

# Learning Sparse Fixed-Structure Gaussian Bayesian Networks

Arnab Bhattacharyya<sup>1</sup> Davin Choo<sup>1</sup> Rishikesh Gajjala<sup>2</sup>  
Sutanu Gayen<sup>1</sup> Yuhao Wang<sup>1</sup>

<sup>1</sup>National University of Singapore

<sup>2</sup>Indian Institute of Science, Bangalore

# Table of Contents

- 1 Background
- 2 High-level approach
- 3 What's new: coefficients recovery
  - Estimators based on least squares
  - Estimator based on Cauchy random variables
- 4 Hardness results
- 5 Experiments

# Bayesian Network

## Background

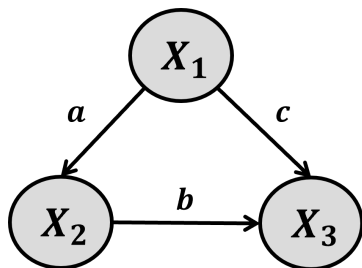


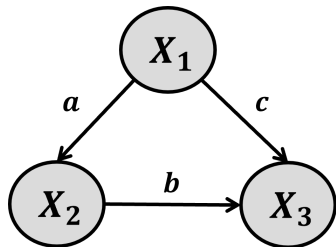
Figure: Bayesian Network

- 1 The DAG  $G = (V, E)$
- 2  $V = \{X_1, \dots, X_n\}$
- 3  $(X_j, X_i) \in E$  whenever  $X_j \rightarrow X_i$
- 4 For variable  $X_i$  with parent indices  $\pi_i \subseteq [n]$
- 5 By Bayesian rule:  
$$P(X_1, \dots, X_n) = \prod_{i=1}^n \Pr_{\mathcal{P}}(X_i \mid \pi_i)$$

**Eg.**  $P(X_1, X_2, X_3) = P(X_1)P(X_2|X_1)p(X_3|X_1, X_2)$

# Structural Equation Models with Gaussian Noise

## Background



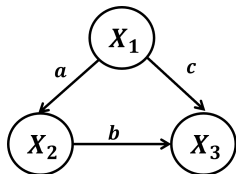
$$\begin{aligned}X_1 &= \eta_1, \\X_2 &= aX_1 + \eta_2, \\X_3 &= bX_2 + cX_3 + \eta_3 \\ \eta_i &\sim \mathcal{N}(0, \sigma_i^2)\end{aligned}$$

- 1  $\hat{M}$  empirical covariance matrix of  $X_1, \dots, X_p$  with Cholesky decomposition  $\hat{M} = \hat{L}\hat{L}^\top$
- 2  $(X_1, \dots, X_p) \sim N(0, M)$  is distributed as a multivariate Gaussian;
- 3 Structural equation model:  $X_i = \eta_i + \sum_{j \in \pi_i} a_{i \leftarrow j} X_j$ ,  $\eta_i \sim N(0, \sigma_i^2)$ .

# Objective

## Background

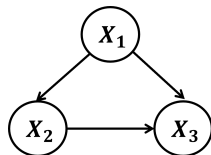
Goal: Parameter learning on a given graph structure



Ground truth with parameters  
 $[a, b, c]$  and  $[\sigma_1, \sigma_2, \sigma_3]$

$X_1$	$X_2$	$X_3$
$X_1^{(1)}$	$X_2^{(1)}$	$X_3^{(1)}$
...	...	...

Draw sample (possibly contaminated)



+ Given graph structure

Objective:

Recover  $[a, b, c]$     $[\sigma_1, \sigma_2, \sigma_3]$    inducing  $\hat{P}$  such that  $d_{TV}(P, \hat{P}) \leq \epsilon$ ;  
Parameter estimation   Variance recovery   Distance measure

# Table of Contents

## High-level approach

- 1 Background
- 2 High-level approach**
- 3 What's new: coefficients recovery
  - Estimators based on least squares
  - Estimator based on Cauchy random variables
- 4 Hardness results
- 5 Experiments

# The distance measure

## High-level approach

- **Total variational (TV) distance:**

$$d_{\text{TV}}(\mathcal{P}, \mathcal{Q}) = \sup_{A \in \mathbb{R}^n} |\mathcal{P}(A) - \mathcal{Q}(A)| = \frac{1}{2} \int_{\mathbb{R}^n} |\mathcal{P}(x) - \mathcal{Q}(x)| dx.$$

- **Kullback–Leibler (KL) divergence:**

$$d_{\text{KL}}(\mathcal{P}, \mathcal{Q}) = \int_{A \in \mathbb{R}^n} \mathcal{P}(A) \log \left( \frac{\mathcal{P}(A)}{\mathcal{Q}(A)} \right) dA.$$

Fact (Pinsker's inequality)

For distributions  $\mathcal{P}$  and  $\mathcal{Q}$ ,  $d_{\text{TV}}(\mathcal{P}, \mathcal{Q}) \leq \sqrt{d_{\text{KL}}(\mathcal{P}, \mathcal{Q})/2}$ .

