
Appendix:

An Adaptive Kernel Approach to Federated Learning of Heterogeneous Causal Effects

Thanh Vinh Vo¹ Arnab Bhattacharyya¹ Young Lee² Tze-Yun Leong¹

¹School of Computing, National University of Singapore

²Roche AG and Harvard University

{votv, arnabb, leongty}@nus.edu.sg

Contents

A	Pre-training step to remove duplicated individuals	2
B	Identification	3
C	Computing CATE, local ATE, and global ATE	3
C.1	Computing the CATE and local ATE	3
C.2	Computing the global ATE from local ATE of each Source	4
D	Comparison metrics	4
E	Derivation of the loss functions	4
E.1	Learning distributions involving latent confounder	4
E.2	Learning auxiliary distributions	5
F	Spectral distribution of some popular kernels	6
G	Proof of Lemma 1	6
H	Proof of Lemma 2	9
H.1	Proof of Part (i)	9
H.2	Proof of Part (ii)	11
I	Further cases of the minimax lower bounds	12
I.1	Further cases of Lemma 1	12
I.2	Further cases of Lemma 2	13
J	Description of IHDP data	14

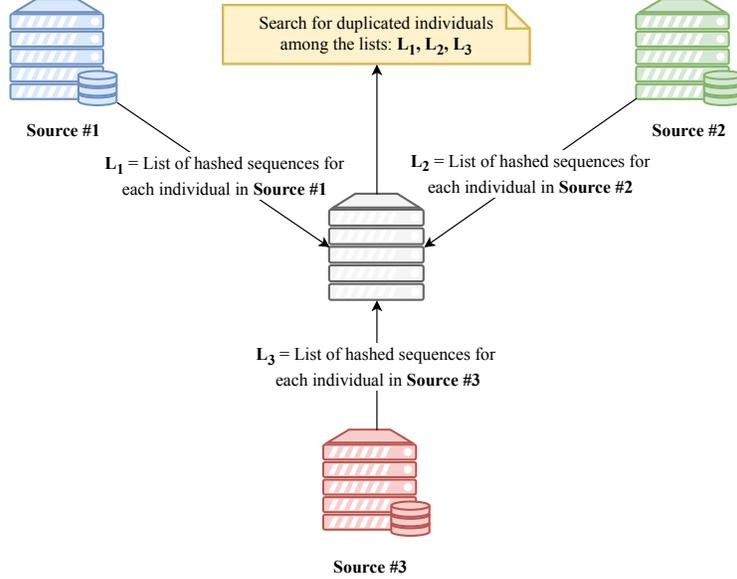


Figure 1: An illustration on how the pre-training step. This step is intended to identify duplicated individuals among the sources. Furthermore, this step preserves privacy since each source sends only their hashed sequences of the individuals.

A Pre-training step to remove duplicated individuals

As mentioned in the main text, we make five assumptions as follows:

- (A1) *Consistency*: $W = w \implies Y(w) = Y$, this follows from the axioms of structural causal model.
- (A2) *No interference*: treatment on one subject does not affect the outcomes of another one. This is because the outcome only has a single node for treatment as a parent.
- (A3) *Positivity* (also known as *Overlap*): every subject has some positive probability to be assigned to every treatment.
- (A4) The individuals in each source must have the same set of *common* covariates.
- (A5) There is no individual whose data exists in more than one source.

Assumptions (A1), (A2) and (A3) are standard in any causal inference algorithm.

Assumption (A4) has been implicitly shown in our setup since all the sources would share the same causal graph. This is a reasonable assumption as we intend to build a unified model on all of the data sources. For example, decentralized data in Choudhury et al. (2019); Vaid et al. (2020); Flores et al. (2020) (to name a few) satisfy this assumption for federated learning.

Assumption (A5) is to ensure that no individuals would dominate the other individuals when training the model. For example, if an individual appears in all of the sources, the trained model would be biased by data of this individual (there is imbalance caused by the use of more data from this particular individual than the others). Hence, this condition would ensure that such bias does not exist.

In practice, Assumption (A5) sometimes does not hold. To address such a problem, we propose a pre-training step to exclude such duplicated individuals. The pre-training step are summarized as follows:

- (1) Suppose that an individual can be uniquely identified via a set of features. For example, a pair of (national identity, nationality) can be used to uniquely identify a person.
- (2) To identify duplicated individuals, we first encode the above features with a hash function such as MD5, SHA256.

- (3) We then send the encoded sequences to a central server.
- (4) The server would collect all encoded sequences from all sources and find among them if an encoded sequence is repeated.
- (5) All of the repeated sequences are associated with duplicated individuals. Thus, we announce the sources to exclude these individual from the training process.

We summarize the pre-training step in Figure 1 with three sources of data.

B Identification

The causal effects are unidentifiable if the confounders are unobserved. However, Louizos et al. (2017) showed that if the joint distribution $p_s(\mathbf{x}^s, y^s, w^s, \mathbf{z}^s)$ can be recovered, then the causal effects are identifiable. In the following, we show how they are identifiable.

Proof. The proof is adapted from Louizos et al. (Theorem 1, 2017). We need to show that the distribution $p_s(y^s | \text{do}(W = w^s), \mathbf{x}^s)$ is identifiable from observational data. We have

$$\begin{aligned} p_s(y^s | \text{do}(W = w^s), \mathbf{x}^s) &= \int p_s(y^s | \text{do}(W = w^s), \mathbf{x}^s, \mathbf{z}^s) p_s(\mathbf{z}^s | \text{do}(W = w^s), \mathbf{x}^s) d\mathbf{z}^s \\ &= \int p_s(y^s | w^s, \mathbf{x}^s, \mathbf{z}^s) p_s(\mathbf{z}^s | \mathbf{x}^s) d\mathbf{z}^s. \end{aligned}$$

where the last equality is obtained by applying the *do*-calculus. The last expression, $\int p_s(y^s | w^s, \mathbf{x}^s, \mathbf{z}^s) p_s(\mathbf{z}^s | \mathbf{x}^s) d\mathbf{z}^s$, can be identified by the joint distribution $p_s(\mathbf{x}^s, y^s, w^s, \mathbf{z}^s)$. In our work, $p_s(\mathbf{x}^s, y^s, w^s, \mathbf{z}^s)$ is recovered by its factorization with the distributions $p_s(w^s | \mathbf{x}^s)$, $p_s(y^s | \mathbf{x}^s, w^s)$, $p_s(\mathbf{z}^s | \mathbf{x}^s, y^s, w^s)$, $p_s(y^s | w^s, \mathbf{z}^s)$, and $p(\mathbf{z}^s)$. Adaptively learning these distributions in a federated setting is the main task of our work. This completes the proof. \square

C Computing CATE, local ATE, and global ATE

This section gives details on how to compute CATE, local ATE and global ATE after training the model.

C.1 Computing the CATE and local ATE

After training the model, each source *can* compute the CATE and the local ATE on for its own source and use it for itself.

$$E[y_i^s | \text{do}(w_i^s = w), \mathbf{x}_i^s] = \int E[y_i^s | w_i^s = w, \mathbf{z}_i^s] p(\mathbf{z}_i^s | \mathbf{x}_i^s) d\mathbf{z}_i^s \simeq \frac{1}{N} \sum_{l=1}^N f_y(w_i^s = w, \mathbf{z}_i^s[l])$$

where $f_y(w_i^s = w, \mathbf{z}_i^s[l])$ is the mean function of $p_s(y_i^s | w_i^s, \mathbf{z}_i^s)$ and $\{\mathbf{z}_i^s[l]\}_{l=1}^N \stackrel{i.i.d.}{\sim} p_s(\mathbf{z}_i^s | \mathbf{x}_i^s)$.

