

# Probabilistic latent variable models for distinguishing between cause and effect

Joris Mooij

joint work with

Oliver Stegle    Dominik Janzing    Kun Zhang    Bernhard Schölkopf

Max Planck Institute for Biological Cybernetics  
Tübingen, Germany  
[joris.mooij@tuebingen.mpg.de](mailto:joris.mooij@tuebingen.mpg.de)



MAX-PLANCK-GESELLSCHAFT

Tübingen, October 7th, 2010



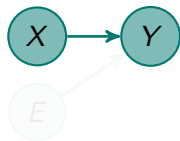
MPI FOR BIOLOGICAL CYBERNETICS

# The core problem

Given  $N$  observations of two variables  $X, Y$

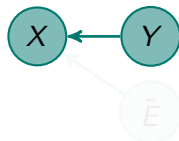
$i$	$x_i$	$y_i$
1	3.4	2.5
2	7.2	5.2
3	2.3	1.3
$\vdots$	$\vdots$	$\vdots$
$N$	4.5	2.1

determine whether:



“X causes Y”

or



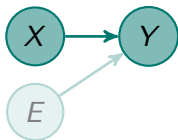
“Y causes X”

# The core problem

Given  $N$  observations of two variables  $X, Y$

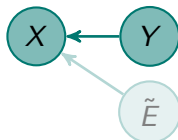
$i$	$x_i$	$y_i$
1	3.4	2.5
2	7.2	5.2
3	2.3	1.3
$\vdots$	$\vdots$	$\vdots$
$N$	4.5	2.1

determine whether:



“ $X$  causes  $Y$ ”

or

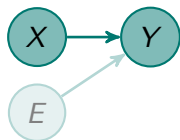


“ $Y$  causes  $X$ ”

# Part I

## The basics

# Basic assumptions



Suppose:

- $X, E$  cause  $Y$ ;
- $X, Y$  are *observed*;
- $E$  is *unobserved* (latent)

- ① **Determinism:** a function  $f$  exists such that:

$$Y = f(X, E)$$

( $f$  is called the *causal mechanism*);

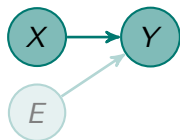
- ② **No common causes** of  $X$  and  $E$  are present:

$$X \perp\!\!\!\perp E,$$

( $X$  and  $E$  are statistically *independent*).

- ③ **Independence of distribution of causes and mechanism:**

$$p(X, E) \perp\!\!\!\perp f$$



Suppose:

- $X, E$  cause  $Y$ ;
- $X, Y$  are *observed*;
- $E$  is *unobserved* (latent)

- ① **Determinism:** a function  $f$  exists such that:

$$Y = f(X, E)$$

( $f$  is called the *causal mechanism*);

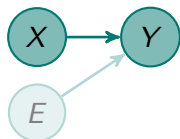
- ② **No common causes** of  $X$  and  $E$  are present:

$$X \perp\!\!\!\perp E,$$

( $X$  and  $E$  are statistically *independent*).

- ③ **Independence of distribution of causes and mechanism:**

$$p(X, E) \perp\!\!\!\perp f$$



Suppose:

- $X, E$  cause  $Y$ ;
- $X, Y$  are *observed*;
- $E$  is *unobserved* (latent)

- 1 **Determinism:** a function  $f$  exists such that:

$$Y = f(X, E)$$

( $f$  is called the *causal mechanism*);

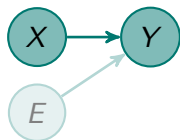
- 2 **No common causes** of  $X$  and  $E$  are present:

$$X \perp\!\!\!\perp E,$$

( $X$  and  $E$  are statistically *independent*).

- 3 **Independence of distribution of causes and mechanism:**

$$p(X, E) \perp\!\!\!\perp f$$



Suppose:

- $X, E$  cause  $Y$ ;
- $X, Y$  are *observed*;
- $E$  is *unobserved* (latent)

- ① **Determinism:** a function  $f$  exists such that:

$$Y = f(X, E)$$

( $f$  is called the *causal mechanism*);

- ② **No common causes** of  $X$  and  $E$  are present:

$$X \perp\!\!\!\perp E,$$

( $X$  and  $E$  are statistically *independent*).