If  $s(\varepsilon)$  samples are needed to ensure  $\underline{d_{\text{KL}}(\mathcal{P}, \mathcal{Q}) \leq \varepsilon}$   
 $\implies s(\varepsilon^2)$  samples are needed to ensure  $\underline{d_{\text{TV}}(\mathcal{P}, \mathcal{Q}) \leq \varepsilon}$ .

# Decomposing KL divergence:

## High-level approach

Decompose KL divergence into  $n$  terms:

$$d_{\text{KL}}(\mathcal{P}, \mathcal{Q}) = \sum_{i=1}^n d_{\text{CP}}(\alpha_i^*, \hat{\alpha}_i)$$

$\implies$  Estimate parameters for each variable independently

$$d_{\text{KL}}(\mathcal{P}, \mathcal{Q}) = \sum_{i=1}^n d_{\text{CP}}(\alpha_i^*, \hat{\alpha}_i) = \sum_{i=1}^n \ln \left( \frac{\hat{\sigma}_i}{\sigma_i} \right) + \frac{\sigma_i^2 - \hat{\sigma}_i^2}{2\hat{\sigma}_i^2} + \frac{\Delta_i^\top M_i \Delta_i}{2\hat{\sigma}_i^2}$$

- $\alpha_i^* = (A_i, \sigma_i)$ : coefficients and variance associated with  $X_i$
- $M_i$ : covariance matrix associated with  $X_i$
- $\hat{\alpha}_i = (\hat{A}_i, \hat{\sigma}_i)$ : estimates for  $\alpha_i^*$
- $\Delta_i = \hat{A}_i - A_i$



# Two-phased recovery approach

## High-level approach

- **First phase:** Estimate the **coefficients**  $\hat{A}_1, \dots, \hat{A}_n$  of the Bayesian network
- **Second phase:** Recover **variances**  $\hat{\sigma}_y^2$  using empirical variances conditioned on our recovered coefficients ( $y$  refers to any arbitrary variable index).
- We are running the same algorithm for **each node**.

# Coefficient estimators overview

## High-level approach

- **There's no explicit sample complexity bound** known for Node-wise least square

Our contributions:

Algorithm	Remarks
LeastSquares	$\tilde{O}(nd/\varepsilon^2)$ samples within TV distance $\varepsilon$ .
BatchAvgLeastSquares	Distributed-friendly generalization.
CauchyEstTree	$\tilde{O}(nd/\varepsilon^2)$ samples within TV distance $\varepsilon$ .
CauchyEst*	Robust when data are contaminated

- All algorithms run in  $\text{poly}(n, d, \log(\delta^{-1}))$  time with success probability at least  $1 - \delta$ .
- \* No theoretical guarantees but empirically robust against contaminated samples

# Table of Contents

What's new: coefficients recovery

- 1 Background
- 2 High-level approach
- 3 What's new: coefficients recovery**
  - Estimators based on least squares
  - Estimator based on Cauchy random variables
- 4 Hardness results
- 5 Experiments

# Vanilla least square

What's new: coefficients recovery Estimators based on least squares

**LeastSquares:** using linear least squares.

$$\begin{pmatrix} X_1^{(1)} & \dots & X_p^{(1)} \\ \vdots & \ddots & \vdots \\ X_1^{(m_1)} & \dots & X_p^{(m_1)} \end{pmatrix} \cdot \begin{pmatrix} a_1 \\ \vdots \\ a_p \end{pmatrix} = \begin{pmatrix} Y^{(1)} \\ \vdots \\ Y^{(m_1)} \end{pmatrix}$$

$X \in \mathbb{R}^{m_1 \times p}$        $A \in \mathbb{R}^p$        $B \in \mathbb{R}^{m_1}$

**Problem:**  $X\hat{A} = B$   
**Solution:**  $\hat{A} = (X^T X)^{-1} X^T B$

No known analysis for the explicit **sample complexity bound**

Conclusion:  $\mathcal{O}(nd_{avg}\varepsilon^{-2} \cdot \log(n\delta^{-1}))$  samples within TV distance  $\varepsilon$ .

