

BIVARIATE CAUSAL DISCOVERY USING BAYESIAN MODEL SELECTION

Mark van der Wilk
ISBA World Meeting 2024



Department of
COMPUTER
SCIENCE

 <https://mvdw.uk>
 @markvanderwilk

3 July 2024

Based on a True Story (ICML 2024)

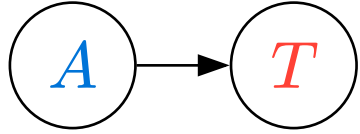
Bivariate Causal Discovery using Bayesian Model Selection

Anish Dhir¹ Samuel Power² Mark van der Wilk³

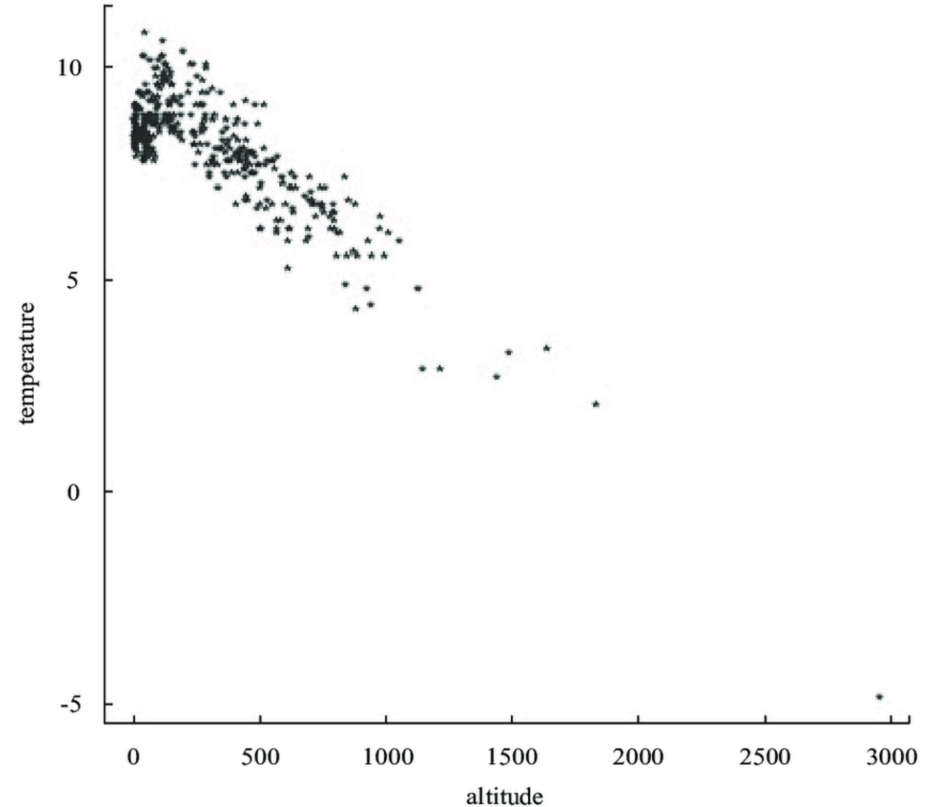


Background & Problem Setting

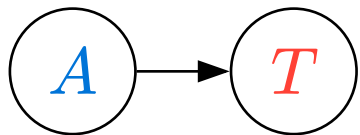
Why is Causality Important?



- Conditional $p(a|t)$ assumes pair (a, t) sampled *jointly* from the same distribution!
- When intervening, **you cannot affect your cause!**
- Intervention breaks links to ancestors, so $p(a|\text{do}(t)) = p(a)$.
- But... $p(t|\text{do}(a)) = p(t|a)$.



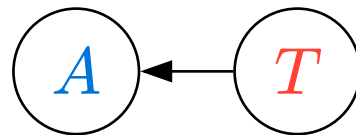
Two Causal Directions, Two Models



$$p(\textcolor{red}{t}, \textcolor{blue}{a} | \varphi, \theta) = p(\textcolor{blue}{a} | \textcolor{red}{t}, \theta) p(\textcolor{red}{t} | \varphi)$$

$p(\textcolor{blue}{a} | \textcolor{red}{t}, \theta)$: cond. density model

$p(\textcolor{red}{t} | \varphi)$: density model

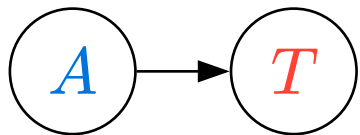


$$p(\textcolor{red}{t}, \textcolor{blue}{a} | \varphi, \theta) = p(\textcolor{red}{t} | \textcolor{blue}{a}, \theta) p(\textcolor{blue}{a} | \varphi)$$

$p(\textcolor{red}{t} | \textcolor{blue}{a}, \theta)$: cond. density model

$p(\textcolor{blue}{a} | \varphi)$: density model

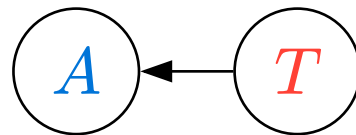
Two Causal Directions, Two Models



$$p(\textcolor{red}{t}, \textcolor{blue}{a} | \varphi, \theta) = p(\textcolor{blue}{a} | \textcolor{red}{t}, \theta) p(\textcolor{red}{t} | \varphi)$$