The problem is to draw $\{\mathbf{z}_i^s[l]\}_{l=1}^N$ from $p_s(\mathbf{z}_i^s | \mathbf{x}_i^s)$. We observe that

$$p_s(\mathbf{z}_i^s | \mathbf{x}_i^s) = \sum_{w_i^s \in \{0,1\}} \int p_s(\mathbf{z}_i^s | \mathbf{x}_i^s, y_i^s, w_i^s) p_s(y_i^s | \mathbf{x}_i^s, w_i^s) p_s(w_i^s | \mathbf{x}_i^s) dy_i^s.$$

Hence, to draw samples, we proceed in the following steps:

- (1) Draw a sample of w_i^s from $p_s(w_i^s | \mathbf{x}_i^s)$.
- (2) Substitute the above sample of w_i^s to $p_s(y_i^s | \mathbf{x}_i^s, w_i^s)$.
- (3) Draw a sample of y_i^s from $p_s(y_i^s | \mathbf{x}_i^s, w_i^s)$.
- (4) Substitute the above sample of y_i^s to $p_s(\mathbf{z}_i^s | \mathbf{x}_i^s, y_i^s, w_i^s)$.
- (5) Draw a sample of \mathbf{z}_i^s from $p_s(\mathbf{z}_i^s | \mathbf{x}_i^s, y_i^s, w_i^s)$.

The density function of $p_s(y_i^s|\mathbf{x}_i^s, w_i^s)$ and $p_s(w_i^s|\mathbf{x}_i^s)$ are available after training the model. As described in the main text, there are two options to draw from $p_s(\mathbf{z}_i^s|\mathbf{x}_i^s, y_i^s, w_i^s)$. The first option is to draw from $q(\mathbf{x}_i^s)$ since it approximates $p_s(\mathbf{z}_i^s|\mathbf{x}_i^s, y_i^s, w_i^s)$. The second option is to use Metropolis-Hastings algorithm with independent sampler (Liu 1996). For the second option, we have that

$$p_s(\mathbf{z}_i^s|\mathbf{x}_i^s, y_i^s, w_i^s) \propto p_s(y_i^s|\mathbf{z}_i^s, w_i^s)p_s(w_i^s|\mathbf{z}_i^s)p_s(\mathbf{x}_i^s|\mathbf{z}_i^s)p(\mathbf{z}_i^s).$$

Hence, it can be used to compute the acceptance probability of interest. Note that the second option would give more exact samples since it further filters the samples based on the exact acceptance probability.

The above would help estimate the CATE given \mathbf{x}_i^s . The local ATE is the average of CATE of individuals in a source s . These quantities can be estimated in a local source’s machine. We show how to compute the global ATE in the next section.

C.2 Computing the global ATE from local ATE of each Source

To compute a global ATE, the server would collect all the local ATE in each source and then compute their weighted average. For example, suppose that we have three sources whose local ATE values are 7.0, 8.5, and 6.8. These local ATEs are averaged over 10, 5, and 12 individuals, in that order. Then, the global ATE is given as follows:

$$\text{global ATE} = \frac{10 \times 7.0 + 5 \times 8.5 + 12 \times 6.8}{10 + 5 + 12} = 7.32.$$

Since each source only shares their local ATE and the number of individuals, it does not leak any sensitive information about the individuals.

D Comparison metrics

We report two error metrics in our experiments:

- Precision in estimation of heterogeneous effects (PEHE):

$$\epsilon_{\text{PEHE}} = \sum_{i=1}^n (\tau(\mathbf{x}_i) - \hat{\tau}(\mathbf{x}_i))^2 / n, \quad (1)$$

- Absolute error:

$$\epsilon_{\text{ATE}} = |\tau - \hat{\tau}|, \quad (2)$$

where $\tau(\mathbf{x}_i)$, τ are the ground truth of ITE and ATE, and $\hat{\tau}(\mathbf{x}_i)$, $\hat{\tau}$ are their estimates. We report the mean and standard error over 10 replicates of the data with different random initializations of the training algorithm.

E Derivation of the loss functions

In this section, we present the loss functions and the form of functions that modulate the desired distributions.

E.1 Learning distributions involving latent confounder

The ELBO of the log marginal likelihood has the following expression

$$\begin{aligned} \log p(\mathbf{x}, \mathbf{y}, \mathbf{w}) &= \log \int p(\mathbf{x}, \mathbf{y}, \mathbf{w}, \mathbf{z}) d\mathbf{z} \\ &\geq \int q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{y}, \mathbf{w}, \mathbf{z})}{q(\mathbf{z})} d\mathbf{z} \\ &= \sum_{s \in \mathcal{S}} \sum_{i=1}^{n_s} \left(E_q \left[\log p_s(y_i^s | w_i^s, \mathbf{z}_i^s) + \log p_s(w_i^s | \mathbf{z}_i^s) + \log p_s(\mathbf{x}_i^s | \mathbf{z}_i^s) \right] - \text{KL}[q(\mathbf{z}_i^s) \| p(\mathbf{z}_i^s)] \right) =: \mathcal{L}. \end{aligned}$$

Using the complete dataset $\tilde{D}^s = \bigcup_{l=1}^M \{(w_i^s, y_i^s, \mathbf{x}_i^s, \mathbf{z}_i^s[l])\}_{i=1}^{n_s}, \forall s \in \mathcal{S}$, we minimize the following loss function J :

$$J = \widehat{\mathcal{L}} + \sum_{c \in \mathcal{A}} R(f_c), \quad \mathcal{A} = \{y_0, y_1, q_0, q_1, x, w\},$$

where $\widehat{\mathcal{L}}$ is the empirical loss function obtained from the negative of \mathcal{L} . In the following, we find the form of f_c based on the representer theorem.

We further define $f_x = [f_{x,1}, \dots, f_{x,d_x}]$, where $f_{x,d}$ is a function taking \mathbf{z}_i^s as input and mapping it to a real value in \mathbb{R} . Similarly, $f_{q_0} = [f_{q_0,1}, \dots, f_{q_0,d_z}]$ and $f_{q_1} = [f_{q_1,1}, \dots, f_{q_1,d_z}]$.

Let \mathcal{H}_c ($c \in \mathcal{A}$) be a reproducing Kernel Hilbert space (RKHS) and $\kappa_c(\cdot, \cdot)$ be kernel function associated with \mathcal{H}_c . We define \mathcal{B}_c as follows:

$$\begin{aligned} \mathcal{B}_{y_0} &= \text{span}\{\kappa_{y_0}(\cdot, \mathbf{z}_i^s[l]), \text{ where } s \in \mathcal{S}; i = 1, \dots, n_s; l = 1, \dots, M\}, \\ \mathcal{B}_{y_1} &= \text{span}\{\kappa_{y_1}(\cdot, \mathbf{z}_i^s[l]), \text{ where } s \in \mathcal{S}; i = 1, \dots, n_s; l = 1, \dots, M\}, \\ \mathcal{B}_x &= \text{span}\{\kappa_x(\cdot, \mathbf{z}_i^s[l]), \text{ where } s \in \mathcal{S}; i = 1, \dots, n_s; l = 1, \dots, M\}, \\ \mathcal{B}_w &= \text{span}\{\kappa_w(\cdot, \mathbf{z}_i^s[l]), \text{ where } s \in \mathcal{S}; i = 1, \dots, n_s; l = 1, \dots, M\}, \\ \mathcal{B}_{q_0} &= \text{span}\{\kappa_{q_0}(\cdot, [\mathbf{x}_i^s, y_i^s]), \text{ where } s \in \mathcal{S}; i = 1, \dots, n_s\}, \\ \mathcal{B}_{q_1} &= \text{span}\{\kappa_{q_1}(\cdot, [\mathbf{x}_i^s, y_i^s]), \text{ where } s \in \mathcal{S}; i = 1, \dots, n_s\}. \end{aligned}$$

We posit the following regularizers:

$$R(f_{y_0}) = \text{reg_factor}_{y_0} \times \|f_{y_0}\|_{\mathcal{H}_{y_0}}^2, \quad R(f_x) = \sum_{d=1}^{d_x} \text{reg_factor}_{x,d} \times \|f_{x,d}\|_{\mathcal{H}_{x,d}}^2 \quad (d = 1, \dots, d_x).$$

The regularizers $R(f_{y_1})$ and $R(f_w)$ are similar to that of $R(f_{y_0})$, and $R(f_{q_0})$, $R(f_{q_1})$ are similar to that of $R(f_x)$.

