

Learning Causal Effect for High-dimensional Observation Data with Unmeasured Confounding

Hedong YAN, Computer Science Department

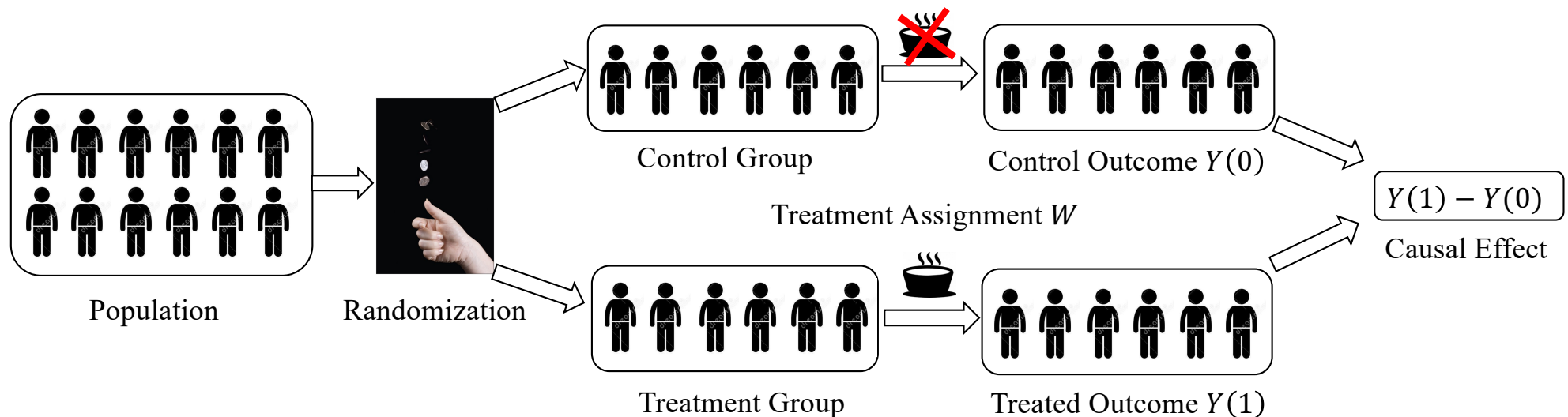
Hong Kong Baptist University

Supervisor: Yiu-ming Cheung

- Background
- Application
- Problem and Challenge
- Related Work
- Our Preliminary work
- Our Approach to Address Challenges
- Futural Plan

Background

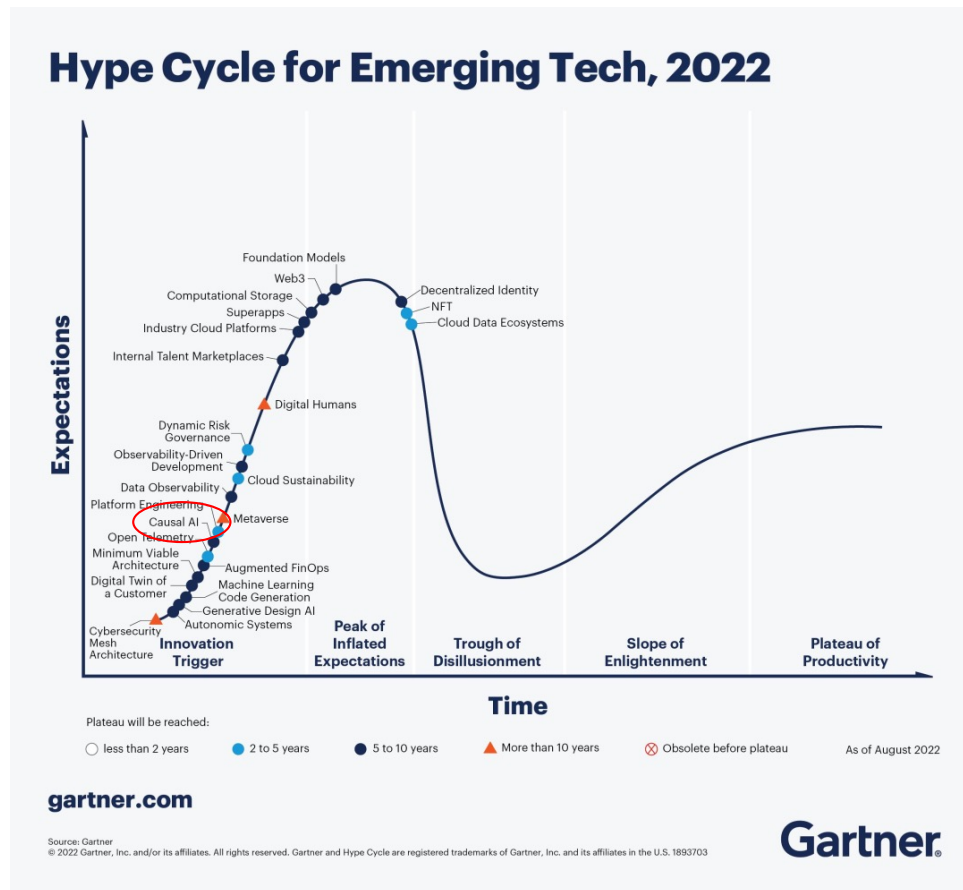
- Causal Effect and Randomized Control Trail
- "Unknown potential yields" of Neyman's agriculture experiment



