\eta} \right)^\alpha \right\}.$$

The slowest convergence rate of this lower bound can be obtained by tuning γ to be converging at the rate of $n^{-\frac{1}{2+\alpha}}$. In particular, by choosing $\gamma = \eta^{\frac{\alpha}{2+\alpha}} \left(\frac{v-1}{n} \right)^{\frac{1}{2+\alpha}}$ assuming $\gamma \leq \frac{1}{4}$, the exponential term can be bounded from below by $\exp\{-2\sqrt{2}\}$, so we obtain the

following lower bound:

$$\frac{1}{2} \eta^{-\frac{\alpha}{2+\alpha}} \left(\frac{v-1}{n} \right)^{\frac{1+\alpha}{2+\alpha}} \exp\{-2\sqrt{2}\}. \quad (\text{A.26})$$

Recall that γ is constrained to $\gamma \leq \min\{\eta, \frac{1}{4}\}$. This implies that the obtained bound is valid for

$$n \geq (\max\{\eta^{-1}, 4\})^{2+\alpha} \eta^\alpha (v-1),$$

whose stronger but simpler form is given by

$$n \geq \max\{\eta^{-2}, 4^{2+\alpha}\} (v-1). \quad (\text{A.27})$$

The lower bound presented in this theorem follows by denormalizing the outcomes, that is, multiply M to (A.26) and substitute η/M for η appearing in (A.26) and (A.27). *Q.E.D.*

A.4. Proof of Theorems 2.5 and 2.6

PROOF OF THEOREM 2.5: Let $W_n^\tau(G)$ be the sample analogue of the welfare criterion (1.2) in the main text that one would construct if the true regression equations were known, $W_n^\tau(G) \equiv E_n(m_0(X_i) + E_n(\tau(X_i) \cdot 1\{X_i \in G\}))$, and $\hat{W}_n^\tau(G)$ be the empirical welfare with the conditional treatment effect estimators $\hat{\tau}^m(\cdot)$ plugged in,

$$\hat{W}_n^\tau(G) \equiv E_n[m_0(X_i) + \hat{\tau}^m(X_i)1\{X_i \in G\}]. \quad (\text{A.28})$$

Since the m -hybrid rule maximizes $\hat{W}_n^\tau(\cdot)$, it holds that $\hat{W}_n^\tau(\hat{G}_{m\text{-hybrid}}) - \hat{W}_n^\tau(\tilde{G}) \geq 0$ for any $\tilde{G} \in \mathcal{G}$. The following inequalities therefore follow:

$$\begin{aligned} W(\tilde{G}) - W(\hat{G}_{m\text{-hybrid}}) &\leq W_n^\tau(\tilde{G}) - \hat{W}_n^\tau(\tilde{G}) - W_n^\tau(\hat{G}_{m\text{-hybrid}}) + \hat{W}_n^\tau(\hat{G}_{m\text{-hybrid}}) \\ &\quad + W(\tilde{G}) - W(\hat{G}_{m\text{-hybrid}}) - W_n^\tau(\tilde{G}) + W_n^\tau(\hat{G}_{m\text{-hybrid}}) \\ &= \frac{1}{n} \sum_{i=1}^n [\tau(X_i) - \hat{\tau}^m(X_i)][1\{X_i \in \tilde{G}\} - 1\{X_i \in \hat{G}_{m\text{-hybrid}}\}] \\ &\quad + W(\tilde{G}) - W_n^\tau(\tilde{G}) + W_n^\tau(\hat{G}_{m\text{-hybrid}}) - W(\hat{G}_{m\text{-hybrid}}) \\ &\leq \frac{1}{n} \sum_{i=1}^n |\hat{\tau}^m(X_i) - \tau(X_i)| + 2 \sup_{G \in \mathcal{G}} |W_n^\tau(G) - W(G)|. \end{aligned} \quad (\text{A.29})$$

This implies that the average welfare loss of the m -hybrid rule can be bounded by

$$\begin{aligned} E_{P^n} [W_G^* - W(\hat{G}_{m\text{-hybrid}})] &\leq E_{P^n} \left[\frac{1}{n} \sum_{i=1}^n |\hat{\tau}^m(X_i) - \tau(X_i)| \right] \\ &\quad + 2E_{P^n} \left[\sup_{G \in \mathcal{G}} |W_n^\tau(G) - W(G)| \right]. \end{aligned} \quad (\text{A.30})$$

For the e -hybrid rule, replacing $W_n^\tau(\cdot)$ and $\hat{W}_n^\tau(\cdot)$ in (A.29) with the empirical welfare $W_n(\cdot)$ defined in (1.7) and $\hat{W}_n(G) \equiv E_n[\frac{Y_i(1-D_i)}{1-e(X_i)} + \hat{\tau}_i^e \cdot 1\{X_i \in G\}]$, respectively, yields a

similar upper bound

$$E_{P^n} [W_{\mathcal{G}}^* - W(\hat{G}_{e\text{-hybrid}})] \leq E_{P^n} \left[\frac{1}{n} \sum_{i=1}^n |\hat{\tau}_i^e - \tau_i| \right] + 2E_{P^n} \left[\sup_{G \in \mathcal{G}} |W_n(G) - W(G)| \right], \quad (\text{A.31})$$

where $\tau_i = \frac{Y_i D_i}{e^{(X_i)} - 1} - \frac{Y_i(1-D_i)}{1-e^{(X_i)}}$. Note that the uniform convergence rate of $E_{P^n}[\sup_{G \in \mathcal{G}} |W_n^\tau(G) - W(G)|]$ is $n^{-1/2}$, same as that of $E_{P^n}[\sup_{G \in \mathcal{G}} |W_n(G) - W(G)|]$, since the proof of Theorem 2.1 can be applied to the following class of functions:

