

Identifiable Energy-based Representations: An Application to Estimating Heterogeneous Causal Effects

AISTATS 2022

Yao Zhang*, Jeroen Berrevoets*, Mihaela van der Schaar



van_der_Schaar
\ LAB

vanderschaar-lab.com



UNIVERSITY OF
CAMBRIDGE



jb2384@cam.ac.uk



[@J_Berrevoets](https://twitter.com/J_Berrevoets)



linkedin.com/in/jeroenberrevoets

We collect many variables to identify the causal effect

Consider the conditional average treatment effect (CATE):

$$\mathbb{E}[Y(1) - Y(0)|X]$$

Where $X \in \mathbb{R}^d$ are an individual's covariates, and $Y(W) \in \mathbb{R}$ is the potential outcome, with $W \in \{0, 1\}$ the treatment.

Before we can predict the CATE, we have to assume that *all* confounding variables are in X . Otherwise, we may be predicting a spurious (or biased) treatment effect.

In practice, people will collect a **lot** of variables to make sure all of the confounders are actually measured

While it does help avoiding spurious estimates, it also results in poor estimates due to *the curse of dimensionality*, especially in low samples.



We could reduce dimensions...

Of course, one could reduce dimensions using: PCA, AEs, ...

... but they have downsides:

- Linear
- Not consistent

With linear dim-reduction, we lose valuable non-linear interactions: **not good**

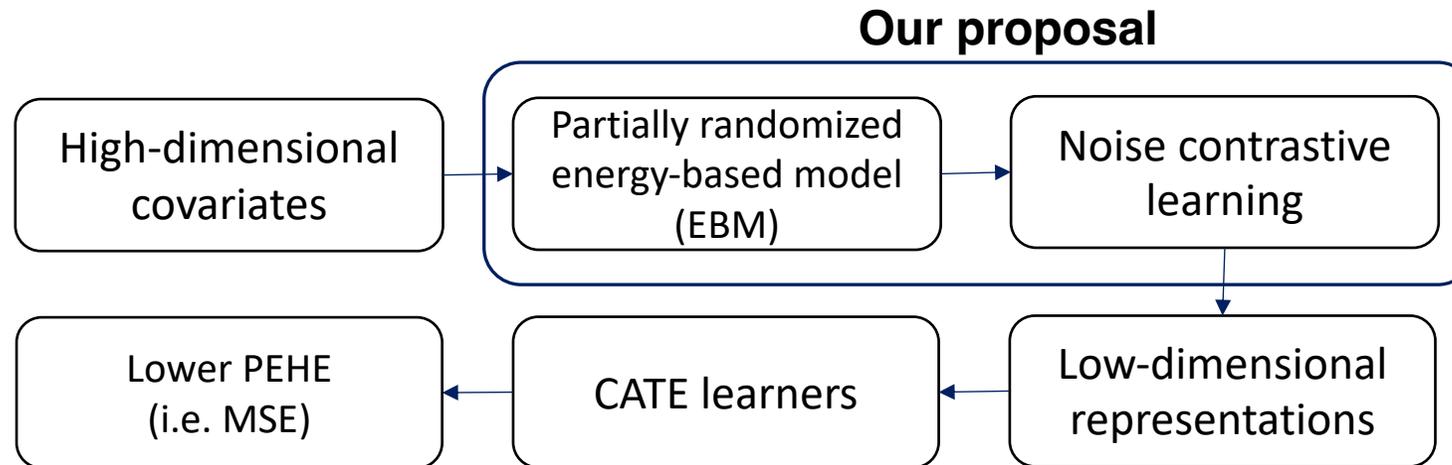
With non-consistent methods, we throw away some of the CATE-literature's advances: **also not good**

We introduce a method that is both non-linear, **and** consistent



We propose an identifiable EBM to reduce dimensions

Let's have a high-level look at our proposal



We propose an identifiable EBM to reduce dimensions

Specifically...

