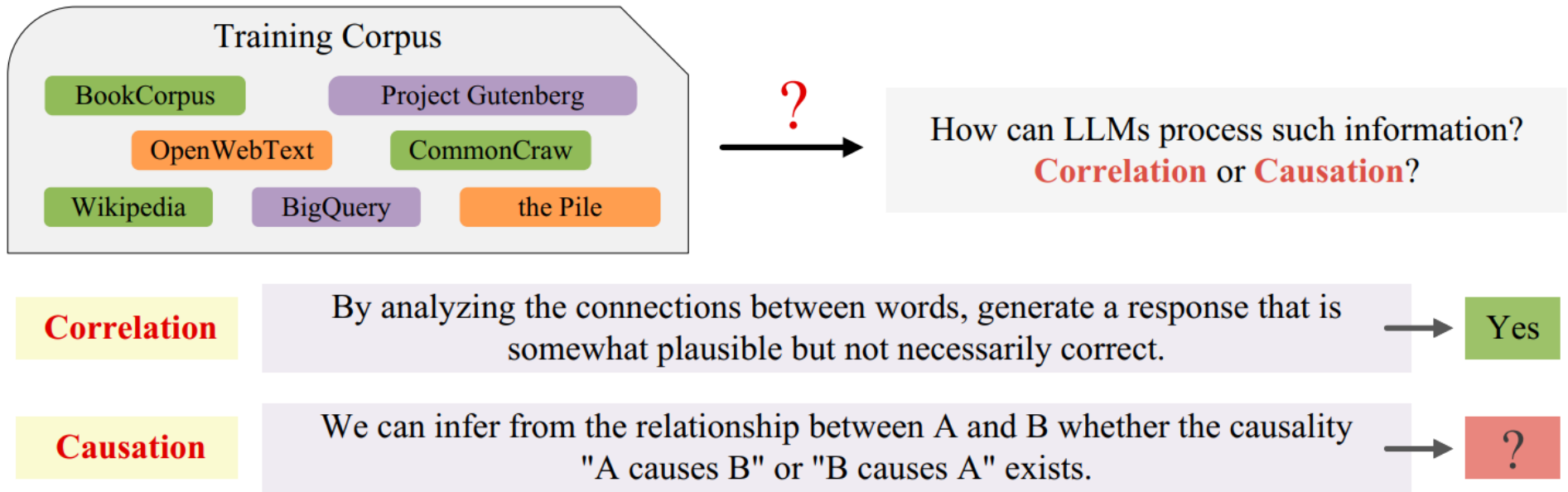


Outlier Robust Causal Discovery

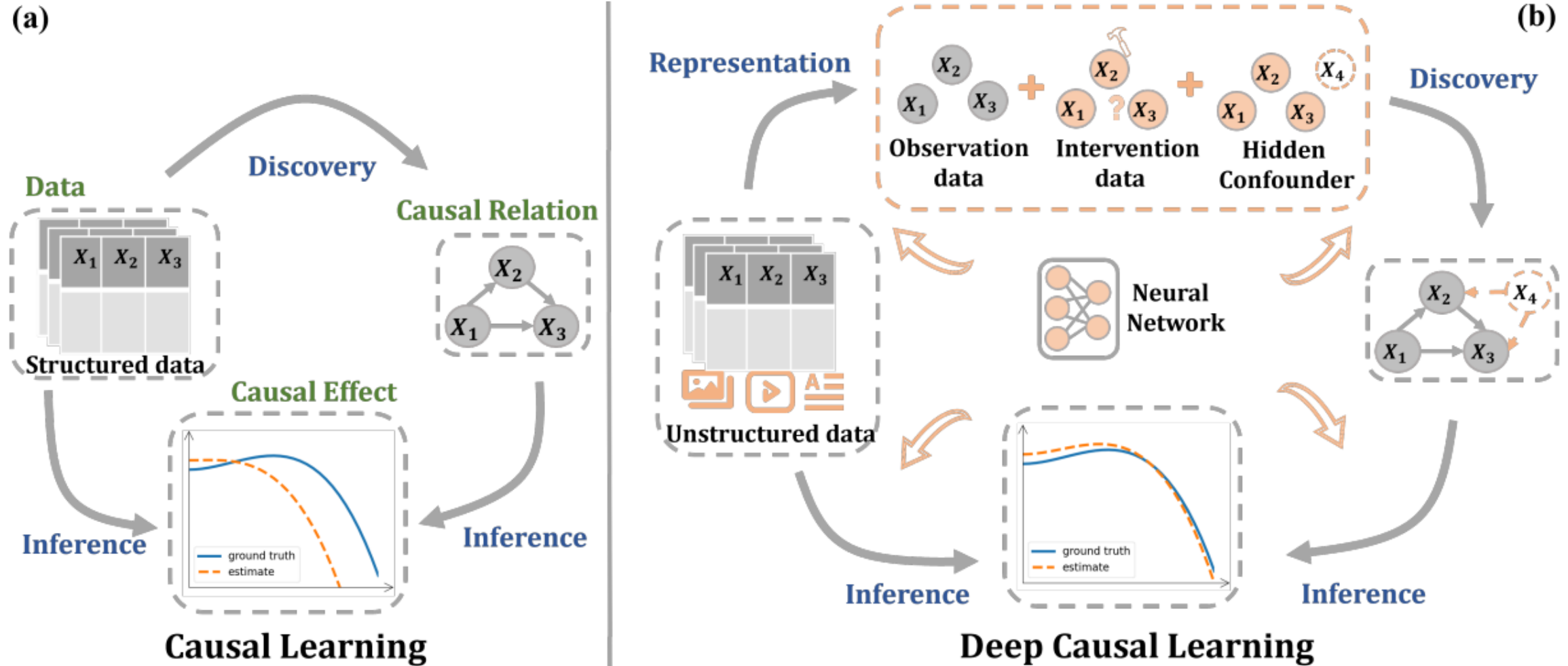
서울대학교 DSML Lab 송의종

Recent Works of Causal Learning in Age of AI

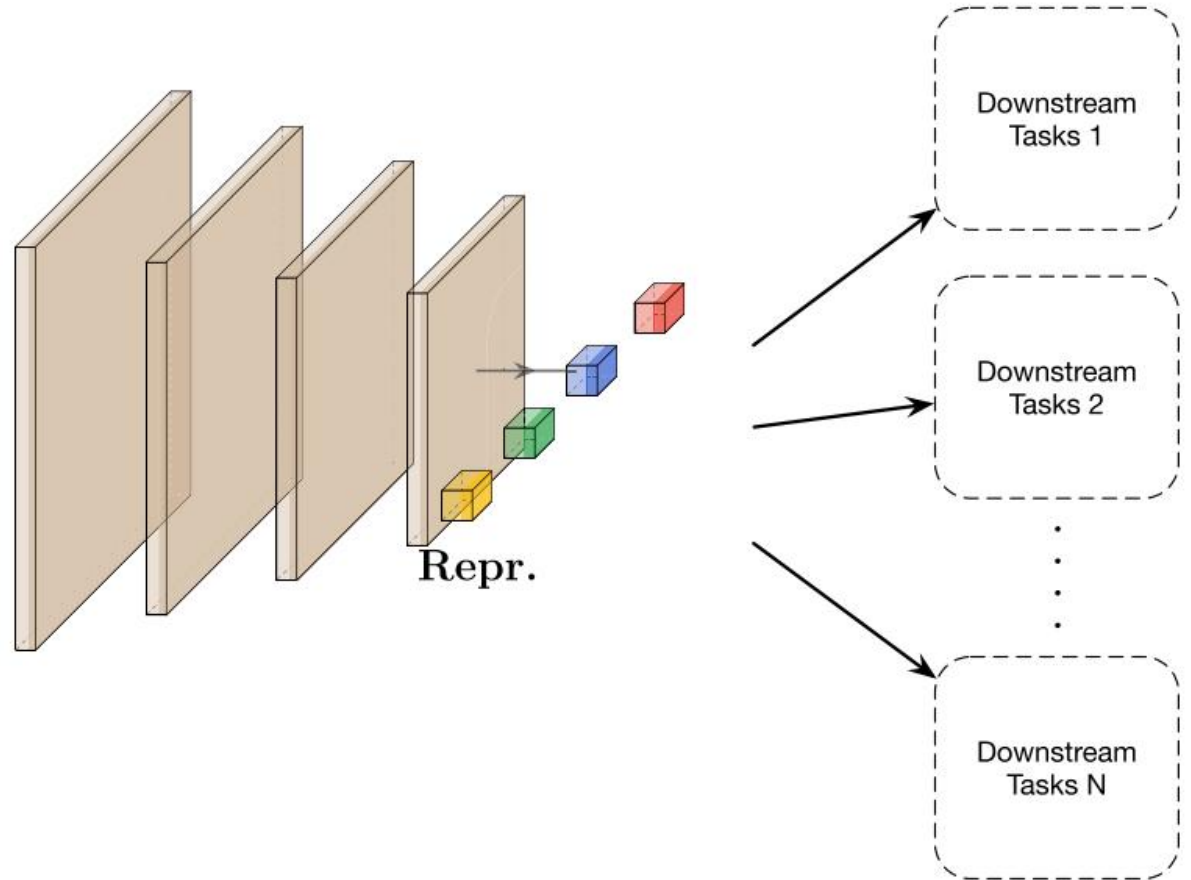
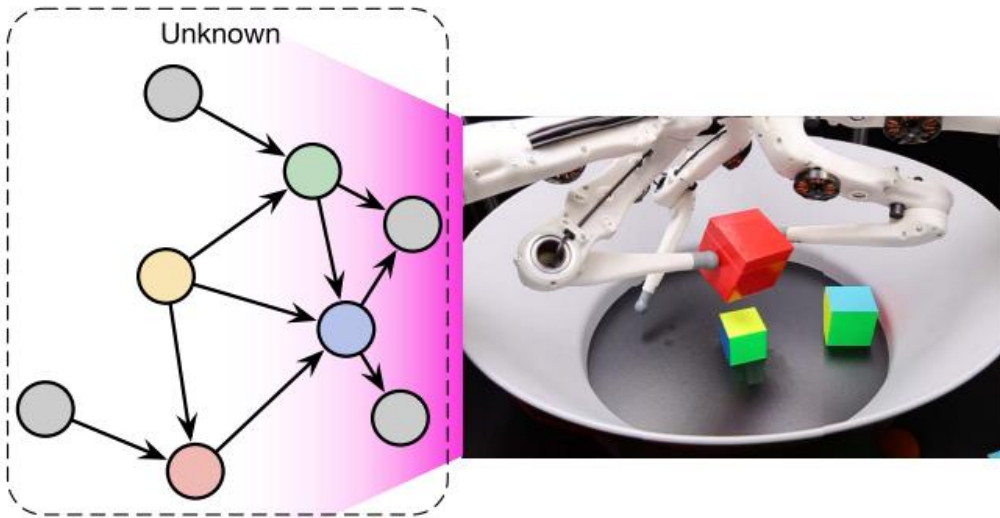
LLM and Causation



