

Appendix:

Bayesian Federated Estimation of Causal Effects from Observational Data

Thanh Vinh Vo¹

Young Lee²

Trong Nghia Hoang³

Tze-Yun Leong¹

¹School of Computing, National University of Singapore

²Harvard University

³School of Electrical Engineering and Computer Science, Washington State University

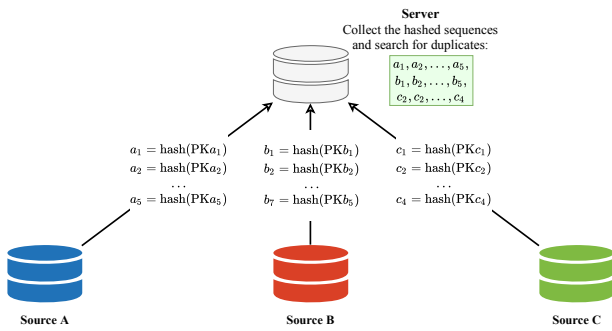


Figure 1: The secure preprocessing procedures to identify duplicated individuals among multiple sources. $\text{PK}a_i$ ($i = 1, \dots, 5$), $\text{PK}b_i$ ($i = 1, \dots, 7$), $\text{PK}c_i$ ($i = 1, \dots, 4$) are the primary keys of each individual in each source. a_i ($i = 1, \dots, 5$), b_i ($i = 1, \dots, 7$), c_i ($i = 1, \dots, 4$) are the hashed sequences of these individuals.

A THE PREPROCESSING PROCEDURE

The assumptions were described briefly in Section 3.2 of the main text. Here we present the preprocessing procedures to remove duplicated individuals.

The preprocessing procedure are summarized as follows. Firstly, each source would use a one-way hash function (such as MD4, MD5, SHA or SHA256) to encrypt each individuals' primary key and then send the hashed sequences to a server. By doing this, the individuals' data are secured. Note that the one-way hash function is agreed among the sources so that they would use the same function. Then, the server collects all hashed sequences from all sources and perform a matching algorithm to see if there exists repeated individuals among different sources. For each repeated individual, the server randomly choose to keep it on a small number (predefined) of sources and inform the other sources to exclude this individual from the training process. The whole procedure is to ensure that an individual does not exists in a huge number of sources, thus prevent learning a biased model. We summarize the procedure in Figure 1.

Assumption 4 and the preprocessing procedure are required for data that are highly repeated in different sources only. For data that are not likely to have a high number of repetitions such as patients from different hospitals of different countries, the above assumption and the preprocessing procedure are not required. Note that the existing methods also need Assumption 4 since they need to combine data and remove repeated individuals.

In this work, we assume that all of the assumptions described in this section are satisfied, and the preprocessing procedure was performed if it is necessary.

B THE FEDERATED EVIDENCE LOWER BOUND

Naively applying variational inference would lead to a non-decomposable ELBO. The proposed ELBO can be decomposed into multiple components, thus enabling federated optimization. We give a full derivation as follows:

$$\begin{aligned} & \log p(\mathbf{y}_{\text{obs}} | \mathbf{X}, \mathbf{w}) \\ &= \log \int p(\mathbf{y}_{\text{obs}}, \mathbf{g}, \Psi, \Sigma | \mathbf{X}, \mathbf{w}) d\mathbf{g} d\Psi d\Sigma \\ &= \log \int p(\mathbf{y}_{\text{obs}} | \mathbf{g}, \Psi, \Sigma, \mathbf{X}, \mathbf{w}) p(\mathbf{g}, \Psi, \Sigma | \mathbf{X}, \mathbf{w}) d\mathbf{g} d\Psi d\Sigma. \end{aligned}$$

From Figure 2, we see that $\mathbf{g}, \Psi, \Sigma \perp\!\!\!\perp \mathbf{X}^s, \mathbf{w}^s$ (for all $s = 1, 2, \dots, m$), i.e. \mathbf{g}, Ψ, Σ are independent with $\mathbf{X}^s, \mathbf{w}^s$ when $\mathbf{y}_{\text{obs}}^s, \mathbf{y}_{\text{mis}}^s$ are not given. Thus, $p(\mathbf{g}, \Psi, \Sigma | \mathbf{X}, \mathbf{w}) = p(\mathbf{g}, \Psi, \Sigma)$.