- Partially randomized EBM p_θ is a combination of multiple EBMs $p_{\theta,j}$ with a shared representation $f_\theta(x)$ and β 's randomized and orthogonalized:

$$p_\theta = (p_{\theta,j} : j \in [k])$$

$$p_{\theta,j}(x) = Z_{\theta,j}^{-1} \exp[-\beta_j^\top f_\theta(x)], \quad \forall j \in [k]$$

- With $f_\theta(x)$ parameterized by a neural network, we show the partially randomized p_θ has two properties:
 - ✓ $f_\theta(x)$ is identifiable up to some universal constants.
 - ✓ universal capacity for approximating any density functions
- Split the clean samples between $p_{\theta,j}$, then generate corrupted samples conditionally,

$$\bar{X}_i = (X_i, \tilde{X}_{i1}, \dots, \tilde{X}_{ib}) \sim p_X(x) \prod_{a=1}^b p_{\tilde{X}|X}(\tilde{x} | x)$$

- Train $p_{\theta,j}$ jointly so the shared f_θ goes through noise contrastive training on full samples.
- We show f_θ converges to limits that are different by some universal constants



Some results on prediction

Methods		X-Learner		DR-Learner		T-Learner		R-Learner	
EBM		\times	\checkmark	\times	\checkmark	\times	\checkmark	\times	\checkmark
d	n	<i>Synth. data with increasing sample size and increasing dimensions</i>							
50	100	2.309 \pm .00	1.994 \pm .02	4.594 \pm .56	2.017 \pm .04	2.441 \pm .00	1.993 \pm .01	3.194 \pm .26	1.982 \pm .04
100	250	2.779 \pm .00	2.018 \pm .01	4.056 \pm .32	2.154 \pm .39	2.838 \pm .00	2.019 \pm .01	3.702 \pm .23	2.018 \pm .01
150	500	2.618 \pm .00	2.000 \pm .01	3.030 \pm .12	2.001 \pm .01	2.641 \pm .00	2.000 \pm .01	2.877 \pm .08	2.000 \pm .01
200	1k	2.185 \pm .00	1.940 \pm .01	2.283 \pm .02	1.941 \pm .01	2.189 \pm .00	1.939 \pm .01	2.271 \pm .01	1.940 \pm .01
250	1.5k	2.267 \pm .00	1.949 \pm .02	2.427 \pm .01	1.976 \pm .00	2.271 \pm .00	1.948 \pm .01	2.436 \pm .02	1.949 \pm .02
n		<i>Synth. data with increasing sample size and dimensions fixed at $d = 100$</i>							
100		2.134 \pm .00	1.927 \pm .01	24.61 \pm 9.9	2.096 \pm .09	2.279 \pm .00	1.929 \pm .01	3.192 \pm .13	1.925 \pm .01
250		2.779 \pm .00	2.018 \pm .01	4.056 \pm .32	2.154 \pm .39	2.838 \pm .00	2.019 \pm .01	3.702 \pm .23	2.018 \pm .01
500		2.155 \pm .00	2.056 \pm .02	2.334 \pm .07	2.273 \pm .67	2.166 \pm .00	2.053 \pm .02	2.271 \pm .05	2.056 \pm .02
1k		2.059 \pm .00	1.964 \pm .02	2.105 \pm .01	2.016 \pm .16	2.061 \pm .00	1.964 \pm .02	2.086 \pm .01	1.965 \pm .02
1.5k		2.013 \pm .00	1.998 \pm .02	2.043 \pm .01	1.998 \pm .02	2.014 \pm .00	1.998 \pm .02	2.024 \pm .01	1.991 \pm .02
n		<i>Twins ($d = 48$) with increasing sample size</i>							
500		0.214 \pm .00	0.144 \pm .00	0.236 \pm .04	0.182 \pm .05	0.221 \pm .00	0.145 \pm .00	0.222 \pm .02	0.145 \pm .00
1k		0.294 \pm .00	0.162 \pm .00	0.348 \pm .12	0.173 \pm .03	0.301 \pm .00	0.162 \pm .01	0.532 \pm .11	0.161 \pm .00
1.5k		0.165 \pm .00	0.154 \pm .00	0.189 \pm .06	0.159 \pm .01	0.165 \pm .00	0.154 \pm .00	0.172 \pm .01	0.154 \pm .00
2k		0.167 \pm .00	0.156 \pm .00	0.197 \pm .03	0.159 \pm .00	0.167 \pm .00	0.156 \pm .00	0.222 \pm .05	0.157 \pm .00
2.5k		0.297 \pm .00	0.153 \pm .00	0.390 \pm .19	0.156 \pm .00	0.297 \pm .00	0.153 \pm .00	0.358 \pm .22	0.153 \pm .00