$$\mathcal{F}^\tau \equiv \{f(X_i; G) = m_0(X_i) + \tau(X_i) \cdot 1\{X_i \in G\} : G \in \mathcal{G}\},$$

which is the VC-subgraph class with the VC-dimension at most v by Lemma A.1. Combined with Condition 2.1 (m), (A.30) implies the uniform convergence rate of the m -hybrid rule given in the current theorem. Similarly, combined with Condition 2.1 (e) and $n^{-1/2}$ -convergence rate of $E_{P^n}[\sup_{G \in \mathcal{G}} |W_n(G) - W(G)|]$, (A.31) leads to the uniform convergence rate of $\phi_n^{-1} \vee n^{-1/2}$ for the e -hybrid rule. *Q.E.D.*

The next lemma gives a linearized solution of a certain polynomial inequality. We owe this lemma to Shin Kanaya (2014, personal communication). The technique of applying the mean value expansion to an implicit function defined as the root of a polynomial equation has been used by Shin Kanaya and Dennis Kristensen in unpublished work on bandwidth choice.

LEMMA A.8: *Let $A \geq 0, B \geq 0$, and $X \geq 0$. For any $\alpha \geq 0$, $X \leq AX^{\frac{\alpha}{1+\alpha}} + B$ implies*

$$X \leq A^{1+\alpha} + (1+\alpha)B.$$

PROOF: When $A = B = 0$, the conclusion trivially holds. When $B > 0$, $X = AX^{\frac{\alpha}{1+\alpha}} + B$ has a unique root, and we denote it by $X^* = g(A, B)$. When $A > 0$ and $B = 0$, we mean by $g(A, 0)$ the nonzero root of $X = AX^{\frac{\alpha}{1+\alpha}}$. Let $f(X, A, B) = X - AX^{\frac{\alpha}{1+\alpha}} - B$. By the form of the inequality, the original inequality can be equivalently written as $X \leq X^* = g(A, B)$, so we aim to verify that X^* is bounded from above by $A^{1+\alpha} + (1+\alpha)B$. Consider the mean value expansion of $g(A, B)$ in B at $B = 0$,

$$X^* = g(A, 0) + \frac{\partial g}{\partial B}(A, \tilde{B}) \times B \quad \text{for some } 0 \leq \tilde{B} \leq B.$$

Note $g(A, 0) = A^{1+\alpha}$. In addition, by the implicit function theorem, we have, with $\tilde{X} = g(A, \tilde{B})$,

$$\begin{aligned} \frac{\partial g}{\partial B}(A, \tilde{B}) &= -\frac{\frac{\partial f}{\partial B}(\tilde{X}, A, \tilde{B})}{\frac{\partial f}{\partial X}(\tilde{X}, A, \tilde{B})} = \frac{1}{1 - \frac{\alpha}{1+\alpha} A \tilde{X}^{-\frac{1}{1+\alpha}}} = \frac{\tilde{X}}{\frac{\tilde{X}}{1+\alpha} + \frac{\alpha}{1+\alpha} (\tilde{X} - A \tilde{X}^{\frac{\alpha}{1+\alpha}})} \\ &= \frac{\tilde{X}}{\frac{\tilde{X}}{1+\alpha} + \frac{\alpha}{1+\alpha} \tilde{B}} \leq 1 + \alpha. \end{aligned}$$

Hence, $X^* \leq A^{1+\alpha} + (1+\alpha)B$ holds.

Q.E.D.

The next lemma provides an exponential tail probability bound of the supremum of the centered empirical processes. This lemma follows from Theorem 2.14.9 in van der Vaart and Wellner (1996) combined with their Theorem 2.6.4.

LEMMA A.9: Assume \mathcal{G} is a VC-class of subsets in \mathcal{X} with VC-dimension $v < \infty$. Let $P_{X,n}(\cdot)$ be the empirical probability distribution on \mathcal{X} constructed upon (X_1, \dots, X_n) generated i.i.d. from $P_X(\cdot)$. Then,

$$P^n \left(\sup_{G \in \mathcal{G}} |P_{X,n}(G) - P_X(G)| > t \right) \leq \left(\frac{C_4 t}{\sqrt{2v}} \right)^{2v} n^v \exp(-nt^2)$$

holds for every $t > 0$, where C_4 is a universal constant.

PROOF OF THEOREM 2.6: We first consider the m -hybrid case. Set $\tilde{G} = G_{\text{FB}}^*$ in (A.29) and rewrite (A.29) in terms of the normalized welfare loss for $\hat{G}_{m\text{-hybrid}}$,

$$\begin{aligned} d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) &\leq \frac{\kappa}{M} [W_n^\tau(G_{\text{FB}}^*) - \hat{W}_n^\tau(G_{\text{FB}}^*) - W_n^\tau(\hat{G}_{m\text{-hybrid}}) + \hat{W}_n^\tau(\hat{G}_{m\text{-hybrid}})] \\ &\quad + d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) - d_n^\tau(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) \\ &\leq \frac{1}{n} \sum_{i=1}^n \frac{\kappa}{M} [\tau(X_i) - \hat{\tau}^m(X_i)] [1\{X_i \in G_{\text{FB}}^*\} - 1\{X_i \in \hat{G}_{m\text{-hybrid}}\}] \quad (\text{A.32}) \\ &\quad + d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) - d_n^\tau(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) \\ &\leq \rho_n + d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) - d_n^\tau(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}), \end{aligned}$$