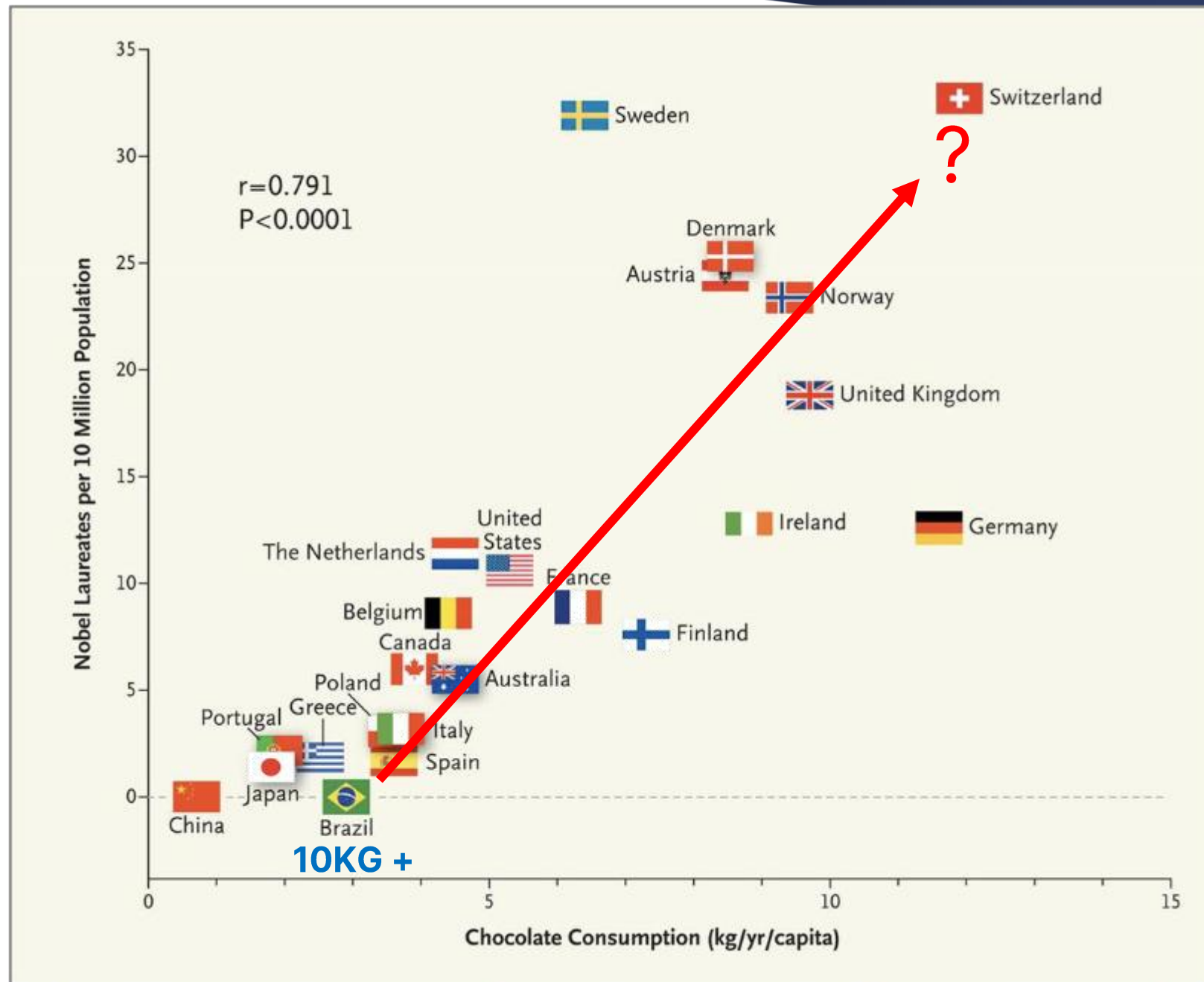
Deep Causal Learning

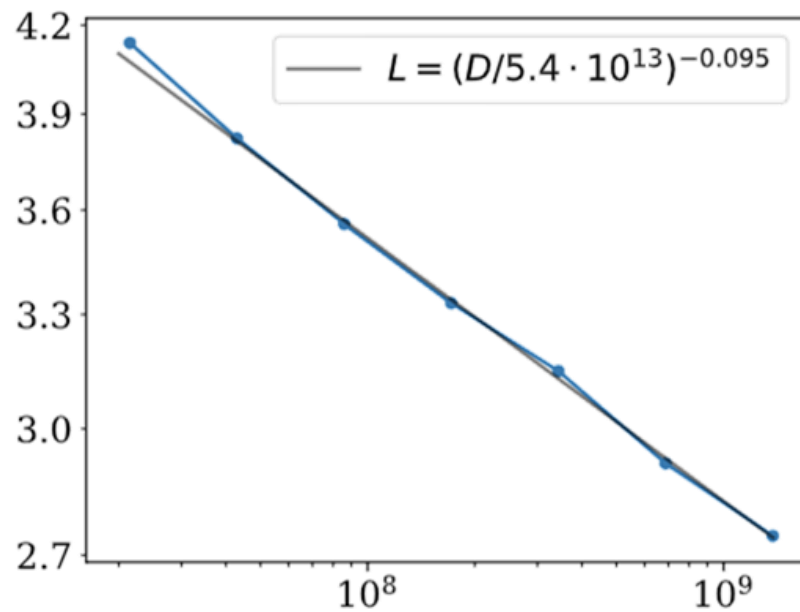


Causal Representation Learning

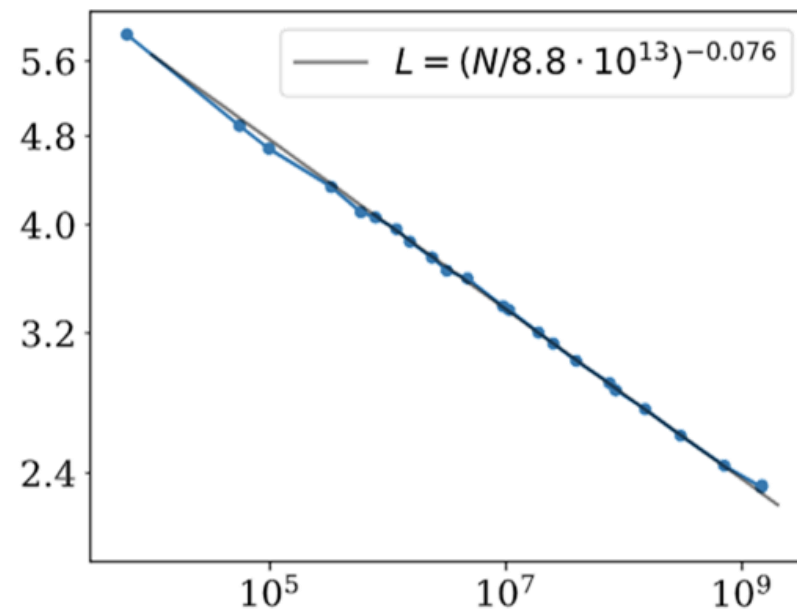


Motivation of Causal Learning

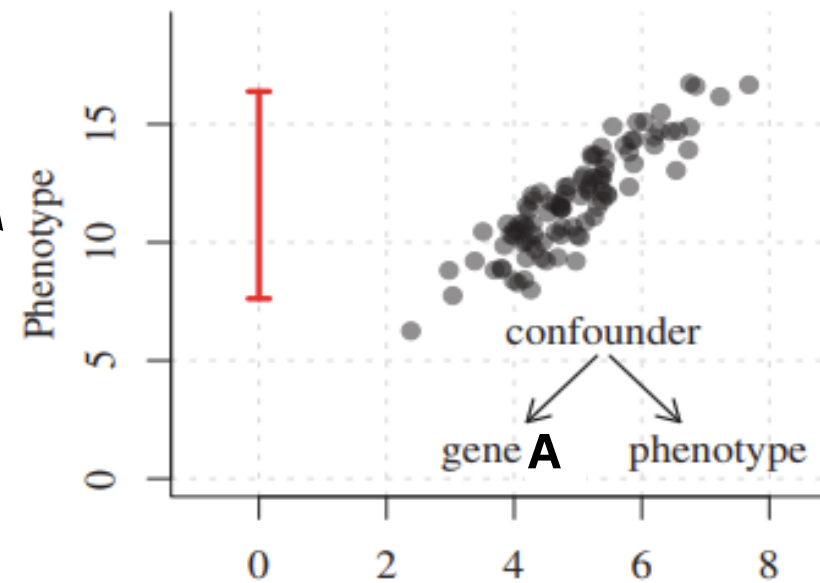
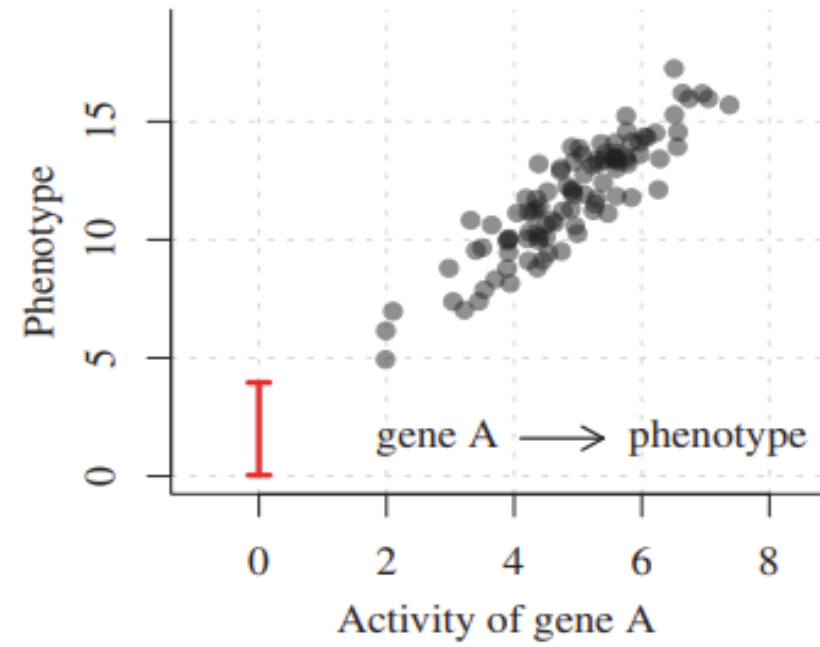
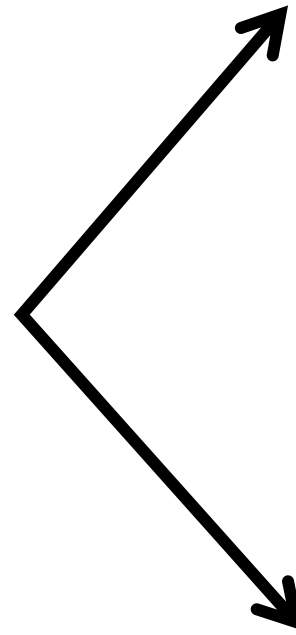
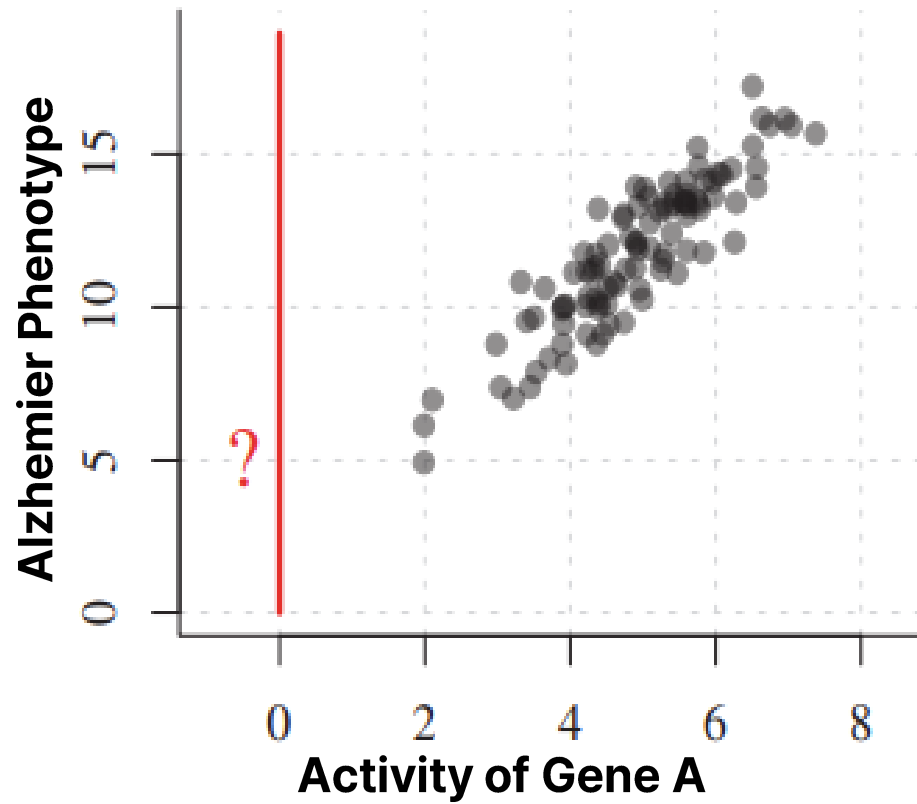