In addition, from Figure 2, we also have

$$p(\mathbf{y}_{\text{obs}} | \mathbf{g}, \Psi, \Sigma, \mathbf{X}, \mathbf{w}) = \prod_{s=1}^m p(\mathbf{y}_{\text{obs}}^s | \mathbf{g}, \Psi, \Sigma, \mathbf{X}^s, \mathbf{w}^s).$$

Thus,

$$\log p(\mathbf{y}_{\text{obs}} | \mathbf{X}, \mathbf{w})$$

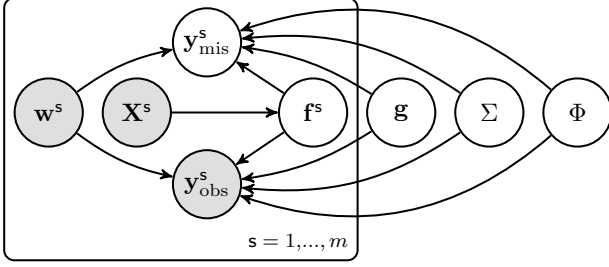


Figure 2: Graphical model that summarizes the proposed framework with treatment w^s , covariate X^s , and the two potential outcomes y^s_{mis} and y^s_{obs} . The quantity f^s is idiosyncratic to the sources and \mathbf{g} contains shared characteristics across all the sources. Σ and Ψ are shared parameters. Note that this is not a causal graph.

$$\begin{aligned}
&= \log \int q(\mathbf{g}, \Psi, \Sigma) \prod_{s=1}^m p(\mathbf{y}^s_{\text{obs}} | \mathbf{g}, \Psi, \Sigma, \mathbf{X}^s, \mathbf{w}^s) \\
&\quad \times \frac{p(\mathbf{g}, \Psi, \Sigma)}{q(\mathbf{g}, \Psi, \Sigma)} d\mathbf{g} d\Psi d\Sigma \\
&\geq \int q(\mathbf{g}, \Psi, \Sigma) \log \left(\prod_{s=1}^m p(\mathbf{y}^s_{\text{obs}} | \mathbf{g}, \Psi, \Sigma, \mathbf{X}^s, \mathbf{w}^s) \right. \\
&\quad \left. \times \frac{p(\mathbf{g}, \Psi, \Sigma)}{q(\mathbf{g}, \Psi, \Sigma)} \right) d\mathbf{g} d\Psi d\Sigma \\
&= \sum_{s=1}^m \mathbb{E}_q[\log p(\mathbf{y}^s_{\text{obs}} | \mathbf{g}, \Psi, \Sigma, \mathbf{X}^s, \mathbf{w}^s)] \\
&\quad - \mathbb{D}_{\text{KL}}[q(\mathbf{g}, \Psi, \Sigma) \| p(\mathbf{g}, \Psi, \Sigma)] \\
&= \sum_{s=1}^m \left(\mathbb{E}_q[\log p(\mathbf{y}^s_{\text{obs}} | \mathbf{g}, \Psi, \Sigma, \mathbf{X}^s, \mathbf{w}^s)] \right. \\
&\quad \left. - \frac{1}{m} \mathbb{D}_{\text{KL}}[q(\mathbf{g}, \Psi, \Sigma) \| p(\mathbf{g}, \Psi, \Sigma)] \right) \\
&= \sum_{s=1}^m \mathbf{L}^s,
\end{aligned}$$

where

$$\begin{aligned}
\mathbf{L}^s &:= \mathbb{E}_q[\log p(\mathbf{y}^s_{\text{obs}} | \mathbf{g}, \Psi, \Sigma, \mathbf{X}^s, \mathbf{w}^s)] \\
&\quad - \frac{1}{m} \mathbb{D}_{\text{KL}}[q(\mathbf{g}, \Psi, \Sigma) \| p(\mathbf{g}, \Psi, \Sigma)].
\end{aligned}$$

Hence, we can divide the ELBO into multiple components, which leads to federated training of the model. Without the proposed model, the ELBO cannot be decomposed into multiple components and hence cannot be trained in a federated setting.

