

# Bayesian model selection allows for less restrictive functional assumptions when doing multivariate causal discovery on observational data

## Bayesian Model Selection for Multivariate Causal Discovery

Anish Dhir<sup>\*1</sup>, Ruby Sedgwick<sup>\*1,2</sup>, Avinash Kori<sup>1</sup>, Ben Glocker<sup>1</sup>, Mark van der Wilk<sup>3</sup>

<sup>1</sup> Department of Computing, Imperial College London, UK    <sup>\*</sup>equal contribution  
<sup>2</sup> Xyme, Oxford, UK  
<sup>3</sup> University of Oxford, UK

paper



### Introduction

- Current causal discovery approaches require **restrictive model assumptions** to ensure identifiability, e.g. additive noise (ANM). In real-world data these restrictive assumptions are commonly violated, **losing identifiability guarantees**.

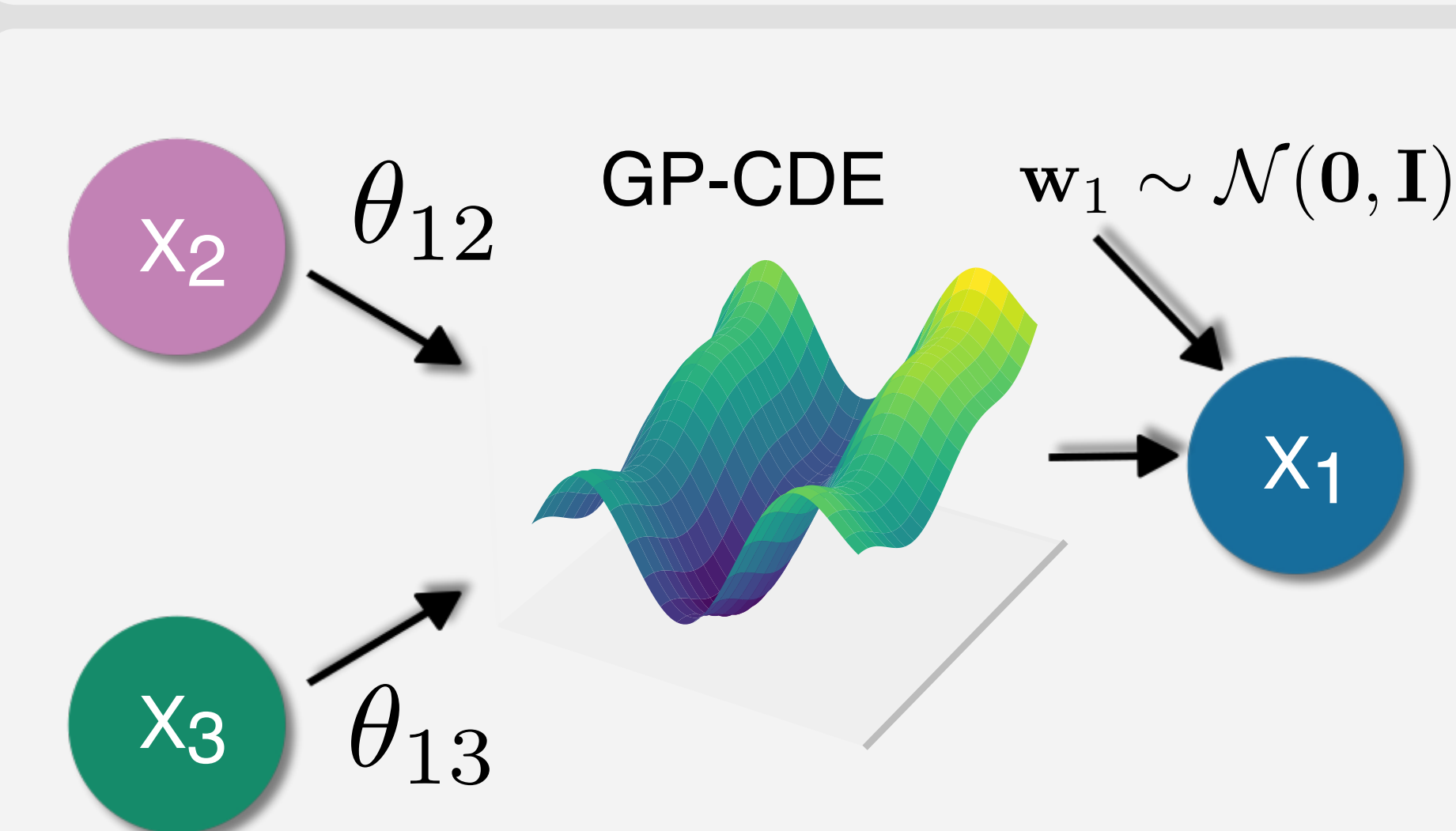
- Bayesian model selection has been proven to improve performance in the bivariate case by allowing **more flexible functional assumptions**.