$p(\textcolor{blue}{a} | \textcolor{red}{t}, \theta)$: cond. density model

$p(\textcolor{red}{t} | \varphi)$: density model



$$p(\textcolor{red}{t}, \textcolor{blue}{a} | \varphi, \theta) = p(\textcolor{red}{t} | \textcolor{blue}{a}, \theta) p(\textcolor{blue}{a} | \varphi)$$

$p(\textcolor{red}{t} | \textcolor{blue}{a}, \theta)$: cond. density model

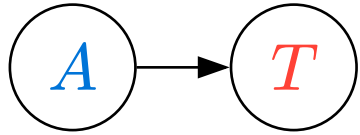
$p(\textcolor{blue}{a} | \varphi)$: density model



Causality should determine model structure

We want sensible results if we apply intervention rules to our *model*!

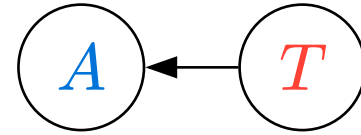
Two Causal Directions, Two Models



$$p(\textcolor{red}{t}, \textcolor{blue}{a} | \varphi, \theta) = p(\textcolor{blue}{a} | \textcolor{red}{t}, \theta) p(\textcolor{red}{t} | \varphi)$$

$p(\textcolor{blue}{a} | \textcolor{red}{t}, \theta)$: cond. density model

$p(\textcolor{red}{t} | \varphi)$: density model



$$p(\textcolor{red}{t}, \textcolor{blue}{a} | \varphi, \theta) = p(\textcolor{red}{t} | \textcolor{blue}{a}, \theta) p(\textcolor{blue}{a} | \varphi)$$

$p(\textcolor{red}{t} | \textcolor{blue}{a}, \theta)$: cond. density model

$p(\textcolor{blue}{a} | \varphi)$: density model



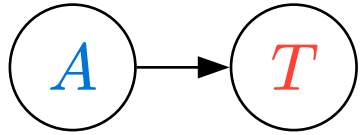
Causality should determine model structure

We want sensible results if we apply intervention rules to our *model*!



Goal: Predict causal structure from observational data.

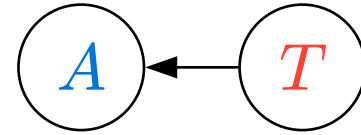
Two Causal Directions, Two Models



$$p(\textcolor{red}{t}, \textcolor{blue}{a} | \varphi, \theta) = p(\textcolor{blue}{a} | \textcolor{red}{t}, \theta) p(\textcolor{red}{t} | \varphi)$$

$p(\textcolor{blue}{a} | \textcolor{red}{t}, \theta)$: cond. density model

$p(\textcolor{red}{t} | \varphi)$: density model



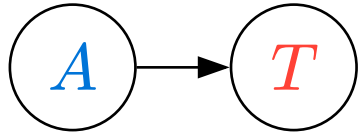
$$p(\textcolor{red}{t}, \textcolor{blue}{a} | \varphi, \theta) = p(\textcolor{red}{t} | \textcolor{blue}{a}, \theta) p(\textcolor{blue}{a} | \varphi)$$

$p(\textcolor{red}{t} | \textcolor{blue}{a}, \theta)$: cond. density model

$p(\textcolor{blue}{a} | \varphi)$: density model

? Could try to fit φ, θ with maximum likelihood..?

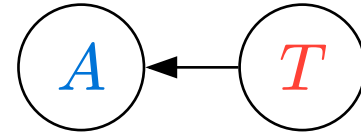
Two Causal Directions, Two Models



$$p(\textcolor{red}{t}, \textcolor{blue}{a} | \varphi, \theta) = p(\textcolor{blue}{a} | \textcolor{red}{t}, \theta) p(\textcolor{red}{t} | \varphi)$$

$p(\textcolor{blue}{a} | \textcolor{red}{t}, \theta)$: cond. density model

$p(\textcolor{red}{t} | \varphi)$: density model



$$p(\textcolor{red}{t}, \textcolor{blue}{a} | \varphi, \theta) = p(\textcolor{red}{t} | \textcolor{blue}{a}, \theta) p(\textcolor{blue}{a} | \varphi)$$

$p(\textcolor{red}{t} | \textcolor{blue}{a}, \theta)$: cond. density model

$p(\textcolor{blue}{a} | \varphi)$: density model

? Could try to fit φ, θ with maximum likelihood..?

🚧 For *flexible* models, both directions give equally good fit! 🤔

Both models are in the same *Markov Equivalence Class*.

