

Causal Discovery in the presence of Latent Confounders

A Stochastic Optimization Approach

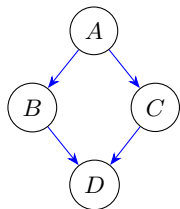
Danish Shakeel Mohammad

Queen Mary, University of London

December 2, 2025

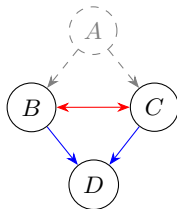
Introduction

Causal Discovery is the automated extraction of cause and effect relationships from data.



DAG

Causally Sufficient



ADMG

Causally Insufficient

Linear Gaussian SCMs

$$\mathbf{V} = \mathbf{\Theta}\mathbf{V} + \mathbf{U}$$

$$\mathbf{U} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Omega})$$

\mathbf{V} : observable variables

\mathbf{U} : unobservable variables

$\mathbf{\Theta}$: structural coefficients matrix

$\mathbf{\Omega}$: error covariances matrix

In the presence of latent confounding, off-diagonal elements of $\mathbf{\Omega}$ are non-zero.

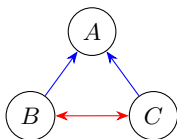
Definition (Causal effect — real world definition)

A variable X has a causal effect on the variable Y if forcing X to take some value x , the distribution of Y explicitly depends on x :

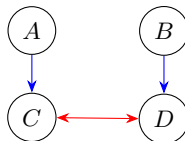
$$\exists x \in \mathcal{X} : P(Y|\text{do}(X = x)) \neq P(Y)$$

Graph Classes and Problem Statement

Bow-free



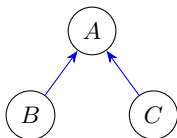
Ancestral



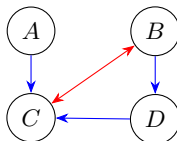
Bow-free ADMGs

- ① In the case of Linear Gaussian SCMs are almost-everywhere identifiable, in the limit of infinite data
- ② Can capture Verma Constraints

Non Bow-free



Non Ancestral

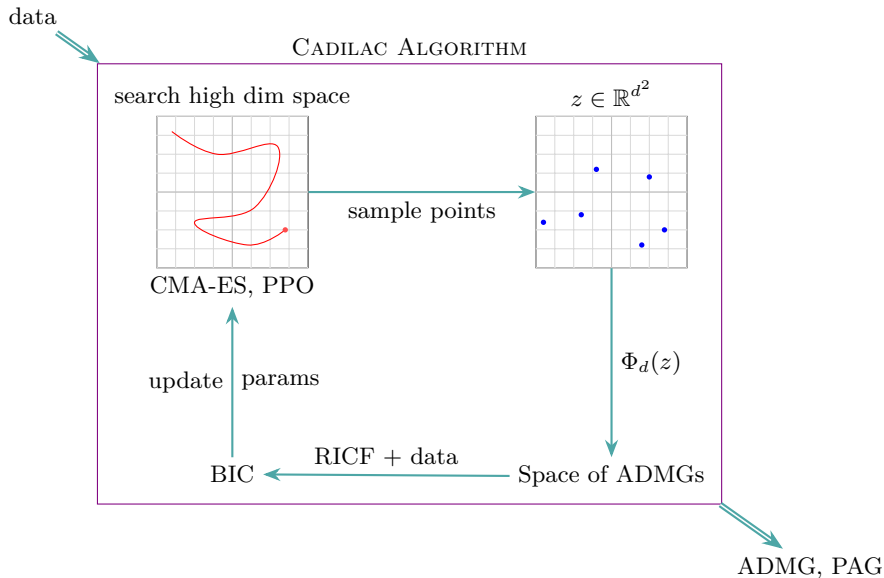


Ancestral ADMGs

- ① In the case of Linear Gaussian SCMs are globally identifiable, in the limit of infinite data
- ② Can't capture Verma Constraints
- ③ Are a subset of Bow-free ADMGs

Problem Statement

Causal Discovery on observational data for Linear Gaussian Structural Causal Models (SCMs), targeting ancestral and bow-free Acyclic Directed Mixed Graphs (ADMGs).