Background

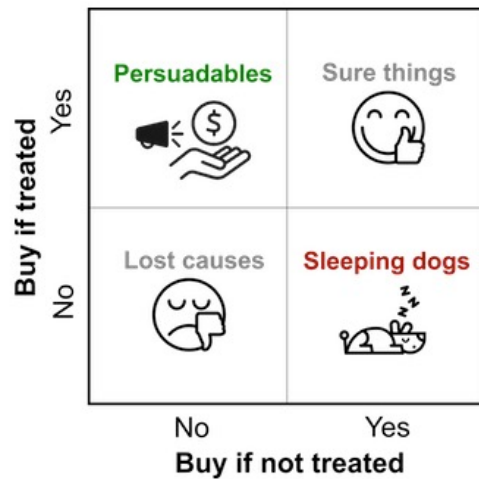
- Limitation and Opportunity
 - RCT can only give a *population-level* conclusion.
 - RCT can not be performed due to *immorality* and *high cost*.
 - Observational *data* upsurges.
- As an alternative, learning causal effects from the observational dataset is not totally impossible.

Background



- Causal AI is still in the innovation trigger stage.

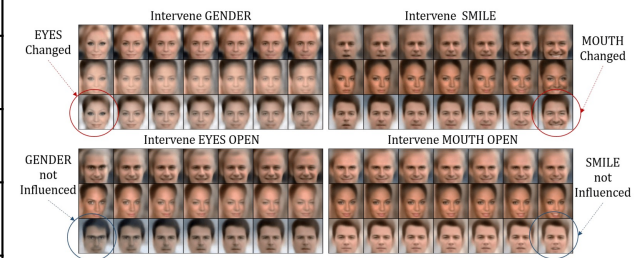
Application



Uplift Marketing

Individuals	$E(Y_i(1))$	$E(Y_i(0))$	Recommendation
u_1	Good	Good	No
u_2	Bad	Bad	No
u_3	Good	Bad	Yes
u_4	Bad	Good	No
u_5	Good	Good	No

Individual Drug Recommendation



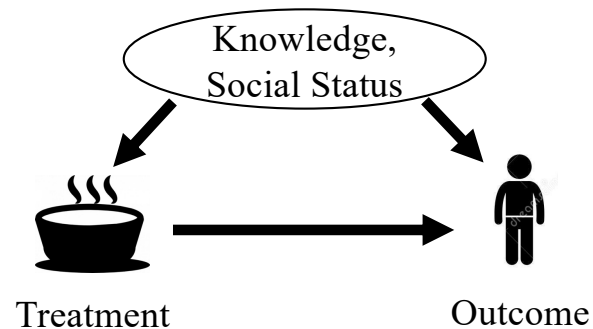
Causal Feature

Problem and Challenge

- Open Problem: Learning causal effects from observational data
- Challenge:
 - Hidden confounding
 - High dimensionality
 - Robustness

Challenge 1: Hidden confounding

- Treatment assignment is unknown and not randomized in observation data. We can not identify causal effects from observational data.



Challenge 2: High dimensionality

- Potential dependency is exponential.
 - For example, the number of acyclic-directed mixed graphs is $O(2^{n^2-n} * n! * 1.3^{n^2})$ where n is the number of variables.