where $d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}})$ is as defined in equation (A.16), $d_n^\tau(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) = W_n^\tau(G_{\text{FB}}^*) - W_n^\tau(\hat{G}_{m\text{-hybrid}})$,

$$\rho_n \equiv \frac{\kappa}{M} \max_{1 \leq i \leq n} |\hat{\tau}^m(X_i) - \tau(X_i)| P_{X,n}(G_{\text{FB}}^* \Delta \hat{G}_{m\text{-hybrid}}),$$

and $P_{X,n}$ is the empirical distribution on \mathcal{X} constructed upon (X_1, \dots, X_n) . Define a class of functions generated by $G \in \mathcal{G}$,

$$\mathcal{H}^\tau \equiv \{h(Z_i; G) : G \in \mathcal{G}\},$$

$$h(Z_i; G) \equiv \frac{\kappa}{M} \tau(X_i) \cdot [1\{X_i \in G\} - 1\{X_i \in G_{\text{FB}}^*\}],$$

which is a VC-subgraph class with the VC-dimension at most v with envelope $\bar{H} = 1$ by Lemma A.1. Let $a = \sqrt{kt} \varepsilon_n$ be as defined in the proof of Theorem 2.3 and $V_a^\tau \equiv \sup_{h \in \mathcal{H}^\tau} \left\{ \frac{E_n(h) - E_P(h)}{-E_P(h) + a^2} \right\}$. By noting

$$d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) - d_n^\tau(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) \leq V_a^\tau (d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) + a^2),$$

inequality (A.32) implies

$$d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) \leq \rho_n + V_a^\tau (d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) + a^2). \quad (\text{A.33})$$

Denote event $\{V_a^\tau < \frac{1}{2}\}$ by Ω_t , which is equivalent to event $\{d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) \leq 2\rho_n + k\varepsilon_n^2 t\}$. The same line of argument that leads to (A.25) in the proof of Theorem 2.3 leads to, for $t \geq 1$,

$$P^n(\Omega_t) = P^n(d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) \leq 2\rho_n + k\varepsilon_n^2 t) \geq 1 - \exp(-t), \quad (\text{A.34})$$

where ε_n is given in (A.23). We bound ρ_n from above by

$$\rho_n \leq \frac{\kappa}{M} \left[\max_{1 \leq i \leq n} |\hat{\tau}^m(X_i) - \tau(X_i)| P_X(G_{\text{FB}}^* \Delta \hat{G}_{m\text{-hybrid}}) + \mathcal{V}_{0,n} \max_{1 \leq i \leq n} |\hat{\tau}^m(X_i) - \tau(X_i)| \right],$$

where

$$\mathcal{V}_{0,n} = \sup_{G \in \mathcal{G}} |P_{X,n}(G_{\text{FB}}^* \Delta G) - P_X(G_{\text{FB}}^* \Delta G)|.$$

Let $\lambda > 0$, that will be chosen properly later. Define events

$$\Lambda_1 = \{\mathcal{V}_{0,n} \leq n^{-\lambda}\}, \quad \Lambda_2 = \{P_X(G_{\text{FB}}^* \Delta \hat{G}_{m\text{-hybrid}}) \geq n^{-\lambda}\}.$$

Then, on $\Lambda_1 \cap \Lambda_2$, it holds that $\mathcal{V}_{0,n} \leq P_X(G_{\text{FB}}^* \Delta \hat{G}_{m\text{-hybrid}})$. Therefore, on $\Lambda_1 \cap \Lambda_2 \cap \Omega_t$, $d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}})$ can be bounded by

$$\begin{aligned} d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) &\leq 4 \frac{\kappa}{M} \max_{1 \leq i \leq n} |\hat{\tau}^m(X_i) - \tau(X_i)| P_X(G_{\text{FB}}^* \Delta \hat{G}_{m\text{-hybrid}}) + k\varepsilon_n^2 t \\ &\leq 4c_1 \frac{\kappa}{M} \max_{1 \leq i \leq n} |\hat{\tau}^m(X_i) - \tau(X_i)| d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}})^{\frac{\alpha}{1+\alpha}} + k\varepsilon_n^2 t, \end{aligned}$$

where the second line follows from Lemma A.7 with the same definition of c_1 given there. By Lemma A.8 and substituting (A.23) to ε_n , we obtain, on event $\Lambda_1 \cap \Lambda_2 \cap \Omega_t$,

$$d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) \leq c_6 \left[\max_{1 \leq i \leq n} |\hat{\tau}^m(X_i) - \tau(X_i)| \right]^{1+\alpha} + c_7 \left(\frac{v}{n} \right)^{\frac{1+\alpha}{2+\alpha}} t, \quad (\text{A.35})$$

where constants c_6 and c_7 depend only on $(M, \kappa, \eta, \alpha)$.

