

# When interpretability met causality: Causal Interpretability<sup>1</sup>

Zirui Yan

Department of Electrical, Computer, and Systems Engineering

Rensselaer Polytechnic Institute

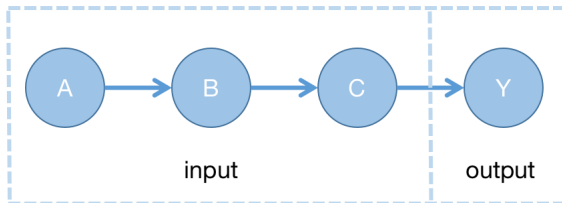
email: `yanz11@rpi.edu`

webpage: `ziruiyan.github.io`

April 25, 2022

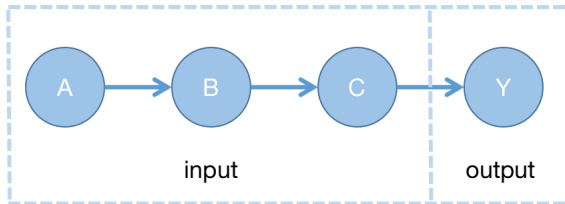
---

<sup>1</sup><https://github.com/ZiruiYan/Causal-Interpretability>



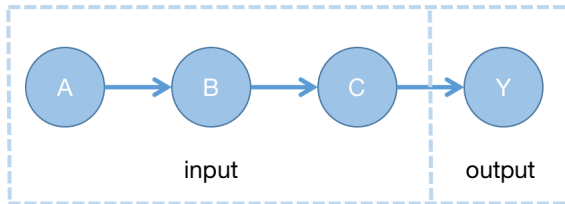
- ▶ **Problem:** Assuming the inputs and output admits some causality, how output  $Y$  will change if we do intervention on variable  $B$ .
- ▶ **Related work:**
  - ▶ There are some interpretability methods such as Quantitative input influence (QII), Accumulated Local Effects (ALE) plot. They may consider the correlation but ignore the causality.
  - ▶ Researchers in causality community has investigate into this problem, but they focus on the binary variable and didn't link it to interpretability.

## Simple Example



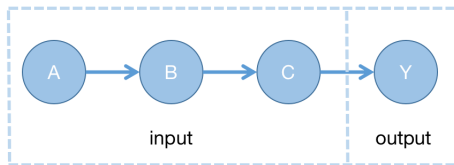
- ▶ A: SAT Score → B: College Addimision → C: Take TML course? → D: get job at IBM?
- ▶ A: 1500 → B: RPI → C: Yes → D: Yes
- ▶ A: 1500 → B: Stanford → C: Yes → D: Yes

## Simple Example



- ▶ A: SAT Score → B: College Addimision → C: Take TML course? → D: get job at IBM?
- ▶ A: 1500 → B: RPI → C: Yes → D: Yes
- ▶ A: 1500 → B: Stanford → C: Yes → D: Yes
- ▶ A: 1500 → B: Stanford → C: No → D: No

- ▶ Assume we know the **Directed acyclic graph (DAG)** between the variables  
Edges represent causal relations between the variables  
Output  $Y$  is terminal nodes



- ▶ Assume that there is **no latent variables**
- ▶  $do(X = x)$ : Change  $X$  to become  $x$ , and affect the descendants of  $X$

- ▶ Assume we interested in the influence of  $\mathbf{X}_s$  on function  $f(\mathbf{X}_s, \mathbf{X}_c)$  and  $\mathbf{X}_c$  are other inputs.
- ▶ Individual Conditional Expectation (ICE) for individual  $i$  with input  $(\mathbf{x}_s^{(i)}, \mathbf{x}_c^{(i)})$ <sup>3</sup>

$$f_{s,\text{ICE}}^{(i)}(\mathbf{x}_s) = f(\mathbf{x}_s, \mathbf{x}_c^{(i)}) . \quad (1)$$

- ▶ Partial Dependence Plot (PDP)<sup>4</sup>

$$f_{s,\text{PDP}}(\mathbf{x}_s) = \mathbb{E}_{\mathbf{X}_c} [f(\mathbf{x}_s, \mathbf{X}_c)] = \int f(\mathbf{x}_s, \mathbf{X}_c) d\mathbb{P}(\mathbf{X}_c) . \quad (2)$$

---

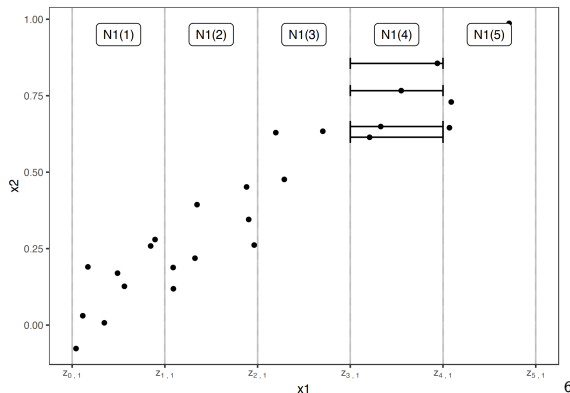
<sup>3</sup>Goldstein, Alex, et al. Peeking inside the black box: Visualizing statistical learning with plots of individual conditional expectation. *Journal of Computational and Graphical Statistics* (2015)