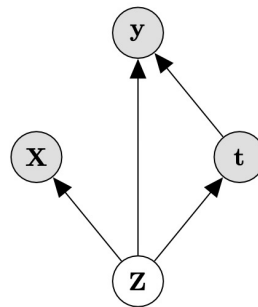
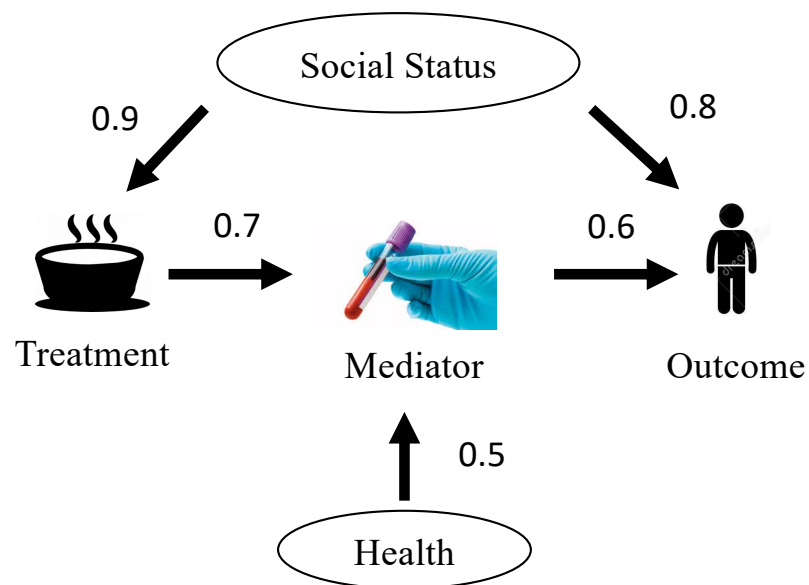


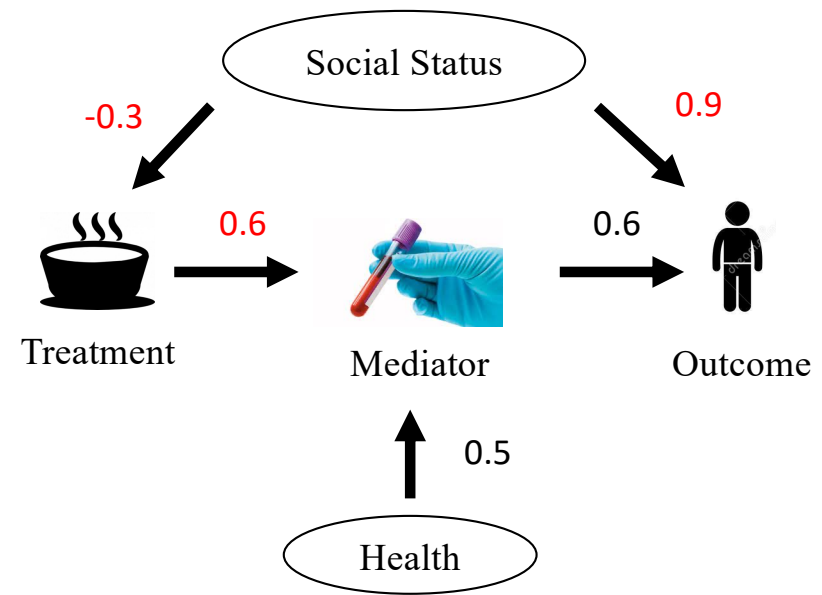
Figure 1: Example of a proxy variable. t is a treatment, e.g. medication; y is an outcome, e.g. mortality. Z is an unobserved confounder, e.g. socio-economic status; and X is noisy views on the hidden confounder Z , say income in the last year and place of residence.

Challenge 3: Robustness

- What if the dependency relationship (structure and parameters) changed?



$$M_1 : E(Y(W)) = E_{X|W=0}(E_W(E(Y|W, X)))$$



M_1 still work well?

Related works: Algorithm

Four Components

Counterfactual Imputation: impute the influence, such as loss value, of counterfactual data on our model.

Balancing Regularization: treatment group and control group are sampled from the same distribution.

Potential Outcome Prediction: learning potential outcomes prediction function for causal effect estimation.

Estimand Modeling: learning a function for the specific causal quantity that we want.

Table 1: Algorithms of causal effect learning from observation data. BLR/BNN: Shalit et al. (2017); TARNet/CFR-MMD/CFR-Wasserstein: Johansson et al. (2016); Dargonet: Shi et al. (2019); X-learner: Künzel et al. (2019); CEVAE: Louizos et al. (2017); Deconfounder: Wang & Blei (2019); GANITE: Yoon et al. (2018); SITE: Yao et al. (2018); DRNets: Schwab et al. (2020); VCNets: Nie et al. (2021).