Dataset Size
tokens



Parameters
non-embedding



Contents

1. Introduction to Causal Inference

Causal Reasoning & Learning

2. Outlier Robust Causal Discovery

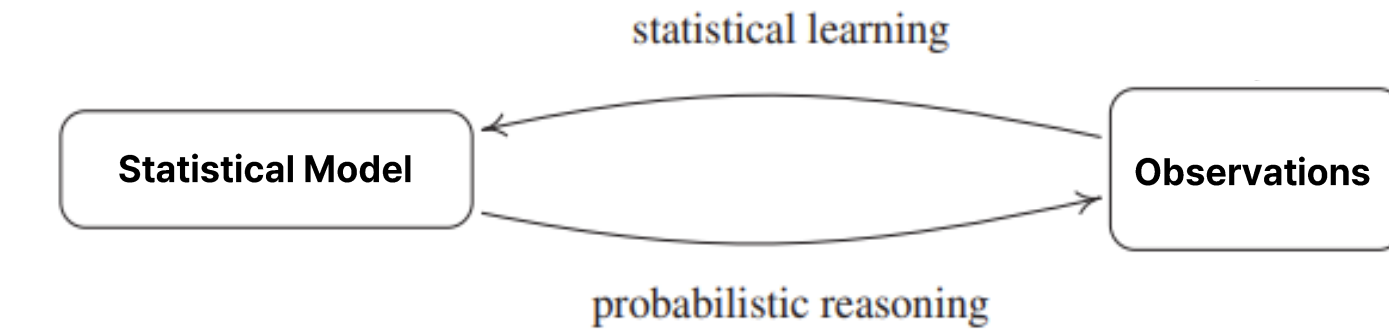
Outlier Model(CCSEM)

Algorithm for Causal Discovery

1. Introduction to Causal Inference

1.1. Causal Reasoning & Learning

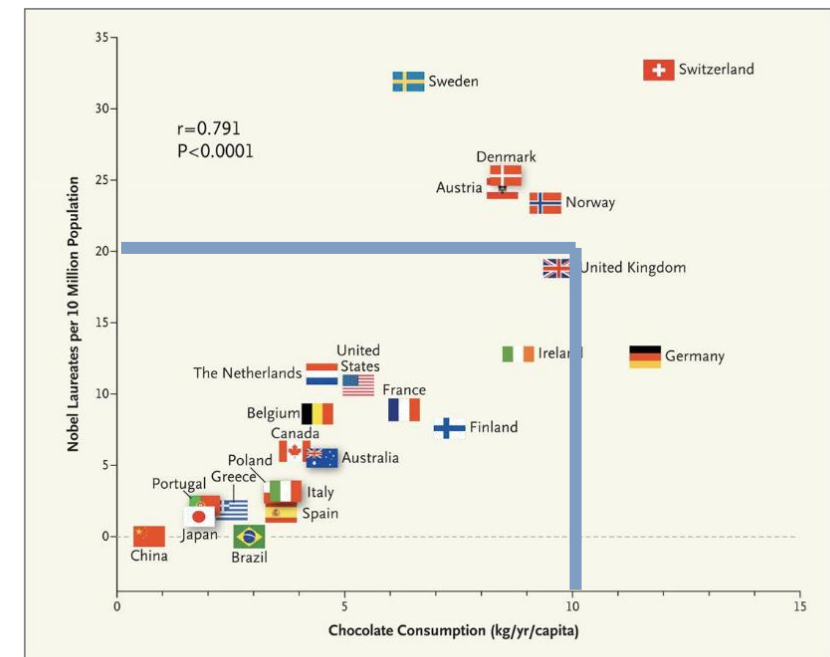
[Statistical Learning & Reasoning]



$$\text{노벨상} = 2\text{초콜릿} + \epsilon$$

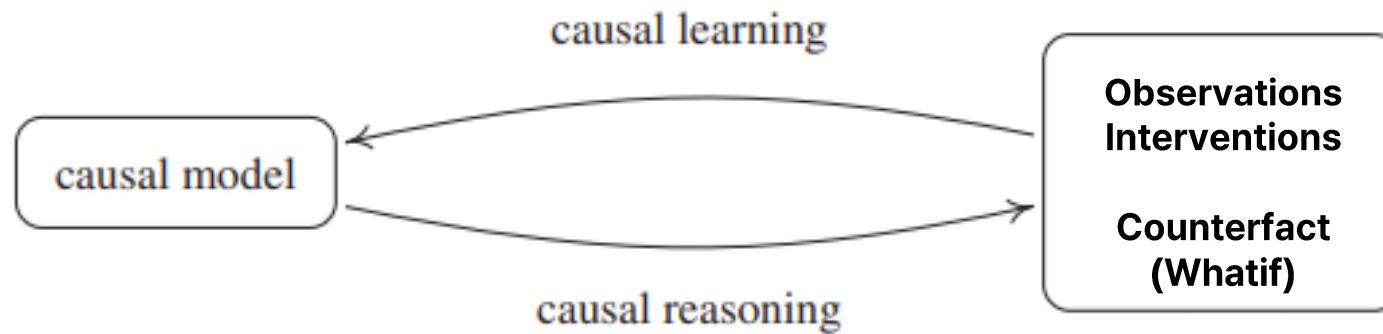
초콜릿 섭취량=10

$$\text{노벨상} = 20$$

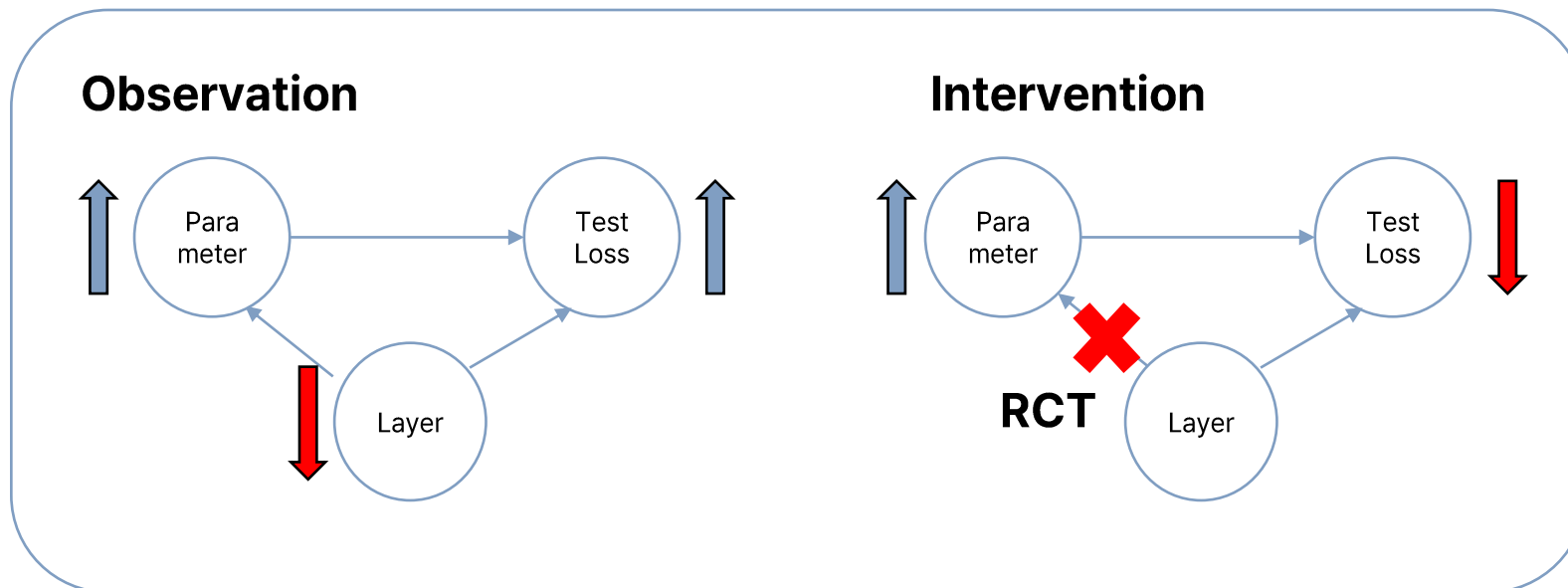


1.1. Causal Reasoning & Learning

[Causal Learning & Reasoning]

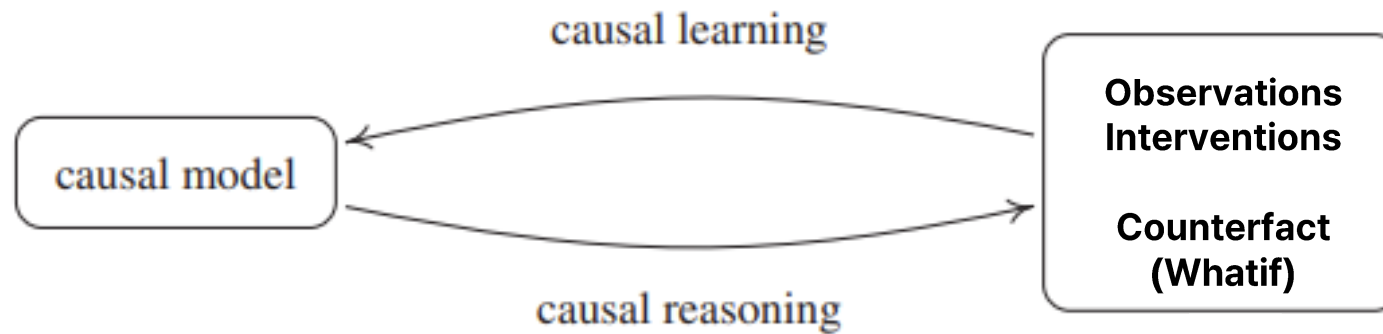


Intervention(개입) Data

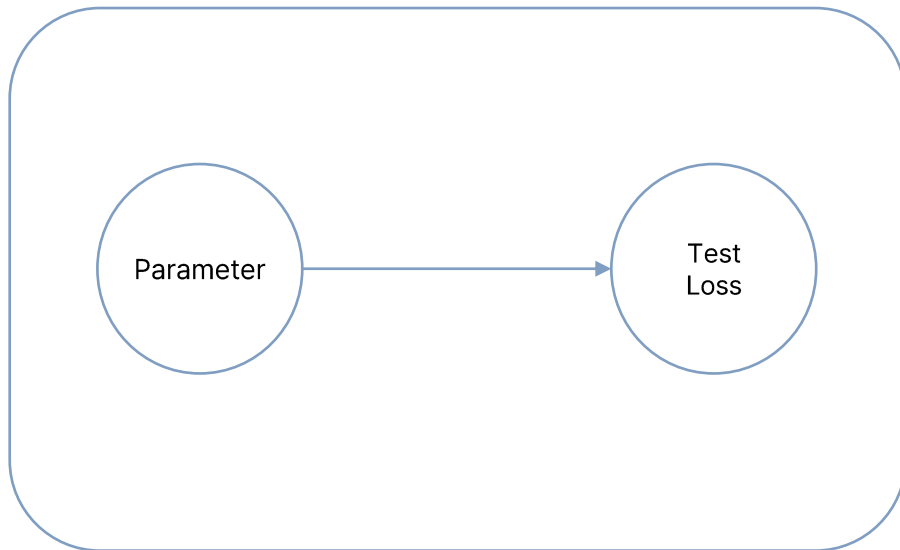


1.1. Causal Reasoning & Learning

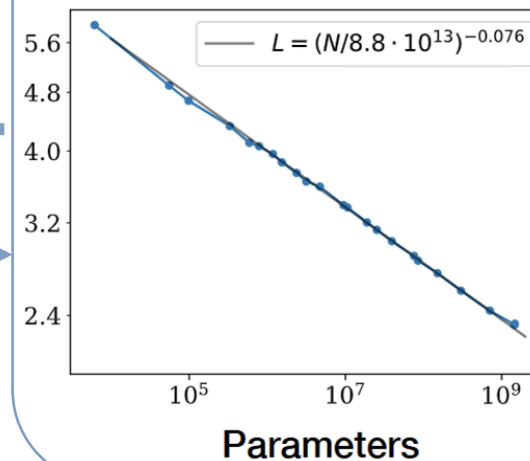
[Causal Learning & Reasoning]



