

Estimation of Individual Treatment Effect in Latent Confounder Models via Adversarial Learning

1. INTRODUCTION

Estimating Individual Treatment Effect (ITE)

- A key challenge with observational data
→ Will treatment A help patient B recover?
- Most previous work relies on the **unconfoundedness assumption**, which posits that all the confounders are measurable; see Figure 1(a)

Latent Confounder Model

- In practice, there are often unmeasurable (latent) confounders; see Figure 1(b)
→ Socio-economic status affects medications available to a patient and her health
- If not appropriately accounted for, the estimated ITE will be subject to **confounding bias**

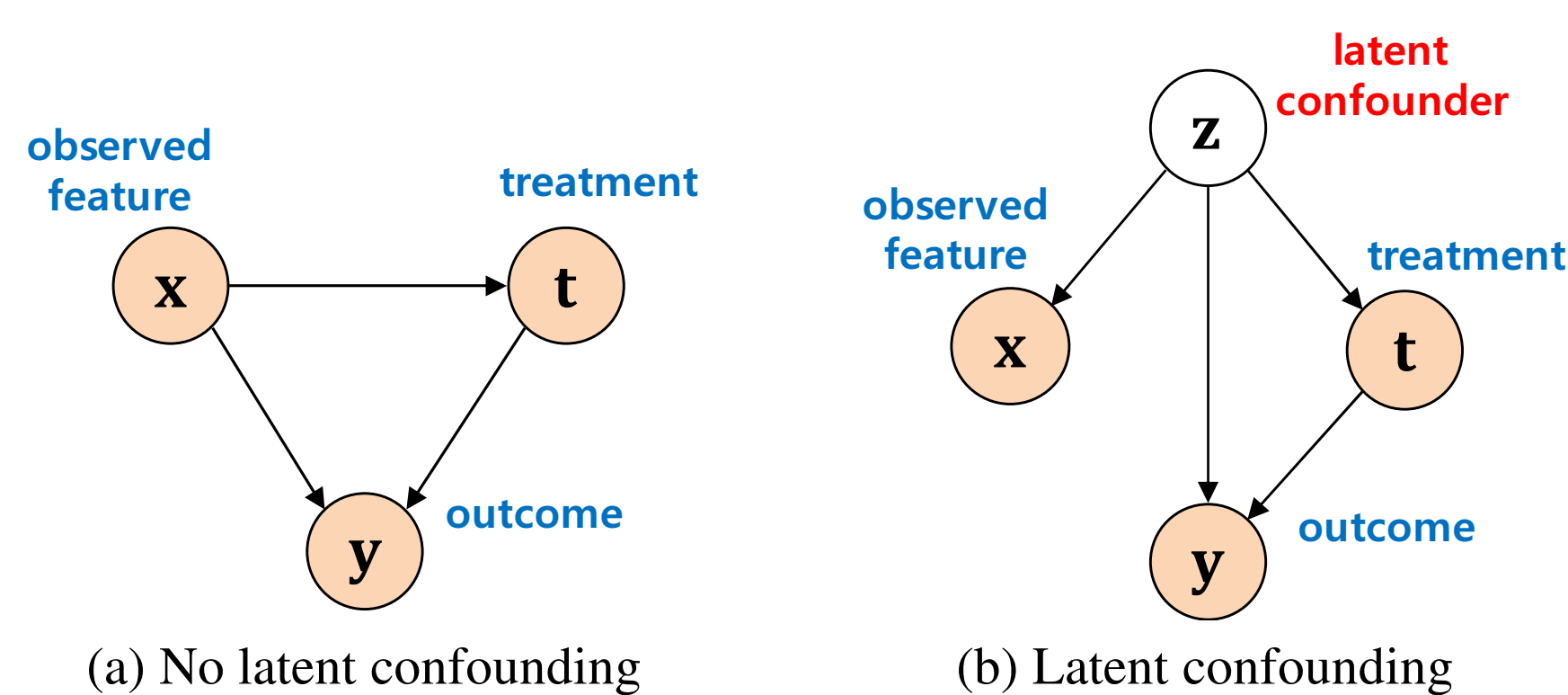


Figure 1: Causal diagrams

2. PROBLEM FORMULATION

Observational Dataset: $\mathcal{D} = \{(x_i, t_i, y_i)\}_{i=1}^N$

- x_i : feature vector
- t_i : treatment (we assume $t \in \{0, 1\}$)
- y_i : outcome vector
- z_i : latent confounder that is not in \mathcal{D}

Objective:

- **Estimate ITE without confounding bias:**

$$ITE(x) = \mathbb{E}[y|x, do(t=1)] - \mathbb{E}[y|x, do(t=0)]$$

How to Account for Latent Confounding?