Using the upper bound derived in (A.35), we obtain, for $t \geq 1$,

$$\begin{aligned} &E_{P^n}(d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}})) \\ &= E_{P^n}(d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) 1\{\Lambda_1 \cap \Lambda_2 \cap \Omega_t\}) + E_{P^n}(d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) 1\{\Lambda_1^c \cup \Lambda_2^c \cup \Omega_t^c\}) \\ &\leq c_6 E_{P^n} \left(\left[\max_{1 \leq i \leq n} |\hat{\tau}^m(X_i) - \tau(X_i)| \right]^{1+\alpha} \right) + c_7 \left(\frac{v}{n} \right)^{\frac{1+\alpha}{2+\alpha}} t + P^n(\Lambda_1^c) \\ &\quad + E_{P^n}(d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) 1\{\Lambda_2^c\}) + P^n(\Omega_t^c) \\ &\leq c_6 \underbrace{\tilde{\psi}_n^{-(1+\alpha)} E_{P^n} \left(\left[\tilde{\psi}_n \max_{1 \leq i \leq n} |\hat{\tau}^m(X_i) - \tau(X_i)| \right]^{1+\alpha} \right)}_{A_{1,n}} + c_7 \underbrace{\left(\frac{v}{n} \right)^{\frac{1+\alpha}{2+\alpha}} t}_{A_{2,n}} \\ &\quad + \underbrace{\left(\frac{C_4}{\sqrt{2v}} \right)^{2v} n^{-2v(\lambda-\frac{1}{2})} \exp(-n^{-2(\lambda-\frac{1}{2})})}_{A_{3,n}} + \underbrace{n^{-\lambda}}_{A_{4,n}} + \underbrace{\exp(-t)}_{A_{5,n}}, \end{aligned}$$

where $\tilde{\psi}_n$ is a sequence as specified in equation (2.10) in the main text. In these inequalities, the third line uses (A.35) and $d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) \leq 1$. In the fourth line, $A_{3,n}$ follows from Lemma A.9, $A_{4,n}$ follows from $d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}) \leq P_X(G_{\text{FB}}^* \Delta \hat{G}_{m\text{-hybrid}})$ and $P_X(G_{\text{FB}}^* \Delta \hat{G}_{m\text{-hybrid}}) < n^{-\lambda}$ on Λ_2^c , and $A_{5,n}$ follows from (A.34).

We now discuss convergence rates of $A_{j,n}$, $j = 1, \dots, 5$, individually with suitable choices of t and λ . Equation (2.10) assumed in this theorem implies

$$\begin{aligned} & \sup_{P \in \mathcal{P}_m} E_{P^n} \left(\left(\tilde{\psi}_n \max_{1 \leq i \leq n} |\hat{\tau}^m(X_i) - \tau(X_i)| \right)^{1+\alpha} \right) \\ &= \sup_{P \in \mathcal{P}_m} E_{P^n} \left(\left[\left(\tilde{\psi}_n \max_{1 \leq i \leq n} |\hat{\tau}^m(X_i) - \tau(X_i)| \right)^2 \right]^{\frac{1+\alpha}{2}} \right) \\ &\leq \left(\left[\sup_{P \in \mathcal{P}_m} E_{P^n} \left(\tilde{\psi}_n \max_{1 \leq i \leq n} |\hat{\tau}^m(X_i) - \tau(X_i)| \right)^2 \right]^{\frac{1+\alpha}{2}} \right) \\ &= O(1), \end{aligned}$$

where the third line follows from Jensen's inequality since $(1 + \alpha)/2 \leq 1$. Hence, $A_{1,n}$ satisfies $\sup_{P \in \mathcal{P}_m} A_{1,n} = O(\tilde{\psi}_n^{-(1+\alpha)})$. By setting $t = (1 + \alpha) \log \tilde{\psi}_n$, we can make the convergence rate of $A_{5,n}$ equal to that of $A_{1,n}$. At the same time, by choosing $\lambda > \frac{1+\alpha}{2+\alpha} \geq \frac{1}{2}$, we can make $A_{3,n}$ and $A_{4,n}$ converge faster than $A_{2,n}$. Hence, the uniform convergence rate of $E_{P^n}(d(G_{\text{FB}}^*, \hat{G}_{m\text{-hybrid}}))$ over $P \in \mathcal{P}_m \cap \mathcal{P}_{\text{FB}}(M, \kappa, \eta, \alpha)$ is bounded by the convergence rates of the $A_{1,n}$ and $A_{2,n}$,

$$O\left(\sup_{P \in \mathcal{P}_m} A_{1,n} \vee \sup_{P \in \mathcal{P}_{\text{FB}}(M, \kappa, \eta, \alpha)} A_{2,n}\right) = O(\tilde{\psi}_n^{-(1+\alpha)} \vee n^{-\frac{1+\alpha}{2+\alpha}} \log \tilde{\psi}_n).$$

This completes the proof for the m -hybrid case.

A proof for the e -hybrid case follows almost identically to the proof of the m -hybrid case. The differences are that ρ_n in inequality (A.32) is given by

$$\rho_n = \frac{\kappa}{M} \max_{1 \leq i \leq n} |\hat{\tau}_i^e - \tau_i| P_{X,n}(G_{\text{FB}}^* \Delta \hat{G}_{e\text{-hybrid}}),$$

and that inequality (A.33) is replaced by

$$d(G_{\text{FB}}^*, \hat{G}_{e\text{-hybrid}}) \leq \rho_n + V_a(d(G_{\text{FB}}^*, \hat{G}_{e\text{-hybrid}}) + a^2), \quad (\text{A.36})$$

where V_a is as defined in the proof of Theorem 2.3. The rest of the proof goes similarly to the proof of the first claim except that the rate $\tilde{\phi}_n$ given in equation (2.11) replaces $\tilde{\psi}_n$ in the first claim. *Q.E.D.*

APPENDIX B: INFERENCE FOR WELFARE GAIN

In the proposed EWM procedure, the maximized empirical welfare $W_n(\hat{G}_{\text{EWM}})$ can be seen as an estimate of $W(\hat{G}_{\text{EWM}})$, the welfare level attained by implementing the estimated treatment rule.¹ In situations where propensity scores are known, this section pro-

¹It is important to note that in finite samples, $W_n(\hat{G}_{\text{EWM}})$ estimates $W(\hat{G}_{\text{EWM}})$ with an upward bias. With fixed n , the size of the bias becomes bigger as \mathcal{G} becomes more complex.

vides a procedure for constructing asymptotically valid confidence intervals for the population welfare gain of implementing the estimated rule.