C PROOF OF LEMMA 1

Proof. We denote $\xi_0^s \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{n_s})$ and $\xi_1^s \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{n_s})$. Then, from the model definition (Eq. (5) in the main text),

we have

$$\begin{aligned}
&\begin{bmatrix} y_1^s(0) & \dots & y_{n_s}^s(0) \\ y_1^s(1) & \dots & y_{n_s}^s(1) \end{bmatrix} \\
&= \Psi^{\frac{1}{2}} \begin{bmatrix} f_1^s(0) + g^s(0) & \dots & f_{n_s}^s(0) + g^s(0) \\ f_1^s(1) + g^s(1) & \dots & f_{n_s}^s(1) + g^s(1) \end{bmatrix} \\
&\quad + \Sigma^{\frac{1}{2}} \begin{bmatrix} \varepsilon_1^s(0) & \dots & \varepsilon_{n_s}^s(0) \\ \varepsilon_1^s(1) & \dots & \varepsilon_{n_s}^s(1) \end{bmatrix}.
\end{aligned}$$

The above equation is equivalent to the following

$$\mathbf{Y}^s = [\boldsymbol{\mu}_0 \quad \boldsymbol{\mu}_1](\Psi^{\frac{1}{2}})^\top + [\varepsilon_0^s \quad \varepsilon_1^s](\Sigma^{\frac{1}{2}})^\top,$$

where

$$\begin{aligned}
\boldsymbol{\mu}_0 &= \boldsymbol{\mu}_0(\mathbf{X}^s) + \mathbf{g}_0^s + (\mathbf{K}^s)^{\frac{1}{2}} \xi_0^s, \\
\boldsymbol{\mu}_1 &= \boldsymbol{\mu}_1(\mathbf{X}^s) + \mathbf{g}_1^s + (\mathbf{K}^s)^{\frac{1}{2}} \xi_1^s.
\end{aligned}$$

Further expanding the right hand side, we have

$$\begin{aligned}
\mathbf{Y}^s &= [\boldsymbol{\mu}_0(\mathbf{X}^s) + \mathbf{g}_0^s \quad \boldsymbol{\mu}_1(\mathbf{X}^s) + \mathbf{g}_1^s] (\Psi^{\frac{1}{2}})^\top \\
&\quad + (\mathbf{K}^s)^{\frac{1}{2}} [\xi_0^s \quad \xi_1^s] (\Psi^{\frac{1}{2}})^\top + [\varepsilon_0^s \quad \varepsilon_1^s] (\Sigma^{\frac{1}{2}})^\top \\
\text{vec}(\mathbf{Y}^s) &= (\Psi^{\frac{1}{2}} \otimes \mathbf{I}_{n_s}) \begin{bmatrix} \boldsymbol{\mu}_0(\mathbf{X}^s) + \mathbf{g}_0^s \\ \boldsymbol{\mu}_1(\mathbf{X}^s) + \mathbf{g}_1^s \end{bmatrix} \\
&\quad + (\Psi^{\frac{1}{2}} \otimes (\mathbf{K}^s)^{\frac{1}{2}}) \begin{bmatrix} \xi_0^s \\ \xi_1^s \end{bmatrix} + (\Sigma^{\frac{1}{2}} \otimes \mathbf{I}_{n_s}) \begin{bmatrix} \varepsilon_0^s \\ \varepsilon_1^s \end{bmatrix},
\end{aligned}$$

where $\text{vec}(\cdot)$ denotes the vectorization of a matrix, which converts a matrix into a column vector.

For the second term on the right hand side of the above equation, note that $\xi_0^s \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{n_s})$ and $\xi_1^s \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{n_s})$, so we have the following

$$\begin{aligned}
&\begin{bmatrix} \xi_0^s \\ \xi_1^s \end{bmatrix} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{2n_s}) \\
&(\Psi^{\frac{1}{2}} \otimes (\mathbf{K}^s)^{\frac{1}{2}}) \begin{bmatrix} \xi_0^s \\ \xi_1^s \end{bmatrix} \\
&\quad \sim \mathcal{N}\left(\mathbf{0}, (\Psi^{\frac{1}{2}} \otimes (\mathbf{K}^s)^{\frac{1}{2}}) \mathbf{I}_{2n} (\Psi^{\frac{1}{2}} \otimes (\mathbf{K}^s)^{\frac{1}{2}})^\top\right) \\
&(\Psi^{\frac{1}{2}} \otimes (\mathbf{K}^s)^{\frac{1}{2}}) \begin{bmatrix} \xi_0^s \\ \xi_1^s \end{bmatrix} \sim \mathcal{N}(\mathbf{0}, \Psi \otimes \mathbf{K}^s).
\end{aligned}$$