- ③ **Independence of distribution of causes and mechanism:**

$$p(X, E) \perp\!\!\!\perp f$$



# Previous work: restricting the model class

Several recent approaches restrict the class of possible causal mechanisms:

LINGAM	$f$ is linear	$f(X, E) = \alpha X + \beta E$
AN	Additive Noise	$f(X, E) = F(X) + E$
PNL	Post-Non-Linear model	$f(X, E) = G(F(X) + E)$
HS	Hetero-Schedastic noise	$f(X, E) = F(X) + E \cdot G(X)$

Causal discovery is possible because of the following *identifiability* results:

Theorem (*al.*)

*Let  $\mathcal{M} \in \{\text{LINGAM, AN, PNL}\}$ : generically, if a model  $X \rightarrow Y$  in  $\mathcal{M}$  exists for  $p(X, Y)$ , no model  $Y \rightarrow X \in \mathcal{M}$  exists for  $p(X, Y)$ .*

The idea is to fit a restricted model in both directions ( $X \rightarrow Y$  and  $Y \rightarrow X$ ) and infer the causal direction to be the one that yields the best fit with the data.

# Previous work: restricting the model class

Several recent approaches restrict the class of possible causal mechanisms:

LINGAM	$f$ is linear	$f(X, E) = \alpha X + \beta E$
AN	Additive Noise	$f(X, E) = F(X) + E$
PNL	Post-Non-Linear model	$f(X, E) = G(F(X) + E)$
HS	Hetero-Schedastic noise	$f(X, E) = F(X) + E \cdot G(X)$

Causal discovery is possible because of the following *identifiability* results:

## Theorem (*al.*)

Let  $\mathcal{M} \in \{\text{LINGAM}, \text{AN}, \text{PNL}\}$ : generically, if a model  $X \rightarrow Y$  in  $\mathcal{M}$  exists for  $p(X, Y)$ , no model  $Y \rightarrow X \in \mathcal{M}$  exists for  $p(X, Y)$ .

The idea is to fit a restricted model in both directions ( $X \rightarrow Y$  and  $Y \rightarrow X$ ) and infer the causal direction to be the one that yields the best fit with the data.

# Previous work: restricting the model class

Several recent approaches restrict the class of possible causal mechanisms:

LINGAM	$f$ is linear	$f(X, E) = \alpha X + \beta E$
AN	Additive Noise	$f(X, E) = F(X) + E$
PNL	Post-Non-Linear model	$f(X, E) = G(F(X) + E)$
HS	Hetero-Schedastic noise	$f(X, E) = F(X) + E \cdot G(X)$

Causal discovery is possible because of the following *identifiability* results:

## Theorem (*al.*)

Let  $\mathcal{M} \in \{\text{LINGAM}, \text{AN}, \text{PNL}\}$ : generically, if a model  $X \rightarrow Y$  in  $\mathcal{M}$  exists for  $p(X, Y)$ , no model  $Y \rightarrow X \in \mathcal{M}$  exists for  $p(X, Y)$ .

The idea is to fit a restricted model in both directions ( $X \rightarrow Y$  and  $Y \rightarrow X$ ) and infer the causal direction to be the one that yields the best fit with the data.

## Previous work: comparing model complexities

A different (recently proposed) approach is based on information theory. It does not restrict the class of possible causal mechanisms.

Theorem (Janzing, Schölkopf)

If  $I(p(X) : p(Y | X)) \stackrel{+}{=} 0$ , then

$$K(p(X)) + K(p(Y | X)) \stackrel{+}{\leq} K(p(Y)) + K(p(X | Y)),$$

where  $K(\cdot)$  is the Kolmogorov complexity and  $I(\cdot : \cdot)$  is the analogue of mutual information based on Kolmogorov complexity.

Unfortunately, Kolmogorov complexity is uncomputable, so this result is not directly useful in practice.

## Previous work: comparing model complexities

A different (recently proposed) approach is based on information theory. It does not restrict the class of possible causal mechanisms.

### Theorem (Janzing, Schölkopf)

If  $I(p(X) : p(Y | X)) \stackrel{+}{=} 0$ , then

$$K(p(X)) + K(p(Y | X)) \stackrel{+}{\leq} K(p(Y)) + K(p(X | Y)),$$

where  $K(\cdot)$  is the Kolmogorov complexity and  $I(\cdot : \cdot)$  is the analogue of mutual information based on Kolmogorov complexity.

Unfortunately, Kolmogorov complexity is uncomputable, so this result is not directly useful in practice.

## Previous work: comparing model complexities

A different (recently proposed) approach is based on information theory. It does not restrict the class of possible causal mechanisms.

### Theorem (Janzing, Schölkopf)

If  $I(p(X) : p(Y | X)) \stackrel{+}{=} 0$ , then

$$K(p(X)) + K(p(Y | X)) \stackrel{+}{\leq} K(p(Y)) + K(p(X | Y)),$$

where  $K(\cdot)$  is the Kolmogorov complexity and  $I(\cdot : \cdot)$  is the analogue of mutual information based on Kolmogorov complexity.