Let $\hat{G} \in \mathcal{G}$ be an estimated treatment rule such as \hat{G}_{EWM} or other data-driven way of selecting G from the set of candidate policies. Define the welfare gain of implementing the estimated treatment rule $\hat{G} \in \mathcal{G}$ by

$$V(\hat{G}) \equiv W(\hat{G}) - W(G_0),$$

where G_0 is a benchmark treatment assignment rule with which the estimated treatment rule \hat{G} is compared in terms of the social welfare. For instance, if the estimated treatment rule \hat{G} is compared with the “no treatment” case, G_0 is the empty set \emptyset . Alternatively, if a benchmark policy is the non-individualized uniform adoption of the treatment, G_0 is set at $G_0 = \mathcal{X}$, and $V(\hat{G})$ is interpreted as the welfare gain of implementing individualized treatment assignment instead of the non-individualized implementation of the treatment.

A construction of one-sided confidence intervals for $V(\hat{G})$ proceeds as follows. Let $\nu_n(G) = \sqrt{n}(V_n(G) - V(G))$, where $V_n(G) \equiv W_n(G) - W_n(G_0)$. If there is a random variable $\tilde{\nu}_n$ such that $\nu_n(\hat{G}) \leq \tilde{\nu}_n$ holds P^n -almost surely, and if $\tilde{\nu}_n$ converges in distribution to a non-degenerate random variable $\tilde{\nu}$, then, with $q_{\tilde{\nu}}(1 - \bar{\alpha})$, the $(1 - \bar{\alpha})$ th quantile of $\tilde{\nu}$, it holds that

$$P^n(\nu_n(\hat{G}) \leq q_{\tilde{\nu}}(1 - \bar{\alpha})) \geq P^n(\tilde{\nu}_n \leq q_{\tilde{\nu}}(1 - \bar{\alpha})) \rightarrow \Pr(\tilde{\nu} \leq q_{\tilde{\nu}}(1 - \bar{\alpha})) = 1 - \bar{\alpha} \quad \text{as } n \rightarrow \infty.$$

Hence, if $\hat{q}_{\tilde{\nu}}(1 - \bar{\alpha})$, a consistent estimator of $q_{\tilde{\nu}}(1 - \bar{\alpha})$, is available, an asymptotically valid one-sided confidence interval for $V(\hat{G})$ with coverage probability $(1 - \bar{\alpha})$ can be given by

$$\left[V_n(\hat{G}) - \frac{\hat{q}_{\tilde{\nu}}(1 - \bar{\alpha})}{\sqrt{n}}, \infty \right). \quad (\text{B.1})$$

Two-sided confidence intervals for $V(\hat{G})$ can be constructed similarly by considering a random variable $\tilde{\nu}_n$ that satisfies $|\nu_n(\hat{G})| \leq \tilde{\nu}_n$, P^n -almost surely, and converges to a nondegenerate random variable $\tilde{\nu}$. With $\hat{q}_{\tilde{\nu}}(1 - \bar{\alpha})$ a consistent estimator for the $(1 - \bar{\alpha})$ th quantile of $\tilde{\nu}$, two-sided confidence interval for $V(\hat{G})$ can be given by

$$\left[V_n(\hat{G}) - \frac{\hat{q}_{\tilde{\nu}}(1 - \bar{\alpha})}{\sqrt{n}}, V_n(\hat{G}) + \frac{\hat{q}_{\tilde{\nu}}(1 - \bar{\alpha})}{\sqrt{n}} \right]. \quad (\text{B.2})$$

In the algorithm summarized below, we specify $\tilde{\nu}_n$ to be $\tilde{\nu}_n = \sqrt{n} \sup_{G \in \mathcal{G}} (V_n(G) - V(G))$ and $\tilde{\nu}_n$ to be $\tilde{\nu}_n = \sqrt{n} \sup_{G \in \mathcal{G}} |V_n(G) - V(G)|$, and estimate the $(1 - \bar{\alpha})$ -quantiles of their asymptotic distributions by bootstrapping the centered empirical processes.²

²The current choices of $\tilde{\nu}_n$ and $\tilde{\nu}_n$ are likely to yield conservative confidence intervals. Keeping the same nominal coverage probability, it is feasible to tighten up the confidence intervals with more sophisticated choices of $\tilde{\nu}_n$ and $\tilde{\nu}_n$, such as $\tilde{\nu}_n = \sqrt{n} \sup_{G \in \hat{\mathcal{G}}} (V_n(G) - V(G))$ and $\tilde{\nu}_n = \sqrt{n} \sup_{G \in \hat{\mathcal{G}}} |V_n(G) - V(G)|$, where $\hat{\mathcal{G}}$ is a data-dependent subclass of \mathcal{G} that contains \hat{G} with probability approaching 1. Such $\hat{\mathcal{G}}$ can be obtained by applying the technique of contact set estimation in the context of stochastic dominance testing. See [Linton, Song, and Whang \(2010\)](#) and [Donald and Hsu \(2016\)](#), as well as the literature on moment inequalities with moment selection ([Andrews and Shi \(2013\)](#), among others).