We see that \mathcal{B}_c is a subspace of \mathcal{H}_c . We project f_{y_0} , f_{y_1} , f_w , $f_{x,d}$ ($d = 1, \dots, d_x$), $f_{q_0,d}$ ($d = 1, \dots, d_z$) and $f_{q_1,d}$ ($d = 1, \dots, d_z$) onto the subspaces \mathcal{B}_{y_0} , \mathcal{B}_{y_1} , \mathcal{B}_w , \mathcal{B}_x , \mathcal{B}_{q_0} and \mathcal{B}_{q_1} , respectively, and obtain f'_{y_0} , f'_{y_1} , f'_w , $f'_{x,d}$, $f'_{q_0,d}$ and $f'_{q_1,d}$. Next, we also project them onto the perpendicular spaces of $\mathcal{B}_{(\cdot)}$ to obtain $f_{y_0}^\perp$, $f_{y_1}^\perp$, f_w^\perp , $f_{x,d}^\perp$, $f_{q_0,d}^\perp$ and $f_{q_1,d}^\perp$.

Note that $f_{(\cdot)} = f'_{(\cdot)} + f_{(\cdot)}^\perp$. Hence, $\|f_{(\cdot)}\|_{\mathcal{H}_{(\cdot)}}^2 = \|f'_{(\cdot)}\|_{\mathcal{H}_{(\cdot)}}^2 + \|f_{(\cdot)}^\perp\|_{\mathcal{H}_{(\cdot)}}^2 \geq \|f'_{(\cdot)}\|_{\mathcal{H}_{(\cdot)}}^2$, which implies that $\text{reg_factor}_{(\cdot)} \times \|f_{(\cdot)}\|_{\mathcal{H}_{(\cdot)}}^2$ is minimized if $f_{(\cdot)}$ is in its subspace $\mathcal{B}_{(\cdot)}$. **(I)**

In addition, due to the reproducing property, we have

$$f_{y_0}(\mathbf{z}_i^s[l]) = \langle f_{y_0}, \kappa_{y_0}(\cdot, \mathbf{z}_i^s[l]) \rangle_{\mathcal{H}_{y_0}} = \langle f'_{y_0}, \kappa_{y_0}(\cdot, \mathbf{z}_i^s[l]) \rangle_{\mathcal{H}_{y_0}} + \langle f_{y_0}^\perp, \kappa_{y_0}(\cdot, \mathbf{z}_i^s[l]) \rangle_{\mathcal{H}_{y_0}} = f'_{y_0}(\mathbf{z}_i^s[l]).$$

Similarly, we also have $f_{y_1}(\mathbf{z}_i^d[l]) = f'_{y_1}(\mathbf{z}_i^d[l])$, $f_w(\mathbf{z}_i^d[l]) = f'_w(\mathbf{z}_i^d[l])$, $f_{x,d}(\mathbf{z}_i^d[l]) = f'_{x,d}(\mathbf{z}_i^d[l])$, $f_{q_0,d}(y_i^d, \mathbf{x}_i^d) = f'_{q_0,d}(y_i^d, \mathbf{x}_i^d)$ and $f_{q_1,d}(y_i^d, \mathbf{x}_i^d) = f'_{q_1,d}(y_i^d, \mathbf{x}_i^d)$. Hence,

$$\widehat{\mathcal{L}}(f_{y_0}, f_{y_1}, f_{q_0}, f_{q_1}, f_x, f_w) = \widehat{\mathcal{L}}(f'_{y_0}, f'_{y_1}, f'_{q_0}, f'_{q_1}, f'_x, f'_w). \quad \textbf{(II)}$$

(I) and **(II)** imply that f_{y_0} , f_{y_1} , $f_{q_0,d}$, $f_{q_1,d}$, $f_{x,d}$, f_w are the weighted sum of elements in their corresponding subspace. Hence,

$$f_c(\mathbf{u}^s) = \sum_{\mathbf{v} \in \mathcal{S}} \sum_{j=1}^{n_v \times M} \kappa(\mathbf{u}^s, \mathbf{u}_j^{\mathbf{v}}) \alpha_j^{\mathbf{v}}.$$

Using this form with the adaptive kernel and Random Fourier Feature described in the main text (Section 4.1), we obtain the desired model.

E.2 Learning auxiliary distributions

The derivation of J_w , J_y and the form of functions modulated the auxiliary distributions are similar to those of J as detailed in Section E.1. The difference is that the empirical loss functions are obtained from the negative log-likelihood instead of the ELBO.

F Spectral distribution of some popular kernels

Table 1 (adopted from Milton et al. (2019)) presents some popular kernels and their associated spectral density $s(\boldsymbol{\omega})$. Those density functions are needed to draw samples of $\boldsymbol{\omega}$ for Random Fourier Features presented in Section 4 of the main text. In our experiments, we used Gaussian kernel.

Table 1: Some popular kernels and their associated spectral density. Note that $K_\nu(\cdot)$ denotes the modified Bessel function of the second kind, $\Gamma(\cdot)$ is the gamma function.

Kernel	Kernel function, $k(\mathbf{x}_1 - \mathbf{x}_2)$	Spectral density, $s(\boldsymbol{\omega})$
Gaussian	$\exp\left(-\frac{\ \mathbf{x}_1 - \mathbf{x}_2\ _2^2}{2\ell^2}\right)$	$\left(\frac{2\pi}{\ell^2}\right)^{-\frac{d}{2}} \exp\left(-\frac{\ell^2\ \boldsymbol{\omega}\ _2^2}{2}\right)$
Laplacian	$\exp\left(-\ell\ \mathbf{x}_1 - \mathbf{x}_2\ _1\right)$	$\left(\frac{2}{\pi}\right)^{\frac{d}{2}} \prod_{i=1}^d \frac{\ell}{\ell^2 + \omega_i^2}$
Matérn	$\frac{2^{1-\nu}}{\Gamma(\nu)} \left(\sqrt{2\nu} \frac{\ \mathbf{x}_1 - \mathbf{x}_2\ _2}{\ell}\right)^\nu K_\nu\left(\sqrt{2\nu} \frac{\ \mathbf{x}_1 - \mathbf{x}_2\ _2}{\ell}\right)$	$\frac{2^d \pi^{\frac{d}{2}} \Gamma(\nu + \frac{d}{2})(2\nu)^\nu}{\Gamma(\nu)\ell^{2\nu}} \left(\frac{2\nu}{\ell^2} + 4\pi^2\ \boldsymbol{\omega}\ _2^2\right)^{-(\nu + \frac{d}{2})}$

G Proof of Lemma 1

We repeat Lemma 1 here for convenience:

Lemma 1 (With presence of latent variables). *Let $\boldsymbol{\theta} = \{\theta^s\}_{s=1}^m$ and $\hat{\boldsymbol{\theta}}$ be its estimate. Let $y_i^s \in \mathbb{R}$ and $\mathbf{x}_i^s \in \mathbb{R}^{d_x}$. Then,*

$$\inf_{\hat{\boldsymbol{\theta}}} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}(P)\|_2 \right] \geq \frac{\sqrt{m(d_x + 3)} \log(2\sqrt{m})}{64\sqrt{B} \sum_{s \in \mathcal{S}} n_s (1 + \sum_{v \in \mathcal{S}, v \neq s} \lambda^{s,v})^2}.$$

Let $\mathcal{S}_{\setminus s} := \mathcal{S} \setminus \{s\}$. The model is summarized as follows:

$$\begin{aligned} p(\mathbf{z}_i^s) &= \mathcal{N}(0, \sigma_z^2 \mathbf{I}_{d_z}), \\ p(w_i^s | \mathbf{z}_i^s) &= \text{Bern}\left(\varphi\left(\left(\theta_w^s + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \theta_w^v\right)^\top \phi(\mathbf{z}_i^s)\right)\right), \\ p(y_i^s | w_i^s, \mathbf{z}_i^s) &= \mathcal{N}\left(\left(w_i^s (\theta_{y1}^s + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \theta_{y1}^v) + (1 - w_i^s) (\theta_{y0}^s + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \theta_{y0}^v)\right)^\top \phi(\mathbf{z}_i^s), \sigma_y^2\right), \\ p(\mathbf{x}_i^s | \mathbf{z}_i^s) &= \mathcal{N}\left(\left(\theta_x^s + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \theta_x^v\right)^\top \phi(\mathbf{z}_i^s), \sigma_x^2 \mathbf{I}_{d_x}\right), \end{aligned}$$

where $\mathbf{z}_i^{(\cdot)} \in \mathbb{R}^{d_z}$, $y_i^{(\cdot)} \in \mathbb{R}$, $w_i^{(\cdot)} \in \{0, 1\}$, $\mathbf{x}_i^{(\cdot)} \in \mathbb{R}^{d_x}$, $\lambda > 0$.

Let $\boldsymbol{\theta} = \{\theta_w^s, \theta_{y0}^s, \theta_{y1}^s, \theta_x^s\}_{s \in \mathcal{S}}$. Let $\mathcal{V}_w, \mathcal{V}_{y0}, \mathcal{V}_{y1}, \mathcal{V}_x$ be $1/(2\sqrt{m})$ -packing of the unit $\|\cdot\|_2$ -balls with cardinality at least $(2\sqrt{m})^{2B}$, $(2\sqrt{m})^{2B}$, $(2\sqrt{m})^{2B}$, $(2\sqrt{m})^{2Bd_x}$, respectively. Let $\mathcal{V}^s = \delta(\mathcal{V}_w \times \mathcal{V}_{y0} \times \mathcal{V}_{y1} \times \mathcal{V}_x)$ and $\mathcal{V} = \mathcal{V}^{s_1} \times \mathcal{V}^{s_2} \times \dots \times \mathcal{V}^{s_m}$. We see that

$$|\mathcal{V}| \geq (2\sqrt{m})^{2mB(d_x+3)}.$$

In the following, we derive the minimax bound:

Proof. We have that

$$\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2 = \sqrt{\sum_{s \in \mathcal{S}} \sum_{c \in \mathcal{A}} \|(\theta_c^s)_1 - (\theta_c^s)_2\|_2^2} \geq \sqrt{\sum_{s \in \mathcal{S}} 4 \left(\frac{\delta}{2\sqrt{m}}\right)^2} = \delta.$$

The marginal distribution

$$p_{\Theta}(w, y, \mathbf{x}) = \int p_{\Theta}(w, y, \mathbf{x}, \mathbf{z}) d\mathbf{z} = \int p_{\Theta}(y|w, \mathbf{z}) p_{\Theta}(w|\mathbf{z}) p_{\Theta}(\mathbf{x}|\mathbf{z}) p(\mathbf{z}) d\mathbf{z}.$$

Moreover, we have that

$$D_{\text{KL}}(p_{\Theta_1}^n \| p_{\Theta_2}^n) = \sum_{s \in \mathcal{S}} D_{\text{KL}}(p_{\Theta_1}^{n_s} \| p_{\Theta_2}^{n_s}).$$

We divide the proof into three parts **(I)**, **(II)**, and **(III)**:

(I) The upper bound of $D_{\text{KL}}(p_{\Theta_1}^{n_s} \| p_{\Theta_2}^{n_s})$

Since the data is independent, we have that

$$\begin{aligned} D_{\text{KL}}(p_{\Theta_1}^{n_s} \| p_{\Theta_2}^{n_s}) &= n_s D_{\text{KL}}(p_{\Theta_1}^1 \| p_{\Theta_2}^1) \\ &\leq n_s \int D_{\text{KL}}\left(p_{\Theta_1}(y|w, \mathbf{z}) p_{\Theta_1}(w|\mathbf{z}) p_{\Theta_1}(\mathbf{x}|\mathbf{z}) \middle\| p_{\Theta_2}(y|w, \mathbf{z}') p_{\Theta_2}(w|\mathbf{z}') p_{\Theta_2}(\mathbf{x}|\mathbf{z}')\right) p(\mathbf{z}) p(\mathbf{z}') d\mathbf{z} d\mathbf{z}' \\ &= n_s \int \left[p_{\Theta_1}(w = 0|\mathbf{z}) D_{\text{KL}}[p_{\Theta_1}(y|w = 0, \mathbf{z}) \| p_{\Theta_2}(y|w = 0, \mathbf{z}')] \right. \\ &\quad + p_{\Theta_1}(w = 1|\mathbf{z}) D_{\text{KL}}[p_{\Theta_1}(y|w = 1, \mathbf{z}) \| p_{\Theta_2}(y|w = 1, \mathbf{z}')] \\ &\quad \left. + D_{\text{KL}}[p_{\Theta_1}(w|\mathbf{z}) \| p_{\Theta_2}(w|\mathbf{z}')] + D_{\text{KL}}[p_{\Theta_1}(\mathbf{x}|\mathbf{z}) \| p_{\Theta_2}(\mathbf{x}|\mathbf{z}')] \right] p(\mathbf{z}) p(\mathbf{z}') d\mathbf{z} d\mathbf{z}'. \end{aligned}$$

In the following, we find the upper bound of each component.

◇ Upper bound of the first and second component

$$\begin{aligned} p_{\Theta_1}(w = 0|\mathbf{z}) D_{\text{KL}}[p_{\Theta_1}(y|w = 0, \mathbf{z}) \| p_{\Theta_2}(y|w = 0, \mathbf{z}')] &\leq \frac{1}{2\sigma_y^2} \left(\left((\theta_{y0}^s)_1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_{y0}^v)_1 \right)^\top \phi(\mathbf{z}) - \left((\theta_{y0}^s)_2 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_{y0}^v)_2 \right)^\top \phi(\mathbf{z}') \right)^2 \\ &\leq \frac{8B^2 \delta^2 (1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v})^2}{\sigma_y^2}. \end{aligned}$$

Similarly, we also have

$$p_{\Theta_1}(w = 1|\mathbf{z}) D_{\text{KL}}[p_{\Theta_1}(y|w = 1, \mathbf{z}) \| p_{\Theta_2}(y|w = 1, \mathbf{z}')] \leq \frac{8B^2 \delta^2 (1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v})^2}{\sigma_y^2}.$$