Definition

$\forall d \in \mathbb{N}^+$ and $z \in \mathbb{R}^{d^2}$, $p \in \mathbb{R}^d$ is the vector formed from the first d elements of z ,
 $E_{\rightarrow}, E_{\leftrightarrow} \in \mathbb{R}^{d \times d}$ are strictly lower triangular matrices formed from the next $\frac{d(d-1)}{2}$ and final $\frac{d(d-1)}{2}$ elements of z , respectively:

$$\Phi_d^{BF}(z)[D] := H(E_{\rightarrow} + E_{\rightarrow}^{\top}) \odot H(\text{grad}(p))$$

$$\Phi_d^{BF}(z)[B] := H(E_{\leftrightarrow} + E_{\leftrightarrow}^{\top}) \odot \Psi(D)$$

$$\Phi_d^{AN}(z)[D] := H(E_{\rightarrow} + E_{\rightarrow}^{\top}) \odot H(\text{grad}(p))$$

$$\Phi_d^{AN}(z)[B] := H(E_{\leftrightarrow} + E_{\leftrightarrow}^{\top}) \odot \Psi(D^+)$$

where $\Psi(M) := (I - M) \odot (I - M^{\top})$,

$H(x) := \begin{cases} 1, & \text{if } x > 0 \\ 0, & \text{if } x \leq 0 \end{cases}$ is the Heaviside step

function, \odot is the element-wise product,
 $\text{grad}(u)_{ij} := u_j - u_i$, and A^+ is the transitive closure of A .

Properties

- ❶ Φ_d^{BF} has time complexity $\mathcal{O}(d^2)$
- ❷ Φ_d^{AN} has time complexity $\mathcal{O}(d^{2.8})$
- ❸ Automatic acyclicity
- ❹ No differentiable constraints needed
- ❺ Surjective
- ❻ Scale and Translation Invariance

PPO

- 1 Deep RL algorithm
- 2 Stochastic Gradient Descent based optimizer
- 3 Tries to find optimal policy π_θ which maximises $J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [\sum_{t=0}^{\infty} \gamma^t r_t]$
- 4 Policy updates are limited to trust region to add stability to training
- 5 $\pi_\theta(z) = \mathcal{N}(z; \mu_\theta, \text{diag}(\sigma_\theta^2))$
- 6 $R(z) = -\frac{1}{n} \text{BIC}(X, \Phi_d(z))$
- 7 We do entropy annealing to help escape local minima
- 8 One-step environment
- 9 Designed to function well in high dimensional spaces
- 10 On-policy algorithm, sample in-efficient

CMA-ES

- 1 Evolutionary algorithm
- 2 Stochastic derivative-free black-box optimizer
- 3 Iterative algorithm: sample, rank, update loop
- 4 Tries to find samples from the search space that maximise an objective function
- 5 We restrict covariance matrix to be diagonal for scalability
- 6 Objective fn
$$f_d(z) = \text{BIC}(X, \Phi_d(z)) - \gamma \Gamma(z)$$
$$\Gamma(z) = \sum_{i < j}^{\{i,j\} < d} \min(|z_i - z_j|, \delta) + \sum_{k, k \geq d} \min(|z_k|, \delta)$$
- 7 Designed to function even in non-convex or ill-behaved landscapes
- 8 Sample efficient

GFCI

- 1 Hybrid Algorithm: Has constraint based and score based phases
- 2 Only outputs a PAG.
- 3 Used as a baseline for comparison.

DCD

- 1 Differentiable constraints for acyclicity and to restrict search to bow-free/arid/ancestral graph classes
- 2 Uses modified RICF algorithm
- 3 Uses augmented Lagrangian to obtain unconstrained optimization problem
- 4 Solves the optimization problem with dual descent

Synthetic Data Generation

- 1 Uses modified version of Erdős-Rényi random graph generation model
- 2 Modification accounts for existence of bidirected edges, requirement for bow-free/ancestral graphs
- 3 Inputs: average degree of graph skeleton $\bar{\rho}$, fraction of directed edges f^{\rightarrow}
- 4 Guarantees: $|\hat{E}| = \bar{\rho}d/2 \pm 5$ and $\widehat{f^{\rightarrow}} = f^{\rightarrow} \pm 0.1$

Table: Θ is the structural coefficients matrix and Ω error covariances matrix

Matrix	Distribution
Θ	$\mathcal{U}(\pm[0.5, 2])$
Off-diag Ω	$\mathcal{U}(\pm[0.4, 0.7])$
Diag Ω	$\mathcal{U}([0.7, 1.2] + \sum(\Omega_{i,-i})$