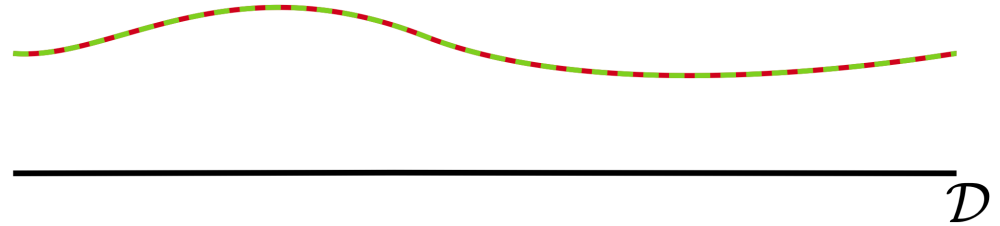
Approach: Restricted Model Classes



For *flexible* models, both directions give equally good fit! 🤔

$$\max_{\theta, \phi} p(\mathcal{D}|\theta, \phi, \mathcal{M}_{X \rightarrow Y})$$

$$\max_{\theta, \phi} p(\mathcal{D}|\theta, \phi, \mathcal{M}_{Y \rightarrow X})$$

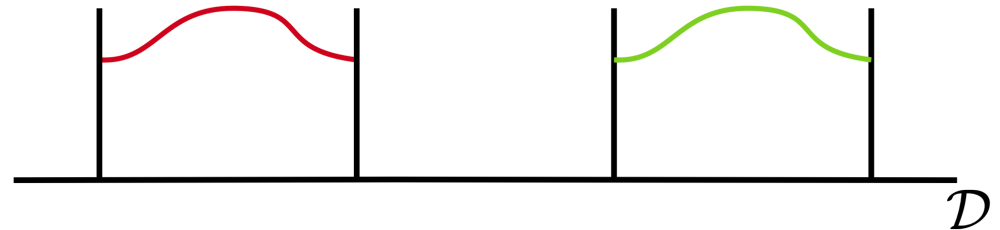


Add restrictions, e.g. ANM

effect = $f(\text{cause}) + \text{noise}$

$$\max_{\theta, \phi} p(\mathcal{D}|\theta, \phi, \mathcal{M}'_{X \rightarrow Y})$$

$$\max_{\theta, \phi} p(\mathcal{D}|\theta, \phi, \mathcal{M}'_{Y \rightarrow X})$$



⇒ Non-overlapping data support

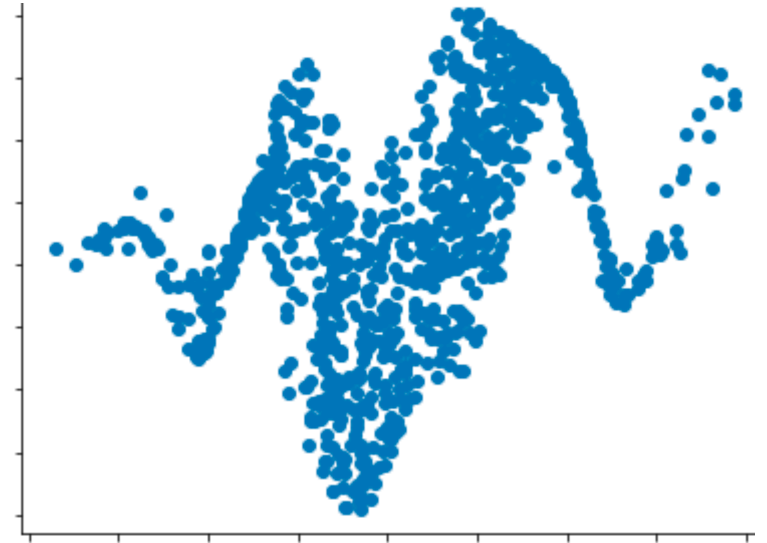
⇒ So... identifiable! (as $N \rightarrow \infty$)

Problem: Restricted Model Classes

But what to do with a dataset like this one?

- Outside datasets covered by ANM!
- Poor fit \Rightarrow bad predictions.
- Loss of identifiability guarantees.

To model realistic datasets, we **want** our model to have support over all datasets!

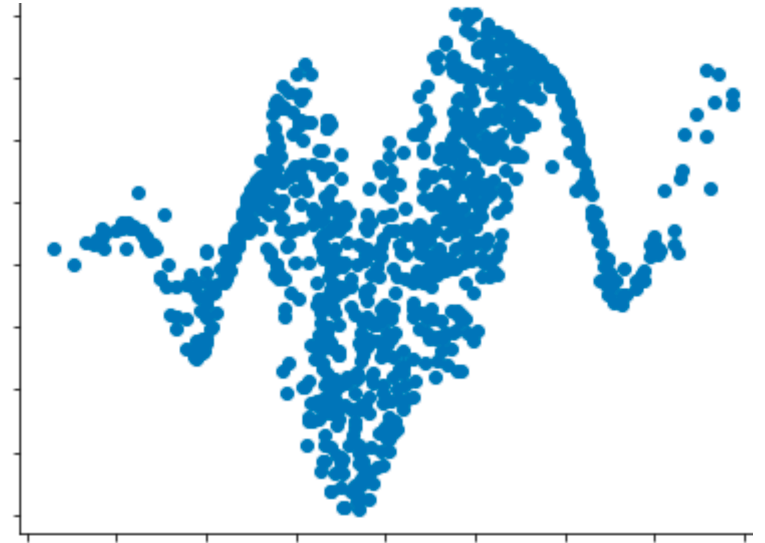


Problem: Restricted Model Classes

But what to do with a dataset like this one?