For the last term, note that $\varepsilon_0^s \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{n_s})$, $\varepsilon_1^s \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{n_s})$, thus

$$\begin{aligned}
&\begin{bmatrix} \varepsilon_0^s \\ \varepsilon_1^s \end{bmatrix} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{2n_s}) \\
&(\Sigma^{\frac{1}{2}} \otimes \mathbf{I}_{n_s}) \begin{bmatrix} \varepsilon_0^s \\ \varepsilon_1^s \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, (\Sigma^{\frac{1}{2}} \otimes \mathbf{I}_{n_s}) \mathbf{I}_{2n} (\Sigma^{\frac{1}{2}} \otimes \mathbf{I}_{n_s})^\top\right) \\
&(\Sigma^{\frac{1}{2}} \otimes \mathbf{I}_{n_s}) \begin{bmatrix} \varepsilon_0^s \\ \varepsilon_1^s \end{bmatrix} \sim \mathcal{N}(\mathbf{0}, \Sigma \otimes \mathbf{I}_{n_s}).
\end{aligned}$$

Consequently,

$$\text{vec}(\mathbf{Y}^s) | \Psi, \Sigma, \mathbf{X}^s, \mathbf{w}^s, \mathbf{g}^s$$

$$\sim N \left(\left(\Psi^{\frac{1}{2}} \otimes \mathbf{I}_{n_s} \right) \begin{bmatrix} \mu_0(\mathbf{X}^s) + \mathbf{g}_0^s \\ \mu_1(\mathbf{X}^s) + \mathbf{g}_1^s \end{bmatrix}, \Psi \otimes \mathbf{K}^s + \Sigma \otimes \mathbf{I}_{n_s} \right),$$

which implies that

$$\begin{aligned} & \begin{bmatrix} \mathbf{y}^s(0) \\ \mathbf{y}^s(1) \end{bmatrix} \mid \Psi, \Sigma, \mathbf{X}^s, \mathbf{w}^s, \mathbf{g}^s \\ & \sim N \left(\left(\Psi^{\frac{1}{2}} \otimes \mathbf{I}_{n_s} \right) \begin{bmatrix} \mu_0(\mathbf{X}^s) + \mathbf{g}_0^s \\ \mu_1(\mathbf{X}^s) + \mathbf{g}_1^s \end{bmatrix}, \Psi \otimes \mathbf{K}^s + \Sigma \otimes \mathbf{I}_{n_s} \right). \end{aligned}$$

This completes the proof. \square

D PROOF OF LEMMA 2

Proof. Following the proof of Lemma 2, we note that if the observed treatment $w_i^s = 0$, then the mean of $p(y_{i,\text{obs}}^s \mid \mathbf{X}^s, \mathbf{w}^s, \Psi, \Sigma, \mathbf{g}^s)$ equals to the mean of $p(y_{i,\text{obs}}^s(0) \mid \Psi, \Sigma, \mathbf{X}^s, \mathbf{w}^s, \mathbf{g}^s)$ and the mean of $p(y_{i,\text{mis}}^s \mid \mathbf{X}^s, \mathbf{w}^s, \Psi, \Sigma, \mathbf{g}^s)$ equals to the mean of $p(y_{i,\text{obs}}^s(1) \mid \Psi, \Sigma, \mathbf{X}^s, \mathbf{w}^s, \mathbf{g}^s)$. If the observed treatment $w_i^s = 1$, then the mean of $p(y_{i,\text{obs}}^s \mid \mathbf{X}^s, \mathbf{w}^s, \Psi, \Sigma, \mathbf{g}^s)$ equals to the mean of $p(y_{i,\text{obs}}^s(1) \mid \Psi, \Sigma, \mathbf{X}^s, \mathbf{w}^s, \mathbf{g}^s)$ and the mean of $p(y_{i,\text{mis}}^s \mid \mathbf{X}^s, \mathbf{w}^s, \Psi, \Sigma, \mathbf{g}^s)$ equals to the mean of $p(y_{i,\text{obs}}^s(0) \mid \Psi, \Sigma, \mathbf{X}^s, \mathbf{w}^s, \mathbf{g}^s)$. Hence, we have