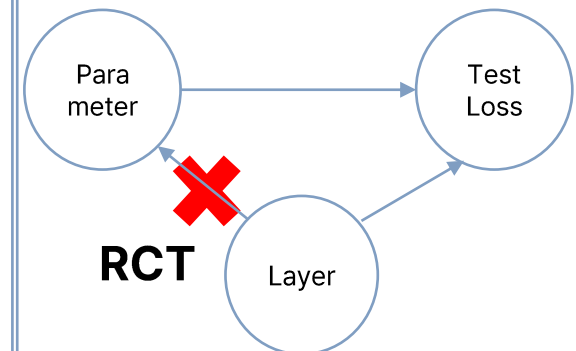
Intervention(개입) Data



Step2: Data form 1

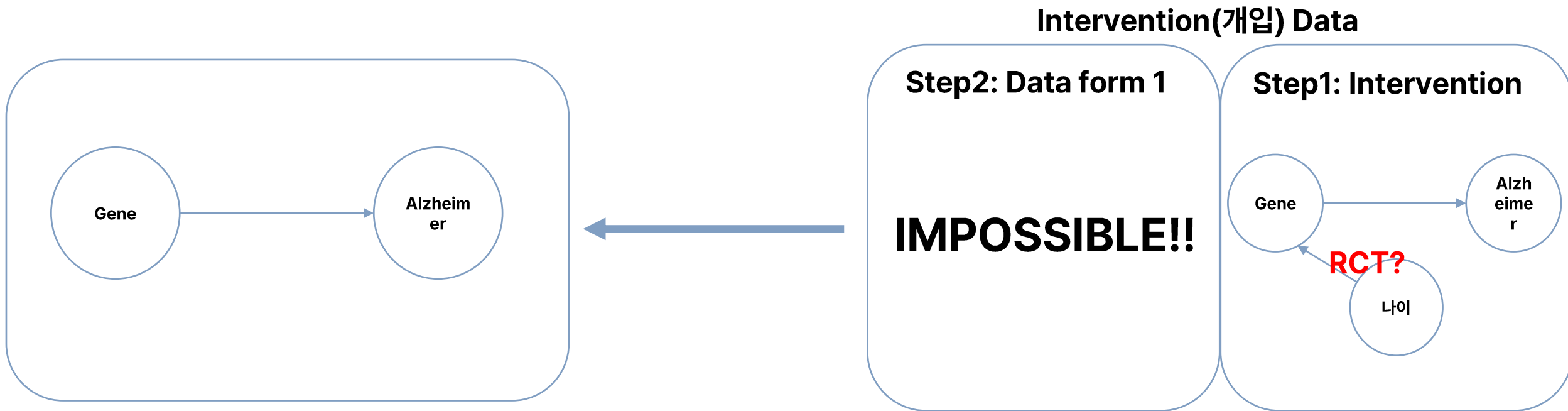


Step1: Intervention



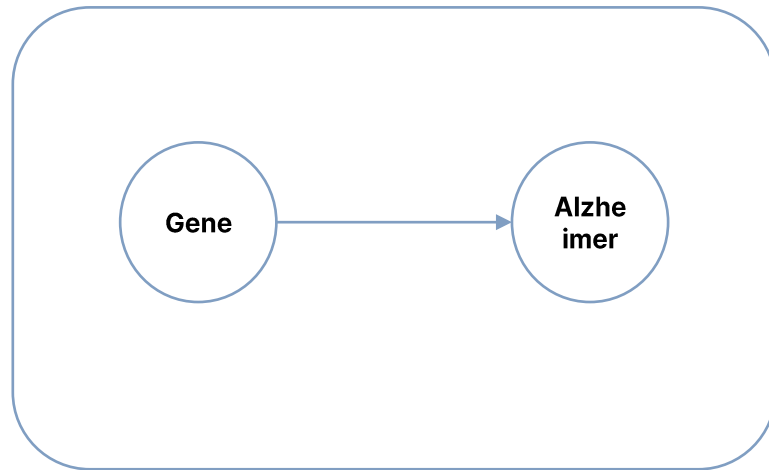
1.1. Causal Reasoning & Learning

[Causal Learning & Reasoning]

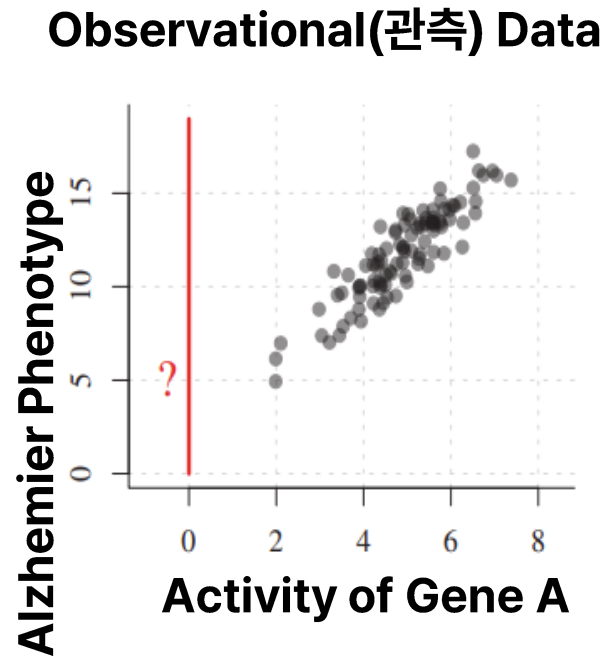


1.1. Causal Reasoning & Learning

[Causal Learning & Reasoning]

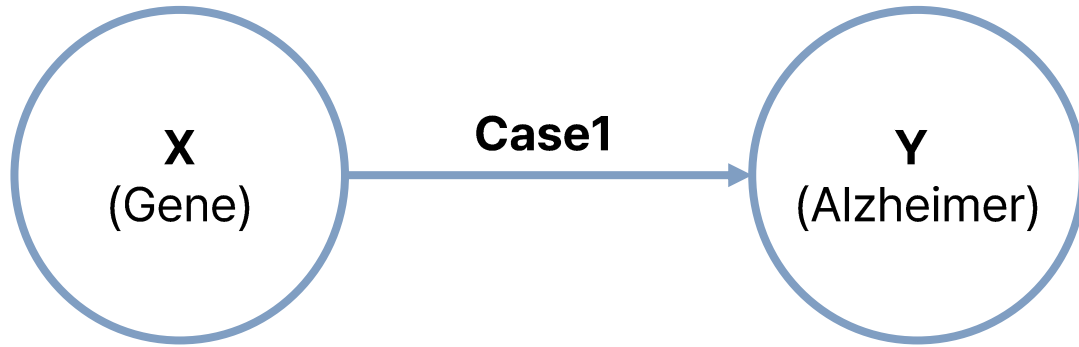


Impossible?

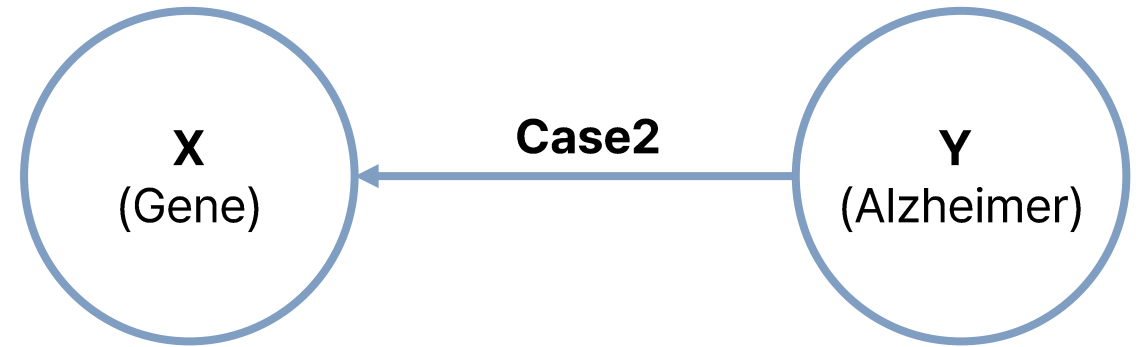


1.1. Causal Reasoning & Learning

[Causal Learning from Observation]



- $X = \varepsilon_X, Y = aX + \varepsilon_Y, \varepsilon_X \perp \varepsilon_Y,$
 - $\text{Var}(\varepsilon_X) = \text{Var}(\varepsilon_Y) = \sigma^2$
- ↓
- $\text{Var}(X) = \sigma^2 \leq \text{Var}(Y) = \sigma^2(1 + a^2)$



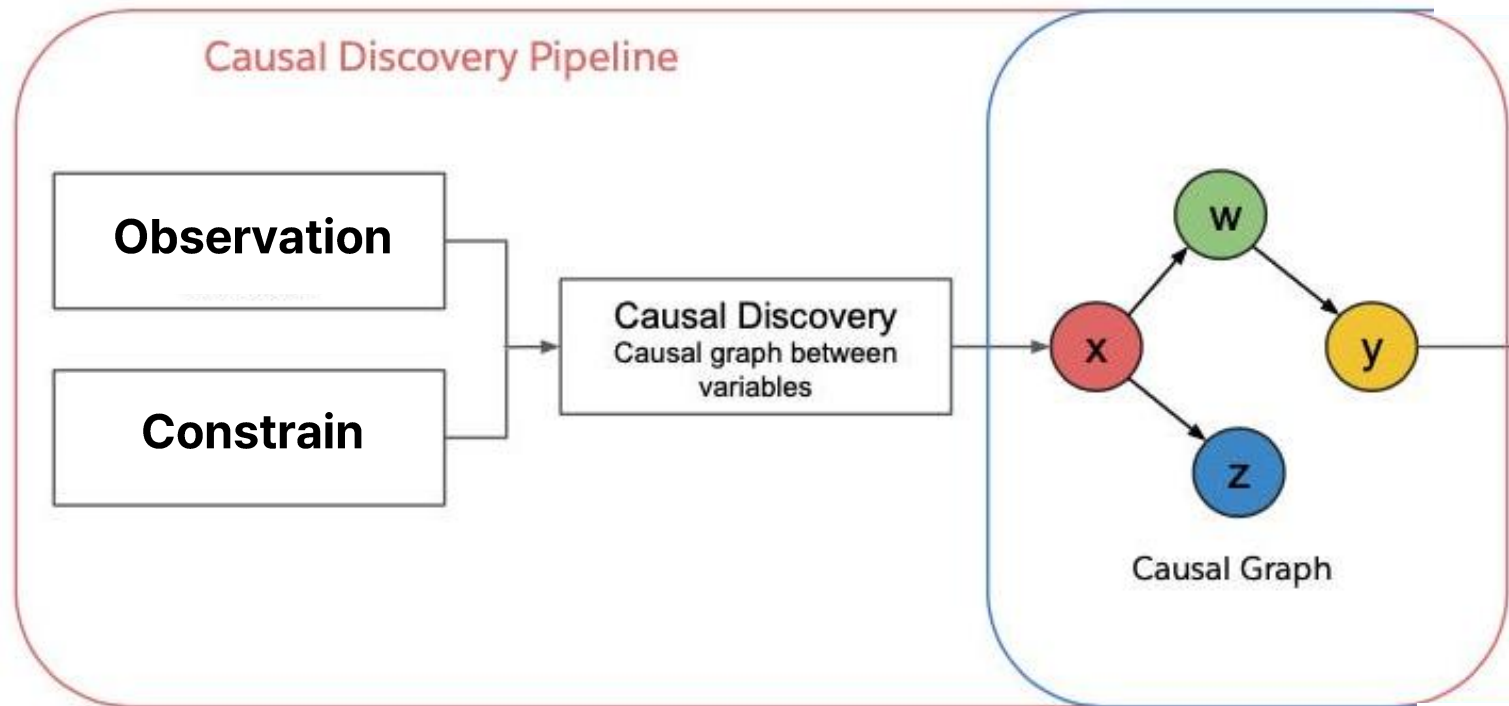
- $X = bY + \varepsilon_X, Y = \varepsilon_Y, \varepsilon_X \perp \varepsilon_Y,$
 - $\text{Var}(\varepsilon_X) = \text{Var}(\varepsilon_Y) = \sigma^2$
- ↓
- $\text{Var}(Y) = \sigma^2 \leq \text{Var}(X) = \sigma^2(1 + b^2)$

Now We can Identify Cause and Effect!!