- We assume the latent confounder model in Figure 1(b); x is treated as a **proxy variable** that provides a noisy view of z
- We can identify $p(y|x, do(t=1))$ (or, similarly, $p(y|x, do(t=0))$) by

$$p(y|x, do(t=1)) = \int_{\mathbf{z}} p(y|z, x, do(t=1))p(z|x, do(t=1))dz$$

$$= \int_{\mathbf{z}} p(y|z, x, t=1)p(z|x)dz,$$

- We adopt an **adversarial learning framework** to learn $p(y|z, x, t)$ and $p(z|x)$

6. PERFORMANCE METRIC

Precision in Estimation of Heterogeneous Effect (PEHE)

- A commonly used metric to quantify the **goodness** of ITE estimation

$$\epsilon_{PEHE} = \frac{1}{N} \sum_{i=1}^N \left((y_i(1) - y_i(0)) - (\hat{y}_i(1) - \hat{y}_i(0)) \right)^2$$

7. BENCHMARKS

- CFR_{WASS}: counterfactual reg. w/ Wasserstein
- CMGP: causal multi-task Gaussian process
- CEVAE: causal effect VAE (CEVAE)

3. CEGAN ARCHITECTURE & COMPONENTS

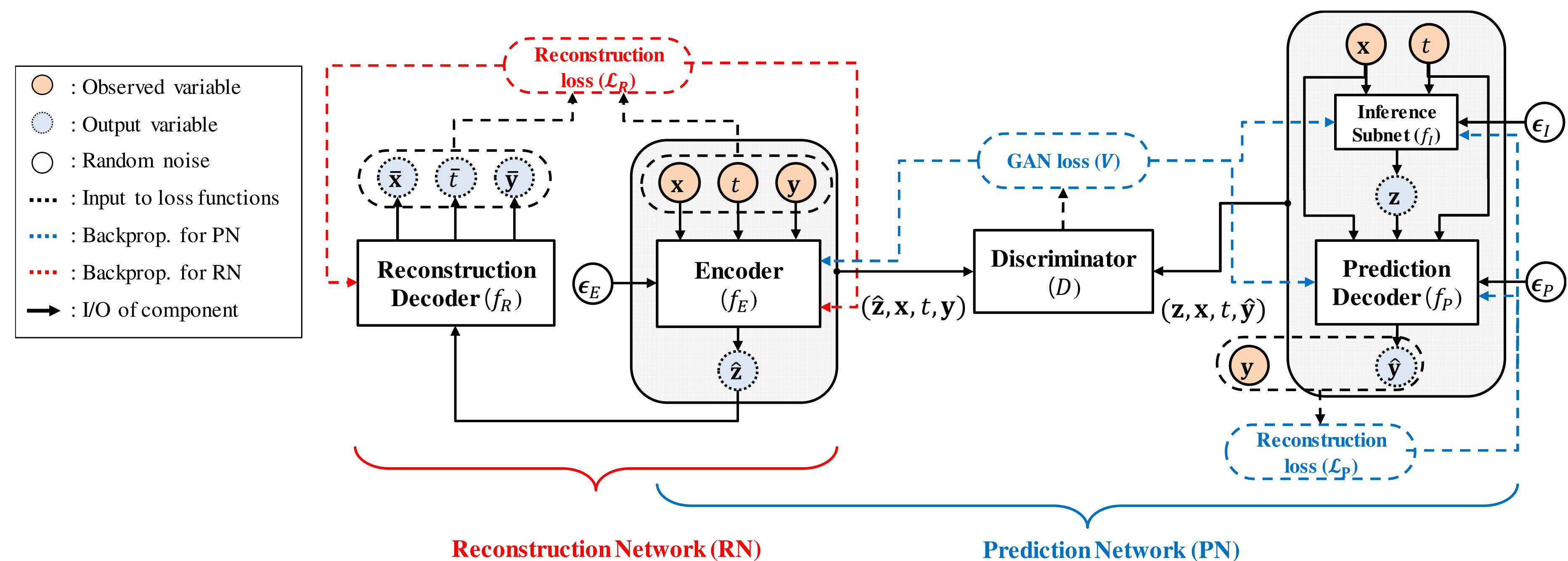


Figure 2: CEGAN architecture

Prediction Network:

- **Generator (G):**
 - Comprises encoder (f_E), inference subnet (f_I), and prediction decoder (f_P) which output $\hat{z} \sim q_E(z|x, t, y)$, $z \sim q_I(z|x, t)$, $\hat{y} \sim q_P(y|z, x, t)$ (via universal approximator technique)
 - Constructs samples of tuples (z, x, t, y) drawn from two joint distributions, i.e., $q_E(z, x, t, y) = p_d(x, t, y)q_E(z|x, t, y)$ and $q_P(z, x, t, y) = p_d(x, t)q_I(z|x, t)q_P(y|z, x, t)$
 - Tries to **fool the discriminator** into believing the tuples are drawn from the same distribution
- **Discriminator (D):**
 - **Distinguishes between tuples** (z, x, t, y) that are drawn from $q_E(z, x, t, y)$ and $q_P(z, x, t, y)$

Reconstruction Network:

- Comprises the same encoder (f_E) and reconstruction decoder (f_R)
- Nudge f_E to learn a meaningful mapping by reconstructing its original input

CEGAN **matches the two distribution** by playing an **adversarial game** between G and D .

4. EXPERIMENTS: SEMI-SYNTHETIC

TWINS Dataset:

- Based on records of twin births in the USA from 1989-1991
- **Artificially create a binary treatment:** $t = 1$ ($t = 0$) denotes being born the heavier (lighter)
- Outcome corresponds to the **mortality** of each of the twins in their first year

Data Generation Process:

- Select GESTAT (i.e. the gestational age in weeks) as the latent confounder z .
- Assign binary treatment $t_i \sim \text{Bern}(\sigma(wz_i))$, where $w \sim \mathcal{N}(10, 0.1^2)$.
- Choose outcome of the heavier twin, $y_i(1)$, if $t_i = 1$ and that of the lighter twin, $y_i(0)$, if $t_i = 0$.

Table 1: Comparison of $\sqrt{\epsilon_{PEHE}}$ (mean \pm std)

Method	no latent confounding		latent confounding	
	In-sample	Out-sample	In-sample	Out-sample
LR-1	0.365 \pm 0.00	0.367 \pm 0.00	0.413 \pm 0.01	0.423 \pm 0.02
LR-2	0.404 \pm 0.02	0.411 \pm 0.02	0.442 \pm 0.02	0.454 \pm 0.02
kNN	0.486 \pm 0.02	0.506 \pm 0.02	0.492 \pm 0.02	0.515 \pm 0.02
CForest	0.356\pm0.01	0.372 \pm 0.01	0.417 \pm 0.02	0.429 \pm 0.02
CMGP	0.367 \pm 0.01	0.365 \pm 0.01	0.430 \pm 0.05	0.438 \pm 0.05
CFR _{WASS}	0.371 \pm 0.03	0.371 \pm 0.03	0.427 \pm 0.05	0.438 \pm 0.05
CEVAE	0.363 \pm 0.00	0.364 \pm 0.00	0.423 \pm 0.00	0.428 \pm 0.00
CEGAN	0.363 \pm 0.00	0.362\pm0.00	0.369\pm0.00	0.369\pm0.00

- “no latent confounding” includes GESTAT in the observational data \mathcal{D}
– Causal model \rightarrow Figure 1(a)
- “latent confounding” excludes GESTAT from the observational data \mathcal{D}
– Causal model \rightarrow Figure 1(b)

6. PERFORMANCE METRIC

Precision in Estimation of Heterogeneous Effect (PEHE)

- A commonly used metric to quantify the **goodness** of ITE estimation

$$\epsilon_{PEHE} = \frac{1}{N} \sum_{i=1}^N \left((y_i(1) - y_i(0)) - (\hat{y}_i(1) - \hat{y}_i(0)) \right)^2$$

7. BENCHMARKS

- CFR_{WASS}: counterfactual reg. w/ Wasserstein
- CMGP: causal multi-task Gaussian process
- CEVAE: causal effect VAE (CEVAE)

5. EXPERIMENTS: SYNTHETIC

Toy Example:

- To assess the robustness of CEGAN to latent confounders (due to noise in the proxy variables)

Data Generation Process:

- Assume latent confounding model in Figure 1(b):

$$z_{ij} \sim \mathcal{N}(3(\mu - 1), 1^2), \quad j = 1, \dots, d_z,$$

$$\mu \sim \text{Bern}(0.5), \quad \mathbf{n} \sim \mathcal{N}(0, \zeta^2 \mathbf{I})$$

$$\mathbf{x}_i | \mathbf{z}_i = \mathbf{z}_i + \mathbf{n},$$

$$t_i | \mathbf{z}_i \sim \text{Bern}(\sigma(0.25 \cdot z_{id_z})),$$

$$y_i | \mathbf{z}_i, t_i = \sigma(\mathbf{1}^T \mathbf{z}_i + (2t_i - 1))$$

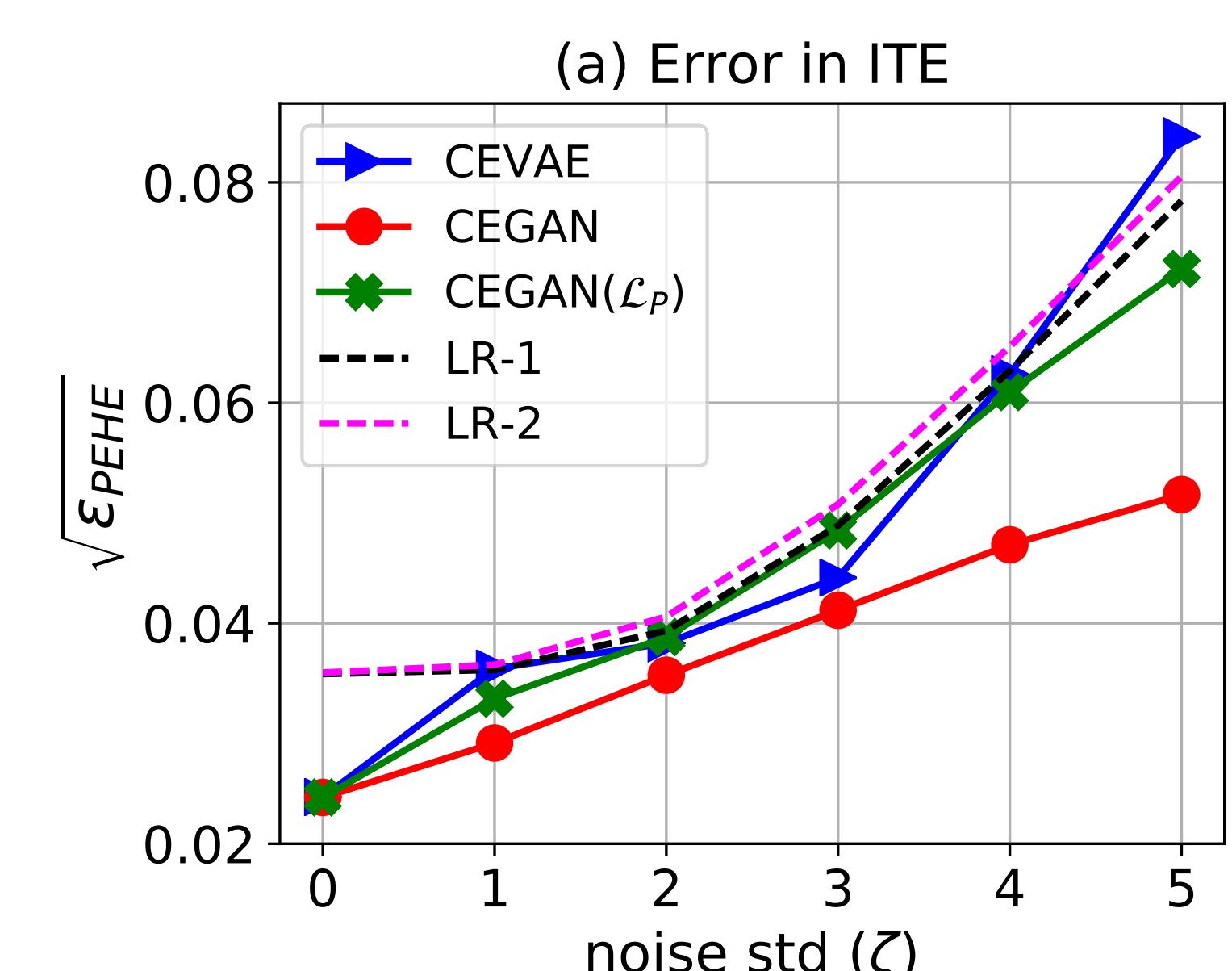


Figure 3: $\sqrt{\epsilon_{PEHE}}$ with respect to noise level ζ