

A PROOFS OF THE THEOREMS

Whenever the context of an expectation operation is not clear, we disambiguate it by specifying the variable that the expectation is taken over and its distribution $\mathbb{E}_{\mathbf{x} \sim f(\mathbf{x})}[\mathbf{x}]$.

A.1 PROOF OF THEOREM 1

Proof. Given that the logarithm function is a strictly increasing function, we can omit it in the optimization; i.e., $\boldsymbol{\lambda}^* = \operatorname{argmin}_{\boldsymbol{\lambda}} \mathbb{E} [\exp(\mathbf{g}^\top \boldsymbol{\lambda} + \ell)]$. Because this is an unconstrained optimization, the optimal solution occurs when the gradient is equal to zero.

$$\begin{aligned} \mathbb{E}[\mathbf{g} \exp(\mathbf{g}^\top \boldsymbol{\lambda}^* + \ell)] &= \mathbf{0}, \\ \mathbb{E}[\mathbf{g} w^*(\mathbf{a}, \mathbf{x})] &= \mathbf{0}, \end{aligned} \quad (7)$$

where the last equation is due to the equation of the weights in the population optimization problem.

Using the definition for the \mathbf{g} vector, Eq. (7) implies that $\mathbb{E}[w^*(\mathbf{a}, \mathbf{x}) \mathbf{a} \phi(\mathbf{x})] = \mathbf{0}$. Thus, we conclude that in the weighted population (with distribution \tilde{F}), the \mathbf{a} and $\phi(\mathbf{x})$ are uncorrelated:

$$\mathbb{E}_{(\mathbf{a}, \mathbf{x}) \sim \tilde{F}}[\mathbf{a} \phi(\mathbf{x})] = \mathbf{0} \quad (8)$$

For every set $\mathcal{B} \subset \mathbb{A} \times \mathbb{X}$, we can write:

$$\tilde{F}(\mathcal{B}) = \int_{\mathcal{B}} w^*(a, \mathbf{x}) dF(a, \mathbf{x}). \quad (9)$$

The Radon-Nikodym theorem implies that $w^*(a, \mathbf{x})$ is the Radon-Nikodym derivative:

$$w^*(\mathbf{x}, a) = \frac{d\tilde{F}(\mathbf{x}, a)}{dF(\mathbf{x}, a)} = \frac{\tilde{f}(\mathbf{x}, a)}{f(\mathbf{x}, a)} \quad (10)$$

$$= \frac{\tilde{f}(\mathbf{x})\tilde{f}(a) + \left\{ \tilde{f}(\mathbf{x}, a) - \tilde{f}(\mathbf{x})\tilde{f}(a) \right\}}{f(\mathbf{x}, a)}, \quad (11)$$

$$= w_{GSW}(a, \mathbf{x}) + \frac{\tilde{f}(\mathbf{x}, a) - \tilde{f}(\mathbf{x})\tilde{f}(a)}{f(\mathbf{x}, a)} \quad (12)$$

Thus, using Eq. (8) and Assumptions 1 and 3 we can write

$$\sup_{a, \mathbf{x}} |w^*(a, \mathbf{x}) - w_{GSW}(a, \mathbf{x})| \leq \delta_{\phi_K}/c.$$

□

A.2 THEOREM 2

Proof. Given that the logarithm function is a strictly increasing function, we can omit it in the optimizations. Thus the sample and population solutions are:

$$\begin{aligned} \hat{\boldsymbol{\lambda}}_n &= \operatorname{argmin}_{\boldsymbol{\lambda}} \frac{1}{n} \sum_{i=1}^n \exp(\mathbf{g}_i^\top \boldsymbol{\lambda} + \ell_i), \\ \boldsymbol{\lambda}^* &= \operatorname{argmin}_{\boldsymbol{\lambda}} \mathbb{E} [\exp(\mathbf{g}^\top \boldsymbol{\lambda} + \ell)]. \end{aligned}$$

The estimator is an M-estimator and given our sample-splitting, the proof follows the asymptotic normality of the estimator (Van der Vaart, 2000, Chapter 5.3).

$$\sqrt{n} \left(\hat{\boldsymbol{\lambda}}_n - \boldsymbol{\lambda}^* \right) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{V}), \quad (13)$$

To obtain the value of \mathbf{V}_1 , note that the optimal sample solution occurs at the solution of the following equation (Z-estimator equation):

$$\sum_{i=1}^n \mathbf{g}_i \exp(\mathbf{g}_i^\top \hat{\boldsymbol{\lambda}}_n) = \mathbf{0}.$$