# Batch average least squares

What's new: coefficients recovery Estimators based on least squares

**BatchAvgLeastSquares**: any interpolation between “batch size” and “number of batches” as long as the **total number of samples** used is sufficiently large.

$$\begin{array}{l} \text{batch 1} \\ \dots \\ \text{batch b} \end{array} \begin{pmatrix} X_1^{(1)} & \dots & X_p^{(1)} \\ \vdots & \ddots & \vdots \\ X_1^{(k)} & \dots & X_p^{(k)} \\ \dots & \dots & \dots \\ X_1^{((b-1)k+1)} & \dots & X_p^{((b-1)k+1)} \\ \vdots & \ddots & \vdots \\ X_1^{(m_1)} & \dots & X_p^{(m_1)} \end{pmatrix} \cdot \begin{pmatrix} a_1 \\ \vdots \\ a_p \end{pmatrix} = \begin{pmatrix} Y^{(1)} \\ \vdots \\ Y^{(k)} \\ \dots \\ Y^{((b-1)k+1)} \\ \vdots \\ Y^{(m_1)} \end{pmatrix}$$

$X \in \mathbb{R}^{k \times p}$                        $A \in \mathbb{R}^p$                        $b \in \mathbb{R}^k$

# CauchyEstTree

What's new: coefficients recovery Estimator based on Cauchy random variables

**Notation:**  $\hat{M}$  empirical covariance matrix of  $X_1, \dots, X_p$  with Cholesky decomposition  $\hat{M} = \hat{L}\hat{L}^\top$

$$d_{\text{KL}}(\mathcal{P}, \mathcal{Q}) = \sum_{i=1}^n d_{\text{CP}}(\alpha_i^*, \hat{\alpha}_i) = \sum_{i=1}^n \ln \left( \frac{\hat{\sigma}_i}{\sigma_i} \right) + \frac{\sigma_i^2 - \hat{\sigma}_i^2}{2\hat{\sigma}_i^2} + \underbrace{\frac{\Delta_i^\top M_i \Delta_i}{2\hat{\sigma}_i^2}}$$

$$\Rightarrow \left| \Delta^\top M \Delta \right| = \left| \Delta^\top L L^\top \Delta \right| = \left\| L^\top \Delta \right\|^2$$

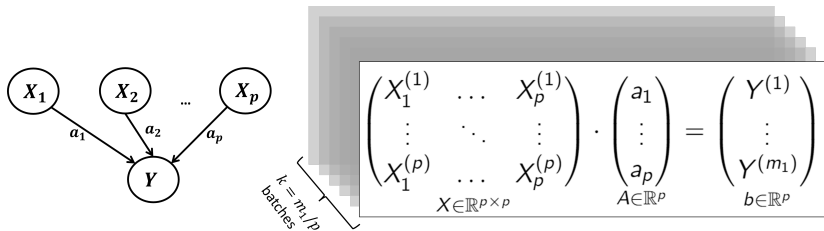
## CauchyEstTree

Consider a batch estimate  $\tilde{A}$  and define  $\Delta = \tilde{A} - A$ . If the Bayesian network is a polytree, then  $\Delta_i = (\tilde{A} - A)_i \sim \frac{\sigma_y}{\sigma_i} \cdot \text{Cauchy}(0, 1)$  for all  $i \in [n]$ .

The study of Cauchy random variables

# CauchyEstTree: Algorithm

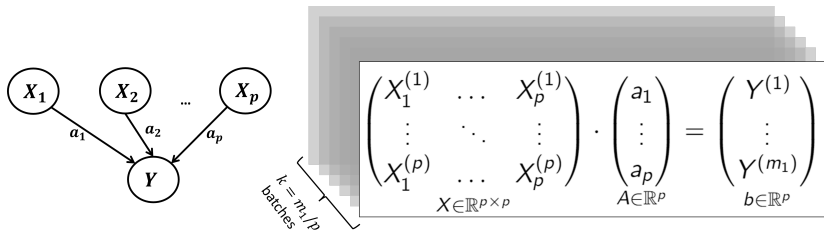
What's new: coefficients recovery Estimator based on Cauchy random variables



- Compute  $\tilde{A} = [\hat{a}_{y \leftarrow 1}, \dots, \hat{a}_{y \leftarrow p}]^T$  as **any solution** to  $X\tilde{A} = [Y^{(1)}, \dots, Y^{(p)}]^T$