Algorithms	Learning Stage	Counterfactual Imputation	Balancing Regularization	Potential Outcome Prediction	Estimand Modeling	Hidden Confounding
BLR BNN	Two-stage	Nearest Neighbor	Moment's Difference	Linear Neural Network	None	None
TARNet CFR-MMD CFR-Wasserstein Dargonet	End-to-end	Perfect Counterfactual	None MMD Wasserstein CrossEntropy	Twin Neural Networks	None	None
X-Learner	Three-stage	Perfect Counterfactual	None	Twin BARTs	Yes	None
CEVAE	End-to-End	Perfect Counterfactual	Bayesian Variational Inference Network	Model Network	None	Proxy variables
Deconfounder	Two-stage	Perfect Counterfactual	Posterior Predictive Check of Factor Model	Linear	None	Proxy variables
GANITE	Two-stage	Counterfactual GAN	None	ITE GAN	None	None
SITE	End-to-end	PDDM Similarity	Middle Point Distance	Neural Network	None	None
DRNets VCNets	End-to-end	Nearest Neighbor	None	Treatment-Dose Networks Varying Coefficient Network	None	None

Related works: Benchmark

Table 3: Causal Dataset. Causeme: 202; JustCause: Hawkins & Kim (2021); e-CARE: Du et al. (2022); IHDP: Hill (2011); News: Johansson et al. (2016); Twins: Louizos et al. (2017); Jobs: Shalit et al. (2017); Movies: Wang & Blei (2019); GWAS: Song et al. (2015).

Type	Name	Introduction	Website
Benchmark	Causeme	time-series	https://causeme.uv.es/
Benchmark	JustCause	support IHDP, ACIC etc.	https://justcause.readthedocs.io/en/latest/
Benchmark	e-CARE	reasoning and explanation for NLP	https://scir-sp.github.io
Dataset	IHDP	home visits and IQ testing	https://github.com/vdorie/npci
Dataset	News	New York Times corpus	https://archive.ics.uci.edu/ml/datasets/Bag+of+Words
Dataset	Twins	birth weight and mortality	http://www.nber.org/data/linked-birth-infant-death-data-vital-statistics-data.html
Dataset	Jobs	labor earnings	https://users.nber.org/~rdehejia/data/.nswdata3.html
Dataset	Movies	Movie income and stars	https://www.kaggle.com/tmdb
Dataset	GWAS	genome-wide association studies	https://github.com/StoreyLab/gctest
Competition	ACIC 2022	conference challenge	https://acic2022.mathematica.org/data
Competition	PCIC 2022	conference challenge	https://pattern.swarma.org/pcic/competition.html

- Benchmarking Difficulty:
 - Lacking randomized interventions and well-matched twins
 - Counterfactual missing
 - High deployment cost

Related works: Dimensionality Reduction

Table 2: Dimensionality reduction assumptions. G: Gaussian; I: independent; nG: non-Gaussian; \perp : orthogonal; \rightarrow : generate; ANN: additive normal noise; DAG: directed acyclic graph.

Method	Mapping	$p(\mathbf{z})$	$p(\mathbf{x})$
PCA	Linear	IG	IG
ICA	Linear	InG	InG+G
t -SNE	Nonlinear	Local continuity	Local continuity
β VAE	Nonlinear	IG with β	\
NGCA	Linear	$G \perp nG$	ANN
LinGAM	Linear	$G \rightarrow nG$	ANN with DAG

Related works: Toolbox

Table 4: Causal Packages. Tetrad: Ramsey et al. (2018); CausalDiscoveryToolbox: Kalainathan & Goudet (2019); Ananke: Nabi et al. (2020), Lee & Shpitser (2020), Bhattacharya et al. (2020); EconML: Keith Battocchi (2019); dowhy: Sharma et al. (2019); causalml: Chen et al. (2020); Causal-Curve: Kobrosly (2020); grf: Athey et al. (2019); dosearch: Tikka et al. (2021); causaleffect: Tikka & Karvanen (2017); dagitty: Textor et al. (2016).

Motivation	Toolbox	Support Team	Introduction
Causal Learning	causal-learn	CMU, DMIR, Gong Mingming team, Shouhei Shimizu team	python version of Tetrad
	Tetrad	CMU	Java
	CausalDiscoveryToolbox	FenTechSolutions	python, DAG/Pair, dataset, independence, structure learning, metrics
	gCastle	Huawei Noah	python, data generation and process, causal structure learning, metrics
Causal Reasoning	tigramite	Jakob Runge	python, learning from time-series data
	Ananke	Ilya Shpitser team	python, support do-calculus
	EconML	Microsoft	python, Econometrics
	dowhy	Microsoft	python
	causalml	Uber	python, campaign target optimization, personalized engagement
	CausalImpact	Google	R, time-series, advertisement and click
	WhyNot	John Miller	python, simulator and environment
	Causal-Curve	Kobrosly, R.W.	python, continuous variable such as price, time and income
	grf	grf-lab of Standford	R
	dosearch	Sanntu Tikka	R
End-to-End	causaleffect	Sanntu Tikka	R
	dagitty	\	R, support adjustment formula
	causalnex	QuantumBlack	python, 0.11.1, structure learning, domain knowledge, estimation
	Y-learn	CSDN	python, June 2022