- $(\text{Var}(X) < \text{Var}(Y) \Leftrightarrow X \rightarrow Y) \wedge (\text{Var}(Y) < \text{Var}(X) \Leftrightarrow Y \rightarrow X)$

1.1. Causal Reasoning & Learning

[Causal Discovery and Inference]

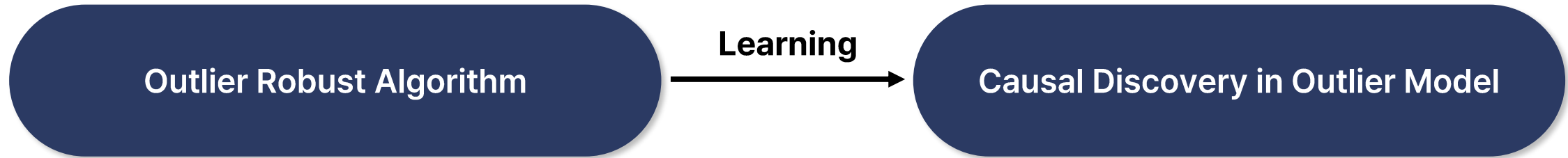


2. Outlier Robust Causal Discovery



2.1 . The Principles of Causal Reasoning

[Table of Contents]



2.1. Outlier Model(CCSEM)

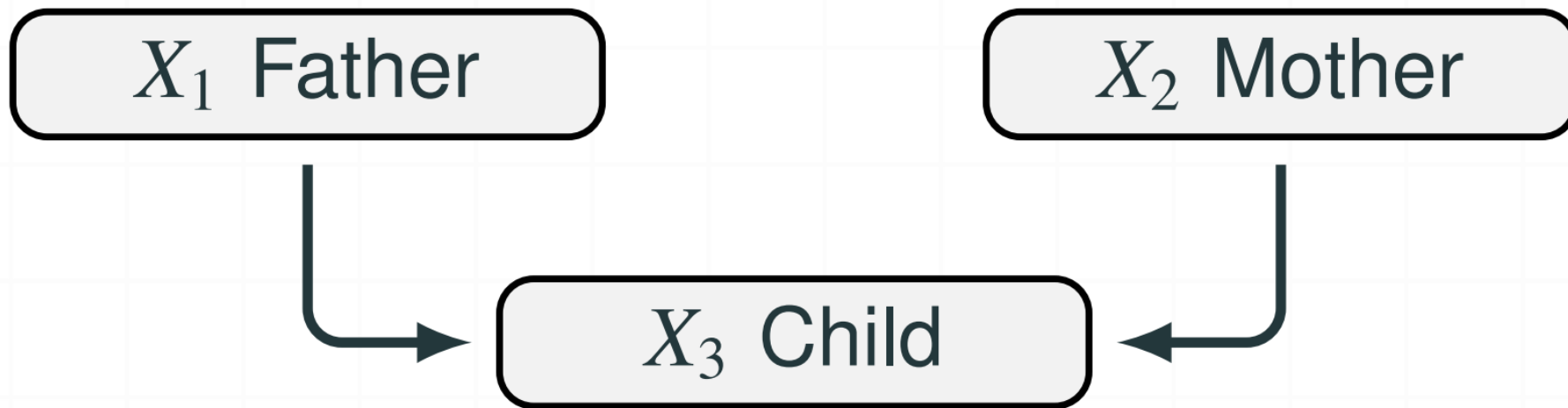
2.2. Outlier Robust Algorithm

2.3. Consistency of Outlier Robust Algorithms

2.4. Simulations and Real-Data

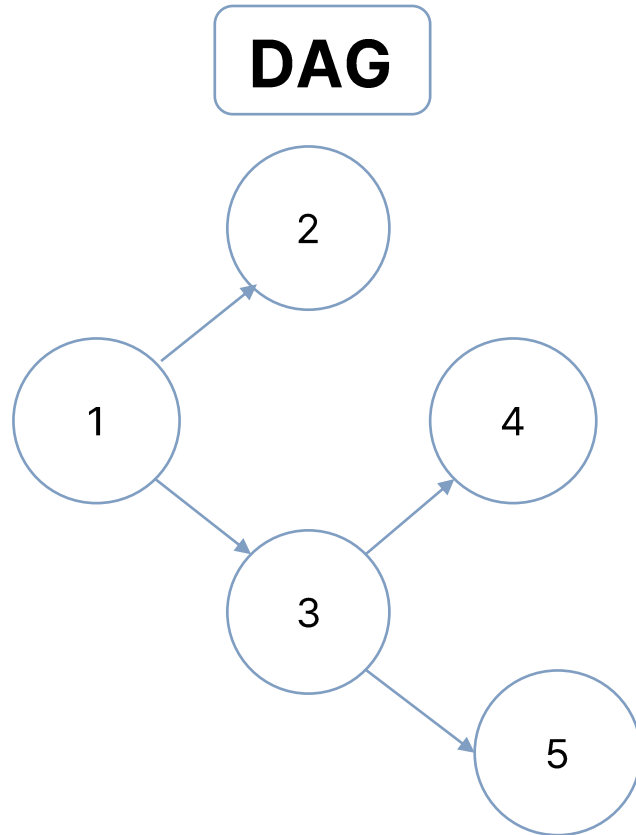
2.0 . Preliminaries

[Definitions]



2.0 . Preliminaries

[Definitions]



Graph Structure

- Parent (Pa): $1 \rightarrow \{2, 3\}$, $Pa(3) = \{1\}$
- Child (Ch): $Ch(3) = \{4, 5\}$
- Ancestor (An): $An(4) = \{1, 3\}$
- Descendant (De): $De(1) = \{2, 3, 4, 5\}$

Orderings

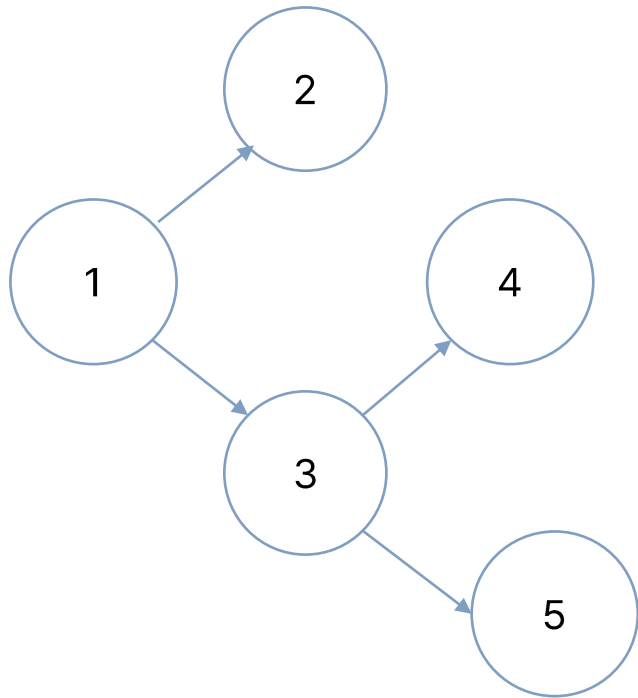
- $\pi = (1, 2, 3, 4, 5)$ or $(1, 3, 2, 5, 4)$

Maximum Indegree

- $d_{\text{in}} = 1$

2.0 . Preliminaries

[Definitions]



Graph Structure

- Parent (Pa): $1 \rightarrow \{2, 3\}$, $Pa(3) = \{1\}$
- Child (Ch): $Ch(3) = \{4, 5\}$
- Ancestor (An): $An(4) = \{1, 3\}$
- Descendant (De): $De(1) = \{2, 3, 4, 5\}$

Orderings

- $\pi = (1, 2, 3, 4, 5)$ or $(1, 3, 2, 5, 4)$

Maximum Indegree

- $d_{\text{in}} = 1$

2.0 . Preliminaries

[Causal Discovery]

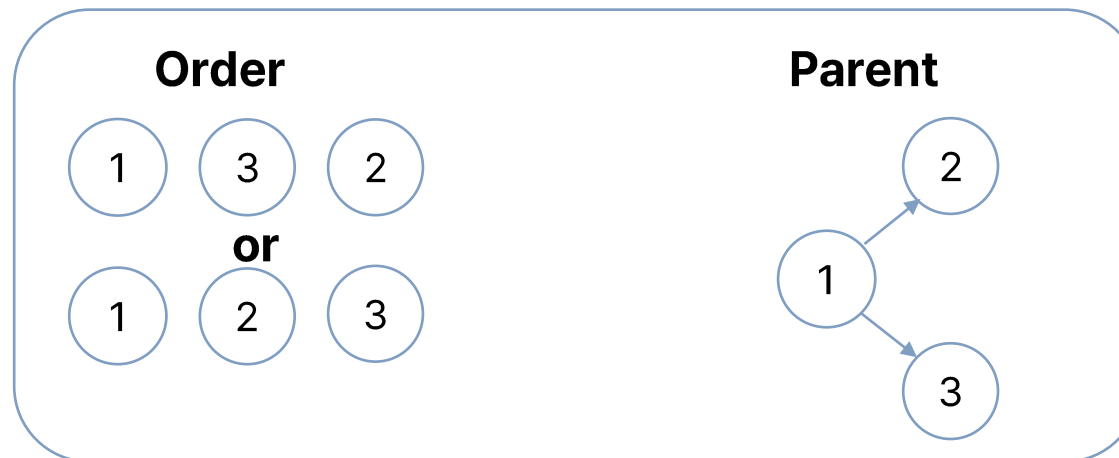
Order-based & Two-Stage (Order \rightarrow Parent)

Definition. First infer a ordering, then select parent set.

Key characteristics. If the order is correct, search is fast.

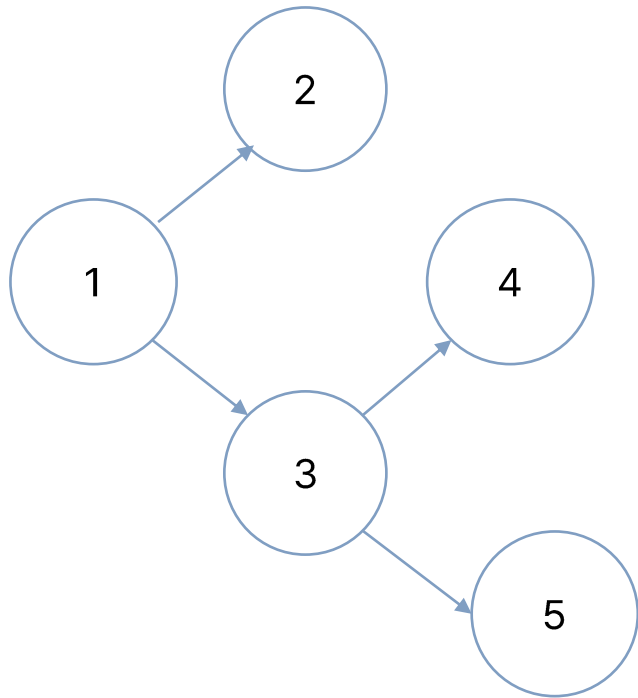
Representative methods (mixed classic+recent).

- *K2* (G.F. Cooper & E. Herskovits, 1992)
- *Ordering-Based Search* (M. Teyssier & D. Koller, 2005)
- *BOSS — Best Order Score Search* (B. Andrews, J. Ramsey, R. Sanchez-Romero, J. Camchong, & E. Kummerfeld, 2023)



2.0 . Preliminaries

[Backward Selection]



Backward Selection: Order then Parents

- **Order fixed:** (1, 2, 3, 4, 5)

Parent selection (check & pick)

- 5: {1, 2, 3, 4} Check $\Rightarrow 5 \leftarrow 3$
- 4: {1, 2, 3} Check $\Rightarrow 4 \leftarrow 3$
- 3: {1, 2} Check $\Rightarrow 3 \leftarrow 1$
- 2: {1} Check $\Rightarrow 2 \leftarrow 1$

2.0 . Preliminaries

[Causal Learning]

Differentiable DAG Learning

Definition. learning as continuous optimization with a smooth acyclicity constraint.

Key characteristics. Works well with deep learning models.

Representative methods (mixed classic+recent).