Unfortunately, Kolmogorov complexity is uncomputable, so this result is not directly useful in practice.

# This work: non-parametric causal discovery

The goal of this work is to avoid the restrictions on the model class for the causal mechanism.

However, in this way we lose identifiability:

## Theorem

*Given random variables  $X, Y, E$  and a function  $f$  such that*

$$Y = f(X, E), \quad X \perp\!\!\!\perp E,$$

*we can always construct a function  $\tilde{f}$  and a random variable  $\tilde{E}$  such that*

$$X = \tilde{f}(Y, \tilde{E}), \quad Y \perp\!\!\!\perp \tilde{E}.$$

The crucial idea is now to compare the *model complexities* and infer the least complex model to be the true causal direction. However, we use other complexity measures than Kolmogorov complexity, so that we can actually compute them.

# This work: non-parametric causal discovery

The goal of this work is to avoid the restrictions on the model class for the causal mechanism.

However, in this way we loose identifiability:

## Theorem

*Given random variables  $X, Y, E$  and a function  $f$  such that*

$$Y = f(X, E), \quad X \perp\!\!\!\perp E,$$

*we can always construct a function  $\tilde{f}$  and a random variable  $\tilde{E}$  such that*

$$X = \tilde{f}(Y, \tilde{E}), \quad Y \perp\!\!\!\perp \tilde{E}.$$

The crucial idea is now to compare the *model complexities* and infer the least complex model to be the true causal direction. However, we use other complexity measures than Kolmogorov complexity, so that we can actually compute them.



# This work: non-parametric causal discovery

The goal of this work is to avoid the restrictions on the model class for the causal mechanism.

However, in this way we lose identifiability:

## Theorem

*Given random variables  $X, Y, E$  and a function  $f$  such that*

$$Y = f(X, E), \quad X \perp\!\!\!\perp E,$$

*we can always construct a function  $\tilde{f}$  and a random variable  $\tilde{E}$  such that*

$$X = \tilde{f}(Y, \tilde{E}), \quad Y \perp\!\!\!\perp \tilde{E}.$$

The crucial idea is now to compare the *model complexities* and infer the least complex model to be the true causal direction. However, we use other complexity measures than Kolmogorov complexity, so that we can actually compute them.

## Bayesian model selection

Prefer the model with the highest *evidence*:

$$p(D | M) = \int p(D | \theta, M) p(\theta | M) d\theta,$$

which is a trade-off between the *likelihood* (goodness-of-fit) and the *prior* (model complexity). ( $D$  = data,  $M$  = model,  $\theta$  = model parameters)

## Basic idea

Causal discovery (for our “core problem”) can be done simply by comparing the evidences for two models ( $X \rightarrow Y$  and  $Y \rightarrow X$ ).

## Bayesian model selection

Prefer the model with the highest *evidence*:

$$p(D | M) = \int p(D | \theta, M) p(\theta | M) d\theta,$$

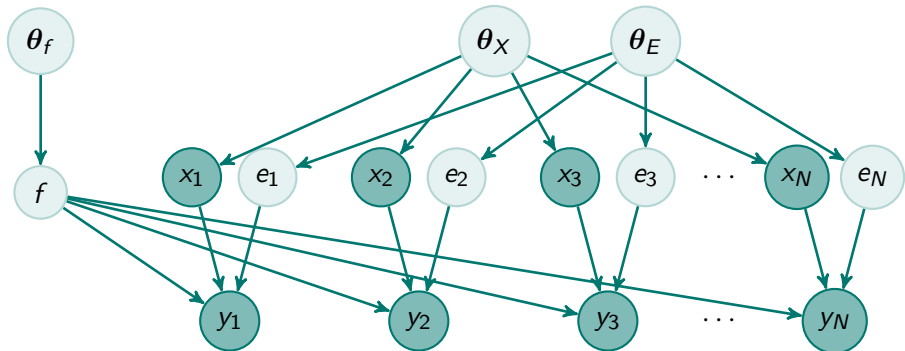
which is a trade-off between the *likelihood* (goodness-of-fit) and the *prior* (model complexity). ( $D = \text{data}$ ,  $M = \text{model}$ ,  $\theta = \text{model parameters}$ )

## Basic idea

Causal discovery (for our “core problem”) can be done simply by comparing the evidences for two models ( $X \rightarrow Y$  and  $Y \rightarrow X$ ).