Our Preliminary works

- Open Package: Identification and Structural Causal Model
- A Rejected Paper (UAI 2022 January): OOD Robustness

Our Preliminary works: Open Package

herdonyan / EstimandIdentification Public

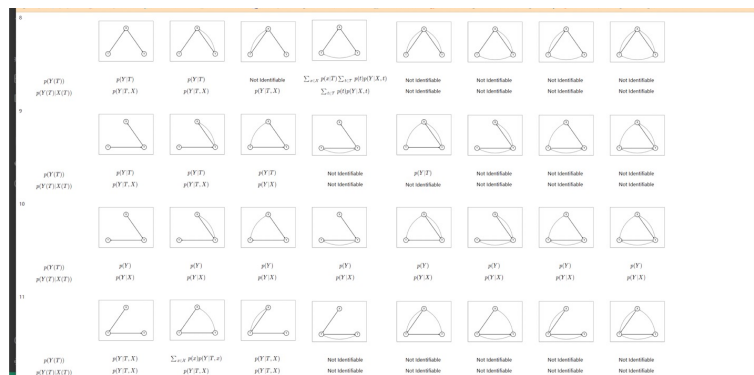
<> Code Issues Pull requests Actions Projects Wiki Security Insights Settings

main 1 branch 0 tags Go to file Add file <> Code

File	Commit Message	Time Ago
3variablesfigs	change address to raw pictures	2 months ago
LICENSE	Create LICENSE	4 days ago
NOTICE	Create NOTICE	4 days ago
README.md	Update README.md	7 months ago
id.py	Update id.py	4 days ago
idc.py	Update idc.py	4 days ago
scm.py	Update scm.py	4 days ago

• Characteristics

- ✓ Automatic Identification
- ✓ Sampling data from given SCM with parameters



Our Preliminary works: Robustness

- Novelty: introduce auto identification into causal effect estimation.

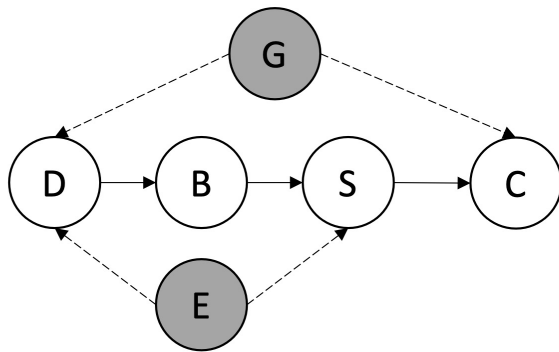
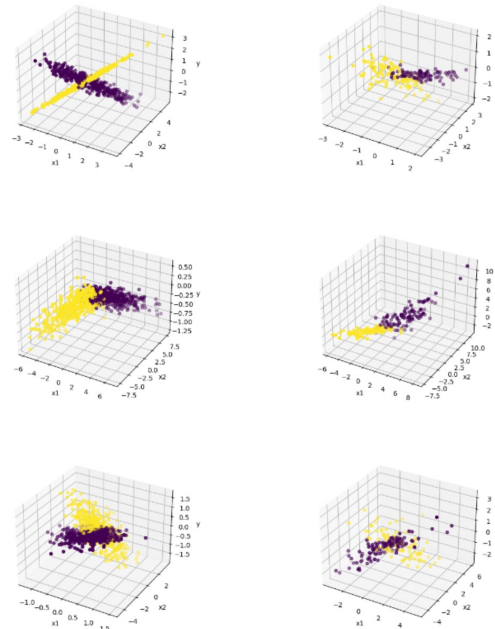