- *NOTEARS* (X. Zheng, B. Aragam, P. K. Ravikumar, & E. P. Xing, 2018)
- *GOLEM* (I. Ng, A. E. Ghassemi, & K. Zhang, 2020)
- *SDCD* — Stable Differentiable Causal Discovery (Achille Nazaret, Justin Hong, Elham Azizi, David M. Blei, 2024)

Discrete

$$\begin{aligned} & \min_{W \in \mathbb{R}^{d \times d}} F(W) \\ & \text{subject to } G(W) \in \text{DAGs} \end{aligned}$$

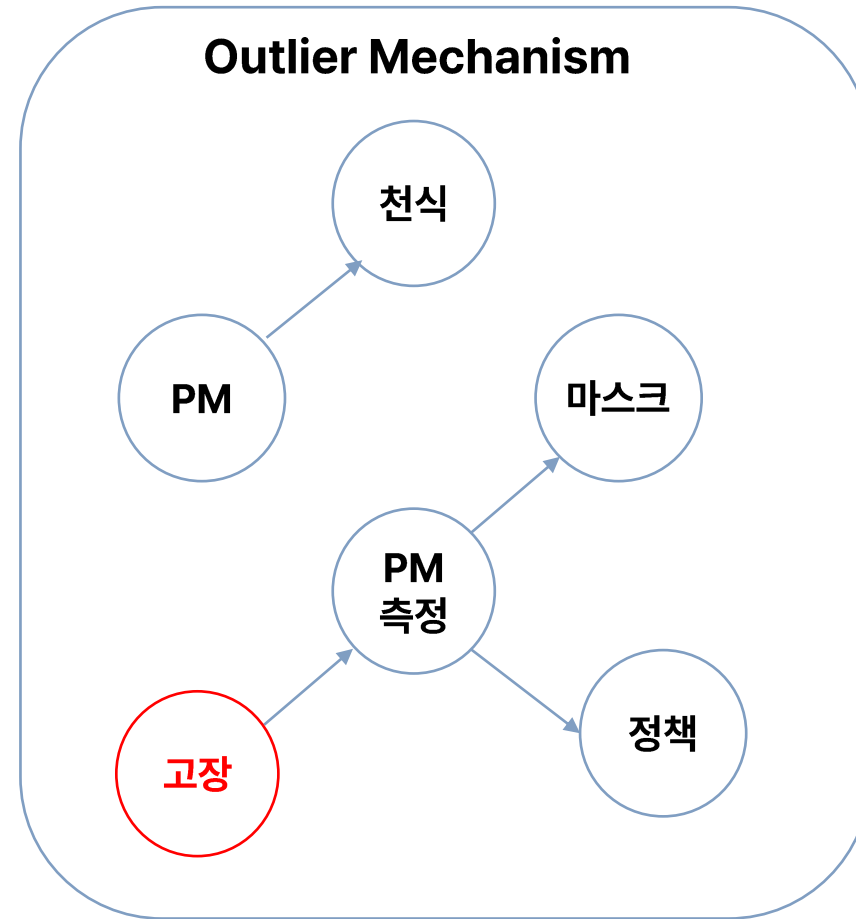
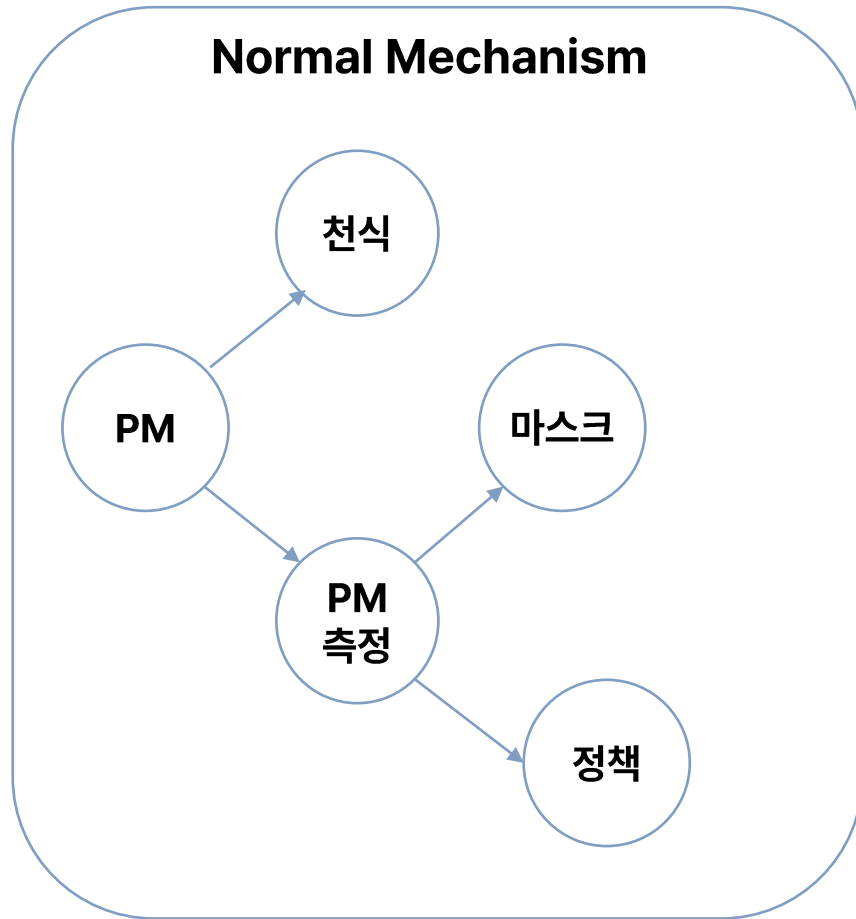
\iff

Continuous!!

$$\begin{aligned} & \min_{W \in \mathbb{R}^{d \times d}} F(W) \\ & \text{subject to } h(W) = 0, \end{aligned}$$

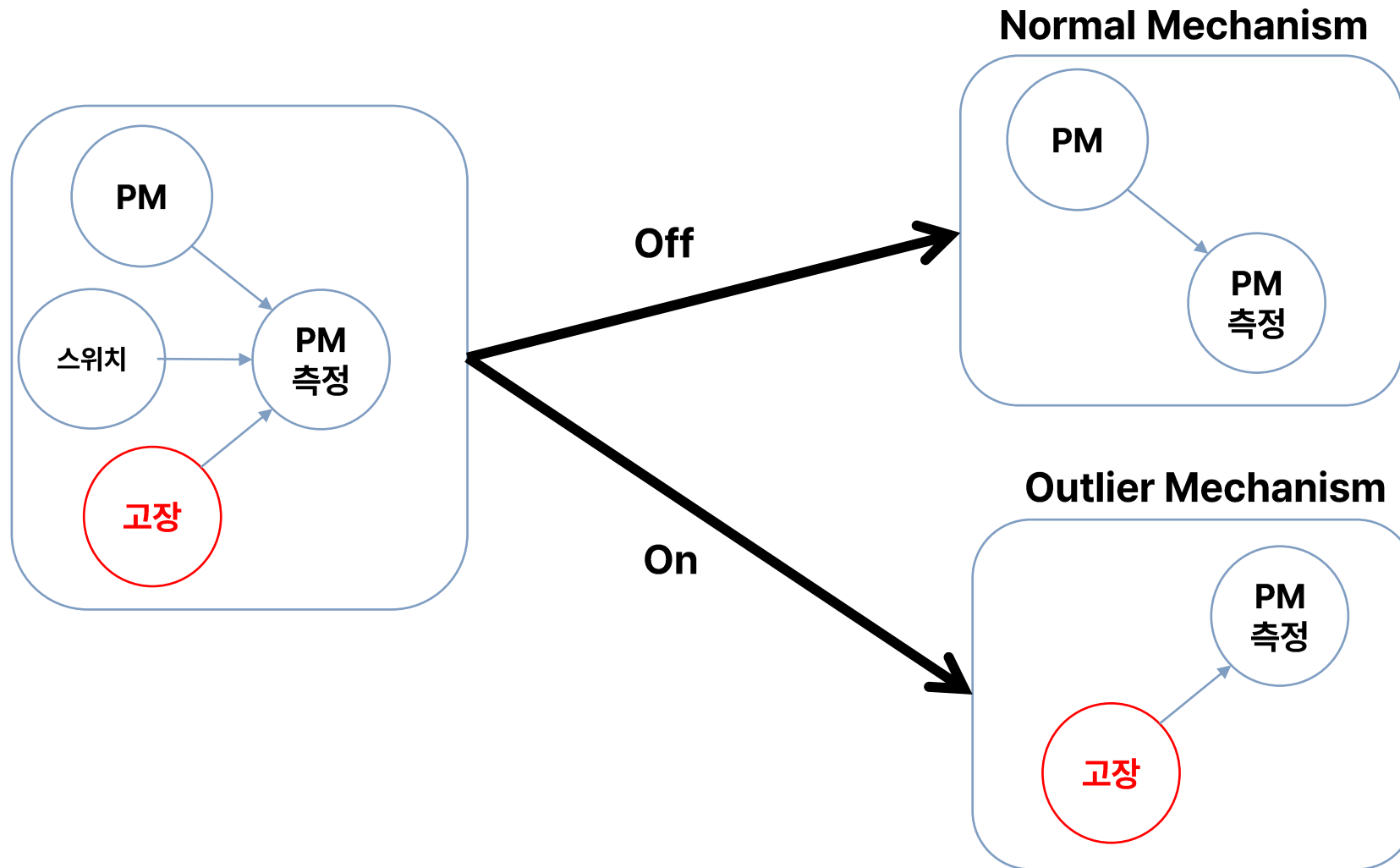
2.0 . Preliminaries

[Outlier]



2.0 . Preliminaries

[Outlier]



2.1 . Outlier Model

[Cellwise Contaminated SEM]

CCSEM

Structural assignments (for $j = 1, \dots, p$):

$$X_j = O_j \zeta_j + (1 - O_j) \left(\sum_{k \in \text{Pa}_G^X(j)} \beta_{k,j} X_k + \varepsilon_j \right).$$

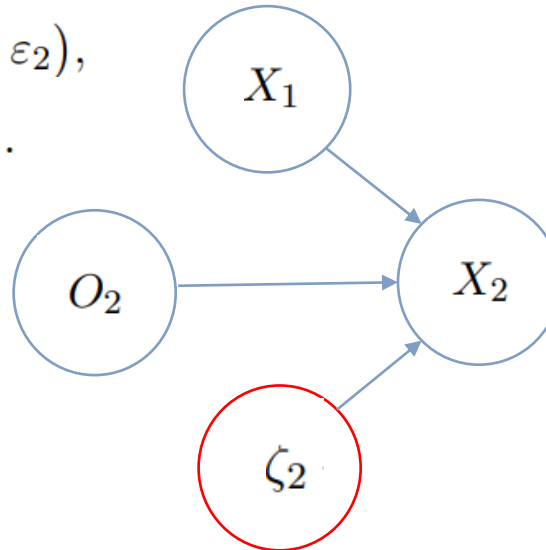
$$O_j = \mathbf{1}\{U_j \leq q_j(O_{\text{Pa}_G^O(j)})\}.$$

$$\zeta_j = g_j(\text{Pa}_G(\zeta_j)) + \xi_j.$$

Errors: $\varepsilon_j \sim N(0, \sigma_j^2)$ and are independent across j .

$$X_2 = O_2 \zeta_2 + (1 - O_2)(\beta X_1 + \varepsilon_2),$$

$$O_2 \sim \text{Bernoulli}(0.2), \quad X_1 = \varepsilon_1.$$



2.2 . Outlier Robust Algorithm

[Structure of Algorithm]

Step1: Order Recovery

$$\widehat{\pi}_r := \arg \max_{j \in V \setminus \{\widehat{\pi}_{p+1-r}, \dots, \widehat{\pi}_p\}} \left[\widehat{\Sigma}_{S'_j(r)}^{\widehat{G}_j(r)} \right]_{j,j}^{-1}$$

- where $\widehat{\Sigma}_{S'_j(r)}^{\widehat{G}_j(r)}$ is the sample covariance of $X_{S'_j(r)}^{\widehat{G}_j(r)}$.
- $X_{S'_j(r)}^{\widehat{G}_j(r)} \in \mathbb{R}^{|\widehat{G}_j(r)| \times |S'_j(r)|}$ (genuine observations).
- $S'_j(r) := \{j\} \cup \text{Supp}(\theta_j^*(r))$.
- $\widehat{G}_j(r) :=$ estimated good indices via residual thresholding.

Step2: Parent Recovery

$$\widehat{\text{Pa}}(j) := \text{Supp}(\widehat{\theta}_j(r)),$$

- $\widehat{\theta}_j(r)$ is the solution of the ℓ_1 -regularized LTS.
- $(\widehat{\theta}_j(r), \widehat{H}_j(r)) := \arg \min_{\theta_j \in \mathbb{R}^{|S_j(r)|}, H \subset \{1, \dots, n\}, |H|=h} \frac{1}{2h} \sum_{i \in H} (X_j^{(i)} - \langle X_{S_j(r)}^{(i)}, \theta_j \rangle)^2 + \lambda_j(r) \|\theta_j\|_1$

2.2 . Outlier Robust Algorithm

[Structure of Algorithm]

Step1: Order Recovery

$$\widehat{\pi}_r := \arg \max_{j \in V \setminus \{\widehat{\pi}_{p+1-r}, \dots, \widehat{\pi}_p\}} \widehat{\text{Var}}(X_j \mid X_{S_j(r)})$$

- where $\widehat{\Sigma}_{S_j'(r)}^{G_j(r)}$ is the sample covariance of $X_{S_j'(r)}^{G_j(r)}$.
- $X_{S_j'(r)}^{G_j(r)} \in \mathbb{R}^{|\widehat{G}_j(r)| \times |S_j'(r)|}$ (genuine observations).
- $S_j'(r) := \{j\} \cup \text{Supp}(\theta_j^*(r))$.
- $\widehat{G}_j(r) :=$ estimated good indices via residual thresholding.