- Outside datasets covered by ANM!
- Poor fit \Rightarrow bad predictions.
- Loss of identifiability guarantees.

To model realistic datasets, we **want** our model to have support over all datasets!

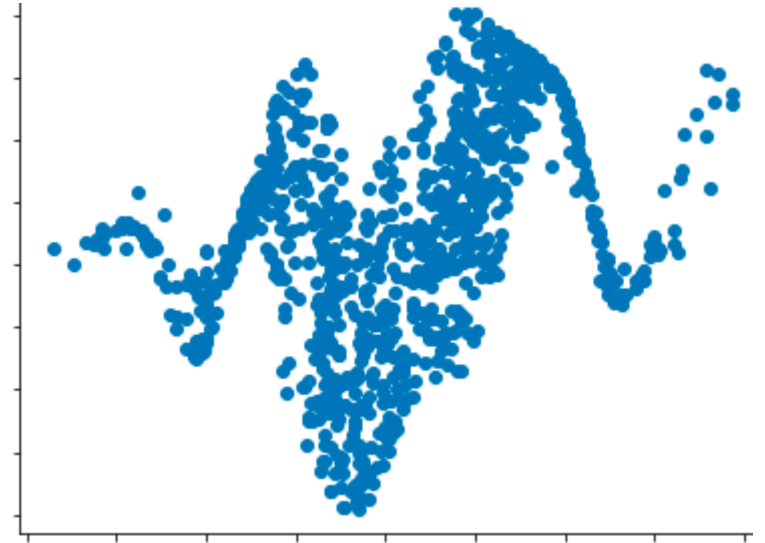


Problem: Restricted Model Classes

But what to do with a dataset like this one?

- Outside datasets covered by ANM!
- Poor fit \Rightarrow bad predictions.
- Loss of identifiability guarantees.

To model realistic datasets, we **want** our model to have support over all datasets!

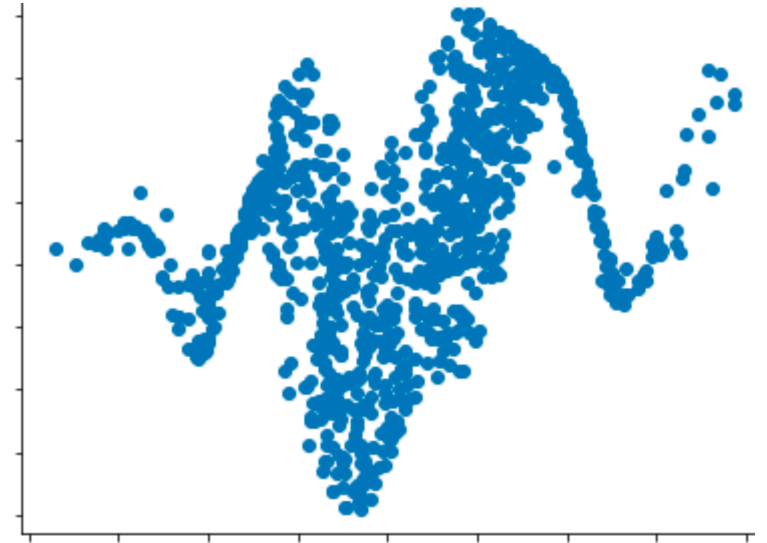


Problem: Restricted Model Classes

But what to do with a dataset like this one?

- Outside datasets covered by ANM!
- Poor fit \Rightarrow bad predictions.
- Loss of identifiability guarantees.

To model realistic datasets, we **want** our model to have support over all datasets!

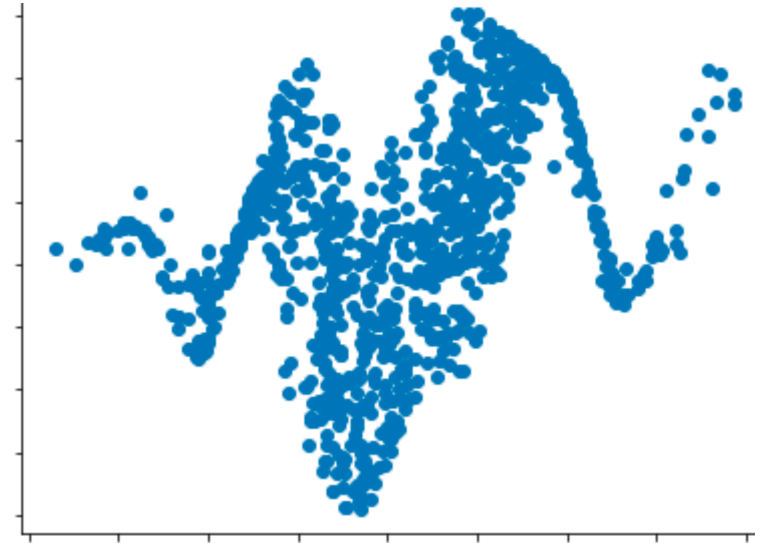


Problem: Restricted Model Classes

But what to do with a dataset like this one?