Thus, the score function is $\psi_{\lambda} = g_i \exp(g_i^{\top} \lambda)$. We denote the matrix of derivatives of the score function by $\dot{\psi}_{\lambda}$ whose elements are defined as $\dot{\psi}_{\lambda, kk'} = \partial \psi_{\lambda, k} / \partial \lambda_{k'}$. Using the theorem in (Van der Vaart, 2000, Chapter 5.3), we can write:

$$V = \mathbb{E}[\dot{\psi}_{\lambda^*}]^{-1} \mathbb{E}[\psi_{\lambda^*} \psi_{\lambda^*}^{\top}] \mathbb{E}[\dot{\psi}_{\lambda^*}]^{-1}. \quad (14)$$

In the above equation we have assumed that $\mathbb{E}[\dot{\psi}_{\lambda^*}]$ matrix is invertible. An unbiased sample estimation of V can be obtained by substituting $\hat{\lambda}_n$ in place of λ^* and taking empirical expectations.

An application of the delta method on Eq. (13) yields:

$$\sqrt{n} \left(\frac{\exp(g_i^{\top} \hat{\lambda}_n + \ell_i)}{\frac{1}{n} \sum_{i=1}^n \exp(g_i^{\top} \hat{\lambda}_n + \ell_i)} - \frac{\exp(g_i^{\top} \lambda^* + \ell_i)}{\frac{1}{n} \sum_{i=1}^n \exp(g_i^{\top} \lambda^* + \ell_i)} \right) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \sigma^2), \quad (15)$$

$$\sqrt{n} \left(\hat{w}_n(a_i, \mathbf{x}_i) - \frac{\exp(g_i^{\top} \lambda^* + \ell_i)}{\mathbb{E}[\exp(g^{\top} \lambda^* + \lambda)]} \right) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \sigma^2), \quad (16)$$

$$\sqrt{n} (\hat{w}_n(a_i, \mathbf{x}_i) - w^*(a_i, \mathbf{x}_i)) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \sigma^2), \quad (17)$$

where Eq. (16) is due to Slutsky's theorem and Eq. (17) is obtained by substitution of the definition for $w^*(a_i, \mathbf{x}_i)$. The variance is obtained by defining the Softmax function $s(\lambda) = \frac{\exp(g_i^{\top} \lambda + \ell_i)}{\frac{1}{n} \sum_{i=1}^n \exp(g_i^{\top} \lambda + \ell_i)}$. We denote the gradient of the Softmax function by $\nabla s(\lambda)$. We can write (Van der Vaart, 2000, Chapter 3):

$$\sigma^2(a_i, \mathbf{x}_i) = \nabla s(\lambda^*)^{\top} V \nabla s(\lambda^*).$$

Substituting the value of V from (14), we conclude:

$$\sigma^2(a_i, \mathbf{x}_i) = \nabla s(\lambda^*)^{\top} \mathbb{E}[\dot{\psi}_{\lambda^*}]^{-1} \mathbb{E}[\psi_{\lambda^*} \psi_{\lambda^*}^{\top}] \mathbb{E}[\dot{\psi}_{\lambda^*}]^{-1} \nabla s(\lambda^*).$$

Note that the value of the softmax function depends on the value of (a_i, \mathbf{x}_i) at each point. \square

B NEURAL NETWORK AND TRAINING DETAILS

B.1 DETAILS OF THE ℓ_{θ} NEURAL NETWORK

The ℓ_{θ} network is defined as follows:

$$\ell_{\theta}(z) = cz + \text{dense3}(\text{elu}(\text{layer_norm}(\text{dense2}(\tanh(\text{dense1}(z))))))$$

The linear term cz acts as a skip connection. The input and output dimensions for the dense linear layers are as follows:

$$\begin{aligned} \text{dense1} : 1 &\mapsto h, \\ \text{dense2} : h &\mapsto h, \\ \text{dense3} : h &\mapsto 1, \end{aligned}$$

where h denotes the hidden dimension. Because the softmax function is invariant to the constant shifts, we do not have any bias terms for dense3 and the skip connection. dense2 also does not have the bias because of the proceeding layer normalization. The dimension h has been tuned as a hyperparameter on a validation data and set to 10.