◇ Upper bound of the third component

$$\begin{aligned} &D_{\text{KL}}[p_{\Theta_1}(w|\mathbf{z}) \| p_{\Theta_2}(w|\mathbf{z}')] \\ &= \varphi \left(\left((\theta_w^s)_1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_w^v)_1 \right)^\top \phi(\mathbf{z}) \right) \log \frac{\varphi \left(\left((\theta_w^s)_1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_w^v)_1 \right)^\top \phi(\mathbf{z}) \right)}{\varphi \left(\left((\theta_w^s)_2 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_w^v)_2 \right)^\top \phi(\mathbf{z}') \right)} \\ &\quad + \varphi \left(- \left((\theta_w^s)_1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_w^v)_1 \right)^\top \phi(\mathbf{z}) \right) \log \frac{\varphi \left(- \left((\theta_w^s)_1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_w^v)_1 \right)^\top \phi(\mathbf{z}) \right)}{\varphi \left(- \left((\theta_w^s)_2 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_w^v)_2 \right)^\top \phi(\mathbf{z}') \right)}. \end{aligned}$$

For the first component,

$$\begin{aligned} &\varphi \left(\left((\theta_w^s)_1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_w^v)_1 \right)^\top \phi(\mathbf{z}) \right) \log \frac{\varphi \left(\left((\theta_w^s)_1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_w^v)_1 \right)^\top \phi(\mathbf{z}) \right)}{\varphi \left(\left((\theta_w^s)_2 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_w^v)_2 \right)^\top \phi(\mathbf{z}') \right)} \\ &\leq \left| \log \left(1 + e^{-\left((\theta_w^s)_2 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_w^v)_2 \right)^\top \phi(\mathbf{z})} \right) - \log \left(1 + e^{-\left((\theta_w^s)_1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_w^v)_1 \right)^\top \phi(\mathbf{z}')} \right) \right| \end{aligned}$$

$$\begin{aligned}
&\leq \left\| (\theta_w^s)_1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_w^v)_1 \right\|_2 \|\phi(\mathbf{z})\|_2 + \left\| (\theta_w^s)_2 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_w^v)_2 \right\|_2 \|\phi(\mathbf{z}')\|_2 \\
&\leq \left(\delta + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \delta \right) \|\phi(\mathbf{z})\|_2 + \left(\delta + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \delta \right) \|\phi(\mathbf{z}')\|_2 \\
&\leq 4B\delta \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \right).
\end{aligned}$$

Similarly, we also have

$$\begin{aligned}
&\varphi \left(- \left((\theta_w^s)_1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_w^v)_1 \right)^\top \phi(\mathbf{z}) \right) \log \frac{\varphi \left(- \left((\theta_w^s)_1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_w^v)_1 \right)^\top \phi(\mathbf{z}) \right)}{\varphi \left(- \left((\theta_w^s)_2 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_w^v)_2 \right)^\top \phi(\mathbf{z}') \right)} \\
&\leq 4B\delta \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \right).
\end{aligned}$$

Thus,

$$D_{\text{KL}}[p_{\theta_1}(w|\mathbf{z}) \| p_{\theta_2}(w|\mathbf{z}')] \leq 8B\delta \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \right).$$

◇ Upper bound of the fourth component

$$\begin{aligned}
&D_{\text{KL}}[p_{\theta_1}(\mathbf{x}|\mathbf{z}) \| p_{\theta_2}(\mathbf{x}|\mathbf{z}')] \\
&= \frac{1}{2\sigma_x^2} \left\| \left((\theta_x^s)_1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_x^v)_1 \right)^\top \phi(\mathbf{z}) - \left((\theta_x^s)_2 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_x^v)_2 \right)^\top \phi(\mathbf{z}') \right\|_2^2 \\
&\leq \frac{1}{2\sigma_x^2} \left(\left\| \left((\theta_x^s)_1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_x^v)_1 \right)^\top \phi(\mathbf{z}) \right\|_2 + \left\| \left((\theta_x^s)_2 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} (\theta_x^v)_2 \right)^\top \phi(\mathbf{z}') \right\|_2 \right)^2 \\
&\leq \frac{8B^2\delta^2 \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \right)^2}{\sigma_x^2}.
\end{aligned}$$

(II) Combining the results

From the above upper bound of each of the components, we obtain

$$\begin{aligned}
D_{\text{KL}}(p_{\theta_1}^{n_s} \| p_{\theta_2}^{n_s}) &\leq n_s \int \left[\frac{16B^2\delta^2(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v})^2}{\sigma_y^2} + 8B\delta \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \right) \right. \\
&\quad \left. + \frac{8B^2\delta^2(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v})^2}{\sigma_x^2} \right] p(\mathbf{z})p(\mathbf{z}') d\mathbf{z}d\mathbf{z}' \\
&= n_s \left[\left(\frac{1}{\sigma_y^2} + \frac{1}{2\sigma_x^2} \right) 16B^2\delta^2 \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \right)^2 + 8B\delta \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \right) \right].
\end{aligned}$$

(III) The minimax lower bound

We have that

$$\begin{aligned}
D_{\text{KL}}(p_{\theta_1}^n \| p_{\theta_2}^n) &= \sum_{s \in \mathcal{S}} D_{\text{KL}}(p_{\theta_1}^{n_s} \| p_{\theta_2}^{n_s}) \\
&\leq \sum_{s \in \mathcal{S}} n_s \left[\left(\frac{1}{\sigma_y^2} + \frac{1}{2\sigma_x^2} \right) 16B^2\delta^2 \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \right)^2 + 8B\delta \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \right) \right].
\end{aligned}$$

Consequently,

$$\inf_{\hat{\theta}_n} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\|\hat{\theta}_n - \theta(P)\|_2 \right]$$

$$\begin{aligned} &\geq \frac{\delta}{2} \left(1 - \frac{\sum_{s \in \mathcal{S}} n_s \left[\left(\frac{1}{\sigma_y^2} + \frac{1}{2\sigma_x^2} \right) 16B^2 \delta^2 \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \right)^2 + 8B\delta \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \right) \right] + \log 2}{\log |\mathcal{V}|} \right) \\ &\geq \frac{\delta}{2} \left(1 - \frac{\sum_{s \in \mathcal{S}} n_s \left[\left(\frac{1}{\sigma_y^2} + \frac{1}{2\sigma_x^2} \right) 16B^2 \delta^2 \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \right)^2 + 8B\delta \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \right) \right] + \log 2}{2mB(d_x + 3) \log(2\sqrt{m})} \right). \end{aligned}$$

We choose $\delta = \frac{\sqrt{mB(d_x+3)} \log(2\sqrt{m})}{4B \sum_{s \in \mathcal{S}} n_s (1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v})^2}$, then

$$\begin{aligned} &1 - \frac{\sum_{s \in \mathcal{S}} n_s \left[\left(\frac{1}{\sigma_y^2} + \frac{1}{2\sigma_x^2} \right) 16B^2 \delta^2 \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \right)^2 + 8B\delta \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \right) \right] + \log 2}{2mB(d_x + 3) \log(2\sqrt{m})} \\ &\geq 1 - \left(\frac{1}{\sigma_y^2} + \frac{1}{2\sigma_x^2} \right) \frac{\log(2\sqrt{m})}{2 \sum_{s \in \mathcal{S}} n_s (1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v})^2} - \frac{1}{\sqrt{mB(d_x + 3)}} - \frac{1}{2mB(d_x + 3)} \\ &\geq 1 - \left(\frac{1}{\sigma_y^2} + \frac{1}{2\sigma_x^2} \right) \frac{\log(2\sqrt{m})}{2 \sum_{s \in \mathcal{S}} n_s (1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v})^2} - \frac{1}{2} - \frac{1}{8}. \end{aligned}$$