# How we model $X \rightarrow Y$

$$\begin{aligned} p(\mathbf{x}, \mathbf{y}) &= p(\mathbf{x})p(\mathbf{y} | \mathbf{x}) \\ &= \int \left( \prod_{i=1}^N p(x_i | \theta_X) \right) p(\theta_X) d\theta_X \\ &\cdot \int \left( \prod_{i=1}^N \delta(y_i - f(x_i, e_i)) p(e_i | \theta_E) \right) p(f | \theta_f) df p(\theta_E) d\theta_E p(\theta_f) d\theta_f \end{aligned}$$



## Part II

# The nasty technical details

# Choosing the priors

Let  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^N$  denote the  $N$  values for  $X, Y$ , and  $\mathbf{e} \in \mathbb{R}^N$  the  $N$  latent variables.

In order to completely specify the generative model  $X \rightarrow Y$ , we need to choose various priors:

- the prior distribution on the inputs  $\mathbf{x}$  (parameterized by  $\theta_X$ )
- the prior distribution on the latents  $\mathbf{e}$  (parameterized by  $\theta_E$ )
- the prior distribution on the function  $f$  (parameterized by  $\theta_f$ )

Note: even in the discrete, finite case (if  $X, Y, E$  can only take a finite number of values), it is not obvious whether the choice of the priors becomes less important as  $N \rightarrow \infty$ : the number of observations ( $2N$ ) is of the same order as the number of unknowns ( $N + K$  for some constant  $K$ ).

# Choosing the priors

Let  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^N$  denote the  $N$  values for  $X, Y$ , and  $\mathbf{e} \in \mathbb{R}^N$  the  $N$  latent variables.

In order to completely specify the generative model  $X \rightarrow Y$ , we need to choose various priors:

- the prior distribution on the inputs  $\mathbf{x}$  (parameterized by  $\theta_X$ )
- the prior distribution on the latents  $\mathbf{e}$  (parameterized by  $\theta_E$ )
- the prior distribution on the function  $f$  (parameterized by  $\theta_f$ )

Note: even in the discrete, finite case (if  $X, Y, E$  can only take a finite number of values), it is not obvious whether the choice of the priors becomes less important as  $N \rightarrow \infty$ : the number of observations ( $2N$ ) is of the same order as the number of unknowns ( $N + K$  for some constant  $K$ ).

# Choosing the priors

Let  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^N$  denote the  $N$  values for  $X, Y$ , and  $\mathbf{e} \in \mathbb{R}^N$  the  $N$  latent variables.

In order to completely specify the generative model  $X \rightarrow Y$ , we need to choose various priors:

- the prior distribution on the inputs  $\mathbf{x}$  (parameterized by  $\theta_X$ )
- the prior distribution on the latents  $\mathbf{e}$  (parameterized by  $\theta_E$ )
- the prior distribution on the function  $f$  (parameterized by  $\theta_f$ )

Note: even in the discrete, finite case (if  $X, Y, E$  can only take a finite number of values), it is not obvious whether the choice of the priors becomes less important as  $N \rightarrow \infty$ : the number of observations ( $2N$ ) is of the same order as the number of unknowns ( $N + K$  for some constant  $K$ ).



# Choosing the priors: the input distribution $p(X | \theta_X)$

For  $p(X | \theta_X)$ , we currently use a Gaussian Mixture Model

$$p(X | \theta_X) = \sum_{i=1}^k \alpha_i \mathcal{N}(\mu_i, \sigma_i^2)$$

with hyperparameters  $\theta_X = (k, \alpha_1, \dots, \alpha_k, \mu_1, \dots, \mu_k, \sigma_1, \dots, \sigma_k)$ , with an improper Dirichlet prior (with parameters  $(-1, -1, \dots, -1)$ ) on the component weights  $\alpha$  and a flat prior on the component parameters  $\mu, \sigma$ .

Instead of integrating over  $\theta_X$ , we maximize over  $\theta_X$ , using a particular penalty for  $k$ , the number of mixture components, which is derived using the MML principle. We use an algorithm proposed by Figueiredo and Jain.<sup>1</sup>

---

<sup>1</sup>Figueiredo & Jain, *Unsupervised learning of finite mixture models*, TPAMI 2002

# Choosing the priors: the input distribution $p(X | \theta_X)$

For  $p(X | \theta_X)$ , we currently use a Gaussian Mixture Model

$$p(X | \theta_X) = \sum_{i=1}^k \alpha_i \mathcal{N}(\mu_i, \sigma_i^2)$$

with hyperparameters  $\theta_X = (k, \alpha_1, \dots, \alpha_k, \mu_1, \dots, \mu_k, \sigma_1, \dots, \sigma_k)$ , with an improper Dirichlet prior (with parameters  $(-1, -1, \dots, -1)$ ) on the component weights  $\alpha$  and a flat prior on the component parameters  $\mu, \sigma$ .

