

Moderately-Balanced Representation Learning for Orthogonal Estimation on Average Treatment Effect

Anonymous Authors

1 Experiments

The final full results for 1000 IHDP and 100 Twins experiments are reported in Table 1. The mark * indicates that the baseline models did not report relevant results.

As a supplementary note, the outcome takes a binary value for Twins experiments. So the factual outcome loss in Eqn (6) will be

$$\mathcal{L}_{fo} = -\frac{1}{N} \sum_{m=1}^N [y_m \log f(\Phi(\mathbf{z}_m)) + (1 - y_m) \log(1 - f(\Phi(\mathbf{z}_m)))].$$

2 Assumptions

Assumption 1 (SUTVA). *The potential outcomes for any individual are not affected by the treatment assignment of other individuals.*

Assumption 2 (Strong Ignorability). *Given the covariates \mathbf{Z} , the potential outcomes are independent of the treatment assignment D : $(Y(0), Y(1)) \perp\!\!\!\perp D \mid \mathbf{Z}$.*

Assumption 3 (Overlap). *The probability of treatment assignment for any unit is positive: $0 < \Pr(D = d \mid \mathbf{Z} = \mathbf{z}) < 1, \forall d \in \{0, 1\}$ and $\mathbf{z} \in \mathcal{Z}$.*

Assumption 4 (Consistency). *The potential outcome for treatment d of each unit is equal to the observed factual outcome if the actual treatment is d : $(Y(d) = Y^F) \mid D = d, \forall d \in \{0, 1\}$.*

3 Proofs

We skip the proofs of Proposition 1 since it can be seen in [Chernozhukov *et al.*, 2018]. In the following, we prove Proposition 2 and Property 1.

3.1 Proof of Proposition 2

Proof. The score functions stated in the Eqn. (2) and Eqn. (3) in the main paper are

$$\begin{aligned} \psi_1(W, \theta^i, \rho) &= \theta^i - g(i, \mathbf{Z}) \\ &\quad - (Y - g(i, \mathbf{Z})) \frac{iD + (1-i)(1-D)}{im(\mathbf{Z}) + (1-i)(1-m(\mathbf{Z}))}; \\ \psi_2(W, \theta^i, \rho) &= \theta^i - g(i, \mathbf{Z}) \\ &\quad - (Y(i) - g(i, \mathbf{Z})) \frac{((D - m(\mathbf{Z})) - \mathbb{E}[\nu \mid \mathbf{Z}])^2}{\mathbb{E}[\nu^2 \mid \mathbf{Z}]}. \end{aligned}$$

We then check if the orthogonal condition (Definition 1) holds for $\psi_1(W, \theta^i, \rho)$.

$$\begin{aligned} \partial_g \psi_1(W, \theta^i, \rho) &= -1 + \frac{iD + (1-i)(1-D)}{im(\mathbf{Z}) + (1-i)(1-m(\mathbf{Z}))}; \\ \partial_m \psi_1(W, \theta^i, \rho) &= (Y - g(i, \mathbf{Z})) \frac{D}{m(\mathbf{Z})^2}, \text{ if } i = 1; \\ \partial_m \psi_1(W, \theta^i, \rho) &= -(Y - g(i, \mathbf{Z})) \frac{1-D}{(1-m(\mathbf{Z}))^2}, \text{ if } i = 0. \end{aligned}$$

If $i = 1$ and under the noise conditions $\mathbb{E}[\nu \mid \mathbf{Z}] = 0$ and $\mathbb{E}[\xi \mid D, \mathbf{Z}] = 0$, then we have

$$\begin{aligned} &\mathbb{E}[\partial_g \psi_1(W, \theta^1, \rho) \mid \mathbf{Z}] \big|_{(g,m)=(g_0,m_0), \theta^1=\theta_0^1} \\ &= -1 + \mathbb{E}\left[\frac{D}{m_0(\mathbf{Z})} \mid \mathbf{Z}\right] \\ &= -1 + \mathbb{E}\left[\frac{D - m_0(\mathbf{Z}) + m_0(\mathbf{Z})}{m_0(\mathbf{Z})} \mid \mathbf{Z}\right] \\ &= -1 + \mathbb{E}\left[\frac{\nu + m_0(\mathbf{Z})}{m_0(\mathbf{Z})} \mid \mathbf{Z}\right] \\ &= -1 + \frac{\mathbb{E}[\nu \mid \mathbf{Z}]}{m_0(\mathbf{Z})} + 1 \\ &= 0. \end{aligned}$$