If $\sum_{s \in \mathcal{S}} n_s (1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v})^2 \geq 2 \left(\frac{1}{\sigma_y^2} + \frac{1}{2\sigma_x^2} \right) \log(2\sqrt{m})$, then

$$\begin{aligned} \inf_{\hat{\theta}_n} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\|\hat{\theta}_n - \theta(P)\|_2 \right] &\geq \frac{1}{2} \times \frac{\sqrt{mB(d_x + 3)} \log(2\sqrt{m})}{4B \sum_{s \in \mathcal{S}} n_s (1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v})^2} \times \left(1 - \frac{1}{4} - \frac{1}{2} - \frac{1}{8} \right) \\ &= \frac{\sqrt{m(d_x + 3)} \log(2\sqrt{m})}{64\sqrt{B} \sum_{s \in \mathcal{S}} n_s (1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v})^2}. \end{aligned}$$

This completes the proof. \square

H Proof of Lemma 2

We repeat Lemma 2 here for convenience:

Lemma 2 (Without the presence of latent variables). *Let $\boldsymbol{\psi} = \{\psi^s\}_{s=1}^m$, $\boldsymbol{\beta} = \{\beta^s\}_{s=1}^m$ and $\hat{\boldsymbol{\psi}}$, $\hat{\boldsymbol{\beta}}$ be their estimates, respectively. Let $y_i^s \in \mathbb{R}$. Then,*

$$\begin{aligned} \text{(i)} \quad \inf_{\hat{\boldsymbol{\psi}}} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\|\hat{\boldsymbol{\psi}} - \boldsymbol{\psi}(P)\|_2 \right] &\geq \frac{m \log(2\sqrt{m})}{256 \sum_{s \in \mathcal{S}} n_s (1 + \sum_{v \in \mathcal{S}, v \neq s} \gamma^{s,v})}, \\ \text{(ii)} \quad \inf_{\hat{\boldsymbol{\beta}}} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}(P)\|_2 \right] &\geq \frac{\sigma}{16\sqrt{2}} \sqrt{\frac{m \log(2\sqrt{m})}{B \sum_{s \in \mathcal{S}} n_s (1 + \sum_{v \in \mathcal{S}, v \neq s} \eta^{s,v})^2}}. \end{aligned}$$

The proof of Lemma 2 is divided into two parts (i) and (ii). We compute them separately:

H.1 Proof of Part (i)

We summarize the model as follows

$$w^s \sim \text{Bern} \left(\varphi \left(\left(\psi^s + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \psi^v \right)^\top \phi(\mathbf{x}^s) \right) \right).$$

Let $\boldsymbol{\psi} = \{\psi^s\}_{s \in \mathcal{S}}$. Let \mathcal{V}_s be $1/(2\sqrt{m})$ -packing of the unit $\|\cdot\|_2$ -balls with cardinality at least $(2\sqrt{m})^{2B}$. We now choose a set $\mathcal{V} = \delta(\mathcal{V}_{s_1} \times \mathcal{V}_{s_2} \times \dots \times \mathcal{V}_{s_m})$. We see that

$$|\mathcal{V}| \geq (2\sqrt{m})^{2mB}.$$

Proof. We have that

$$\|\boldsymbol{\psi}_1 - \boldsymbol{\psi}_2\|_2 = \sqrt{\sum_{s \in \mathcal{S}} \|\psi_1^s - \psi_2^s\|_2^2} \geq \delta/2.$$

Moreover,

$$D_{\text{KL}}(p_{\boldsymbol{\psi}_1}^n \| p_{\boldsymbol{\psi}_2}^n) = \sum_{s \in \mathcal{S}} D_{\text{KL}}(p_{\psi_1^s}^{n_s} \| p_{\psi_2^s}^{n_s}).$$

We first find upper bound of $D_{\text{KL}}(p_{\psi_1^s}^{n_s} \| p_{\psi_2^s}^{n_s})$. Since the data is independent, we have that

$$\begin{aligned} D_{\text{KL}}(p_{\psi_1^s}^{n_s} \| p_{\psi_2^s}^{n_s}) &= n_s D_{\text{KL}}(p_{\psi_1^s}^1 \| p_{\psi_2^s}^1) \\ &= n_s \left[\varphi \left(\left(\psi_1^s + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \psi_1^v \right)^\top \phi(\mathbf{x}^s) \right) \log \frac{\varphi \left(\left(\psi_1^s + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \psi_1^v \right)^\top \phi(\mathbf{x}^s) \right)}{\varphi \left(\left(\psi_2^s + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \psi_2^v \right)^\top \phi(\mathbf{x}^s) \right)} \right. \\ &\quad \left. + \varphi \left(- \left(\psi_1^s + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \psi_1^v \right)^\top \phi(\mathbf{x}^s) \right) \log \frac{\varphi \left(- \left(\psi_1^s + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \psi_1^v \right)^\top \phi(\mathbf{x}^s) \right)}{\varphi \left(- \left(\psi_2^s + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \psi_2^v \right)^\top \phi(\mathbf{x}^s) \right)} \right]. \end{aligned}$$

The first component:

$$\begin{aligned} &\varphi \left(\left(\psi_1^s + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \psi_1^v \right)^\top \phi(\mathbf{x}^s) \right) \log \frac{\varphi \left(\left(\psi_1^s + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \psi_1^v \right)^\top \phi(\mathbf{x}^s) \right)}{\varphi \left(\left(\psi_2^s + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \psi_2^v \right)^\top \phi(\mathbf{x}^s) \right)} \\ &\leq \left| \log \left(1 + e^{-\left(\psi_2^s + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \psi_2^v \right)^\top \phi(\mathbf{x}^s)} \right) - \log \left(1 + e^{-\left(\psi_1^s + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \psi_1^v \right)^\top \phi(\mathbf{x}^s)} \right) \right| \\ &\stackrel{(\star)}{\leq} \left| \left(\psi_2^s + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \psi_2^v \right)^\top \phi(\mathbf{x}^s) - \left(\psi_1^s + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \psi_1^v \right)^\top \phi(\mathbf{x}^s) \right| \\ &\leq 4B\delta \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \right), \end{aligned}$$

where (\star) follows from the fact that the SoftPlus function $\log(1 + e^x)$ is 1-Lipschitz. In particular,

$$\left| \log(1 + e^{x_1}) - \log(1 + e^{x_2}) \right| = \left| \int_{x_1}^{x_2} \frac{e^x}{1 + e^x} dx \right| \leq \left| \int_{x_1}^{x_2} 1 dx \right| = |x_1 - x_2|.$$

Similarly, for the second component, we also have

$$\begin{aligned} &\varphi \left(- \left(\psi_1^s + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \psi_1^v \right)^\top \phi(\mathbf{x}^s) \right) \log \frac{\varphi \left(- \left(\psi_1^s + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \psi_1^v \right)^\top \phi(\mathbf{x}^s) \right)}{\varphi \left(- \left(\psi_2^s + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \psi_2^v \right)^\top \phi(\mathbf{x}^s) \right)} \\ &\leq 4B\delta \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \right). \end{aligned}$$

Thus,

$$D_{\text{KL}}(p_{\boldsymbol{\psi}_1}^{n_s} \| p_{\boldsymbol{\psi}_2}^{n_s}) \leq 8B\delta \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \right) n_s.$$

Consequently,

$$D_{\text{KL}}(p_{\hat{\Psi}_1}^n \| p_{\hat{\Psi}_2}^n) \leq 8B\delta \sum_{s \in \mathcal{S}} n_s \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v}\right).$$