<sup>4</sup>Friedman, Jerome H. Greedy function approximation: a gradient boosting machine. *Annals of statistics* (2001)



## ► Accumulated Local Effects (ALE)<sup>5</sup>

$$f_{s,\text{ALE}}(\mathbf{x}_s) = \int_{x_0}^{x_s} \mathbb{E}_{\mathbf{x}_c | \mathbf{x}_s = z_s} [f(\mathbf{x}_s, \mathbf{x}_c) | \mathbf{x}_s = z_s] dz_s - c, \quad (3)$$



<sup>5</sup>Apley, Daniel W., and Jingyu Zhu. Visualizing the effects of predictor variables in black box supervised learning models. Journal of the Royal Statistical Society (2020)

<sup>6</sup>Sec 8.2 in Molnar, Christoph. Interpretable machine learning. Lulu.com (2020)

- ▶ Individual Intervention Expectation (IIE) for individual  $i$  with input  $(\mathbf{x}_s^{(i)}, \mathbf{x}_c^{(i)})$

$$f_{s,\text{IIE}}^{(i)}(\mathbf{x}_s) = f(\mathbf{X}_s, \mathbf{X}_c | (\mathbf{X}_s, \mathbf{X}_c) = (\mathbf{x}_s^{(i)}, \mathbf{x}_c^{(i)}), do(\mathbf{X}_s = \mathbf{x}_s)) , \quad (4)$$

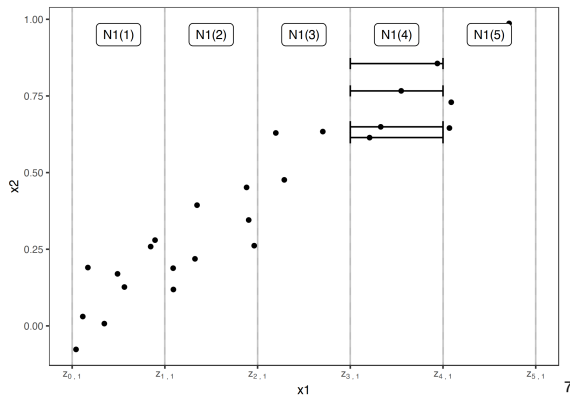
- ▶ Causal Partial Dependence Plot (CPDP)

$$f_{s,\text{CPDP}}^{(i)}(\mathbf{x}_s) = \mathbb{E}_{(\mathbf{x}_s, \mathbf{x}_c)} f(\mathbf{X}_s, \mathbf{X}_c | do(\mathbf{X}_s = \mathbf{x}_s)) , \quad (5)$$



## ► Accumulated Local Casual Effects (ALCE)

$$f_{S,ALCE}(x_s) = \int_{x_0}^{x_s} \mathbb{E}_{X_c|X_s=z_s} [f(X_s, X_c) | do(X_s = z_s)] dz_s - c. \quad (6)$$

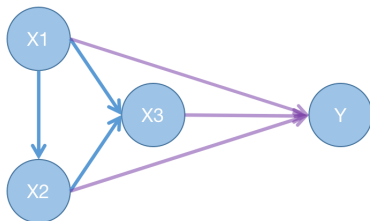


## Experiment setting: synthetic

- ▶ Assume we only know the causal DAG and the class of function, where  $\epsilon_i \sim \mathcal{N}(0, 0.01)$

$$\begin{aligned}X_1 &\leftarrow U(-10, 10), \\X_2 &\leftarrow 10\sigma(X_1) - 5 + \epsilon_2, \\X_3 &\leftarrow 10\sigma(-X_1 - X_2) - 5 + \epsilon_3, \\Y &\leftarrow 10\sigma(X_1 + X_2 - X_3) - 5 + \epsilon_4,\end{aligned}\tag{7}$$

where  $\sigma$  is the sigmoid function.