$$\begin{aligned} &\mathbb{E}[\partial_m \psi_1(W, \theta^1, \rho) \mid \mathbf{Z}] \big|_{(g,m)=(g_0,m_0), \theta^1=\theta_0^1} \\ &= \mathbb{E}\left[(Y - g_0(i, \mathbf{Z})) \frac{D}{m_0(\mathbf{Z})^2} \mid \mathbf{Z}\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[(Y - g_0(i, \mathbf{Z})) \frac{D}{m_0(\mathbf{Z})^2} \mid D, \mathbf{Z}\right] \mid \mathbf{Z}\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[\xi \frac{D}{m_0(\mathbf{Z})^2} \mid D, \mathbf{Z}\right] \mid \mathbf{Z}\right] \\ &= \mathbb{E}\left[\mathbb{E}[\xi \mid D, \mathbf{Z}] \frac{D}{m_0(\mathbf{Z})^2} \mid \mathbf{Z}\right] \\ &= 0. \end{aligned}$$

Table 1: The full results of 1000 IHDP and 100 Twins experiments.

Method	IHDP In-sample		IHDP Out-of-sample		Twins In-sample		Twins Out-of-sample	
	$\sqrt{\epsilon_{PEHE}}$	ϵ_{ATE}	$\sqrt{\epsilon_{PEHE}}$	ϵ_{ATE}	AUC	ϵ_{ATE}	AUC	ϵ_{ATE}
OLS/LR ₁	5.8 ± .3	.73 ± .04	5.8 ± .3	.94 ± .06	.660 ± .005	.004 ± .003	.500 ± .028	.007 ± .006
OLS/LR ₂	2.4 ± .1	.14 ± .01	2.5 ± .1	.31 ± .02	.660 ± .004	.004 ± .003	.500 ± .016	.007 ± .006
k-NN	2.1 ± .1	.14 ± .01	4.1 ± .2	.79 ± .05	.609 ± .010	.003 ± .002	.492 ± .012	.005 ± .004
BART	2.1 ± .1	.23 ± .01	2.3 ± .1	.34 ± .02	.506 ± .014	.121 ± .024	.500 ± .011	.127 ± .024
CF	3.8 ± .2	.18 ± .01	3.8 ± .2	.40 ± .03	*	.029 ± .004	*	.034 ± .008
CEVAE	2.7 ± .1	.34 ± .01	2.6 ± .1	.46 ± .02	.845 ± .003	.022 ± .002	.841 ± .004	.032 ± .003
SITE	.69 ± .0	.22 ± .01	.75 ± .0	.24 ± .01	.862 ± .002	.016 ± .001	.853 ± .006	.020 ± .002
GANITE	1.9 ± .4	.43 ± .05	2.4 ± .4	.49 ± .05	*	.006 ± .002	*	.009 ± .008
BLR	5.8 ± .3	.72 ± .04	5.8 ± .3	.93 ± .05	.611 ± .009	.006 ± .004	.510 ± .018	.033 ± .009
BNN	2.2 ± .1	.37 ± .03	2.1 ± .1	.42 ± .03	.690 ± .008	.006 ± .003	.676 ± .008	.020 ± .007
TARNet	.88 ± .0	.26 ± .01	.95 ± .0	.28 ± .01	.849 ± .002	.011 ± .002	.840 ± .006	.015 ± .002
CFR-WASS	.71 ± .0	.25 ± .01	.76 ± .0	.27 ± .01	.850 ± .002	.011 ± .002	.842 ± .005	.028 ± .003
Dragonnet	1.3 ± .4	.14 ± .01	1.3 ± .5	.20 ± .05	*	.006 ± .005	*	.006 ± .005
MBRL	.522 ± .007	.121 ± .005	.565 ± .008	.133 ± .005	.879 ± .000	.003 ± .000	.874 ± .001	.007 ± .001
MBRL+ θ_1^i	.522 ± .007	.102 ± .004	.565 ± .008	.166 ± .007	.879 ± .000	.003 ± .000	.874 ± .001	.008 ± .000
MBRL+ θ_2^i	.522 ± .007	.114 ± .005	.565 ± .008	.204 ± .008	.879 ± .000	.003 ± .000	.874 ± .001	.006 ± .001