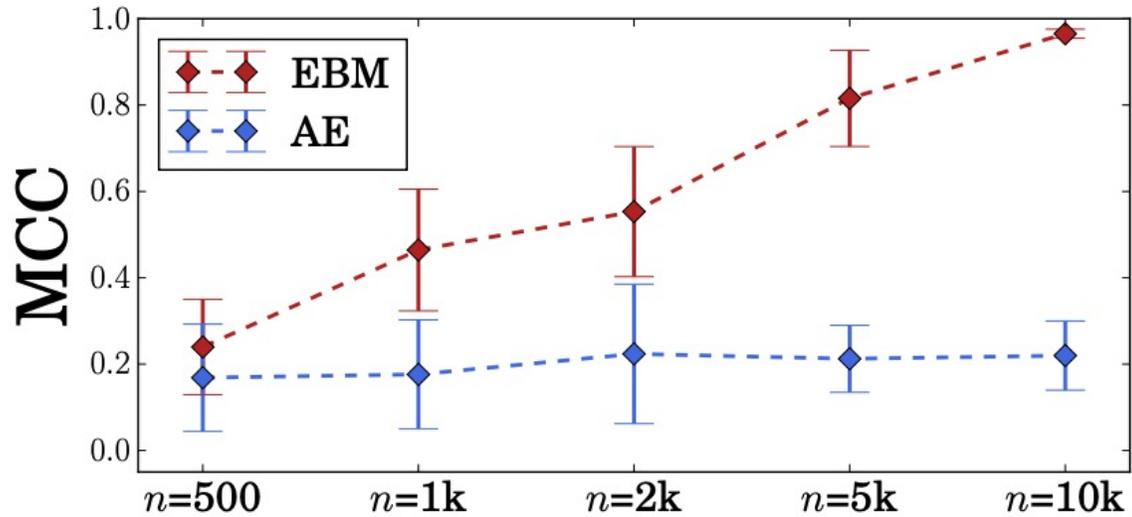
Lower PEHE across learners

Methods		PCA	FA	SE	Isomap	K-PCA	AE	EBM
n		<i>Twins ($d = 48$) with increasing sample size</i>						
500		1.092 \pm .11	1.758 \pm 1.1	1.011 \pm .00	1.006 \pm .00	1.015 \pm .00	0.580 \pm .03	0.145 \pm .00
1k		1.015 \pm .00	0.963 \pm .00	1.010 \pm .00	1.004 \pm .00	1.010 \pm .00	0.549 \pm .04	0.161 \pm .00
1.5k		1.014 \pm .00	0.965 \pm .00	1.005 \pm .00	1.006 \pm .00	1.012 \pm .00	0.546 \pm .04	0.154 \pm .00
2k		1.013 \pm .00	0.957 \pm .00	1.009 \pm .00	1.007 \pm .00	1.013 \pm .00	0.579 \pm .03	0.157 \pm .00
2.5k		1.007 \pm .00	0.951 \pm .00	1.002 \pm .00	1.006 \pm .00	1.006 \pm .00	0.542 \pm .04	0.153 \pm .00

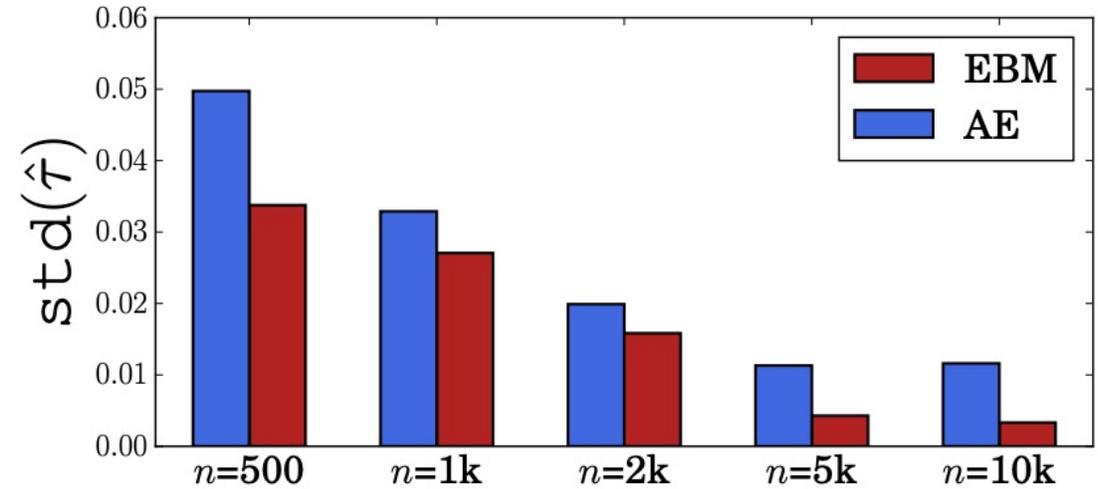
Lower PEHE than other dimensionality reduction methods

Some results on Identifiability

Converging Representations and CATE estimates



Increasing Mean Correlation Coefficient (MCC) with sample size



Decreasing standard error on CATE-estimates