So, we have that

$$\begin{aligned} \inf_{\hat{\Psi}_n} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\|\hat{\Psi}_n - \Psi(P)\|_2 \right] &\geq \frac{\delta}{4} \left(1 - \frac{8B\delta \sum_{s \in \mathcal{S}} n_s \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v}\right) + \log 2}{\log |\mathcal{V}|} \right) \\ &\geq \frac{\delta}{4} \left(1 - \frac{8B\delta \sum_{s \in \mathcal{S}} n_s \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v}\right) + \log 2}{2mB \log(2\sqrt{m})} \right). \end{aligned}$$

We choose $\delta = \frac{m \log(2\sqrt{m})}{16 \sum_{s \in \mathcal{S}} n_s (1 + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v})}$, then

$$1 - \frac{8B\delta \sum_{s \in \mathcal{S}} n_s \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v}\right) + \log 2}{2mB \log(2\sqrt{m})} \geq \frac{1}{4}.$$

Thus,

$$\begin{aligned} \inf_{\hat{\Psi}_n} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\|\hat{\Psi}_n - \Psi(P)\|_2 \right] &\geq \frac{1}{4} \times \frac{mB \log(2\sqrt{m})}{16B \sum_{s \in \mathcal{S}} n_s (1 + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v})} \times \frac{1}{4} \\ &= \frac{m \log(2\sqrt{m})}{256 \sum_{s \in \mathcal{S}} n_s (1 + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v})}. \end{aligned}$$

This completes the proof of part (i). \square

H.2 Proof of Part (ii)

Proof. We summarize the model as follows

$$y^s = \left((1 - w^s) (\beta_0^s + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} \beta_0^v) + w^s (\beta_1^s + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} \beta_1^v) \right)^\top \phi(\mathbf{x}^s) + \epsilon_s, \quad \epsilon_s \sim \mathcal{N}(0, \sigma^2).$$

Let $\beta = \{\beta_0^s, \beta_1^s\}_{s \in \mathcal{S}}$. Let \mathcal{V}_{0s} and \mathcal{V}_{1s} be $1/(2\sqrt{m})$ -packing of the unit $\|\cdot\|_2$ -balls with cardinality at least $(2\sqrt{m})^{2B}$. Let $\mathcal{V}_s = \mathcal{V}_{0s} \times \mathcal{V}_{1s}$. We now choose a set $\mathcal{V} = \delta(\mathcal{V}_{s_1} \times \mathcal{V}_{s_2} \times \dots \times \mathcal{V}_{s_m})$. We see that

$$|\mathcal{V}| \geq (2\sqrt{m})^{4mB}.$$

We have that

$$\|\beta_1 - \beta_2\|_2 = \sqrt{\sum_{s \in \mathcal{S}} \left(\|(\beta_0^s)_1 - (\beta_0^s)_2\|_2^2 + \|(\beta_1^s)_1 - (\beta_1^s)_2\|_2^2 \right)} \geq \delta/\sqrt{2}.$$

Moreover,

$$D_{\text{KL}}(p_{\beta_1}^n \| p_{\beta_2}^n) = \sum_{s \in \mathcal{S}} D_{\text{KL}}(p_{\beta_1}^{n_s} \| p_{\beta_2}^{n_s}) = \sum_{s \in \mathcal{S}} n_s D_{\text{KL}}(p_{\beta_1}^1 \| p_{\beta_2}^1).$$

In addition,

$$\begin{aligned} &D_{\text{KL}}(p_{\beta_1}^1 \| p_{\beta_2}^1) \\ &= \frac{1}{2\sigma^2} \left(\left((1 - w^s) ((\beta_0^s)_1 + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} (\beta_0^v)_1) + w^s ((\beta_1^s)_1 + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} (\beta_1^v)_1) \right)^\top \phi(\mathbf{x}^s) \right. \\ &\quad \left. - \left((1 - w^s) ((\beta_0^s)_2 + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} (\beta_0^v)_2) + w^s ((\beta_1^s)_2 + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} (\beta_1^v)_2) \right)^\top \phi(\mathbf{x}^s) \right)^2 \end{aligned}$$

$$\begin{aligned}
&\leq \frac{1}{2\sigma^2} \left(\left((1-w^s)(2\delta + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} 2\delta) + w^s(2\delta + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} 2\delta) \right) \|\phi(\mathbf{x}^s)\|_2 \right)^2 \\
&\leq \frac{8B^2\delta^2}{\sigma^2} \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} \right)^2,
\end{aligned}$$

Thus,

$$D_{\text{KL}}(p_{\beta_1}^n \| p_{\beta_2}^n) \leq \frac{8B^2\delta^2}{\sigma^2} \sum_{s \in \mathcal{S}} n_s \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} \right)^2.$$

Consequently,

$$\begin{aligned}
\inf_{\hat{\beta}_n} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\|\hat{\beta}_n - \beta(P)\|_2 \right] &\geq \frac{\delta}{2\sqrt{2}} \left(1 - \frac{\frac{8B^2\delta^2}{\sigma^2} \sum_{s \in \mathcal{S}} n_s \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} \right)^2 + \log 2}{\log |\mathcal{V}|} \right) \\
&\geq \frac{\delta}{2\sqrt{2}} \left(1 - \frac{\frac{8B^2\delta^2}{\sigma^2} \sum_{s \in \mathcal{S}} n_s \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} \right)^2 + \log 2}{4mB \log(2\sqrt{m})} \right).
\end{aligned}$$

We choose $\delta^2 = \frac{mB \log(2\sqrt{m})}{4 \frac{B^2}{\sigma^2} \sum_{s \in \mathcal{S}} n_s \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} \right)^2}$, then

$$1 - \frac{\frac{8B^2\delta^2}{\sigma^2} \sum_{s \in \mathcal{S}} n_s \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} \right)^2 + \log 2}{4mB \log(2\sqrt{m})} = 1 - \frac{2mB \log(2\sqrt{m}) + \log 2}{4mB \log(2\sqrt{m})} \geq \frac{1}{4}.$$

Thus,

$$\begin{aligned}
\inf_{\hat{\beta}_n} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\|\hat{\beta}_n - \beta(P)\|_2 \right] &\geq \frac{1}{2\sqrt{2}} \sqrt{\frac{4mB \log(2\sqrt{m})}{2 \frac{8B^2}{\sigma^2} \sum_{s \in \mathcal{S}} n_s \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} \right)^2}} \times \frac{1}{4} \\
&= \frac{\sigma}{16\sqrt{2}} \sqrt{\frac{m \log(2\sqrt{m})}{B \sum_{s \in \mathcal{S}} n_s \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} \right)^2}}.
\end{aligned}$$

This completes the proof of part (ii). \square

I Further cases of the minimax lower bounds

In Lemma 1 and 2, we have presented the minimax lower bounds when $y_i^s \in \mathbb{R}$ and $\mathbf{x}_i^s \in \mathbb{R}^{d_x}$. Here, we briefly describe the other cases.