If $i = 0$ and under the noise conditions $\mathbb{E}[\nu | \mathbf{Z}] = 0$ and $\mathbb{E}[\xi | D, \mathbf{Z}] = 0$, then we have

$$\begin{aligned}
& \mathbb{E}[\partial_g \psi_1(W, \theta^0, \rho) | \mathbf{Z}] |_{(g,m)=(g_0,m_0), \theta^0=\theta_0^0} \\
&= -1 + \mathbb{E}\left[\frac{1-D}{1-m_0(\mathbf{Z})} | \mathbf{Z}\right] \\
&= -1 + \mathbb{E}\left[\frac{1-D-m_0(\mathbf{Z})+m_0(\mathbf{Z})}{1-m_0(\mathbf{Z})} | \mathbf{Z}\right] \\
&= -1 + \mathbb{E}\left[\frac{-\nu+1-m_0(\mathbf{Z})}{1-m_0(\mathbf{Z})} | \mathbf{Z}\right] \\
&= -1 - \frac{\mathbb{E}[\nu | \mathbf{Z}]}{1-m_0(\mathbf{Z})} + 1 \\
&= 0.
\end{aligned}$$

$$\begin{aligned}
& \mathbb{E}[\partial_m \psi_1(W, \theta^0, \rho) | \mathbf{Z}] |_{(g,m)=(g_0,m_0), \theta^0=\theta_0^0} \\
&= \mathbb{E}\left[-(Y-g_0(i, \mathbf{Z})) \frac{1-D}{(1-m_0(\mathbf{Z}))^2} | \mathbf{Z}\right] \\
&= \mathbb{E}\left[\mathbb{E}\left[-(Y-g_0(i, \mathbf{Z})) \frac{1-D}{(1-m_0(\mathbf{Z}))^2} | D, \mathbf{Z}\right] | \mathbf{Z}\right] \\
&= \mathbb{E}\left[\mathbb{E}\left[-\xi \frac{1-D}{(1-m_0(\mathbf{Z}))^2} | D, \mathbf{Z}\right] | \mathbf{Z}\right] \\
&= \mathbb{E}\left[-\mathbb{E}[\xi | D, \mathbf{Z}] \frac{1-D}{(1-m_0(\mathbf{Z}))^2} | \mathbf{Z}\right] \\
&= 0.
\end{aligned}$$

We then check if the orthogonal condition (Definition 1) holds for $\psi_2(W, \theta^i, \rho)$.

$$\partial_m \psi_2(W, \theta^i, \rho) = (Y(i) - g(i, \mathbf{Z})) \frac{2((D - m(\mathbf{Z})) - \mathbb{E}[\nu | \mathbf{Z}])}{\mathbb{E}[\nu^2 | \mathbf{Z}]}$$

$$\partial_g \psi_2(W, \theta^i, \rho) = -1 + \frac{((D - m(\mathbf{Z})) - \mathbb{E}[\nu | \mathbf{Z}])^2}{\mathbb{E}[\nu^2 | \mathbf{Z}]}$$