- However, the naive discrete Bayesian model selection approach isn't feasible for the multivariate case where the number of possible graphs scales super-exponentially.

- We propose a **continuous Bayesian model selection** approach that **scales well to large numbers of variables** while still allowing more flexible functional assumptions.

### Causal Gaussian Process Conditional Estimator (CGP-CDE)

Aim: given a set of variables, find the causal graph.  $X_1$   $X_2$   $X_3$

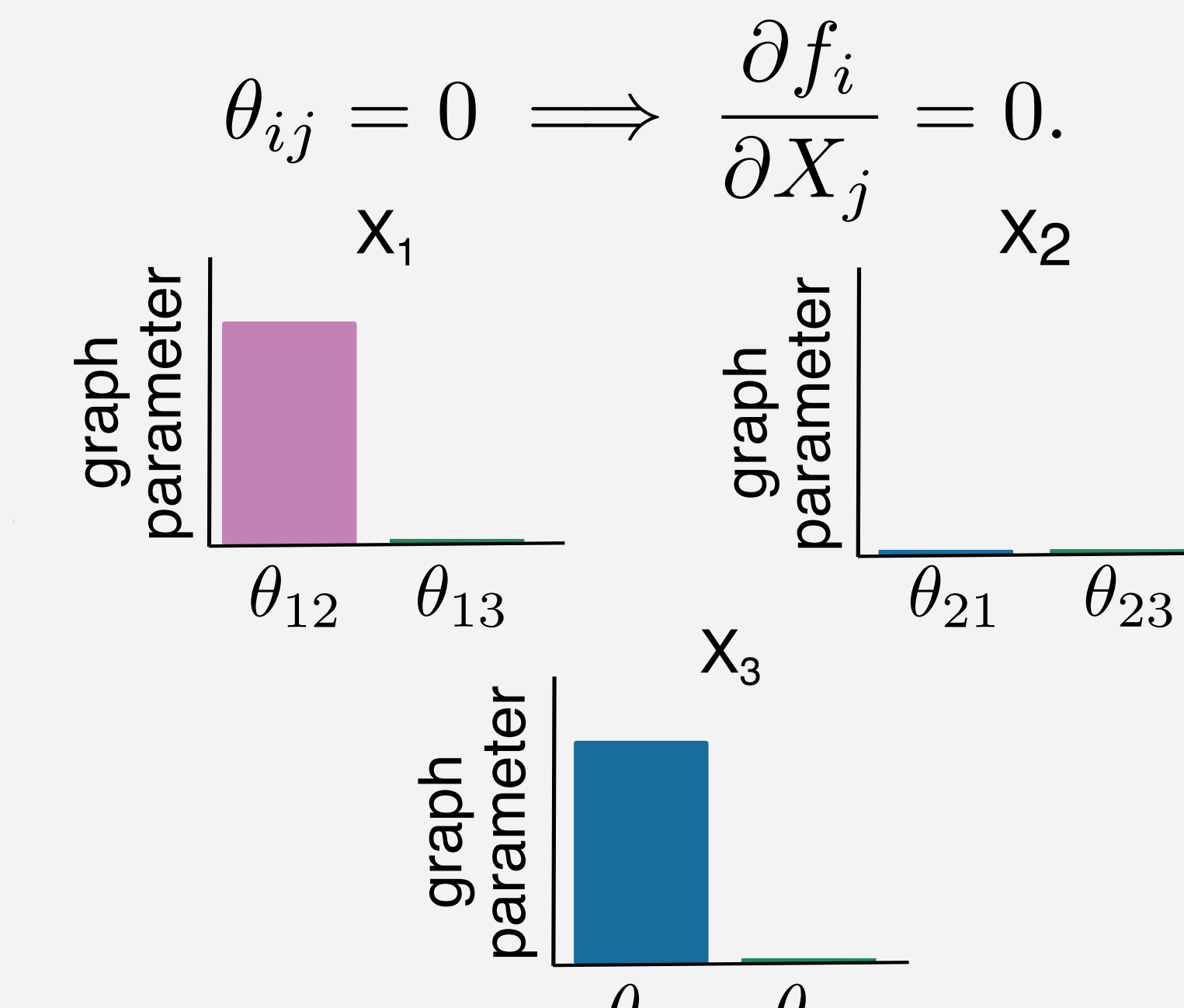


Fit a Gaussian process conditional estimator (GP-CDE) for each variable. Latent variable  $w_1$  allows non-Gaussian and heteroskedastic densities.

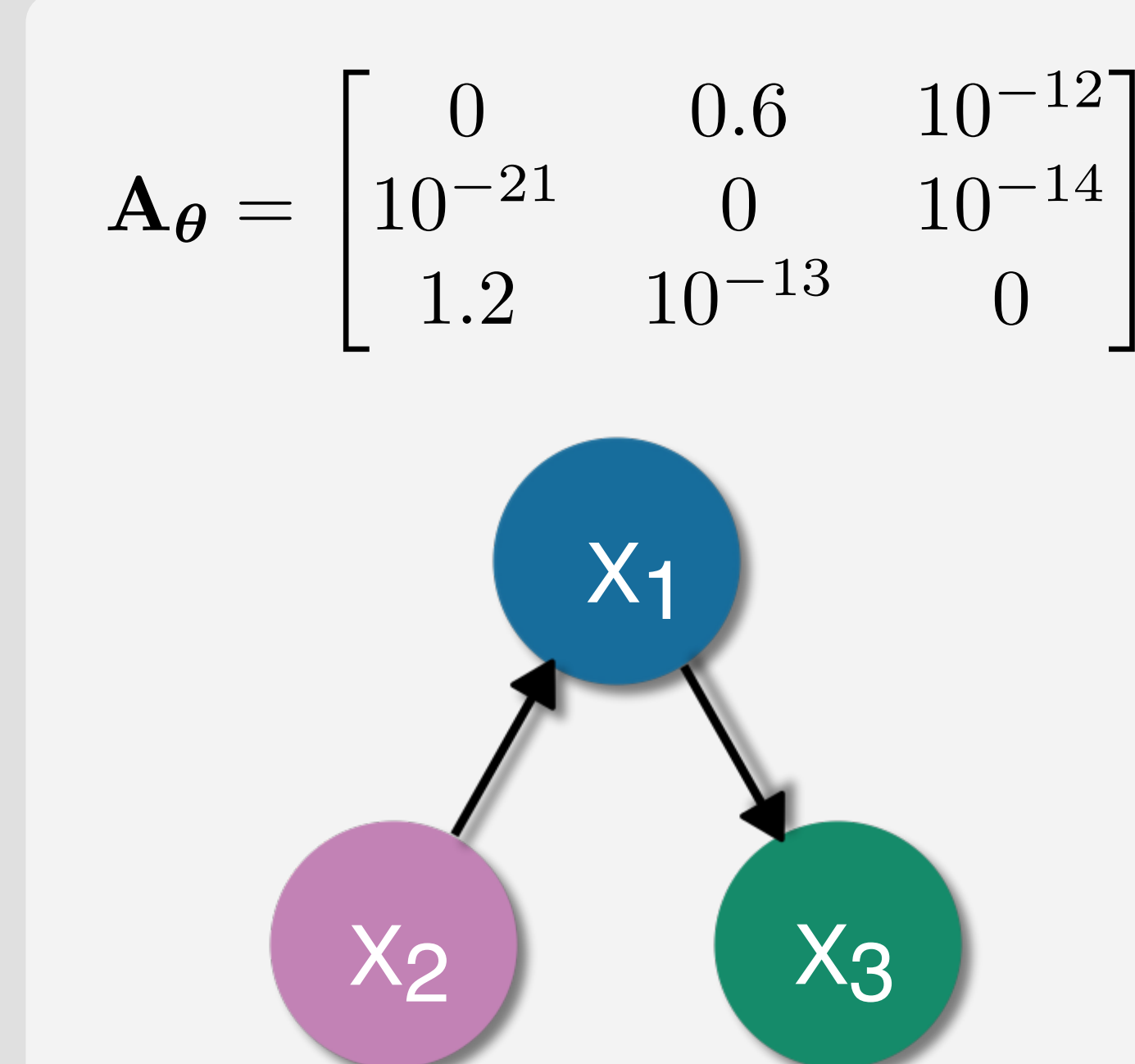
Loss:  $\mathcal{L}_{\text{ELBO}}(q, \theta, \sigma, \phi) + \log p(\theta) - \gamma_t h(\mathbf{A}_\theta)$   
marginal likelihood    acyclicity constraint

$$A_{ij} = \begin{cases} \theta_{ij} & \text{if } j \neq i, \\ 0 & \text{otherwise.} \end{cases}$$