- Outside datasets covered by ANM!
- Poor fit \Rightarrow bad predictions.
- Loss of identifiability guarantees.

To model realistic datasets, we **want** our model to have support over all datasets!

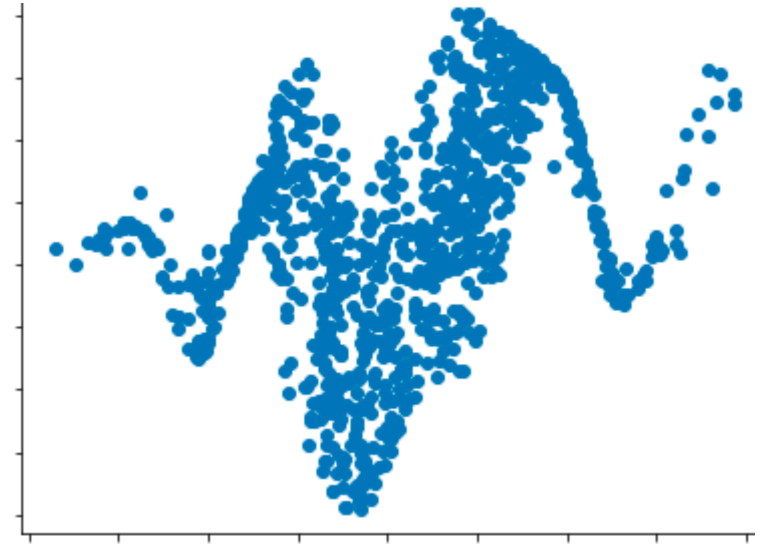



Problem: Restricted Model Classes

But what to do with a dataset like this one?

- Outside datasets covered by ANM!
- Poor fit \Rightarrow bad predictions.
- Loss of identifiability guarantees.

To model realistic datasets, we **want** our model to have support over all datasets!



 **Predict causal structure from observational data with flexible models with realistic assumptions.**

Bayesian Perspective

Model Selection

- We have two models, with different causal assumptions.
- Each model has its own unknown parameters.
- We want to determine which model is appropriate.

 **Is this not just a hierarchical Bayesian inference problem?**

Just find the posterior over the models, using the marginal likelihood:

$$p(\mathcal{M}_{X \rightarrow Y} | \mathbf{x}, \mathbf{y}) \propto p(\mathbf{x}, \mathbf{y} | \mathcal{M}_{X \rightarrow Y}) p(\mathcal{M}_{X \rightarrow Y})$$

$$p(\mathbf{x}, \mathbf{y} | \mathcal{M}_{X \rightarrow Y}) = \iint p(\mathbf{x} | \varphi) p(\mathbf{y} | \mathbf{x}, \theta) p(\varphi, \theta) d\varphi d\theta$$

Has been investigated before, but didn't get it quite right (see paper).

Causal Assumptions in Bayesian Models

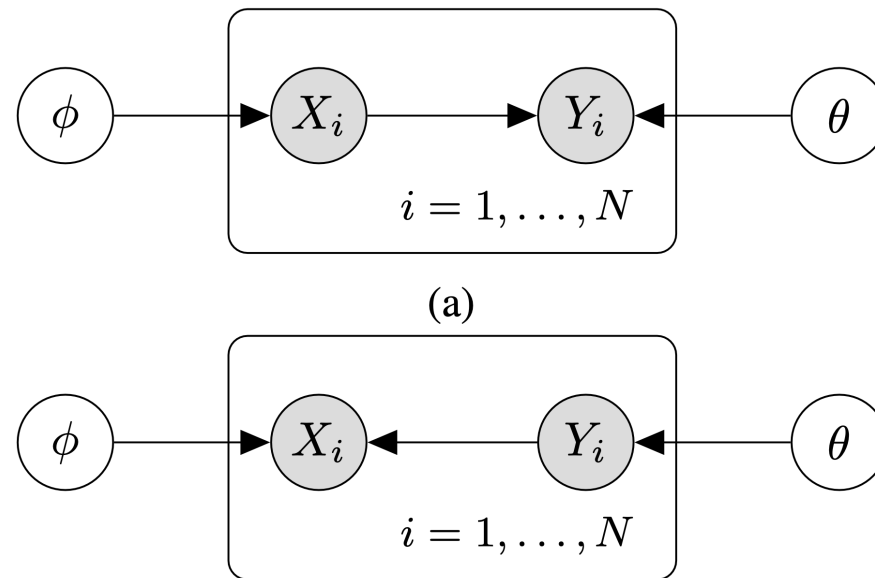
Observational data, so causality enters only through model assumptions.