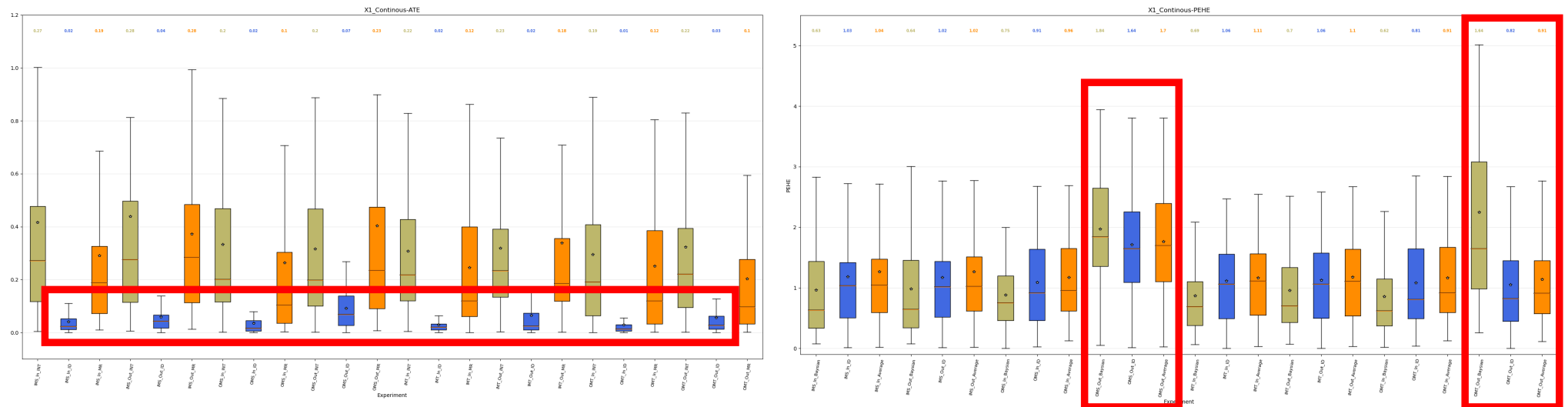
Figure 1: Example of four variables. D means dopamine; B means senior brain activity (frontal lobe); G means unobserved gene/physique; E means social environment not easy to measure. S means smoking behaviour, and C means cancer. For example, $E \rightarrow D$ may represent some life pressures, and $E \rightarrow S$ may be unconscious mimic nature.



- The left column is train data, and the right column is test data. Yellow and purple indicate smoking or not.
- X_1 and X_2 are variable D and B . Y is variable C .

- In our simulation, we want to calculate the causal effect of smoking on cancer.
- We use $p(c|do(s)) = \frac{\sum_D p(D)p(s, c|D, b)}{\sum_D p(D)p(s|D, b)}$ and maximum likelihood to estimate $E(c|do(s))$ for all individuals.

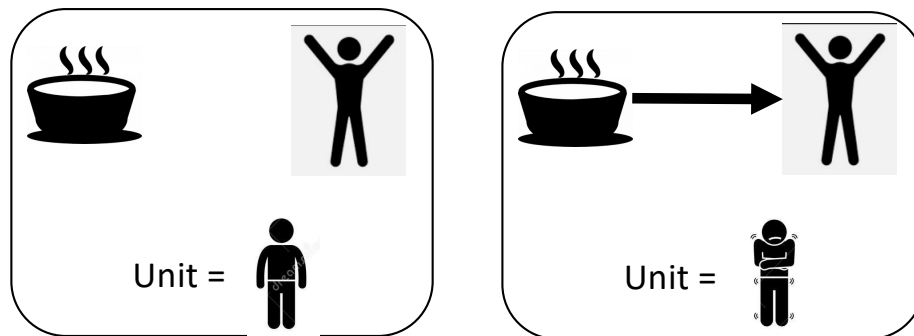
Our Preliminary works: Robustness



- In unbiasedness testing, estimations after identification are more unbiased than MR Freedman (2008) and INT Lin (2013) from ATE estimation results in both discrete and continuous cases. Considering estimation variance, it got better performance when outer mechanisms (dashed line) are changed.

Our Approach to Address Challenges

- Hidden Confounding: Individual Diagram
- Novelty: Will be the **first** to learn **Individual** Structural Causal Model for causal effect estimation.



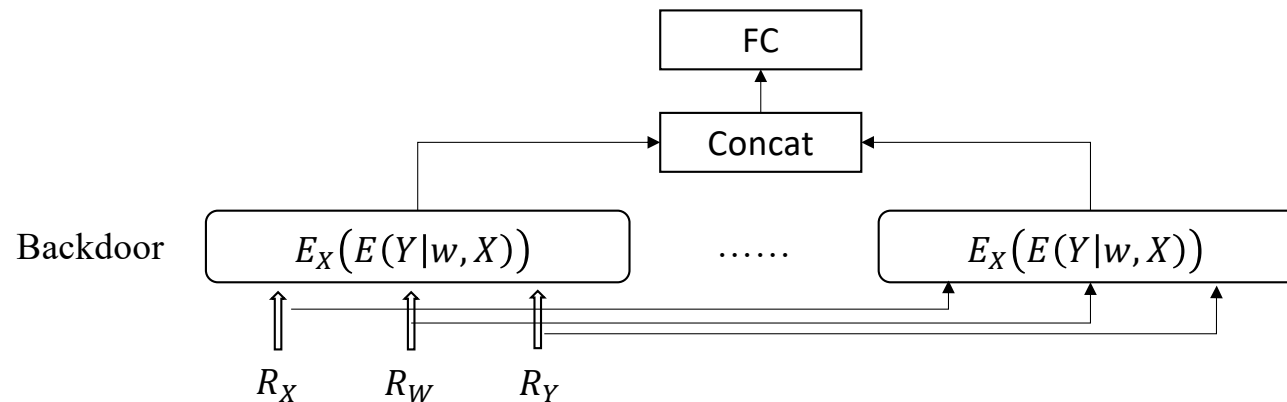
Diagrams often assume non-parametric dependency among variables for all units.
But different units may have different dependencies and parameters.