Instead of integrating over  $\theta_X$ , we maximize over  $\theta_X$ , using a particular penalty for  $k$ , the number of mixture components, which is derived using the MML principle. We use an algorithm proposed by Figueiredo and Jain.<sup>1</sup>

---

<sup>1</sup>Figueiredo & Jain, *Unsupervised learning of finite mixture models*, TPAMI 2002

# Choosing the priors: the noise distribution $p(E | \theta_E)$

For  $p(E | \theta_E)$ , we simply use a standard-normal distribution:

$$p(E | \theta_E) = \mathcal{N}(0, 1),$$

so there are no hyperparameters.

This may look like a severe restriction on the model, but is not as bad as it seems: in general, there exists a function  $g$  such that

$$E = g(\bar{E}), \bar{E} \sim \mathcal{N}(0, 1).$$

Then  $Y = f(X, E)$  corresponds with  $Y = \bar{f}(X, \bar{E})$ , where  $\bar{f} := f(\cdot, g(\cdot))$ .

However, this assumption does introduce a dependency between  $p(E)$  and  $f$ , which apparently violates our basic assumption  $p(X, E) \perp\!\!\!\perp f$ .

# Choosing the priors: the noise distribution $p(E | \theta_E)$

For  $p(E | \theta_E)$ , we simply use a standard-normal distribution:

$$p(E | \theta_E) = \mathcal{N}(0, 1),$$

so there are no hyperparameters.

This may look like a severe restriction on the model, but is not as bad as it seems: in general, there exists a function  $g$  such that

$$E = g(\bar{E}), \bar{E} \sim \mathcal{N}(0, 1).$$

Then  $Y = f(X, E)$  corresponds with  $Y = \bar{f}(X, \bar{E})$ , where  $\bar{f} := f(\cdot, g(\cdot))$ .

However, this assumption does introduce a dependency between  $p(E)$  and  $f$ , which apparently violates our basic assumption  $p(X, E) \perp\!\!\!\perp f$ .

# Choosing the priors: the noise distribution $p(E | \theta_E)$

For  $p(E | \theta_E)$ , we simply use a standard-normal distribution:

$$p(E | \theta_E) = \mathcal{N}(0, 1),$$

so there are no hyperparameters.

This may look like a severe restriction on the model, but is not as bad as it seems: in general, there exists a function  $g$  such that

$$E = g(\bar{E}), \bar{E} \sim \mathcal{N}(0, 1).$$

Then  $Y = f(X, E)$  corresponds with  $Y = \bar{f}(X, \bar{E})$ , where  $\bar{f} := f(\cdot, g(\cdot))$ .

However, this assumption does introduce a dependency between  $p(E)$  and  $f$ , which apparently violates our basic assumption  $p(X, E) \perp\!\!\!\perp f$ .

# Choosing the priors: the function prior $p(f | \theta_f)$

For  $p(f | \theta_f)$ , we currently take a Gaussian Process with zero mean function and squared-exponential covariance function:

$$k((x, e), (x', e')) = \lambda_Y^2 \exp\left(-\frac{(x - x')^2}{2\lambda_X^2}\right) \exp\left(-\frac{(e - e')^2}{2\lambda_E^2}\right)$$

where  $\theta_f = (\lambda_X, \lambda_Y, \lambda_E)$  are length-scale parameters.

We currently preprocess the data  $\mathbf{x}$ ,  $\mathbf{y}$  such that they have mean 0 and variance 1. For the prior on the length-scale parameters, we use a broad Gamma distribution:  $\lambda \sim \Gamma(30, 0.5)$ . The only reason for doing this is a numerical one: if the length scales become too large, the kernel matrix  $\mathbf{K}_{ij} = k((x_i, e_i), (x_j, e_j))$  will become difficult to handle numerically.

Instead of integrating over  $\theta_f$ , we maximize over  $\theta_f$ .

# Choosing the priors: the function prior $p(f | \theta_f)$

For  $p(f | \theta_f)$ , we currently take a Gaussian Process with zero mean function and squared-exponential covariance function:

$$k((x, e), (x', e')) = \lambda_Y^2 \exp\left(-\frac{(x - x')^2}{2\lambda_X^2}\right) \exp\left(-\frac{(e - e')^2}{2\lambda_E^2}\right)$$

where  $\theta_f = (\lambda_X, \lambda_Y, \lambda_E)$  are length-scale parameters.

We currently preprocess the data  $\mathbf{x}$ ,  $\mathbf{y}$  such that they have mean 0 and variance 1. For the prior on the length-scale parameters, we use a broad Gamma distribution:  $\lambda \sim \Gamma(30, 0.5)$ . The only reason for doing this is a numerical one: if the length scales become too large, the kernel matrix  $\mathbf{K}_{ij} = k((x_i, e_i), (x_j, e_j))$  will become difficult to handle numerically.