Symmetry implies that:

- $p(\mathcal{M}_{X \rightarrow Y}) = p(\mathcal{M}_{Y \rightarrow X})$
- We want the same prior on $p(y_i|x_i, \theta, \mathcal{M}_{X \rightarrow Y})$ as on $p(x_i|y_i, \varphi, \mathcal{M}_{Y \rightarrow X})$.
- And similarly for $p(x_i|\varphi, \mathcal{M}_{X \rightarrow Y})$ and $p(y_i|\theta, \mathcal{M}_{Y \rightarrow X})$.

ICM implies independent priors.

Causal direction is encoded in graph.



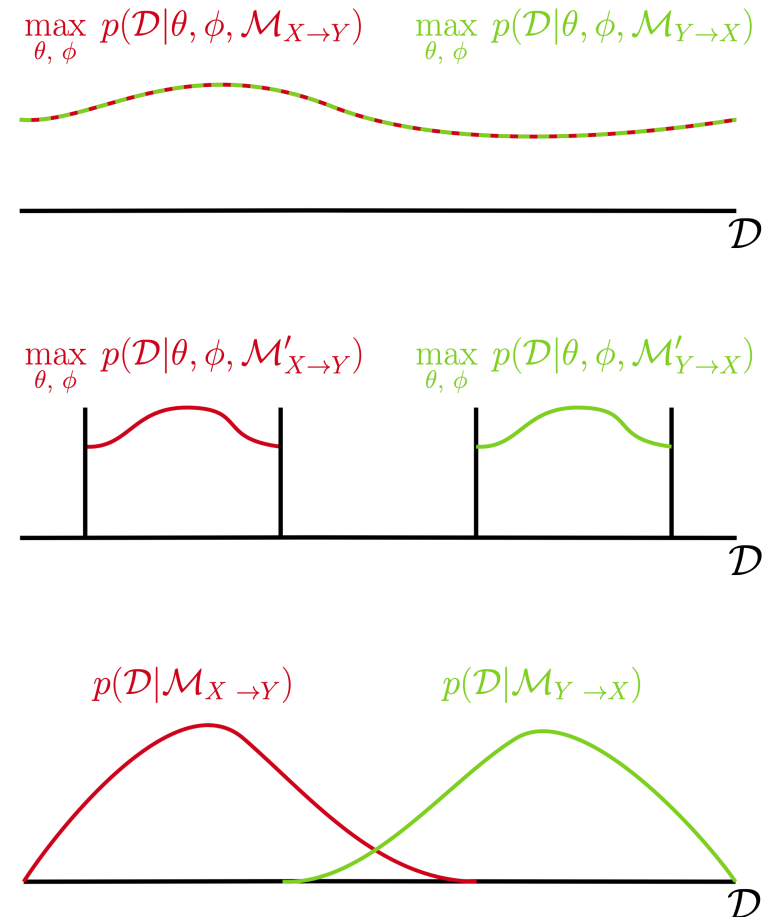
Guarantees (or lack thereof)

- Priors **gives Bayes an opinion** on causal direction, where MaxLik does not.
- Even for flexible models with wide support!
- Price you pay: Overlap in distributions. So no perfect identifiability. Even if data sampled exactly from prior!

$$P(E) = \frac{1}{2}(1 - \text{TV}[P_{\mathcal{D}}(\cdot | \mathcal{M}_{X \rightarrow Y}),$$

$$P_{\mathcal{D}}(\cdot | \mathcal{M}_{Y \rightarrow X})])$$

- Is this so different from existing approach?



Putting this Into Practice

A Practical Model

A conditional GPLVM (Bayesian VAE) for the conditional density:

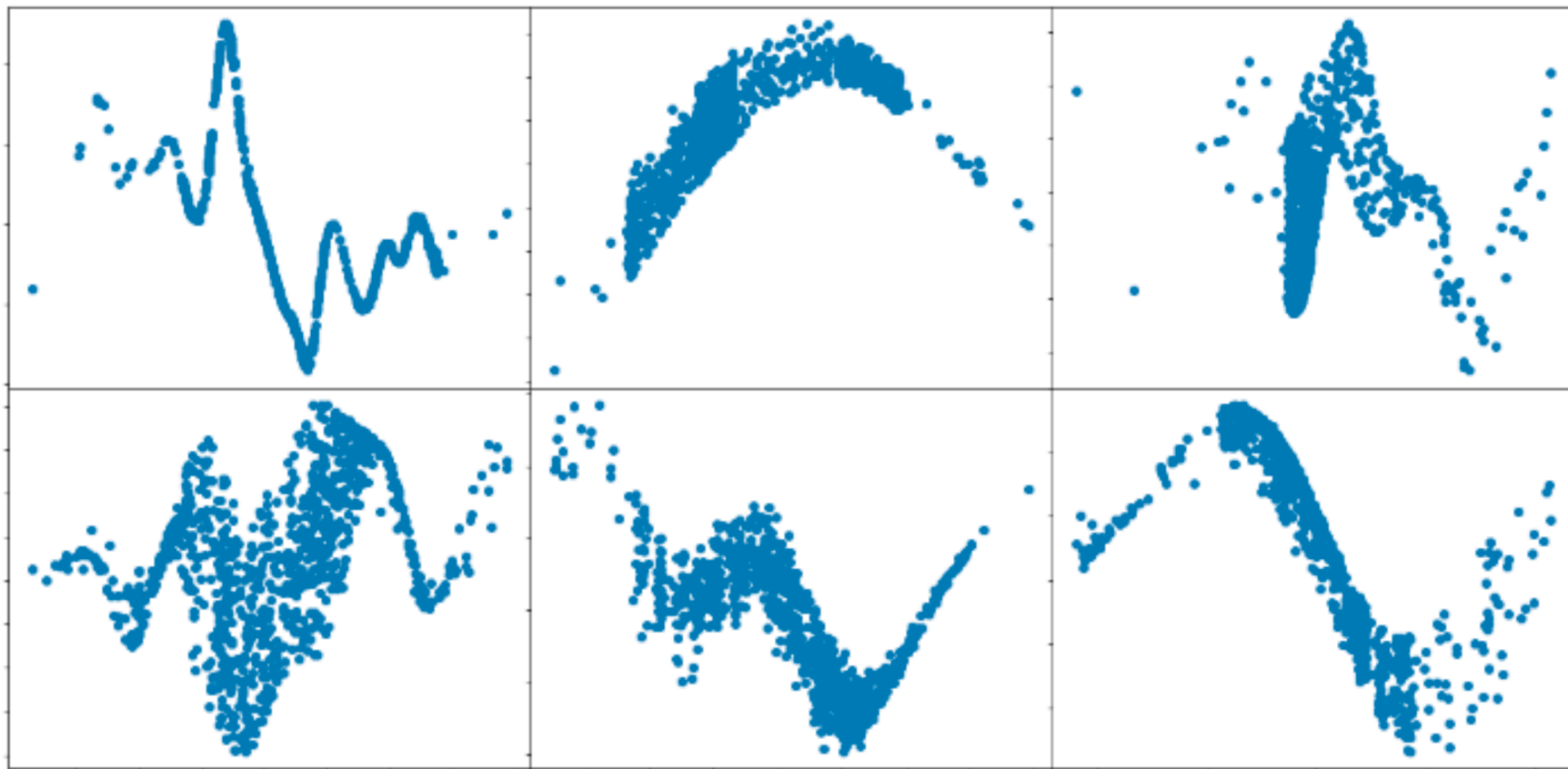
$$p(y_i|x_i, f, \mathcal{M}_{X \rightarrow Y}) = \int \mathcal{N}(y_i; f(x_i, w_i), \sigma^2) \mathcal{N}(w_i) dw_i$$

$$f \sim \mathcal{GP}(0, k)$$

- Flexible (non-parametric) model over many conditional densities.
- Similar GPLVM prior on $p(x_i|g, \mathcal{M}_{X \rightarrow Y})$.
- Relatively standard variational approximation to perform inference.