An adjacency matrix is constructed from the graph parameters (inverse lengthscales).



The ARD properties of the kernel mean edges that aren't evidenced by the data are "switched off" by making the graph parameter small.

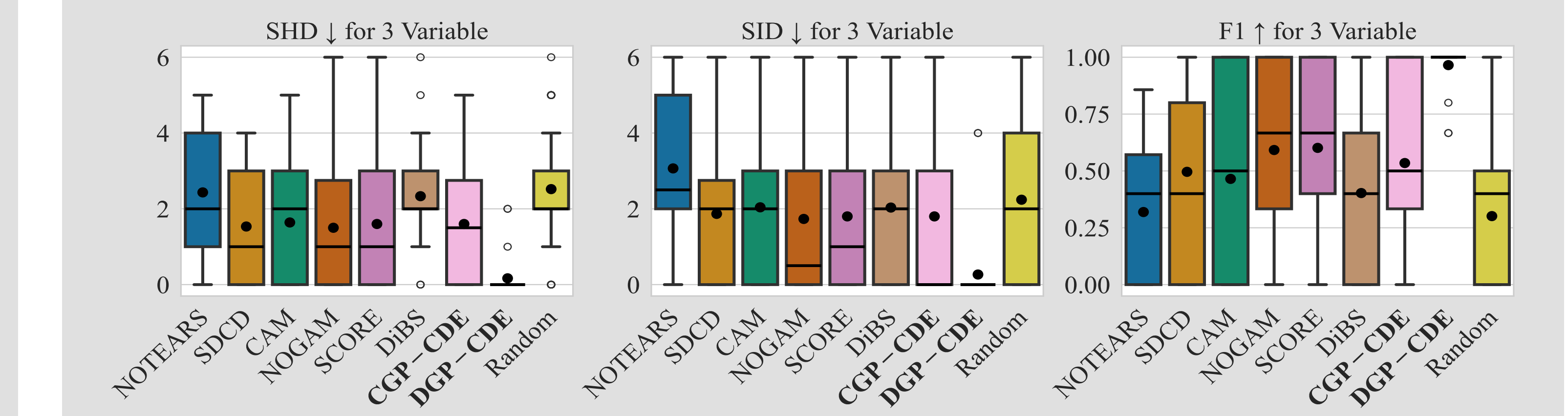


The causal graph is extracted from the graph parameters via thresholding.

### How bad is the continuous approximation?

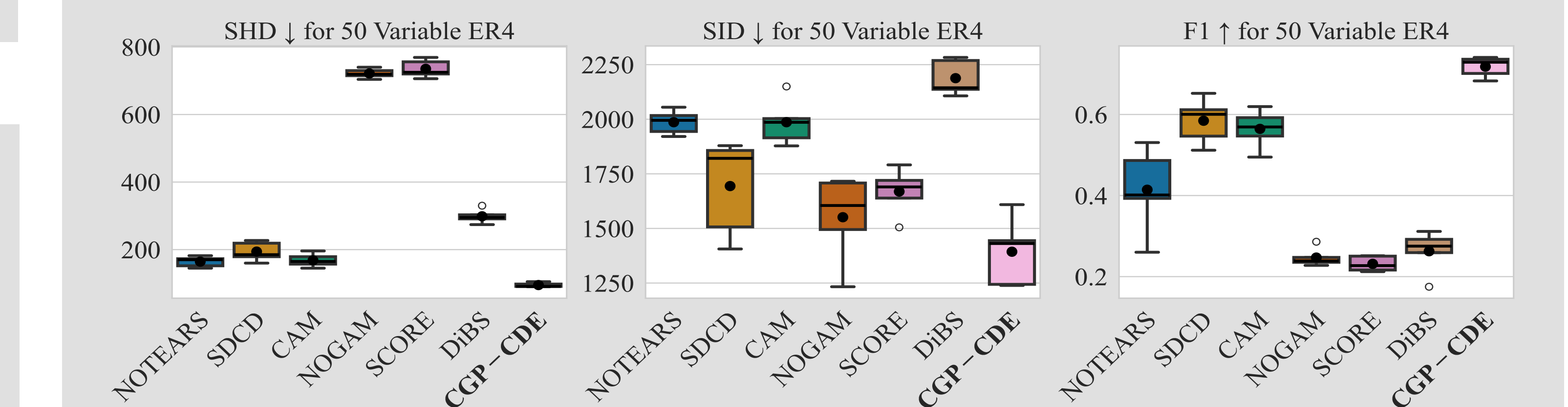
- Excellent performance of discrete Bayesian model selection (DGP-CDE) shows the **probability of error is small**.