Instead of integrating over  $\theta_f$ , we maximize over  $\theta_f$ .

# Back to our integral

By maximizing over the hyperparameters instead of integrating over them, we have reduced the computational problem to solving:

$$p(\mathbf{x}, \mathbf{y} | X \rightarrow Y) \approx \max_{\theta_X} \left( \prod_{i=1}^N p(x_i | \theta_X) \right) p(\theta_X) \\ \cdot \max_{\theta_f} p(\theta_f) \int \left( \prod_{i=1}^N \delta(y_i - f(x_i, e_i)) p(e_i) \right) p(f | \theta_f) \mathbf{d}e \mathbf{d}f$$

The first maximization problem is solved numerically by using a modified EM algorithm written by Figueiredo and Jain, and gives

$$\max_{\theta_X} \left( \prod_{i=1}^N p(x_i | \theta_X) \right) p(\theta_X) = \exp(-\mathcal{L}_{\text{MML}}(\mathbf{x}))$$

From now on, we focus on the second part.



By maximizing over the hyperparameters instead of integrating over them, we have reduced the computational problem to solving:

$$p(\mathbf{x}, \mathbf{y} | X \rightarrow Y) \approx \max_{\boldsymbol{\theta}_X} \left( \prod_{i=1}^N p(x_i | \boldsymbol{\theta}_X) \right) p(\boldsymbol{\theta}_X) \\ \cdot \max_{\boldsymbol{\theta}_f} p(\boldsymbol{\theta}_f) \int \left( \prod_{i=1}^N \delta(y_i - f(x_i, e_i)) p(e_i) \right) p(f | \boldsymbol{\theta}_f) d\mathbf{e} df$$

The first maximization problem is solved numerically by using a modified EM algorithm written by Figueiredo and Jain, and gives

$$\max_{\boldsymbol{\theta}_X} \left( \prod_{i=1}^N p(x_i | \boldsymbol{\theta}_X) \right) p(\boldsymbol{\theta}_X) = \exp(-\mathcal{L}_{\text{MML}}(\mathbf{x}))$$

From now on, we focus on the second part.

# Integrating over the latent variables $\mathbf{e}$

We first integrate out the latent variables  $\mathbf{e}$ , using the Dirac delta function calculus:

$$\int \left( \prod_{i=1}^N \delta(y_i - f(x_i, e_i)) p(e_i) \right) p(f | \boldsymbol{\theta}_f) d\mathbf{e} df = \int \frac{p(\mathbf{e}_0(f))}{J(\mathbf{e}_0(f))} p(f | \boldsymbol{\theta}_f) df$$

where

$$J(\mathbf{e}_0(f)) = \det \left| \frac{\partial f}{\partial \mathbf{e}}(\mathbf{x}, \mathbf{e}_0(f)) \right| = \prod_{i=1}^N \left| \frac{\partial f}{\partial e} (x_i, (\mathbf{e}_0(f))_i) \right|$$

is the absolute value of the determinant of the Jacobian and  $\mathbf{e}_0(f)$  is the unique vector satisfying  $f(x_i, (\mathbf{e}_0(f))_i) = y_i$ .

Here we assumed that for each  $x$ , the function  $f_x : e \mapsto f(x, e)$  is invertible. Although this is not a restriction on the model class, this assumption is not compatible with our Gaussian Process prior on  $f$ .

# Integrating over the latent variables $\mathbf{e}$

We first integrate out the latent variables  $\mathbf{e}$ , using the Dirac delta function calculus:

$$\int \left( \prod_{i=1}^N \delta(y_i - f(x_i, e_i)) p(e_i) \right) p(f | \boldsymbol{\theta}_f) d\mathbf{e} df = \int \frac{p(\mathbf{e}_0(f))}{J(\mathbf{e}_0(f))} p(f | \boldsymbol{\theta}_f) df$$

where

$$J(\mathbf{e}_0(f)) = \det \left| \frac{\partial f}{\partial \mathbf{e}}(\mathbf{x}, \mathbf{e}_0(f)) \right| = \prod_{i=1}^N \left| \frac{\partial f}{\partial e}(x_i, (\mathbf{e}_0(f))_i) \right|$$

is the absolute value of the determinant of the Jacobian and  $\mathbf{e}_0(f)$  is the unique vector satisfying  $f(x_i, (\mathbf{e}_0(f))_i) = y_i$ .

Here we assumed that for each  $x$ , the function  $f_x : e \mapsto f(x, e)$  is invertible. Although this is not a restriction on the model class, this assumption is not compatible with our Gaussian Process prior on  $f$ .