Using the noise condition $\mathbb{E}[\nu | \mathbf{Z}] = 0$, we have

$$\begin{aligned}
& \mathbb{E}[\partial_g \psi_2(W, \theta^i, \rho) | \mathbf{Z}] |_{(g,m)=(g_0,m_0), \theta^i=\theta_0^i} \\
&= -1 + \mathbb{E}\left[\frac{((D - m_0(\mathbf{Z})) - \mathbb{E}[\nu | \mathbf{Z}])^2}{\mathbb{E}[\nu^2 | \mathbf{Z}]} | \mathbf{Z}\right] \\
&= -1 + \frac{1}{\mathbb{E}[\nu^2 | \mathbf{Z}]} \mathbb{E}\left[((D - m_0(\mathbf{Z})) - \mathbb{E}[\nu | \mathbf{Z}])^2 | \mathbf{Z}\right] \\
&= -1 + \frac{1}{\mathbb{E}[\nu^2 | \mathbf{Z}]} \mathbb{E}[(D - m_0(\mathbf{Z}))^2 + (\mathbb{E}[\nu | \mathbf{Z}])^2 \\
&\quad - 2(D - m_0(\mathbf{Z}))\mathbb{E}[\nu | \mathbf{Z}] | \mathbf{Z}] \\
&= -1 + \frac{1}{\mathbb{E}[\nu^2 | \mathbf{Z}]} [\mathbb{E}[\nu^2 | \mathbf{Z}] + (\mathbb{E}[\nu | \mathbf{Z}])^2 - 2(\mathbb{E}[\nu | \mathbf{Z}])^2] \\
&= -1 + \frac{1}{\mathbb{E}[\nu^2 | \mathbf{Z}]} [\mathbb{E}[\nu^2 | \mathbf{Z}] - (\mathbb{E}[\nu | \mathbf{Z}])^2] \\
&= -1 + \frac{\mathbb{E}[\nu^2 | \mathbf{Z}]}{\mathbb{E}[\nu^2 | \mathbf{Z}]} = 0.
\end{aligned}$$

By the model setup $Y = g_0(D, \mathbf{Z}) + \xi$, we have the underlying relation for the potential outcome $Y(i)$ that $Y(i) = g_0(i, \mathbf{Z}) + \xi$. Using the noise condition $\mathbb{E}[\xi | D, \mathbf{Z}] = 0$, we have

$$\begin{aligned}
& \mathbb{E}[\partial_m \psi_2(W, \theta^i, \rho) | \mathbf{Z}] |_{(g,m)=(g_0,m_0), \theta^i=\theta_0^i} \\
&= \mathbb{E}\left[\mathbb{E}\left[\xi \frac{2((D - m_0(\mathbf{Z})) - \mathbb{E}[\nu | \mathbf{Z}])}{\mathbb{E}[\nu^2 | \mathbf{Z}]} | D, \mathbf{Z}\right] | \mathbf{Z}\right] \\
&= \mathbb{E}\left[\frac{2((D - m_0(\mathbf{Z})) - \mathbb{E}[\nu | \mathbf{Z}])}{\mathbb{E}[\nu^2 | \mathbf{Z}]} \mathbb{E}[\xi | D, \mathbf{Z}] | \mathbf{Z}\right] \\
&= 0.
\end{aligned}$$

Therefore, the noise conditions $\mathbb{E}[\xi | D, \mathbf{Z}] = 0$ and $\mathbb{E}[\nu | \mathbf{Z}] = 0$ are sufficient for the score functions ψ_1 and ψ_2 satisfying the orthogonal condition. \square

3.2 Proof of Property 1

Proof. Using the noise condition $\mathbb{E}[\xi \mid D, \mathbf{Z}] = 0$, we have

$$\begin{aligned} & \mathbb{E}[(Y - g_0(D, \mathbf{Z}))(D - m_0(\mathbf{Z}))] \\ &= \mathbb{E}[\mathbb{E}[(Y - g_0(D, \mathbf{Z}))(D - m_0(\mathbf{Z})) \mid D, \mathbf{Z}]] \\ &= \mathbb{E}[(D - m_0(\mathbf{Z}))\mathbb{E}[(Y - g_0(D, \mathbf{Z})) \mid D, \mathbf{Z}]] \\ &= \mathbb{E}[(D - m_0(\mathbf{Z}))\mathbb{E}[\xi \mid D, \mathbf{Z}]] \\ &= 0. \end{aligned}$$

□

References

- [Chernozhukov *et al.*, 2018] Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins. Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1):C1–C68, 2018.
- [Huang *et al.*, 2021] Yiyan Huang, Cheuk Hang Leung, Xing Yan, and Qi Wu. Higher-order orthogonal causal learning for treatment effect. *arXiv preprint arXiv:2103.11869*, 2021.