Our Approach to Address Challenges

- High Dimensionality:
 - Balancing Representation
 - Variables Grouping (the dependencies between variables within each group can be learned separately)
 - Multi Diagram Identification
 - Model Learning
- Novelty: The dependency among variables will be simplified by representation learning and variable grouping while preserving causal effect estimation performance.

Our Approach to Address Challenges

- Robustness:
 - Use multi-head techniques in the diagram identification stage to improve robustness for causal effect learning
- Novelty: Will be the first to introduce multi-head identification modules in causal effect learning.



Futural plan

Table 5: Ph.D. Program Timeline with Publication Goals

Year	Activities
1	Coursework, literature review, research proposal
2	Data collection, preliminary analysis, conference paper 1, conference presentation, QE
3	Advanced analysis, paper writing, conference paper 2, journal paper, Candidature
4	Finalize dissertation, defend dissertation, conference paper 3, graduation

Thanks!

IHDP (semi-synthetic)

- The causal effect of a home visit on IQ test result

I use experimental data from the Infant Health and Development Program (IHDP), a randomized experiment that began in 1985, targeted low-birth-weight, premature infants, and provided the treatment group with both intensive high-quality child care and home visits from a trained provider. The program was highly successful at significantly raising cognitive test scores of the treated children relative to controls at the end of the intervention (Brooks-Gunn, Liaw, and Klebanov 1991). The study collected data on many pretreatment variables. I use measurements on the child—birth weight, head circumference, weeks born preterm, birth order, first born, neonatal health index (see Scott and Bauer 1989), sex, twin status—as well as behaviors engaged in during the pregnancy—smoked cigarettes, drank alcohol, took drugs—and measurements on the mother at the time she gave birth—age, marital status, educational attainment (did not graduate from high school, graduated from high school, attended some college but did not graduate, graduated from college), whether she worked during pregnancy, whether she received prenatal care—and the site (8 total) in which the family resided at the start of the intervention. There are 6 continuous covariates and 19 binary covariates.

Twins

- The causal effect of birth weight on mortality

We introduce a new benchmark task that utilizes data from twin births in the USA between 1989-1991 [3]³. The treatment $t = 1$ is being born the heavier twin whereas, the outcome corresponds to the mortality of each of the twins in their first year of life. Since we have records for both twins, their outcomes could be considered as the two potential outcomes with respect to the treatment of being born heavier. We only chose twins which are the same sex. Since the outcome is thankfully quite rare (3.5% first-year mortality), we further focused on twins such that both were born weighing less than $2kg$. We thus have a dataset of 11984 pairs of twins. The mortality rate for the lighter twin is 18.9%, and for the heavier 16.4%, for an average treatment effect of -2.5% . For each twin-pair we obtained 46 covariates relating to the parents, the pregnancy and birth: mother and father education, marital status, race and residence; number of previous births; pregnancy risk factors such as diabetes, renal disease, smoking and alcohol use; quality of care during pregnancy; whether the birth was at a hospital, clinic or home; and number of gestation weeks prior to birth.

In this setting, for each twin pair we observed both the case $t = 0$ (lighter twin) and $t = 1$ (heavier twin). In order to simulate an observational study, we selectively hide one of the two twins; if we were to choose at random this would be akin to a randomized trial. In order to simulate the case of hidden confounding with proxies, we based the treatment assignment on a single variable which is highly correlated with the outcome: GESTAT10, the number of gestation weeks prior to birth. It is ordinal with values from 0 to 9 indicating birth before 20 weeks gestation, birth after 20-27 weeks of gestation and so on⁴. We then set $t_i | \mathbf{x}_i, \mathbf{z}_i \sim \text{Bern}(\sigma(w_o^\top \mathbf{x} + w_h(\mathbf{z}/10 - 0.1)))$, $w_o \sim \mathcal{N}(0, 0.1 \cdot I)$, $w_h \sim \mathcal{N}(5, 0.1)$, where \mathbf{z} is GESTAT10 and \mathbf{x} are the 45 other features.