- ▶ We are interested in variable  $X_2$
- ▶ Number of samples is 50000

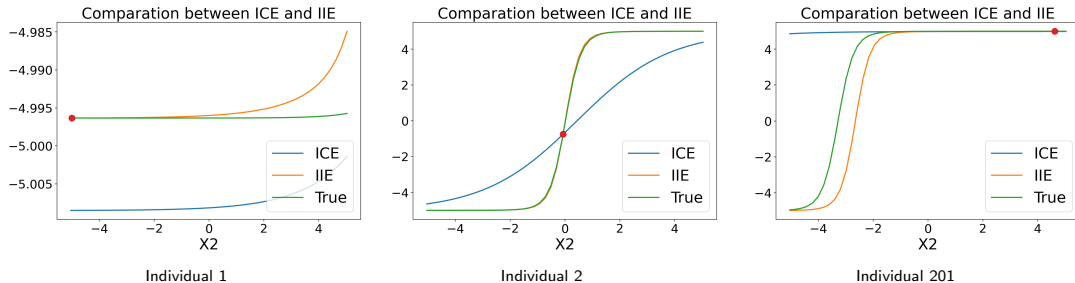
- For individual  $i$

$$\begin{aligned}x_2^{(i)} &\leftarrow 10\sigma(x_1^{(i)}) - 5 + \epsilon_2^{(i)}, \\x_3^{(i)} &\leftarrow 10\sigma(-x_1^{(i)} - x_2^{(i)}) - 5 + \epsilon_3^{(i)}, \\y^{(i)} &\leftarrow 10\sigma(x_1^{(i)} + x_2^{(i)} - x_3^{(i)}) - 5 + \epsilon_4^{(i)},\end{aligned}\tag{8}$$

- For machine learning

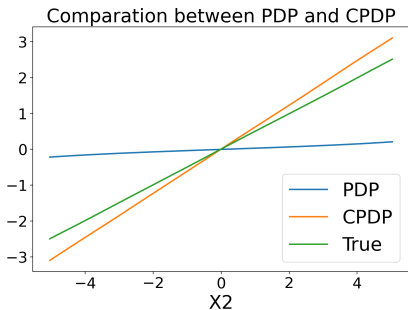
$$\begin{aligned}\hat{\epsilon}_2^{(i)} &= x_2^{(i)} - \hat{x}_2^{(i)}, \\ \hat{\epsilon}_3^{(i)} &= x_3^{(i)} - \hat{x}_3^{(i)}, \\ \hat{\epsilon}_4^{(i)} &= y^{(i)} - \hat{y}^{(i)},\end{aligned}\tag{9}$$

## Experiment results: synthetic

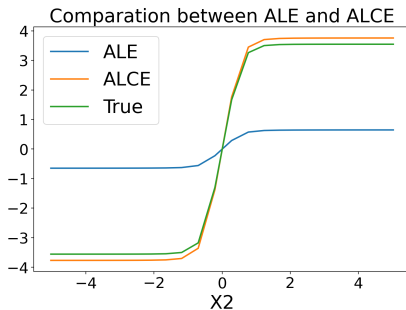


**Figure.** Comparison between ICE and IIE

## Experiment results: synthetic



Comparison between PDP and CPDP

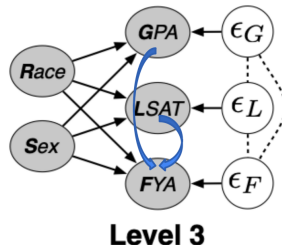


Comparison between ALE and ALCE

**Figure.** PDP and ALCE plots

## Experiment: real dataset

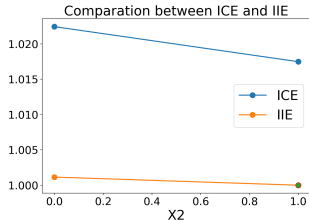
### ► Law School Success<sup>8</sup>



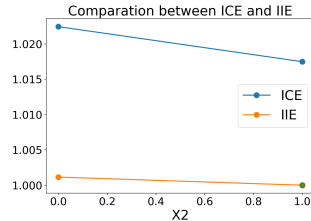
$$\begin{cases} \text{GPA} = b_G + w_G^R R + w_G^S S + \epsilon_G, & \epsilon_G \sim p(\epsilon_G) \\ \text{LSAT} = b_L + w_L^R R + w_L^S S + \epsilon_L, & \epsilon_L \sim p(\epsilon_L) \\ \text{FYA} = b_F + w_F^R R + w_F^S S + w_F^G \text{GPA} + w_F^L \text{LSAT} + \epsilon_F, & \epsilon_F \sim p(\epsilon_F) \end{cases} \quad (10)$$

<sup>8</sup>similar setting as in Kusner, Matt J., et al. Counterfactual fairness. NeurIPS (2017).

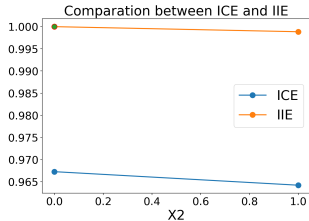
# Experiment: real dataset



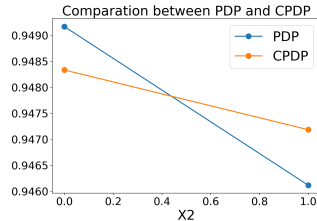
ICE and IIE for female 1



ICE and IIE for female 2



ICE and IIE for male 1



PDP and CPDP plots

## Conclusion

- ▶ Traditional interpretability methods will meet problem when we are interested in causality
- ▶ The proposed causal interpretable methods works better in this scenario

## Problem and Future work

- ▶ DAG is hard to identify based on observational dataset
  1. how to use the intervention dataset
  2. causal interpretable methods can help to identify the DAG
- ▶ Deep neural network have bad local minimum
  1. causal interpretable methods can only interpret the model itself
  2. causal interpretable methods can help to evaluation the model
- ▶ analysis of the performance of causal interpretable methods (e.g. in high dimension)