$$\begin{aligned} \mu_{\text{obs}}(\mathbf{X}^s) &= (\mathbf{1} - \mathbf{w}^s) \odot \mathbf{m}_0 + \mathbf{w}^s \odot \mathbf{m}_1, \\ \mu_{\text{mis}}(\mathbf{X}^s) &= \mathbf{w}^s \odot \mathbf{m}_0 + (\mathbf{1} - \mathbf{w}^s) \odot \mathbf{m}_1, \end{aligned}$$

Similarly, for the covariance matrix, each element in \mathbf{K}_{obs} , \mathbf{K}_{mis} , and \mathbf{K}_{om} also depends on whether $w_i^s = 0$ or $w_i^s = 1$. So each element in these matrices is computed by the following kernel function

$$\begin{aligned} k_{\text{obs}}(\mathbf{x}_i, \mathbf{x}_j) &= [(1-w_i)(1-w_j)\psi_{11} + w_i w_j \psi_{22} \\ &\quad + (1-w_i)w_j \psi_{12} + w_i(1-w_j)\psi_{21}] k(\mathbf{x}_i, \mathbf{x}_j) \\ &\quad + [(1-w_i)\sigma_{11} + w_i \sigma_{22}] \mathbb{1}_{i=j}, \\ k_{\text{mis}}(\mathbf{x}_i, \mathbf{x}_j) &= [w_i w_j \psi_{11} + (1-w_i)(1-w_j)\psi_{22} \\ &\quad + (1-w_i)w_j \psi_{21} + w_i(1-w_j)\psi_{12}] k(\mathbf{x}_i, \mathbf{x}_j) \\ &\quad + [w_i \sigma_{11} + (1-w_i)\sigma_{22}] \mathbb{1}_{i=j}, \\ k_{\text{om}}(\mathbf{x}_i, \mathbf{x}_j) &= [(1-w_i)(1-w_j)\psi_{21} + w_i w_j \psi_{12} \\ &\quad + (1-w_i)w_j \psi_{22} + w_i(1-w_j)\psi_{11}] k(\mathbf{x}_i, \mathbf{x}_j) \\ &\quad + [(1-w_i)\sigma_{21} + w_i \sigma_{12}] \mathbb{1}_{i=j}, \end{aligned}$$

where ψ_{ab} and σ_{ab} are the (a, b) -th elements of Ψ and Σ , respectively.

This completes the proof. \square

E EVALUATION METRICS

The two evaluation metrics reported in our experiments are defined as follows: (i) precision in estimation of heterogeneous effects (PEHE):

$$\epsilon_{\text{PEHE}} := \sum_{s=1}^m \sum_{i=1}^{n_s} (\tau_i^s - \hat{\tau}_i^s)^2 / (mn_s)$$

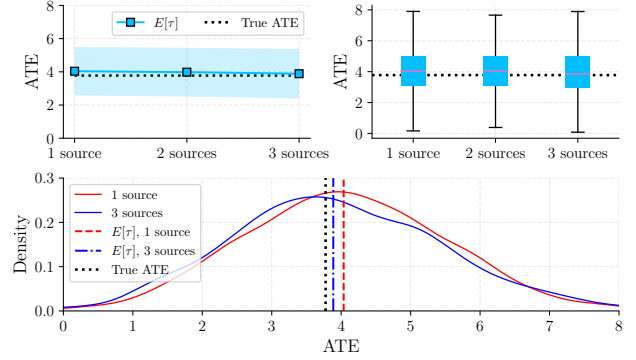


Figure 3: The estimated ATE distribution on source #1 of IHDP dataset. The dotted black lines represent the true ATE.

for evaluating ITE, and (ii) absolute error:

$$\epsilon_{\text{ATE}} := |\tau - \hat{\tau}|$$

for evaluating ATE, where τ_i^s and τ are the *true* ITE and *true* ATE, respectively, and $\hat{\tau}_i^s, \hat{\tau}$ are their estimates.