- The difference in performance between CGP-CDE and DGP-CDE is due to the continuous relaxation in CGP-CDE.



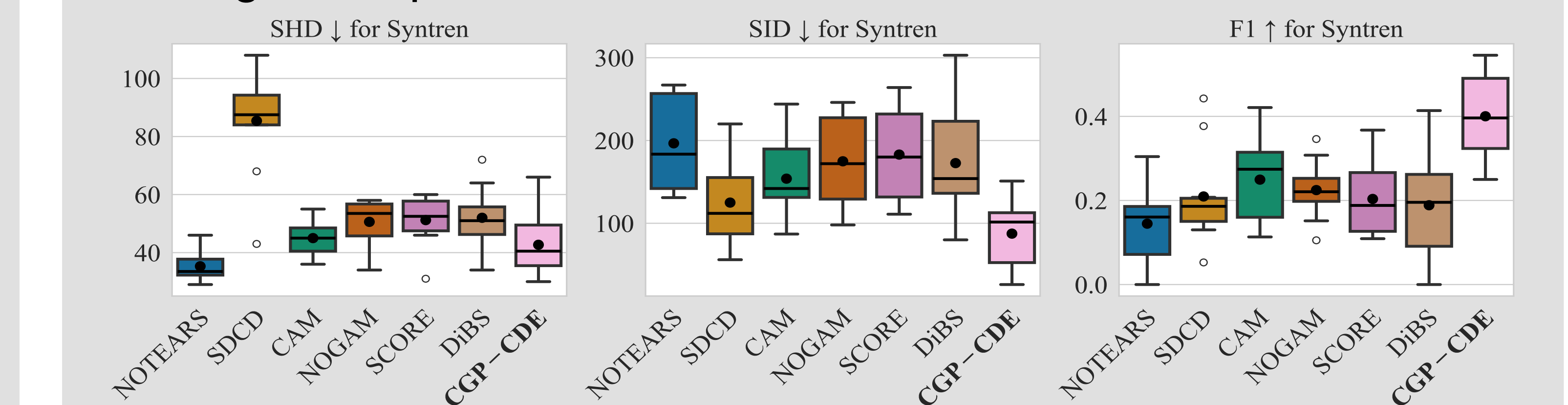
### Does it scale to large numbers of variables?

- Experiments on ER graphs with 4 expected edges per variable and 50 variables. The CGP-CDE outperforms the other models.



### How well does it perform on semi-synthetic data?

- Experiments on Syntren dataset, derived from a gene regulatory network simulator which has 10 datasets of 20 nodes. The CGP-CDE again outperforms the other models.



### Bayesian model selection solves this problem

$$\begin{aligned} P(\mathcal{M}_G|X) &\propto \underbrace{P(X|\mathcal{M}_G)}_{\text{marginal likelihood}} \underbrace{P(\mathcal{M}_G)}_{\text{prior}} \\ P(X|\mathcal{M}_G) &= \int \prod_{i=1}^D (P_i(X_i|X_{\text{PA}_{\mathcal{G}}(i)}) \pi_i(dP_i)) \end{aligned}$$

Causal model:  $\mathcal{M}_G$     Data:  $X \in \mathbb{R}^D$

- Factorised prior on distributions encodes the **independent mechanism (ICM) assumption**.

- The **marginal likelihood** will prefer models whose ICM assumption aligns with the properties of the data generating process.