Step2: Parent Recovery

$$\widehat{\text{Pa}}(j) := \text{Supp}(\widehat{\theta}_j(r)),$$

- $\widehat{\theta}_j(r)$ is the solution of the ℓ_1 -regularized LTS.
- $(\widehat{\theta}_j(r), \widehat{H}_j(r)) := \arg \min_{\theta_j \in \mathbb{R}^{|S_j(r)|}, H \subset \{1, \dots, n\}, |H|=h} \frac{1}{2h} \sum_{i \in H} (X_j^{(i)} - \langle X_{S_j(r)}^{(i)}, \theta_j \rangle)^2 + \lambda_j(r) \|\theta_j\|_1$

2.2 . Outlier Robust Algorithm

[Pseudo Code]

Algorithm 1: Robust Gaussian Linear SEM Learning Algorithm

Input : n independent samples, $X^{1:n}$

Output: Estimated graph structure, $\hat{G} = (V, \hat{E})$

Set $\hat{\pi}_{p+1} = \emptyset$;

for $r = \{1, 2, \dots, p-1\}$ **do**

for $j \in V \setminus \{\hat{\pi}_{p+1}, \dots, \hat{\pi}_{p+2-r}\}$ **do**

$S_j(r) = V \setminus (\{j\} \cup \{\hat{\pi}_{p+1}, \dots, \hat{\pi}_{p+2-r}\})$;

 Estimate $\hat{\theta}_j(r)$ for ℓ_1 -regularized LTS in Equation (2);

 Estimate truncated conditional variances $\widehat{\text{Var}}(X_j \mid X_{S_j(r)})$ using Equation (3);

end

 Determine $\hat{\pi}_{p+1-r} = \arg \max_j \widehat{\text{Var}}(X_j \mid X_{S_j(r)})$;

 Determine $\widehat{\text{Pa}}(\hat{\pi}_{p+1-r}) = \{k \in S_{\hat{\pi}_{p+1-r}}(r) : [\hat{\theta}_{\hat{\pi}_{p+1-r}}(r)]_k \neq 0\}$;

end

Return: Estimate an edge set, $\hat{E} = \cup_{r \in \{1, 2, \dots, p-1\}} \{(k, \hat{\pi}_{p+1-r}) : k \in \widehat{\text{Pa}}(\hat{\pi}_{p+1-r})\}$

2.3 . Consistency of Algorithm

[Order Recovery]

Assumption (Outlier–distance separation)

- For $r \in \{1, \dots, p-1\}$, $j \in \{\pi_1, \dots, \pi_{p+1-r}\}$, $T_j(r) := \{\pi_1, \dots, \pi_{p+1-r}\} \setminus \{j\}$.
- There exists $\eta_{\min} > 0$ such that

$$\min_{i \in B_j(r, o)} \left| X_j^{(i)} - \mathbb{E}[X_j^{(i)} \mid X_{T_j(r)}^{(i)}] \right| > \eta_{\min}.$$

- Bad–row set: $B_j(r, o) := \{i \in \{1, \dots, n\} : \exists k \in \{j\} \cup T_j(r) \text{ with } o_k^{(i)} = 1\}$.

Assumption (Truncated conditional–variance gap)

- For $r \in \{1, \dots, p-1\}$, $j \in \{\pi_1, \dots, \pi_{p+1-r}\}$, $T_j(r) := \{\pi_1, \dots, \pi_{p+1-r}\} \setminus \{j\}$.
- For any $\eta \in (0, \eta_{\min})$, $\exists \tau_{\min} > 0$ such that for every $\ell \in \text{An}(j)$,

$$\begin{aligned} & \text{Var}\left(X_j \mid X_{T_j(r)}, |X_j - \mathbb{E}(X_j \mid X_{T_j(r)})| < \eta\right) \\ & - \text{Var}\left(X_\ell \mid X_{T_j(r)}, |X_\ell - \mathbb{E}(X_\ell \mid X_{T_j(r)})| < \eta\right) > \tau_{\min}. \end{aligned}$$

2.3 . Consistency of Algorithm

[Order Recovery]

Theorem (Recovery of the Ordering)

- Gaussian linear CCSEM with cellwise contamination. Assume (Outlier–distance separation) and (Truncated conditional–variance gap) hold.
- Suppose, for each $r \in \{1, \dots, p-1\}$, $j \in \{\pi_1, \dots, \pi_{p+1-r}\}$, the supports of $\theta_j^*(r)$ are provided, and choose $\eta \in (0, \eta_{\min})$.
- If the truncated–variance sample size $h' := \min_{j,r} |\widehat{G}_j(r)| \geq C'_\epsilon d^2 \log p$,
 $\Rightarrow \Pr(\widehat{\pi} \in \Pi^*) \geq 1 - \epsilon$, for any $\epsilon > 0$, with some constant $C'_\epsilon > 0$.

2.3 . Consistency of Algorithm

[Parent Recovery]

Assumption (Clean-row fraction under CCSEM)

- Fix iteration $r \in \{1, \dots, p-1\}$ and $j \in \{\pi_1, \dots, \pi_{p+1-r}\}$.
- Let $S_j(r)$ be the design set and define the bad-row set under config. o by
$$B_j(r, o) := \{ i \in \{1, \dots, n\} : \exists k \in \{j\} \cup S_j(r) \text{ with } o_k^{(i)} = 1 \}.$$
- There exists $\alpha_{\min} \in (0, 1]$ such that $|B_j(r, o)| \leq \min\{(1 - \alpha_{\min})h, n - h\}$.

2.3 . Consistency of Algorithm

[Parent Recovery]

Assumption (Spectral bound)

- $X_{S_j(r)}^{B_j(r,o)} \in \mathbb{R}^{|B_j(r,o)| \times (p-1)}$ collects rows $B_j(r,o)$ and columns $S_j(r)$.
- There exists $c_{\max} > 0$ such that

$$\max_{j,r} \max_{\substack{K_j \subset V \setminus \{j\} \\ |K_j| = |B_j(r,o)|}} \|X_{K_j}^{B_j(r,o)}\| \leq c_{\max} \sqrt{|B_j(r,o)| \log p},$$

where $\|A\| := \sup_{\|v\|_2 \leq 1} \|Av\|_2$ (spectral norm).

Assumption (Parameter bound)

- $\forall r \in \{1, \dots, p-1\}, j \in \{\pi_1, \dots, \pi_{p+1-r}\}, \exists c_1 > 0$ such that
$$\theta_{\min} < \min_{j,r} \|\theta_j^*(r)\|_{\min} \leq \max_{j,r} \|\theta_j^*(r)\|_1 \leq \rho \leq \frac{c_1}{2} \sqrt{\frac{h}{\log p}},$$
- ρ is the tuning parameter used in the ℓ_1 -regularized LTS step.

2.3 . Consistency of Algorithm

[Parent Recovery]

Theorem (Sign Recovery)

- Consider the ℓ_1 -regularized LTS. with solution $\hat{\theta}_j(r)$.
- Assume (Clean-row fraction), (Spectral bound), and (Parameter bound) hold.
- Then for any $\epsilon > 0$, $\exists c_\epsilon, \kappa_1 > 0$ such that

$$\lambda_j(r) = c_\epsilon \sqrt{\frac{\log p}{h}},$$

$$h \geq \left(\frac{c_\epsilon}{\kappa_1 \alpha_{\min} \theta_{\min}^2} \right)^2 \left(\frac{9}{4} d + 4 |B_j(r, o)| \right) \log p,$$

$$\Rightarrow \Pr\left(\text{sign}(\hat{\theta}_j(r)) = \text{sign}(\theta_j^*(r))\right) \geq 1 - \frac{\epsilon}{p^2}.$$

2.3 . Consistency of Algorithm

[Consistency]

Corollary (Consistency of the Algorithm)

- Setting: Gaussian linear CCSEM with cellwise contamination.
- Assume (Clean-row fraction), (Spectral bound), (Parameter bound), and (Truncated conditional-variance gap) hold, with appropriate tuning (λ, η) .
- If $h = \Omega((d + |B|) \log p)$ and $h' = \Omega(d^2 \log p)$,
where $|B| := \max_{j,r} |B_j(r, o)|$, $h' := \min_{j,r} |\hat{G}_j(r)|$,
 $\Rightarrow \Pr(\hat{G} = G) \rightarrow 1 \quad (n \rightarrow \infty)$.

2.4 . Simulations and Real-Data

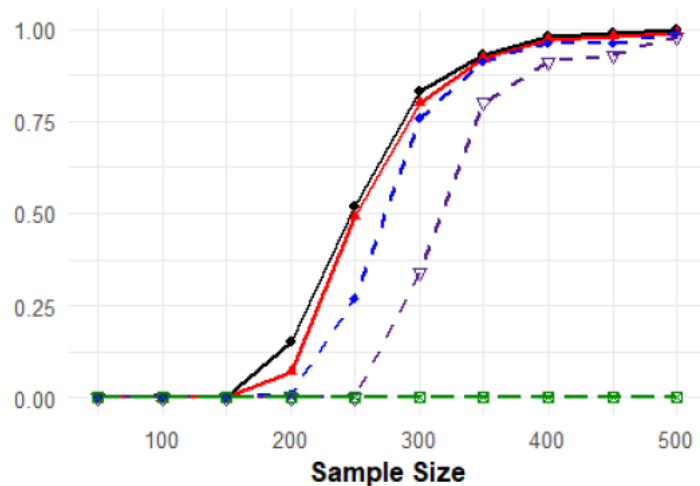
[Simulations]

Simulation Settings (CCSEM)

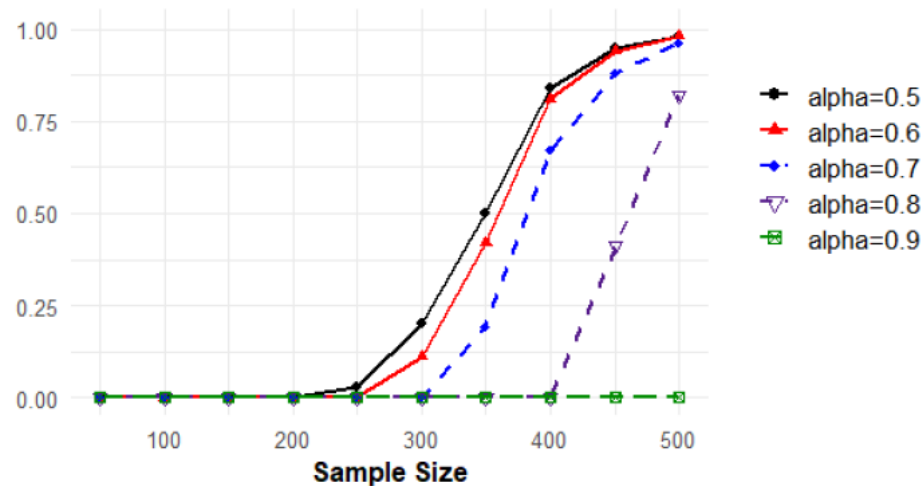
- **Graph size:** $p \in \{5, 10, \dots, 25\}$; 100 realizations per p .
- **Sparsity:** max degree $d \leq 4$, min indegree = 1; no isolated nodes.
- **Edge weights:** $\beta_{k,j}$ i.i.d. uniform in $(-1, -0.75) \cup (0.75, 1)$.
- **Noise:** $\sigma_j^2 = 0.75$ for all j .
- **Outliers:** $B \in \{1, 30, 60, 90\}$,
values generated from $N(100, \sigma_j^2)$.