F ADDITIONAL EXPERIMENTAL RESULTS

In this section, we present some additional results which was skipped in the main text due to limited space.

F.1 SYNTHETIC DATA: DATA-2

In this section, we present additional experimental results on DATA-2. Again, those results were skipped in the main text due to limited space. In Table 1, we present additional results of the baselines trained locally (loc) and the baselines trained with bootstrap aggregating (agg). Similar to the experiments on DATA-1 presented in the main text, the results on DATA-2 also show that FedCI achieves much lower errors, especially the error in predicting ITE.

F.2 IHDP DATASET

In this section, we present additional experimental results on the IHDP dataset. The results here were not presented in the main text due to limited space. In Table 2, we present additional results of the baselines trained locally (loc) and the baselines trained with bootstrap aggregating (agg). Similar to the experiments on synthetic data, the results presented here show that FedCI achieves much smaller errors. The reason is because FedCI has access to all the data sources in a federated fashion while the ‘baselines trained locally’ (loc) and the ‘baselines trained with bootstrap aggregating’ (agg) only have access to a local data source.

Similar to the experiment on synthetic data, the estimated distribution of ATE in the first source ($s = 1$) is presented

Table 1: Out-of-sample errors on DATA-2 where top-3 performances are highlighted in bold (lower is better). The dashes (—) in ‘loc’ and ‘agg’ indicate that the numbers are the same as those of ‘com’.

Method	The error of ITE ($\sqrt{\epsilon_{PEHE}}$)			The error of ATE (ϵ_{ATE})		
	1 source	3 sources	5 sources	1 source	3 sources	5 sources
BART _{loc}	—	18.4±0.3	18.3±0.2	—	3.37±0.7	2.90±0.6
X-Learner _{loc}	—	22.7±0.5	22.8±0.5	—	3.55±1.3	3.09±0.8
R-Learner _{loc}	—	26.3±0.2	26.1±0.2	—	19.7±0.3	19.5±0.3
OthoRF _{loc}	—	38.3±1.4	40.0±0.9	—	4.09±0.9	4.40±1.2
TARNet _{loc}	—	37.6±0.6	37.1±0.4	—	7.31±0.4	7.25±0.3
CFR Wass _{loc}	—	37.2±0.7	37.0±0.5	—	7.24±0.3	7.12±0.2
CFR MMD _{loc}	—	37.2±0.6	36.8±0.4	—	7.21±0.4	7.11±0.3
CEVAE _{loc}	—	21.4±0.7	19.8±0.6	—	2.11±0.4	1.97±0.2
BART _{agg}	—	17.9±0.2	17.7±0.2	—	3.91±0.8	3.15±0.7
X-Learner _{agg}	—	18.2±0.4	17.1±0.2	—	3.43±1.3	3.07±0.8
R-Learner _{agg}	—	26.2±0.3	26.1±0.2	—	19.7±0.4	19.6±0.3
OthoRF _{agg}	—	25.0±1.3	17.3±0.6	—	4.56±1.1	1.30±0.4
TARNet _{agg}	—	36.5±0.3	36.1±0.3	—	7.26±0.3	7.18±0.3
CFR Wass _{agg}	—	35.2±0.5	35.0±0.3	—	7.13±0.3	6.97±0.2
CFR MMD _{agg}	—	35.2±0.5	35.1±0.4	—	7.10±0.4	7.05±0.2
CEVAE _{agg}	—	19.2±0.8	18.3±0.7	—	2.02±0.3	1.91±0.4
BART _{com}	18.0±0.4	17.7±0.2	17.4±0.1	3.54±1.3	2.94±0.8	1.84±0.5
X-Learner _{com}	21.1±0.9	17.9±0.4	16.2±0.2	4.55±1.4	3.29±1.0	2.37±0.8
R-Learner _{com}	25.9±0.6	23.5±0.5	21.3±0.4	19.0±0.8	15.6±0.7	12.3±0.6
OthoRF _{com}	37.8±2.7	10.7±0.5	9.83±0.5	7.88±2.2	1.99±0.4	2.36±0.6
TARNet _{com}	36.1±0.4	35.5±0.2	35.0±0.2	7.11±0.4	7.10±0.3	7.08±0.2
CFR Wass _{com}	35.1±0.4	34.5±0.2	34.1±0.2	7.10±0.4	7.01±0.3	6.90±0.2
CFR MMD _{com}	35.1±0.4	35.0±0.2	34.9±0.2	7.12±0.4	7.02±0.3	7.01±0.2
CEVAE _{com}	20.1±0.5	18.4±0.6	16.6±0.6	1.50±0.3	1.38±0.4	1.89±0.2
FedCI	9.28±0.4	6.34±0.2	5.53±0.1	2.37±0.5	1.47±0.4	0.74±0.2