ALGORITHM B.1: 1. Let $\hat{G} \in \mathcal{G}$ be an estimated treatment assignment rule (e.g., EWM rule), and $V_n(\cdot) = W_n(\cdot) - W_n(G_0)$ be the empirical welfare gain obtained from the original sample.

2. Resample n -observations of $Z_i = (Y_i, D_i, X_i)$ randomly with replacement from the original sample and construct the bootstrap analogue of the welfare gain, $V_n^*(\cdot) = W_n^*(\cdot) - W_n^*(G_0)$, where $W_n^*(\cdot)$ is the empirical welfare of the bootstrap sample.

3. For one-sided confidence intervals, compute $\tilde{v}_n^* = \sqrt{n} \sup_{G \in \mathcal{G}} (V_n^*(G) - V_n(G))$. For two-sided confidence intervals, compute $\tilde{v}_n^* = \sqrt{n} \sup_{G \in \mathcal{G}} |V_n^*(G) - V_n(G)|$.

4. Let $\bar{\alpha} \in (0, 1/2)$. Repeat step 2 and 3 many times. For one-sided (two-sided) confidence intervals, obtain $\hat{q}_{\tilde{v}}(1 - \bar{\alpha})$ ($\hat{q}_{\tilde{v}}(1 - \bar{\alpha})$) by the empirical $(1 - \bar{\alpha})$ th quantile of the bootstrap realizations of \tilde{v}_n^* (\tilde{v}_n^*).

Given Assumption 2.1, the uniform central limit theorem for empirical processes assures that \tilde{v}_n and \tilde{v}_n^* converge in distribution to the supremum of mean-zero Brownian bridge processes and the supremum of their absolute values, respectively. Furthermore, by the well-known result on the asymptotic validity of the bootstrap empirical processes (see, e.g., Section 3.6 of van der Vaart and Wellner (1996)), the bootstrap critical values $\hat{q}_{\tilde{v}}(1 - \bar{\alpha})$ and $\hat{q}_{\tilde{v}^*}(1 - \bar{\alpha})$ consistently estimate the corresponding quantiles of the limiting distributions of \tilde{v}_n and \tilde{v}_n^* , respectively. We can therefore assure that the confidence intervals constructed in (B.1) and (B.2) have the desired asymptotic coverage probability.

The same inference procedure is valid for the welfare gain estimated with demeaned outcomes $V_n^{\text{dm}}(\hat{G}) \equiv W_n^{\text{dm}}(\hat{G}) - W_n^{\text{dm}}(G_0)$. Resampling in this case is from observations $Z_i^{\text{dm}} = (Y_i^{\text{dm}}, D_i, X_i)$, with outcomes $Y_i^{\text{dm}} = Y_i - E_n[Y_i]$ demeaned by the outcome mean in the original sample.

APPENDIX C: COMPUTING EWM TREATMENT RULES

The Empirical Welfare Maximization rule \hat{G}_{EWM} , as well as hybrid rules $\hat{G}_{m\text{-hybrid}}$, and $\hat{G}_{e\text{-hybrid}}$, share the same structure

$$\hat{G} \in \arg \max_{G \in \mathcal{G}} \sum_{i=1}^n g_i \cdot 1\{X_i \in G\}, \quad (\text{C.1})$$

where each g_i is a function of the data, that is, for the EWM rule \hat{G}_{EWM} , $g_i = \frac{1}{n} \left(\frac{Y_i D_i}{e(X_i)} - \frac{Y_i(1-D_i)}{1-e(X_i)} \right)$, for the e -hybrid rule $\hat{G}_{e\text{-hybrid}}$, $g_i = \hat{\tau}_i^e/n$, and for the m -hybrid rule $\hat{G}_{m\text{-hybrid}}$, $g_i = \hat{\tau}^m(X_i)/n$. The objective function in (C.1) is non-convex and discontinuous in G ; thus finding \hat{G} could be computationally challenging. In this section, we propose a set of convenient tools that permit solving this optimization problem and performing inference using widely available software for practically important classes of sets \mathcal{G} defined by linear eligibility scores.³

C.1. Single Linear Index Rules

We start with the problem of computing optimal treatment rules that assign treatments based on a linear index (linear eligibility score; LES, see Examples 2.1 and 2.2). To reduce notational complexity, we include a constant in the covariate vector X throughout

³For the empirical illustration, we used IBM ILOG CPLEX Optimization Studio, which is available free for academic use through the IBM Academic Initiative.

the exposition of this section. An LES rule can be expressed as $1\{X^T\beta \geq 0\}$. This type of treatment rule is commonly used in practice because it offers a simple way to reduce the dimension of observable characteristics. Furthermore, it is easy to enforce monotonicity of treatment assignment in specific covariates by imposing sign restrictions on the components of β .

Let \mathcal{G}_{LES} be a collection of half-spaces of the covariate space \mathcal{X} , which are the upper contour sets of linear functions:

$$\begin{aligned}\mathcal{G}_{\text{LES}} &= \{G_\beta : \beta \in \mathbf{B} \subset \mathbb{R}^{d_x+1}\}, \\ G_\beta &= \{x : x^T\beta \geq 0\}.\end{aligned}$$

Then the optimization problem (C.1) becomes

$$\max_{\beta \in \mathbf{B}} \sum_{i=1}^n g_i \cdot 1\{X_i^T\beta \geq 0\}. \quad (\text{C.2})$$

This problem is similar to the maximum weighted score problem analyzed in Florios and Skouras (2008). They observed that the maximum score objective function could be rewritten as a Mixed Integer Linear Programming problem with additional binary parameters (z_1, \dots, z_n) that replace the indicator functions $1\{X_i^T\beta \geq 0\}$. The equality $z_i = 1\{X_i^T\beta \geq 0\}$ is imposed by a combination of linear inequality constraints and the restriction that z_i 's are binary. The advantage of a MILP representation is that it is a standard optimization problem that could be solved by multiple commercial and open-source solvers. The branch-and-cut algorithms implemented in these solvers are faster than brute force combinatorial optimization.