### Probability of error

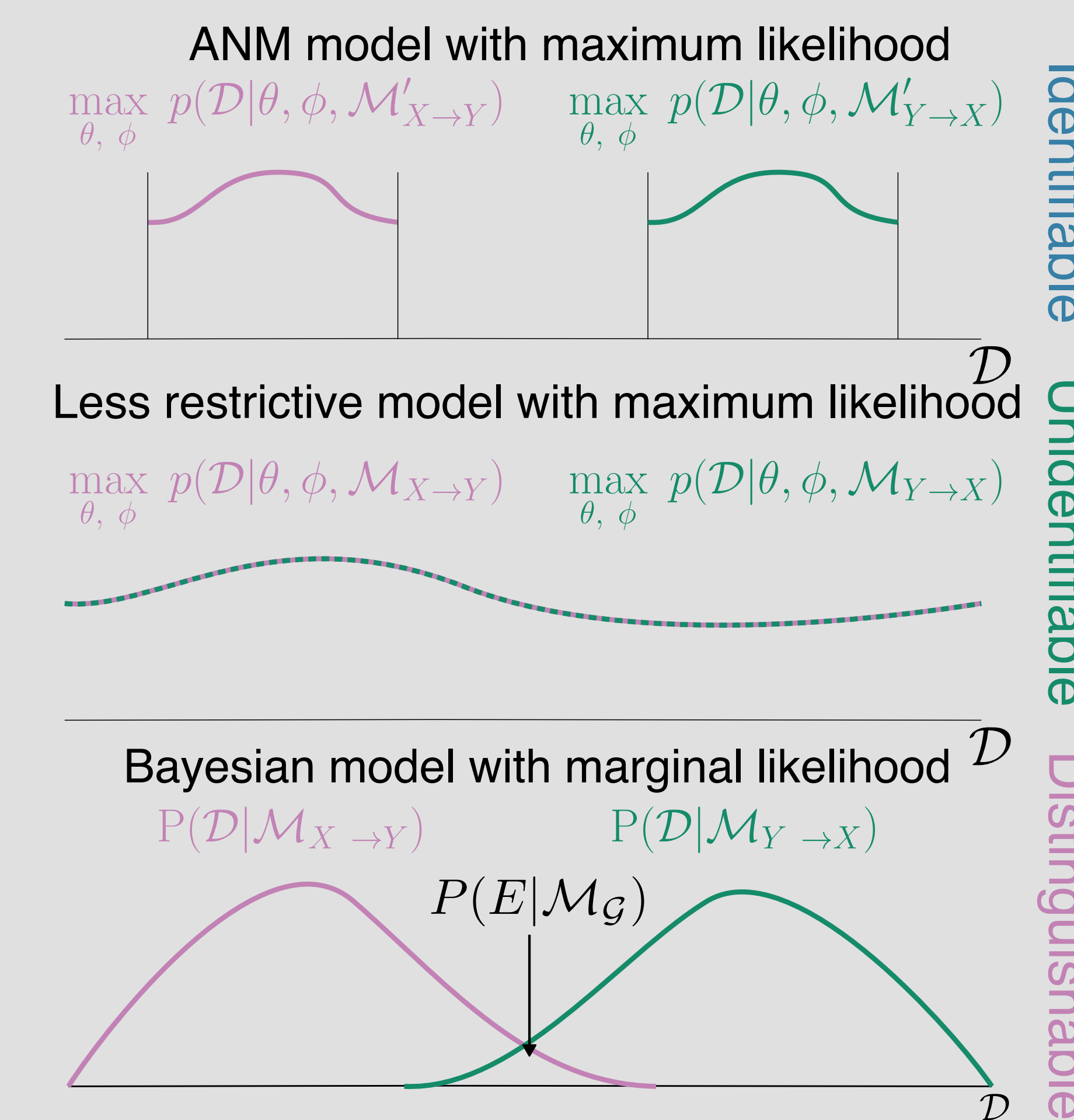
$$P(E|\mathcal{M}_G) = \int_{\mathcal{R}} p(X|\mathcal{M}_G) dX, \quad \mathcal{R} = \{X : \underbrace{S(X|\mathcal{M}_{\mathcal{H}})}_{\text{score}} > S(X|\mathcal{M}_G) \text{ for } \mathcal{H} \neq \mathcal{G}\}$$