# “Integrating over $f$ ”

We then evaluate the remaining integral:

$$\int \frac{p(\mathbf{e}_0(f))p(f | \boldsymbol{\theta}_f)}{J(\mathbf{e}_0(f))} df = \int \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K}) \mathcal{N}(\mathbf{e}_0(f) | \mathbf{0}, \mathbf{I}) J^{-1}(\mathbf{e}_0(f)) d\mathbf{e}_0(f)$$

where  $\mathbf{K} = k(\mathbf{x}, \mathbf{e}_0(f))$ .

We then approximate this integral by maximizing over  $\mathbf{e}_0(f)$  (henceforth simply denoted  $\mathbf{e}$ ):

$$\dots \approx \max_{\mathbf{e}} \left( \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K}) \mathcal{N}(\mathbf{e} | \mathbf{0}, \mathbf{I}) J^{-1}(\mathbf{e}) \right)$$

and using the mean predicted partial derivatives of the GP  $f$ :

$$\frac{\partial f}{\partial \mathbf{e}}(x_i, e_i) = \frac{\partial k}{\partial \mathbf{e}}((x_i, e_i), (\mathbf{x}, \mathbf{e})) \mathbf{K}^{-1} \mathbf{y}$$

for approximating the Jacobian.

# “Integrating over $f$ ”

We then evaluate the remaining integral:

$$\int \frac{p(\mathbf{e}_0(f))p(f | \boldsymbol{\theta}_f)}{J(\mathbf{e}_0(f))} df = \int \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K})\mathcal{N}(\mathbf{e}_0(f) | \mathbf{0}, \mathbf{I})J^{-1}(\mathbf{e}_0(f)) d\mathbf{e}_0(f)$$

where  $\mathbf{K} = k(\mathbf{x}, \mathbf{e}_0(f))$ .

We then approximate this integral by maximizing over  $\mathbf{e}_0(f)$  (henceforth simply denoted  $\mathbf{e}$ ):

$$\dots \approx \max_{\mathbf{e}} \left( \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K})\mathcal{N}(\mathbf{e} | \mathbf{0}, \mathbf{I})J^{-1}(\mathbf{e}) \right)$$

and using the mean predicted partial derivatives of the GP  $f$ :

$$\frac{\partial f}{\partial \mathbf{e}}(x_i, e_i) = \frac{\partial k}{\partial \mathbf{e}}((x_i, e_i), (\mathbf{x}, \mathbf{e}))\mathbf{K}^{-1}\mathbf{y}$$

for approximating the Jacobian.

# Integrating or optimizing latent variables?

Summarizing, we first integrated out the latent variables and then optimized over control points of the GP.

Alternatively, one might start with integrating out the GP exactly, and then optimize over the latent variables, similarly to what is usually done in GPLVMs (Gaussian Process Latent Variable Models).

However, we believe that for the purpose of causal discovery, the latter approach would not work well. The reason is that when optimizing over  $\mathbf{e}$ , the result is often quite dependent on  $\mathbf{x}$ , which violates our basic assumption that  $X \perp\!\!\!\perp E$ . Indeed, our approach seems more related to (nonlinear) ICA than to PCA (which is related to GPLVMs).

# Integrating or optimizing latent variables?

Summarizing, we first integrated out the latent variables and then optimized over control points of the GP.

Alternatively, one might start with integrating out the GP exactly, and then optimize over the latent variables, similarly to what is usually done in GPLVMs (Gaussian Process Latent Variable Models).

However, we believe that for the purpose of causal discovery, the latter approach would not work well. The reason is that when optimizing over  $\mathbf{e}$ , the result is often quite dependent on  $\mathbf{x}$ , which violates our basic assumption that  $X \perp\!\!\!\perp E$ . Indeed, our approach seems more related to (nonlinear) ICA than to PCA (which is related to GPLVMs).

# Integrating or optimizing latent variables?

Summarizing, we first integrated out the latent variables and then optimized over control points of the GP.

Alternatively, one might start with integrating out the GP exactly, and then optimize over the latent variables, similarly to what is usually done in GPLVMs (Gaussian Process Latent Variable Models).

However, we believe that for the purpose of causal discovery, the latter approach would not work well. The reason is that when optimizing over  $\mathbf{e}$ , the result is often quite dependent on  $\mathbf{x}$ , which violates our basic assumption that  $X \perp\!\!\!\perp E$ . Indeed, our approach seems more related to (nonlinear) ICA than to PCA (which is related to GPLVMs).