We propose replacing (C.2) by its equivalent problem:

$$\max_{\substack{\beta \in \mathbf{B}, \\ z_1, \dots, z_n \in \mathbb{R}}} \sum_{i=1}^n g_i \cdot z_i \quad (\text{C.3})$$

$$\begin{aligned}\text{s.t. } & \frac{X_i^T\beta}{C_i} < z_i \leq 1 + \frac{X_i^T\beta}{C_i} \quad \text{for } i = 1, \dots, n, \\ & z_i \in \{0, 1\},\end{aligned} \quad (\text{C.4})$$

where constants C_i should satisfy $C_i > \sup_{\beta \in \mathbf{B}} |X_i^T\beta|$. Then the inequality constraints (C.4) and the restriction that z_i 's are binary imply that $z_i = 1$ if and only if $X_i^T\beta \geq 0$. It follows that the maximum value of (C.4) for each value of β is the same as the value of (C.2).

The problem (C.3) is a linear optimization problem with linear inequality constraints and integer constraints on z_i 's if the set \mathbf{B} is defined by linear inequalities that could be passed to any MILP solver. Florios and Skouras (2008) imposed only one side of the inequality constraint (C.4) for each i . For $g_i > 0$, it is sufficient to impose only the upper bound on z_i , and for $g_i < 0$, only the lower bound. The other side of the bound is always satisfied by the solution due to the direction of the objective function.

Our formulation has significant advantages. Despite a larger number of inequalities, it reduces the computation time in our applications by a factor of 10–40. Furthermore, it is not sufficient to impose only one side of the inequalities on z_i 's for optimization with a capacity constraint considered further below.

Our data contain large sets of observations that differ from each other in only one covariate. Suppose that m observations i_1, \dots, i_m differ only in the value of the last covariate: $X_{i_1} = (1, \tilde{x}_1, \dots, \tilde{x}_{d_x-1}, x_{d_x, i_1}), \dots, X_{i_m} = (1, \tilde{x}_1, \dots, \tilde{x}_{d_x-1}, x_{d_x, i_m})$, and are ordered with $x_{d_x, i_1} \leq x_{d_x, i_2} \leq \dots \leq x_{d_x, i_m}$. Then the solution must satisfy either $z_{i_1} \leq z_{i_2} \leq \dots \leq z_{i_m}$ or $z_{i_1} \geq z_{i_2} \geq \dots \geq z_{i_m}$. We found it advantageous to split the optimization problem in our empirical application into two: one explicitly imposing $z_{i_1} \leq \dots \leq z_{i_m}$ and one explicitly imposing $z_{i_1} \geq \dots \geq z_{i_m}$ for sets of observations that have the same values of education, but different values of prior earnings.

Inference on the welfare gain $V(\hat{G}_{\text{EWM}})$ of the empirical welfare-maximizing policy requires computing $\bar{v}_n^* = \sup_{G \in \mathcal{G}} \sqrt{n}(V_n^*(G) - V_n(G))$ in each bootstrap sample. Denoting the bootstrap weights by $\{w_i^*\}$, $\sum_{i=1}^n w_i^* = n$, \bar{v}_n^* could be expressed as

$$\bar{v}_n^* = \sqrt{n} \sup_{G \in \mathcal{G}} \sum_{i=1}^n (w_i^* - 1) g_i \cdot 1\{X_i^T \beta \geq 0\}. \quad (\text{C.5})$$

The optimization problem for \bar{v}_n^* is analogous to the optimization problem for \hat{G}_{EWM} . Furthermore, solving it does not require the knowledge of \hat{G}_{EWM} ; hence all bootstrap computations could be performed in parallel with the main EWM problem.

C.2. Multiple Linear Index Rules

We extend this method to compute treatment rules based on multiple linear scores. These rules construct J scores that are linear in covariates (or in their functions) and assign an individual to treatment if each score exceeds a specific threshold. An example of a multiple index treatment rule with three indices is when an individual is assigned to a job training program if $(25 \leq \text{age} \leq 35)$ AND (wage at the previous job $< \$15$). The results are easily extended to treatment rules that apply if any of the indices exceeds its threshold, for example, $(\text{age} \geq 40)$ OR (length of unemployment ≥ 2 years).

Let the treatment assignment set G be defined as an intersection of upper contour sets of J linear functions:

$$\begin{aligned} \mathcal{G} &= \{G_{\beta^1, \dots, \beta^J}, \beta^1, \dots, \beta^J \in \mathbf{B}\}, \\ G_{\beta^1, \dots, \beta^J} &= \{x : x^T \beta^1 \geq 0, \dots, x^T \beta^J \geq 0\}. \end{aligned}$$

Then the optimization problem (C.1) becomes

$$\max_{\beta^1, \dots, \beta^J \in \mathbf{B}} \sum_{i=1}^n g_i \cdot 1\{X_i^T \beta^1 \geq 0, \dots, X_i^T \beta^J \geq 0\}. \quad (\text{C.6})$$