I.1 Further cases of Lemma 1

In this section, we further detail the lower bound for binary outcomes and binary proxy variables. In this case, we need to re-derive the upper bound of

$$p_{\theta_1}(w = j|z) D_{\text{KL}}[p_{\theta_1}(y|w = j, z) \| p_{\theta_2}(y|w = j, z')] \quad \text{and} \quad D_{\text{KL}}[p_{\theta_1}(\mathbf{x}|z) \| p_{\theta_2}(\mathbf{x}|z')],$$

where $j = 1, 2$. Using similar derivations as before for the quantity $D_{\text{KL}}[p_{\theta_1}(w|z) \| p_{\theta_2}(w|z')]$, we have that

$$p_{\theta_1}(w = j|z) D_{\text{KL}}[p_{\theta_1}(y|w = j, z) \| p_{\theta_2}(y|w = j, z')] \leq 8B\delta \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \right),$$

and

$$D_{\text{KL}}[p_{\theta_1}(\mathbf{x}|z) \| p_{\theta_2}(\mathbf{x}|z')] \leq d_x 8B\delta \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v} \right).$$

Combining the results, we have

$$D_{\text{KL}}(p_{\hat{\theta}_1}^n \| p_{\hat{\theta}_2}^n) = \sum_{s \in \mathcal{S}} D_{\text{KL}}(p_{\hat{\theta}_1}^{n_s} \| p_{\hat{\theta}_2}^{n_s}) \leq \sum_{s \in \mathcal{S}} n_s 8(d_x + 3)B\delta \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v}\right).$$

Consequently, we have that

$$\inf_{\hat{\theta}_n} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\|\hat{\theta}_n - \theta(P)\|_2 \right] \geq \frac{\delta}{2} \left(1 - \frac{\sum_{s \in \mathcal{S}} n_s 8(d_x + 3)B\delta \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v}\right) + \log 2}{2mB(d_x + 3) \log(2\sqrt{m})} \right).$$

We choose $\delta = \frac{m \log(2\sqrt{m})}{8 \sum_{s \in \mathcal{S}} n_s \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v}\right)}$, then

$$1 - \frac{\sum_{s \in \mathcal{S}} n_s 8(d_x + 3)B\delta \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v}\right) + \log 2}{2mB(d_x + 3) \log(2\sqrt{m})} \geq \frac{3}{8}.$$

Thus,

$$\boxed{\inf_{\hat{\theta}_n} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\|\hat{\theta}_n - \theta(P)\|_2 \right] \geq \frac{3mB \log(2\sqrt{m})}{128 \sum_{s \in \mathcal{S}} n_s B \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \lambda^{s,v}\right)}.$$

Remark 1. Note that the derivation in this Section and in Section H.1 give us enough tools to compute the minimax lower bounds for any further case, i.e., any combination of the outcomes and proxy variables (binary or continuous). The key is to initially find the upper bound of $D_{\text{KL}}(p_{\hat{\theta}_1}^n \| p_{\hat{\theta}_2}^n)$ based on the constructed packing. Then, using Fano's method to obtain the minimax lower bounds.

I.2 Further cases of Lemma 2

Note that the lower bound of Lemma 2, part (i) has only one case since we only focus on binary treatment, and it is presented in the main text. For part (ii), consider $y_i^s \in \{0, 1\}$, then the model of the outcomes would follow a Bernoulli distribution. Reusing the scheme in Section H.2, we need to find the new upper bound of $D_{\text{KL}}(p_{\hat{\beta}_1}^n \| p_{\hat{\beta}_2}^n)$. In particular,

$$D_{\text{KL}}(p_{\hat{\beta}_1}^n \| p_{\hat{\beta}_2}^n) = \sum_{s \in \mathcal{S}} n_s \left[\varphi(v_1) \log \frac{\varphi(v_1)}{\varphi(v_2)} + \varphi(-v_1) \log \frac{\varphi(-v_1)}{\varphi(-v_2)} \right],$$

where $v_j = \left((1 - w^s) \left((\beta_0^s)_j + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} (\beta_0^v)_j \right) + w^s \left((\beta_1^s)_j + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} (\beta_1^v)_j \right) \right)^\top \phi(\mathbf{x}^s)$. We have that

$$\begin{aligned} \varphi(v_1) \log \frac{\varphi(v_1)}{\varphi(v_2)} &\leq \left\| \left((1 - w^s) \left((\beta_0^s)_1 - (\beta_0^s)_2 + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} [(\beta_0^v)_1 - (\beta_0^v)_2] \right) \right. \right. \\ &\quad \left. \left. + w^s \left((\beta_1^s)_1 - (\beta_1^s)_2 + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} [(\beta_1^v)_1 - (\beta_1^v)_2] \right) \right\|_2 \|\phi(\mathbf{x}^s)\|_2 \\ &\leq 4B\delta \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \right), \end{aligned}$$

Similarly, $\varphi(-v_1) \log \frac{\varphi(-v_1)}{\varphi(-v_2)} \leq 4B\delta \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \gamma^{s,v} \right)$. Hence,

$$D_{\text{KL}}(p_{\hat{\beta}_1}^n \| p_{\hat{\beta}_2}^n) \leq 8B\delta \sum_{s \in \mathcal{S}} n_s \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} \right).$$

Using similar technique in Section H.2, we obtain

$$\boxed{\inf_{\hat{\beta}_n} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\|\hat{\beta}_n - \beta(P)\|_2 \right] \geq \frac{m \log(2\sqrt{m})}{32\sqrt{2} \sum_{s \in \mathcal{S}} n_s \left(1 + \sum_{v \in \mathcal{S}_{\setminus s}} \eta^{s,v} \right)}.$$

We observe that the lower bound is similar to that of Lemma 2, part (i) since they are both lower bounds of a binary response variable. The constant in this bound is larger ($1/(32\sqrt{2})$) than that of Lemma 2, part (i) ($1/256$). This is expected since there are more parameters in this model, i.e., $\{\beta_0^s, \beta_1^s\}_{s \in \mathcal{S}}$, as compared to the model in Lemma 2, part (i) ($\{\psi^s\}_{s \in \mathcal{S}}$).

J Description of IHDP data

This section describe details of the IHDP data, which was skipped in the main text due to limited space.

The Infant Health and Development Program (IHDP) is a randomized study on the impact of specialist visits (the treatment) on the cognitive development of children (the outcome). The dataset consists of 747 records with 25 covariates describing properties of the children and their mothers. The treatment group includes children who received specialist visits and control group includes children who did not receive. Further details are presented in Appendix. For each child, a treated and a control outcome are simulated using the numerical schemes provided in the NPCI package (Dorie 2016), thus allowing us to know the *true* individual treatment effect. We use 10 replicates of the dataset in this experiment. For each replicate, we divide into three sources, each consists of 249 data points. For each source, we use the first 50 data points for training, the next 100 for testing and the rest 99 for validating. We report the mean and standard error of the evaluation metrics over 10 replicates of the data.

References

- Choudhury, O., Park, Y., Salonidis, T., Gkoulalas-Divanis, A., Sylla, I., et al. (2019). Predicting adverse drug reactions on distributed health data using federated learning. In *AMIA Annual Symposium Proceedings*, volume 2019, page 313. American Medical Informatics Association.
- Dorie, V. (2016). Npci: Non-parametrics for causal inference. URL: <https://github.com/vdorie/npci>.
- Flores, M., Dayan, I., Roth, H., Zhong, A., Harouni, A., Gentili, A., Abidin, A., Liu, A., Costa, A., Wood, B., et al. (2020). Federated learning used for predicting outcomes in SARS-COV-2 patients. *Preprint. medRxiv. 2020;2020.08.11.20172809*.
- Liu, J. S. (1996). Metropolized independent sampling with comparisons to rejection sampling and importance sampling. *Statistics and Computing*, 6(2):113–119.
- Louizos, C., Shalit, U., Mooij, J. M., Sontag, D., Zemel, R., and Welling, M. (2017). Causal effect inference with deep latent-variable models. In *Advances in Neural Information Processing Systems*, pages 6446–6456.
- Milton, P., Coupland, H., Giorgi, E., and Bhatt, S. (2019). Spatial analysis made easy with linear regression and kernels. *Epidemics*, 29:100362.
- Vaid, A., Jaladanki, S. K., Xu, J., Teng, S., Kumar, A., and Lee, S. (2020). Federated learning of electronic health records improves mortality prediction in patients. *Ethnicity*, 52(77.6):0–001.

— END —