# The final objective function

Thus our final optimization problem is:

$$\min_{\lambda, \mathbf{e}} \left( -\log p(\lambda) - \log \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K}) - \log \mathcal{N}(\mathbf{e} | \mathbf{0}, \mathbf{I}) + \sum_{i=1}^N \log \left| \frac{\partial f}{\partial \mathbf{e}}(x_i, e_i) \right| \right)$$

where

$$\mathbf{K} = k(\mathbf{x}, \mathbf{e})$$

and

$$\frac{\partial f}{\partial \mathbf{e}}(x_i, e_i) = \frac{\partial k}{\partial \mathbf{e}}((x_i, e_i), (\mathbf{x}, \mathbf{e})) \mathbf{K}^{-1} \mathbf{y}$$

## Problems

There are still two major issues here:

- 1 if  $\frac{\partial f}{\partial \mathbf{e}}$  becomes 0 for some  $(x_i, e_i)$ , the objective function becomes  $-\infty$ ;
- 2 the kernel matrix  $\mathbf{K}$  is extremely ill-conditioned.

# The final objective function

Thus our final optimization problem is:

$$\min_{\lambda, \mathbf{e}} \left( -\log p(\lambda) - \log \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K}) - \log \mathcal{N}(\mathbf{e} | \mathbf{0}, \mathbf{I}) + \sum_{i=1}^N \log \left| \frac{\partial f}{\partial \mathbf{e}}(x_i, e_i) \right| \right)$$

where

$$\mathbf{K} = k(\mathbf{x}, \mathbf{e})$$

and

$$\frac{\partial f}{\partial \mathbf{e}}(x_i, e_i) = \frac{\partial k}{\partial \mathbf{e}}((x_i, e_i), (\mathbf{x}, \mathbf{e})) \mathbf{K}^{-1} \mathbf{y}$$

## Problems

There are still two major issues here:

- 1 if  $\frac{\partial f}{\partial \mathbf{e}}$  becomes 0 for some  $(x_i, e_i)$ , the objective function becomes  $-\infty$ ;
- 2 the kernel matrix  $\mathbf{K}$  is extremely ill-conditioned.

To deal with the first problem, we regularize our objective function as follows:

- 1 We approximate  $|\log| x \approx \log \sqrt{x^2 + \epsilon}$  with  $\epsilon \ll 1$  (using  $\epsilon = 10^{-3}$ ).
- 2 We implemented a log barrier that heavily penalized negative values of  $\frac{\partial f}{\partial \mathbf{e}}(\hat{x}_i, \mathbf{e}_i)$ . This was done to avoid sign flips of these terms that would violate the invertability assumption.
- 3 We added a tiny amount of additive  $\mathcal{N}(0, \sigma^2)$ -noise to each  $y_i$ -value, which is equivalent to replacing  $\mathbf{K}$  by  $\mathbf{K} + \sigma^2 \mathbf{I}$  (using  $\sigma = 10^{-5}$ ).

Further, we initialize  $\mathbf{e}$  with additive noise models. The main reason is that in an additive noise model, the  $\frac{\partial f}{\partial \mathbf{e}}(x_i, \mathbf{e}_i)$  are all positive and constant. This initialization effectively leads to a solution that satisfies the invertability assumption that we made.

To deal with the first problem, we regularize our objective function as follows:

- 1 We approximate  $|\log| x \approx \log \sqrt{x^2 + \epsilon}$  with  $\epsilon \ll 1$  (using  $\epsilon = 10^{-3}$ ).
- 2 We implemented a log barrier that heavily penalized negative values of  $\frac{\partial f}{\partial \mathbf{e}}(\hat{x}_i, e_j)$ . This was done to avoid sign flips of these terms that would violate the invertability assumption.
- 3 We added a tiny amount of additive  $\mathcal{N}(0, \sigma^2)$ -noise to each  $y_i$ -value, which is equivalent to replacing  $\mathbf{K}$  by  $\mathbf{K} + \sigma^2 \mathbf{I}$  (using  $\sigma = 10^{-5}$ ).

Further, we initialize  $\mathbf{e}$  with additive noise models. The main reason is that in an additive noise model, the  $\frac{\partial f}{\partial \mathbf{e}}(x_i, e_j)$  are all positive and constant. This initialization effectively leads to a solution that satisfies the invertability assumption that we made.

- See also our NIPS paper *Probabilistic latent variable models for distinguishing between cause and effect*,  
J. M. Mooij, O. Stegle, D. Janzing, K. Zhang, B. Schölkopf,  
Advances in Neural Information Processing Systems 23 (NIPS\*2010)
- See also our code package (also distributed as part of the Causality Challenge):  
<http://webdav.tuebingen.mpg.de/causality/nips2010-gpi-code.tar.gz>