We propose its equivalent formulation as a MILP problem with auxiliary binary variables $\{(z_i^1, \dots, z_i^J, z_i^*), i = 1, \dots, n\}$:

$$\max_{\substack{\beta^1, \dots, \beta^J \in \mathbf{B}, \\ z_i^1, \dots, z_i^J, z_i^* \in \mathbb{R}}} \sum_{i=1}^n g_i \cdot z_i^* \quad (\text{C.7})$$

$$\text{s.t. } \frac{X_i^T \beta^j}{C_i} < z_i^j \leq 1 + \frac{X_i^T \beta^j}{C_i} \quad \text{for } 1 \leq i \leq n, 1 \leq j \leq J, \quad (\text{C.8})$$

$$1 - J + \sum_{j=1}^J z_i^j \leq z_i^* \leq J^{-1} \sum_{j=1}^J z_i^j \quad \text{for } 1 \leq i \leq n, \quad (\text{C.9})$$

$$z_i^1, \dots, z_i^J, z_i^* \in \{0, 1\} \quad \text{for } 1 \leq i \leq n.$$

Similarly to the single index problem, the inequalities (C.8) and the constraint that z_i^j 's are binary imply together that $z_i^j = 1\{X_i^T \beta^j \geq 0\}$. Linear inequalities (C.9) and the binary constraints imply together that

$$z_i^* = z_i^1 \cdot \dots \cdot z_i^J = 1\{X_i^T \beta^1 \geq 0\} \cdot \dots \cdot 1\{X_i^T \beta^J \geq 0\}.$$

The problem for a collection of sets defined by the union of linear inequalities

$$G_{\beta^1, \dots, \beta^J} = \{X : X^T \beta^1 \geq 0 \text{ or } \dots \text{ or } X^T \beta^J \geq 0\}$$

could also be written as a MILP problem with the inequality constraint (C.9) replaced by

$$J^{-1} \sum_{j=1}^J z_i^j \leq z_i^* \leq \sum_{j=1}^J z_i^j \quad \text{for } i = 1, \dots, n. \quad (\text{C.10})$$

C.3. Optimization With a Capacity Constraint

When there is a capacity constraint K on the proportion of population that could be assigned to treatment 1, Empirical Welfare Maximization problem (2.4) on a set \mathcal{G} of half-spaces becomes

$$\max_{\beta \in \mathbf{B}} \left[\min \left\{ 1, \frac{Kn}{\sum_{i=1}^n 1\{X_i^T \beta \geq 0\}} \right\} \sum_{i=1}^n g_i \cdot 1\{X_i^T \beta \geq 0\} \right]. \quad (\text{C.11})$$

This problem cannot be rewritten as a linear optimization problem in the same way as (C.3) because the factor $\min\{1, \frac{Kn}{\sum_{i=1}^n 1\{X_i^T \beta \geq 0\}}\}$ varies with β . This factor could take fewer than n different values and the maximum of (C.11) could be obtained by solving a sequence of optimization problems each of which holds this factor constant:

For $k = \lfloor Kn \rfloor, \dots, n$

$$\begin{aligned} & \max_{\substack{\beta \in \mathbf{B}, \\ z_1, \dots, z_n \in \mathbb{R}}} \min \left\{ 1, \frac{Kn}{k} \right\} \sum_{i=1}^n g_i \cdot z_i \\ & \text{s.t.} \quad \frac{X_i^T \beta}{C_i} < z_i \leq 1 + \frac{X_i^T \beta}{C_i} \quad \text{for } 1 \leq i \leq n, \\ & \quad z_i \in \{0, 1\}, \\ & \quad \sum_{i=1}^n z_i \leq k. \end{aligned}$$

The capacity constrained problem with multiple indexes could be solved similarly.

REFERENCES

- ANDREWS, D., AND X. SHI (2013): "Inference Based on Conditional Moment Inequalities," *Econometrica*, 81, 609–666. [22]
- BOUSQUET, O. (2002): "A Bennet Concentration Inequality and Its Application to Suprema of Empirical Processes," *Comptes Rendus de l'Académie des Sciences—Series I*, 334, 495–500. [10]
- DONALD, S. G., AND Y.-C. HSU (2016): "Improving the Power of Tests of Stochastic Dominance," *Econometric Reviews*, 35 (4), 553–585. [22]
- DUDLEY, R. M. (1999): *Uniform Central Limit Theorems*. Cambridge: Cambridge University Press. [2]
- FLORIOS, K., AND S. SKOURAS (2008): "Exact Computation of Max Weighted Score Estimators," *Journal of Econometrics*, 146 (1), 86–91. [24]
- LINTON, O., K. SONG, AND Y. WHANG (2010): "An Improved Bootstrap Test of Stochastic Dominance," *Journal of Econometrics*, 154, 186–202. [22]
- LUGOSI, G. (2002): "Pattern Classification and Learning Theory," in *Principles of Nonparametric Learning*, ed. by L. Györfi. Vienna: Springer, 1–56. [2]
- MASSART, P., AND É. NÉDÉLEC (2006): "Risk Bounds for Statistical Learning," *The Annals of Statistics*, 34 (5), 2326–2366. [2]
- TSYBAKOV, A. B. (2004): "Optimal Aggregation of Classifiers in Statistical Learning," *The Annals of Statistics*, 32 (1), 135–166. [11]
- VAN DER VAART, A. W., AND J. A. WELLNER (1996): *Weak Convergence and Empirical Processes*. New York: Springer. [2,5,19,23]

Co-editor Liran Einav handled this manuscript.

Manuscript received 9 March, 2015; final version accepted 29 November, 2017; available online 29 November, 2017.