TARNet

- Model

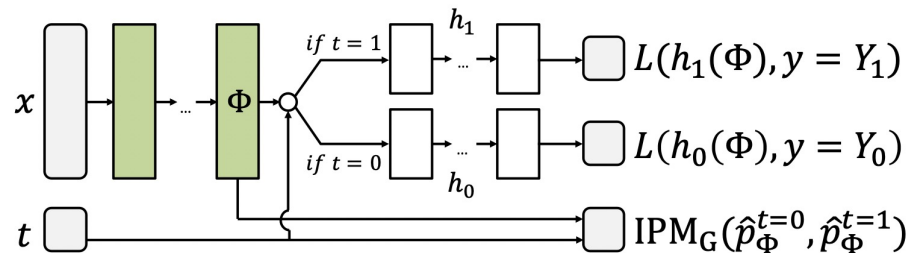


Figure 1. Neural network architecture for ITE estimation. L is a loss function, IPM_G is an integral probability metric. Note that only one of h_0 and h_1 is updated for each sample during training.

- Loss

$$\begin{aligned}
 & \min_{h, \Phi} \quad \frac{1}{n} \sum_{i=1}^n w_i \cdot L(h(\Phi(x_i), t_i), y_i) + \lambda \cdot \mathfrak{R}(h) \\
 & \quad + \alpha \cdot \text{IPM}_G(\{\Phi(x_i)\}_{i:t_i=0}, \{\Phi(x_i)\}_{i:t_i=1}), \\
 & \text{with} \quad w_i = \frac{t_i}{2u} + \frac{1-t_i}{2(1-u)}, \quad \text{where } u = \frac{1}{n} \sum_{i=1}^n t_i, \\
 & \text{and} \quad \mathfrak{R} \text{ is a model complexity term.} \quad (3)
 \end{aligned}$$

Note that $u = p(t = 1)$ is simply the proportion of treated units in the population. The weights w_i compensate for the difference in treatment group size in our sample, see Theorem 1. $\text{IPM}_G(\cdot, \cdot)$ is the (empirical) integral probability metric w.r.t. G . For most IPMs, we cannot compute the factor B_ϕ in (2), but treat it as part of the hyperparameter α . This makes our objective sensitive to the scaling of Φ , even for a constant α . We therefore normalize Φ through either projection or batch-normalization with fixed scale.

CEVAE

- Model

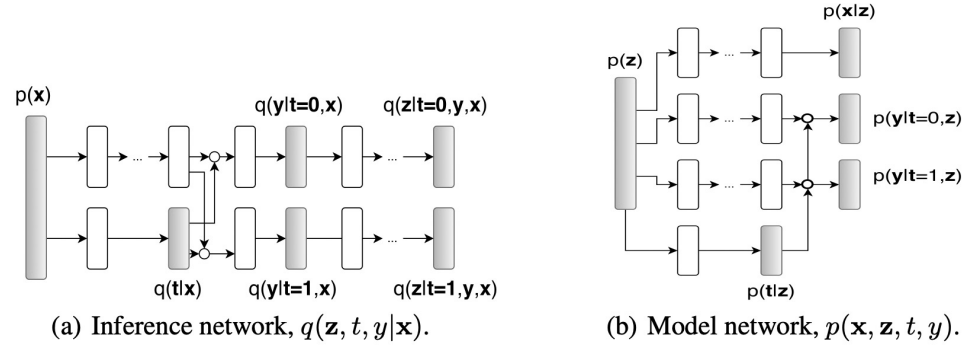


Figure 2: Overall architecture of the model and inference networks for the Causal Effect Variational Autoencoder (CEVAE). White nodes correspond to parametrized deterministic neural network transitions, gray nodes correspond to drawing samples from the respective distribution and white circles correspond to switching paths according to the treatment t .

- Loss

$$\mathcal{L} = \sum_{i=1}^N \mathbb{E}_{q(\mathbf{z}_i|\mathbf{x}_i, t_i, y_i)} [\log p(\mathbf{x}_i, t_i|\mathbf{z}_i) + \log p(y_i|t_i, \mathbf{z}_i) + \log p(\mathbf{z}_i) - \log q(\mathbf{z}_i|\mathbf{x}_i, t_i, y_i)]. \quad (6)$$

$$\mathcal{F}_{\text{CEVAE}} = \mathcal{L} + \sum_{i=1}^N (\log q(t_i = t_i^*|\mathbf{x}_i^*) + \log q(y_i = y_i^*|\mathbf{x}_i^*, t_i^*)), \quad (10)$$