Identifiable:  $P(E|\mathcal{M}_G) = 0$

Unidentifiable:  $P(E|\mathcal{M}_G) = P(E|U(\mathcal{G}))$

Distinguishable:  $0 < P(E|\mathcal{M}_G) < P(E|U(\mathcal{G}))$

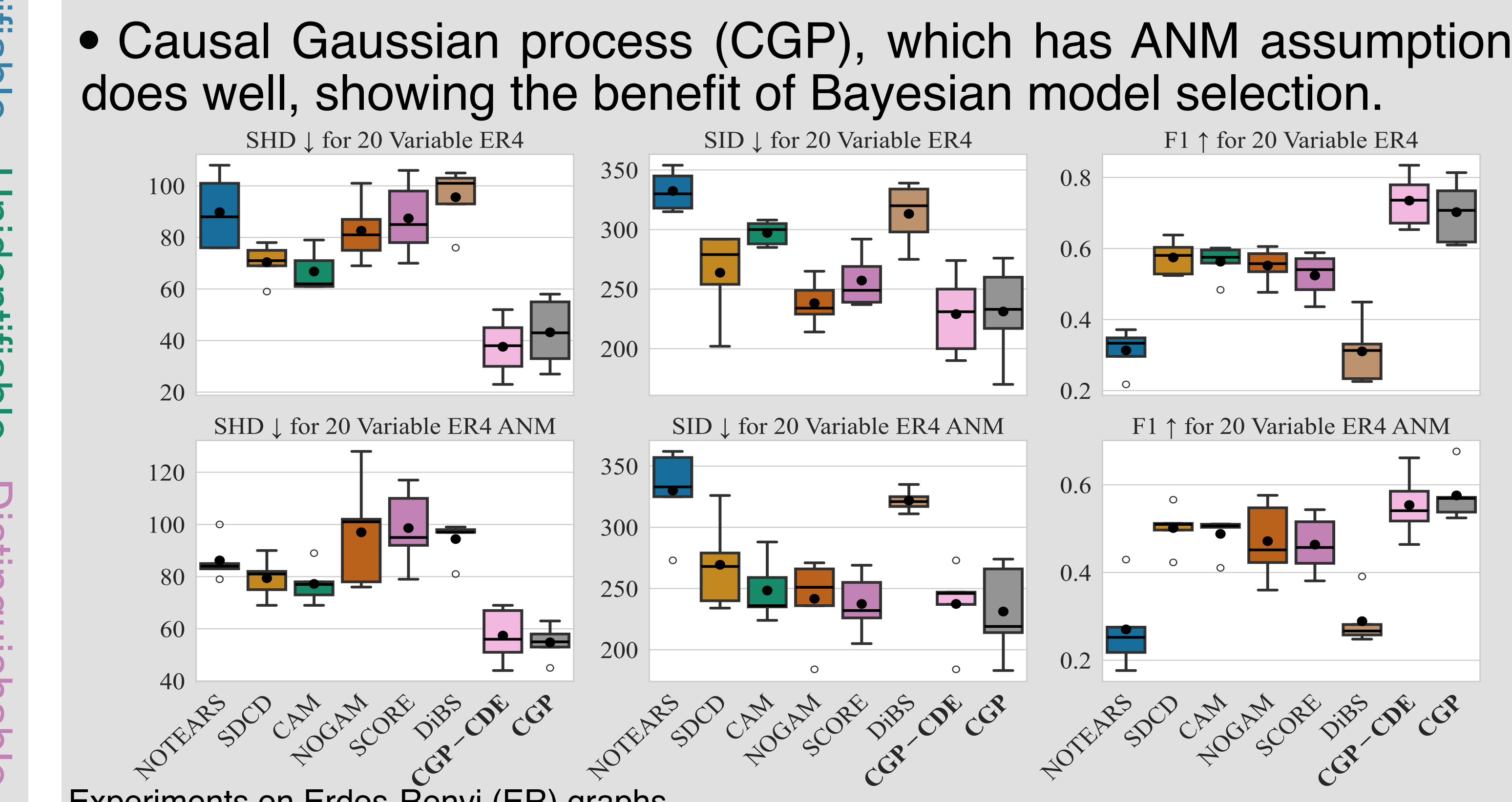
uniformly random selected graph



Figures from Dhir et al. 2024

### Is performance on identifiable data worse?

- CGP-CDE outperforms ANM methods on ANM and non-ANM data.



Experiments on Erdos-Renyi (ER) graphs.