Table 2: Out-of-sample errors on IHDP dataset where top-3 performances are highlighted in bold (lower is better). The dashes (—) in ‘loc’ and ‘agg’ indicate that the numbers are the same as those of ‘com’.

Method	The error of ITE ($\sqrt{\epsilon_{PEHE}}$)			The error of ATE (ϵ_{ATE})		
	1 source	2 sources	3 sources	1 source	2 sources	3 sources
BART _{loc}	—	5.83±2.6	6.56±3.3	—	2.09±0.9	1.38±0.5
X-Learner _{loc}	—	4.14±1.5	4.54±1.9	—	1.51±0.7	0.77±0.5
R-Learner _{loc}	—	6.35±1.9	6.16±2.0	—	2.13±0.7	1.44±0.3
OthoRF _{loc}	—	4.33±1.6	4.59±1.9	—	1.10±0.6	0.75±0.3
TARNet _{loc}	—	3.71±1.0	3.83±1.1	—	1.31±0.5	0.98±0.4
CFR Wass _{loc}	—	3.35±0.8	3.12±0.7	—	0.87±0.5	0.82±0.4
CFR MMD _{loc}	—	3.40±0.9	3.15±1.2	—	1.17±0.5	0.63±0.3
CEVAE _{loc}	—	3.78±0.7	3.93±0.8	—	1.91±0.3	2.37±0.2
BART _{agg}	—	4.05±1.9	3.69±1.8	—	2.09±1.0	1.30±0.5
X-Learner _{agg}	—	3.98±1.5	4.28±1.9	—	1.51±0.7	0.83±0.5
R-Learner _{agg}	—	4.76±1.3	4.46±1.6	—	1.92±0.5	1.41±0.2
OthoRF _{agg}	—	3.40±1.1	4.26±1.9	—	0.87±0.3	1.20±0.6
TARNet _{agg}	—	3.52±0.9	3.81±1.2	—	1.23±0.4	0.95±0.4
CFR Wass _{agg}	—	3.21±0.7	2.93±0.9	—	0.80±0.3	0.71±0.2
CFR MMD _{agg}	—	3.17±0.8	2.91±1.3	—	1.12±0.5	0.57±0.3
CEVAE _{agg}	—	3.63±0.7	3.73±0.5	—	0.92±0.2	0.84±0.5
BART _{com}	5.98±2.7	4.32±2.1	4.04±2.0	1.80±1.1	2.09±1.1	1.21±0.6
X-Learner _{com}	4.22±1.6	4.15±1.5	4.06±1.8	1.64±0.7	1.93±0.8	0.84±0.4
R-Learner _{com}	6.97±2.1	4.43±1.4	4.47±1.7	3.15±0.5	1.34±0.5	1.10±0.3
OthoRF _{com}	4.49±1.9	3.81±1.3	3.75±1.5	1.86±0.8	1.61±0.6	1.56±0.8
TARNet _{com}	4.50±1.4	3.15±0.8	3.79±1.1	1.52±0.5	1.18±0.4	0.91±0.3
CFR Wass _{com}	4.37±1.2	2.93±0.6	2.85±0.9	1.18±0.7	0.72±0.2	0.67±0.1
CFR MMD _{com}	4.43±1.3	2.85±0.6	2.83±1.1	2.32±0.8	0.63±0.2	0.54±0.2
CEVAE _{com}	3.16±0.6	2.34±0.6	2.31±0.7	2.02±0.4	0.53±0.1	0.48±0.2
FedCI	2.88±0.8	2.36±0.5	2.35±0.6	1.43±0.7	1.03±0.4	0.51±0.2

in Figure 3. Again, the figures show that the true ATE is inside the estimated interval and the estimated mean ATE shifts towards its true value (dotted lines) when more data sources are used.