# CauchyEstTree: Algorithm

What's new: coefficients recovery Estimator based on Cauchy random variables

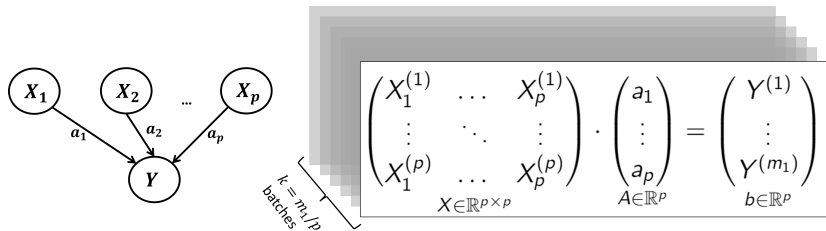


- Compute  $\tilde{A} = [\hat{a}_{y \leftarrow 1}, \dots, \hat{a}_{y \leftarrow p}]^\top$  as **any solution** to  $X\tilde{A} = [Y^{(1)}, \dots, Y^{(p)}]^\top$
- **MED** <sub>$i$</sub>  as the **median** of  $(\hat{L}^\top \tilde{A}^{(1)})_i, \dots, (\hat{L}^\top \tilde{A}^{(\lfloor m/p \rfloor)})_i$ ;



# CauchyEstTree: Algorithm

What's new: coefficients recovery Estimator based on Cauchy random variables



- Compute  $\tilde{A} = [\hat{a}_{y \leftarrow 1}, \dots, \hat{a}_{y \leftarrow p}]^\top$  as **any solution** to  $X\tilde{A} = [Y^{(1)}, \dots, Y^{(p)}]^\top$
- **MED**<sub>*i*</sub> as the **median** of  $(\hat{L}^\top \tilde{A}^{(1)})_i, \dots, (\hat{L}^\top \tilde{A}^{(l[m/p])})_i$ ;
- $\hat{A}^\top = (\hat{L}^\top)^{-1} [\text{MED}_1, \dots, \text{MED}_n]^\top$

# Cauchy-based estimators

What's new: coefficients recovery Estimator based on Cauchy random variables

**Q:** Why take median instead of mean?

**A:** The variance of a Cauchy variable is unbounded

**Conclusion:**  $\mathcal{O}(nd_{avg}d\varepsilon^{-1} \cdot \log(n\delta^{-1}))$  samples within TV distance  $\varepsilon$ .

**Generalization** to random DAGs  $\implies$  **CauchyEst** estimator.

**Limitation:** No guarantees when DAG is not a tree.

# Table of Contents

- 1 Background
- 2 High-level approach
- 3 What's new: coefficients recovery
  - Estimators based on least squares
  - Estimator based on Cauchy random variables
- 4 Hardness results**
- 5 Experiments

# Hardness results

## Hardness results

### Learning Gaussian product distributions

Given samples from a  $n$ -fold Gaussian product distribution  $P$ , learning a  $\hat{P}$  such that  $\text{d}_{\text{TV}}(P, \hat{P}) = O(\varepsilon)$  with success probability  $2/3$  needs  $\Omega(n\varepsilon^{-2})$  samples in general.

### Learning Gaussian Bayesian networks

For any  $0 < \varepsilon < 1$  and  $n, d$  such that  $d \leq n/2$ , there exists a DAG  $G$  over  $[n]$  of in-degree  $d$  such that learning a Gaussian Bayesian network  $\hat{P}$  on  $G$  such that  $\text{d}_{\text{TV}}(P, \hat{P}) \leq \varepsilon$  with success probability  $2/3$  needs  $\Omega(nd\varepsilon^{-2})$  samples in general.

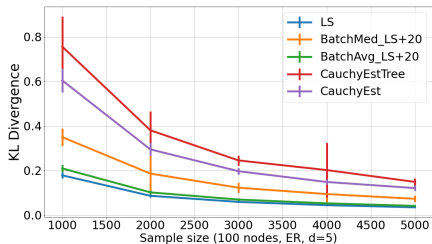
# Table of Contents

## Experiments

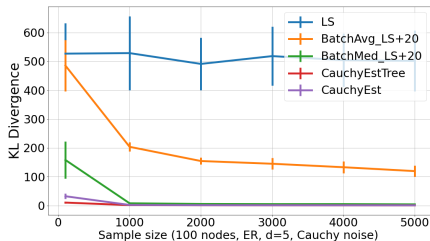
- 1 Background
- 2 High-level approach
- 3 What's new: coefficients recovery
  - Estimators based on least squares
  - Estimator based on Cauchy random variables
- 4 Hardness results
- 5 Experiments**

# Experiment results

## Experiments



Algorithms evaluated on ER graph with  $d = 5$  on *uncontaminated* data



Algorithms evaluated on ER graph with  $d = 5$  on *contaminated* data