B.2 DETAILS OF THE PROPENSITY SCORE COMPUTATION FOR IPW

We model both $f(a)$ and $f(a|\mathbf{x})$ as univariate normal distributions. This is the correct assumption in our synthetic data. The marginal distribution $f(a)$ is estimated by simply finding the mean and standard deviation of the observed treatment values. For the conditional distribution, we write $a|\mathbf{x} \sim \mathcal{N}(\mu_a(\mathbf{x}), \sigma_{a|\mathbf{x}}^2)$, where $\mu_a(\mathbf{x})$ is modeled using a feedforward neural network with two layers and $\sigma_{a|\mathbf{x}}^2$ is estimated using the residuals of the neural network predictions. The dimension of the neural network has been tuned as a hyperparameter on validation data and set to 30.

B.3 FURTHER TRAINING DETAILS

We used PyTorch to implement E2B. For reproducibility purposes, we provide the final settings used for training:

- Learning algorithm: Adam with learning rate 0.001, no AMSGrad.
- Batch size: 100
- Max epochs: 400
- Weight decay: 2.5×10^{-5} .
- Validation on a dataset of size 400, every 10 steps.

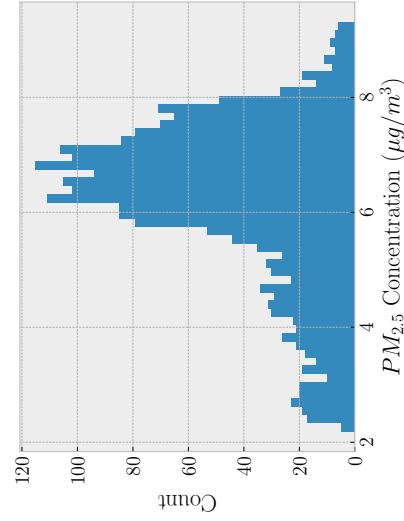
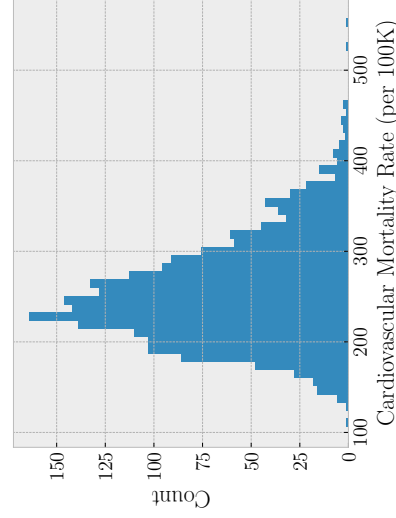
B.4 DETAIL OF PERMUTATION WEIGHTING

We created a stacked data by stacking $\{(\mathbf{x}_i, a_i, a_i \odot \mathbf{x}_i)\}_{i=1}^n$ and $\{(\mathbf{x}_i, \tilde{a}_i, \tilde{a}_i \odot \mathbf{x}_i)\}_{i=1}^n$, where \tilde{a} are permutations of the original treatments. We trained a random forest classifier to predict whether each data is from the permuted or the original set. We tried both random forests and neural networks and obtained better results with the former. We also calibrated the predicted probabilities of the classifier before computation of the weights.

C DATA AND PREPROCESSING DESCRIPTION

Table 2: NSAPH Data Description

	PM2.5	CMR	healthfac	population	ses	unemploy	HH_inc	femaleHH	vacant	owner_occ	eduattain	pctfam_pover
count	2132.0	2132.0	2132.0	2132.0	2132.0	2132.0	2132.0	2132.0	2132.0	2132.0	2132.0	2132.0
mean	6.17	255.25	0.18	10.78	0.0	7.85	10.69	11.92	14.25	71.44	35.03	11.25
std	1.45	56.76	0.5	1.26	0.96	2.83	0.24	3.94	8.71	7.76	7.07	5.2
min	2.19	106.14	-2.85	6.2	-1.84	0.0	9.91	2.1	3.8	19.3	9.4	0.0
25%	5.51	215.38	0.0	10.04	-0.67	6.0	10.54	9.3	8.8	67.7	30.4	7.6
50%	6.43	248.16	0.14	10.62	-0.14	7.6	10.67	11.2	11.65	72.7	35.4	10.55
75%	7.15	288.77	0.34	11.46	0.47	9.3	10.83	13.6	16.6	76.7	39.9	13.82
max	9.3	557.43	3.33	16.07	6.46	30.9	11.66	38.0	74.0	89.7	54.6	44.9

(a) Histogram of $PM_{2.5}$ 

(b) Histogram of Cardiovascular Mortality Rate

Figure 3: The Histograms of Data