2.4 . Simulations and Real-Data

[Number of Trimmed Sample]



(c) $|B| = 60$

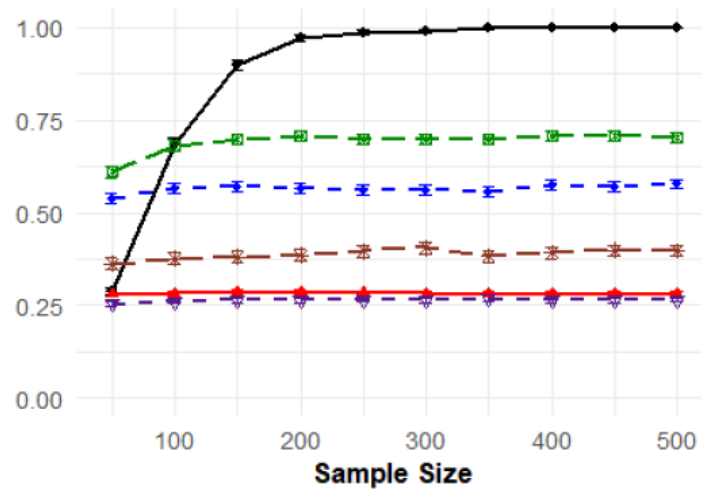


(d) $|B| = 90$

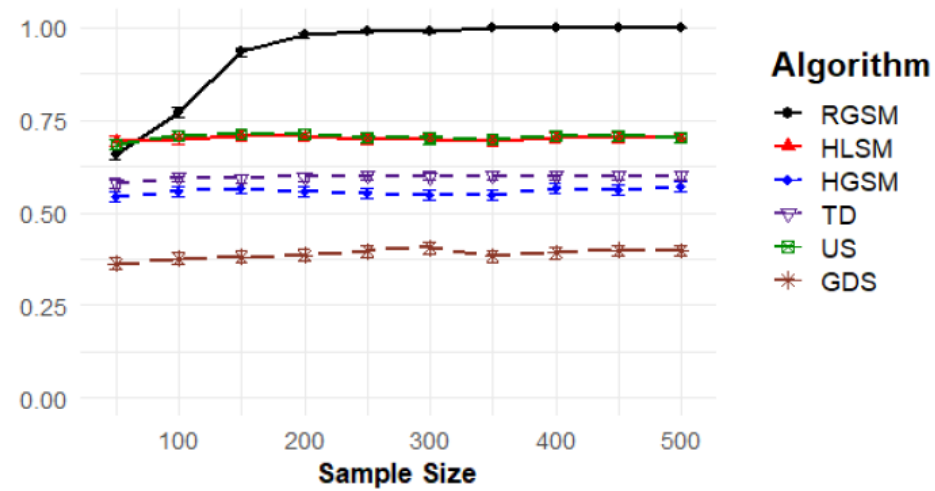
Figure 2: Performance of the proposed algorithm with various fractions of the sample size, $\alpha \in \{0.5, 0.6, \dots, 0.9\}$, for learning 10-node corrupted Gaussian linear SEMs with the different maximum number of bad samples ($|B| \in \{1, 30, 60, 90\}$) on the first element of the ordering. The empirical probability of successful graph recovery is shown versus sample size.

2.4 . Simulations and Real-Data

[Comparison with Other Algorithms]



(c) Precision: $|B| = 30$

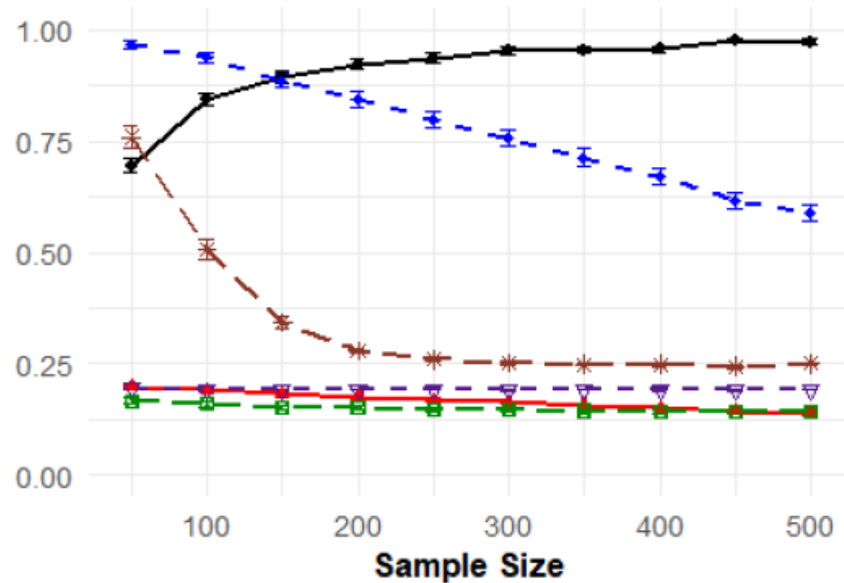


(d) Recall: $|B| = 30$

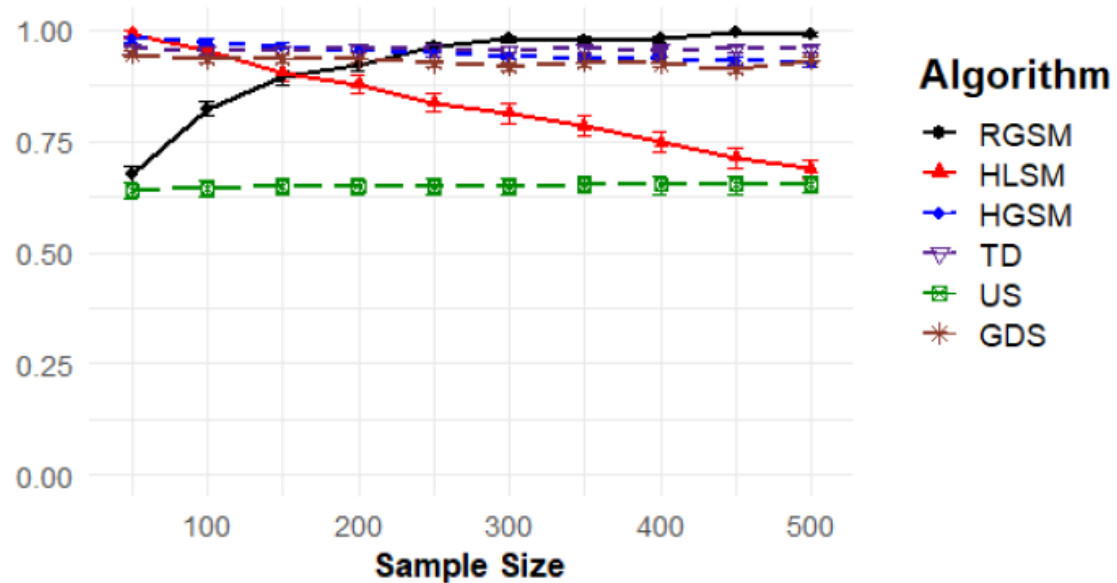
Figure 6: Comparison of the proposed algorithm (RGSM) against the HLSM, HGSM, TD, and US algorithms in terms of the average precision and recall for learning 10-node corrupted Gaussian linear SEMs with $|B| \in \{1, 30\}$.

2.4 . Simulations and Real-Data

[All Outliers]



(a) Precision



(b) Recall

Figure 10: Comparison of the proposed algorithm (RGSM) against the HLSM, HGSM, TD, and US algorithms in terms of average Hamming distance for learning 10-node corrupted Gaussian linear SEMs when all observation are outliers.

2.4 . Simulations and Real-Data

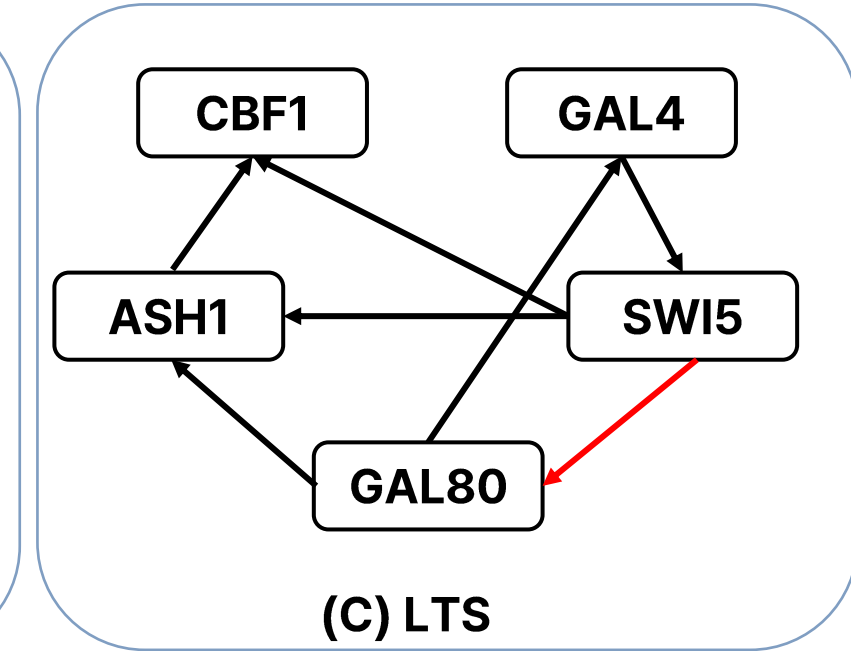
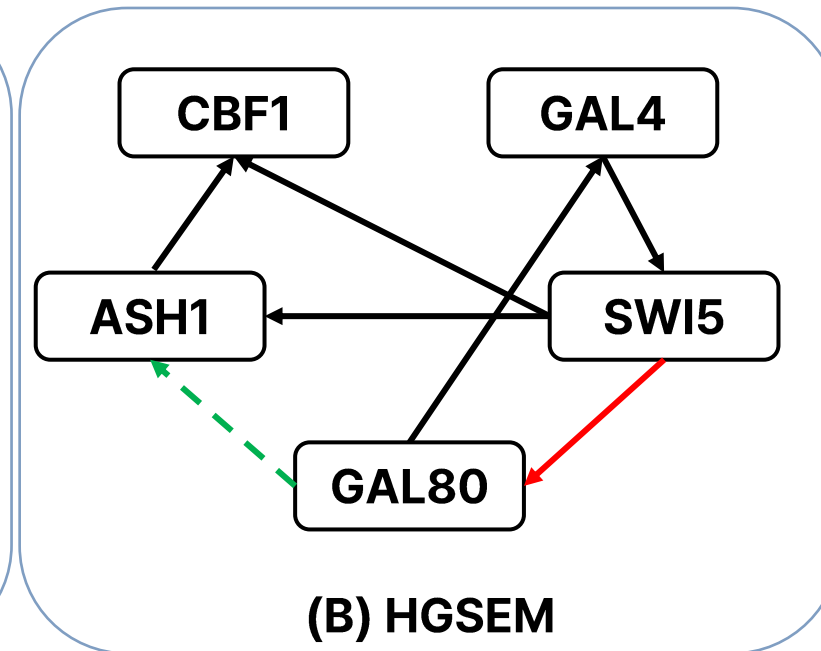
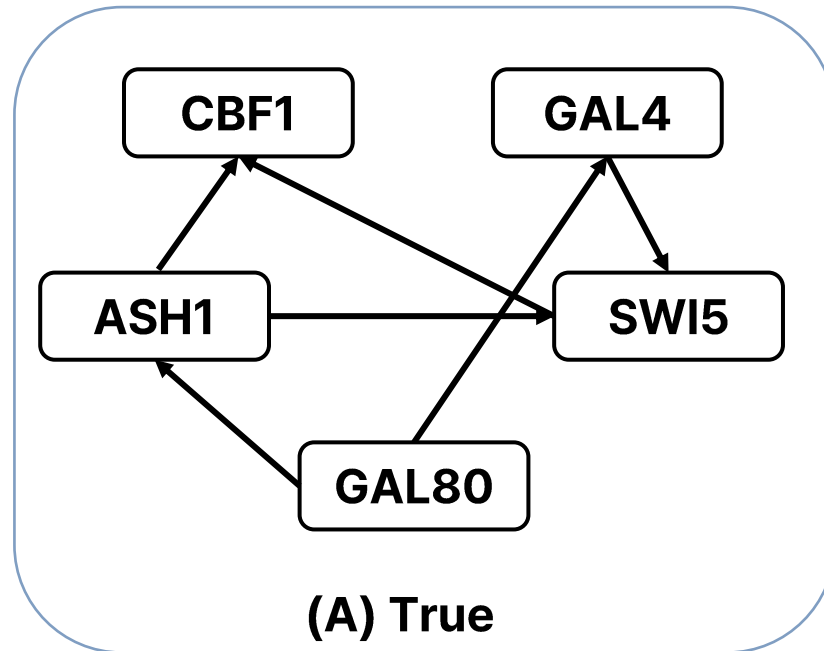
[Real-Data]

IRMA Yeast Real Data: Dataset

- **IRMA Yeast Network:**
 - 5 genes (CBF1, GAL4, SWI5, GAL80, ASH1)
 - Ground-truth DAG with 8 directed edges
- **Experimental Conditions:**
 - “Off” (glucose medium, network repressed)
 - “On” (galactose medium, network induced)
- **Measurements:**
 - 21 time-points per condition ($n = 42$ samples)
 - Steady-state mRNA expression (DNA microarrays)

2.4 . Simulations and Real-Data

[Real-Data]



감사합니다.