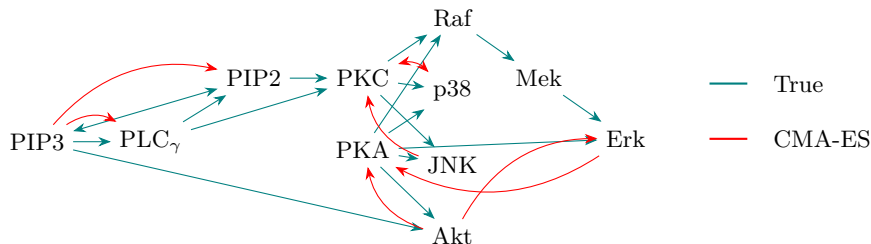
Table: Performance of Algorithms on Sachs Dataset

	SHD (\downarrow)	$ E \cap \hat{E} / \hat{E} ^1(\uparrow)$	PAG F_1 ³ (\uparrow)	$\tau^2(\downarrow)$
DCD	53	3 / 43	0.47	249.5
CMA-ES	22	1 / 8	0.58	105.7
Relcadilac	23	1 / 9	0.48	1029.2

¹ $|E \cap \hat{E}|$ is the number of correct predicted edges, and $|\hat{E}|$ is the total number of predicted edges.

² τ is the runtime of the algorithm in seconds.

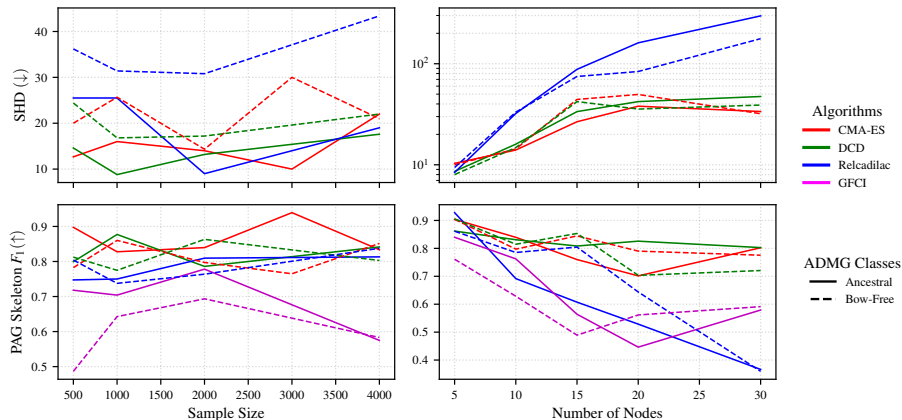
³ The F_1 score is computed on the skeleton of the PAG.



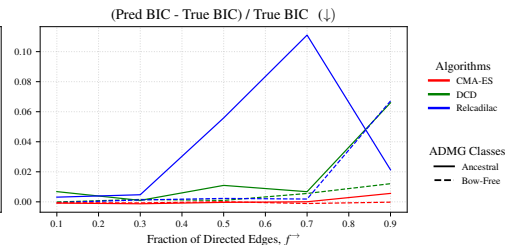
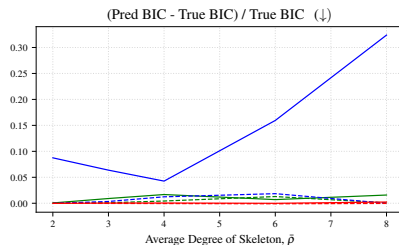
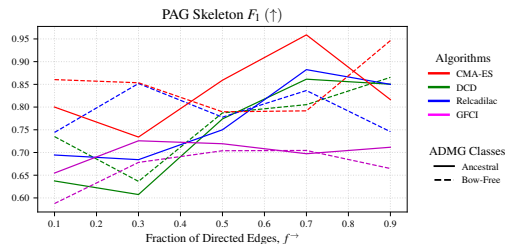
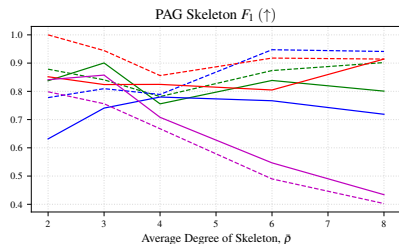
Performance Plots I

Varying Parameters: Sample Size, Number of Nodes, Average Degree of Skeleton, Fraction of Directed Edges

Metrics Captured: Structural Hamming Distance (lower is better), F_1 score of the edges of the PAG Skeleton (higher is better), Fractional Excess BIC score (lower is better)



Performance Plots II



- 1 Explore the SAC RL algorithm for potentially better sample efficiency than PPO
- 2 RICE algorithm bi-directed connected components-based decomposition and caching for potential speedup in sparse graphs
- 3 Extend the Vec2ADMG mappings to include other graph classes like arid graphs