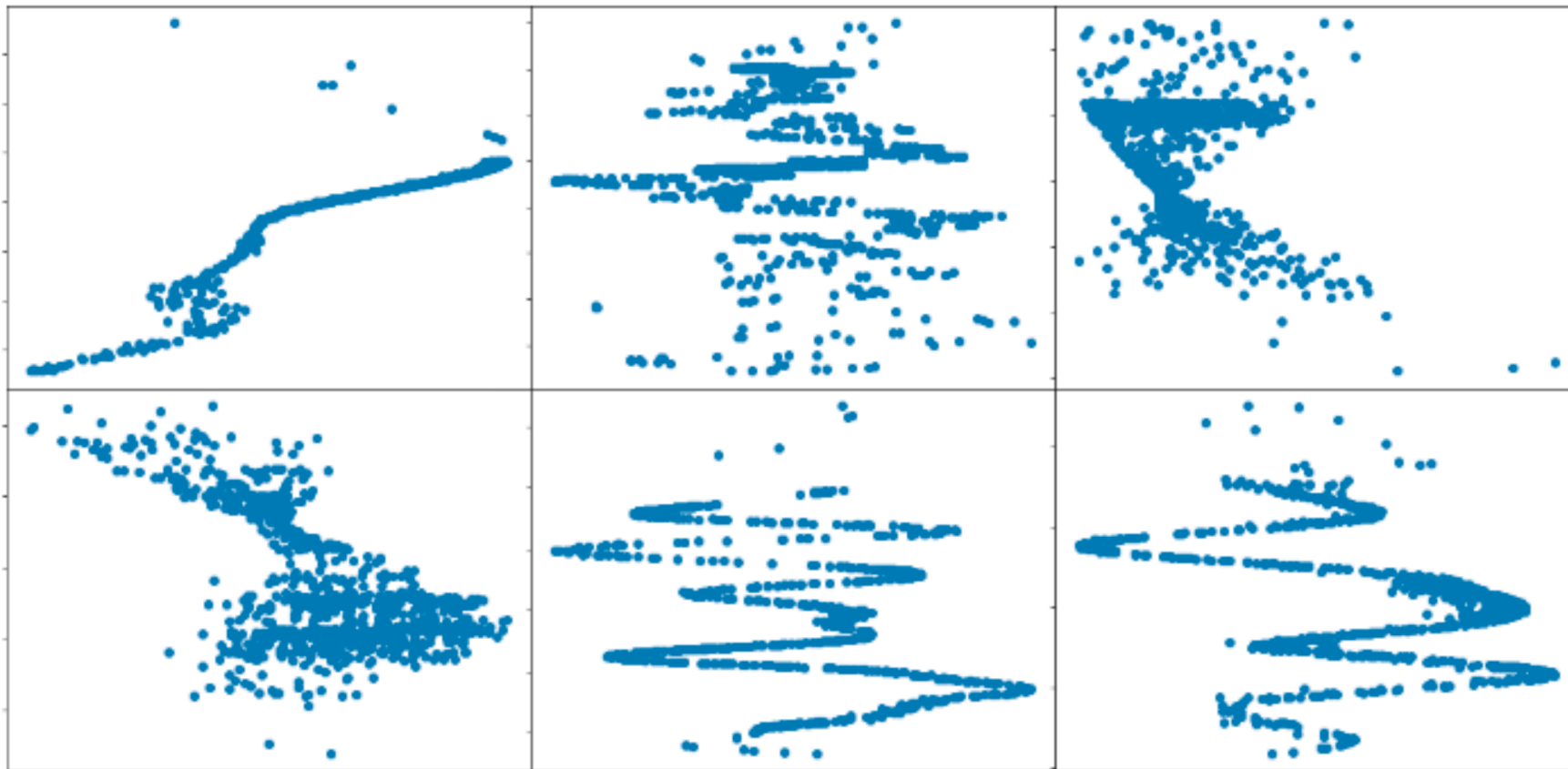
Overlap in Priors

$$\mathcal{M}_{\mathbf{x} \rightarrow \gamma}$$



Overlap in Priors

$$\mathcal{M}_{X \leftarrow Y}$$



Experimental Results


Are our prior assumptions good?

- For identifiable ANM data, GPLVM gets 100% accuracy.
- For real data: Can **only** determine this experimentally, as in other approaches where theoretical assumptions are broken in practice.

| Methods | CE-Cha | CE-Multi | CE-Net | CE-Gauss | CE-Tueb |
|--------------|-------------|-------------|-------------|--------------------|-------------|
| CGNN | <u>76.2</u> | 94.7 | 86.3 | 89.3 | <u>76.6</u> |
| GPI | 71.5 | 73.8 | 88.1 | 90.2 | 70.6 |
| PNL | 78.6 | 51.7 | 75.6 | 84.7 | 73.8 |
| ANM | 43.7 | 25.5 | 87.8 | 90.7 | 63.9 |
| IGCI | 55.6 | 77.8 | 57.4 | 16.0 | 63.1 |
| LiNGAM | 57.8 | 62.3 | 3.3 | 72.2 | 31.1 |
| RECI | 59.0 | 94.7 | 66.0 | 71.0 | 70.5 |
| CCS | 69.3 | <u>96.0</u> | 89.7 | 90.5 | N/A |
| CHD | 72.0 | <u>97.6</u> | 90.5 | 91.4 | N/A |
| CKL | 69.8 | 95.5 | 89.3 | 91.0 | N/A |
| CKM | 69.7 | 90.6 | <u>94.3</u> | 91.6 | N/A |
| CTV | 72.2 | 95.8 | 91.9 | <u>91.8</u> | N/A |
| GPLVM | 82.1 | 97.7 | 98.8 | 90.2 | 78.3 |

Summary

- Causal discovery from observational data is naturally a Bayesian Model Selection problem.
- Bayes allows specifying *realistic* assumptions, without artificial/unverifiable restrictions.

 **A Bayesian method with realistic assumptions without strict guarantees
outperforms methods with unrealistic assumptions that do provide guarantees.**

Future Work & Links to Deep Learning

- Can we express causal assumptions in neural network architecture, and discover them?

Learning Layer-wise Equivariances Automatically using Gradients

Tycho F.A. van der Ouderaa¹

Alexander Immer^{2,3}

Mark van der Wilk^{1,4}

- Can we scale this to multiple variables?
- Can we use deep generative models as meta-learners to replace explicit Bayesian approximate inference? (Everything Bayes can do, meta-learning can do with simulated data.)