

# Probabilistic Active Learning of Structural Causal Models

Causality Workshop

Uncertainty in Artificial Intelligence, August 2017

---

Paul K. Rubenstein, Ilya Tolstikhin, Philipp Hennig, Bernhard Schölkopf

Max Planck Institute for Intelligent Systems, Tübingen

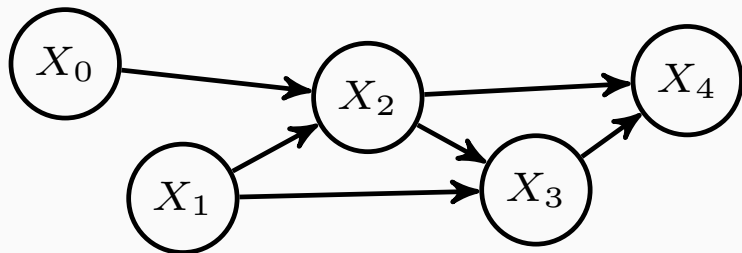
# Motivation

What should we do after causal discovery? Understand **relationships** between variables for control.

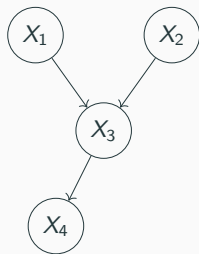


Would it be better to reduce each smoker's consumption by 50%, or stop 50% of smokers completely?

# Motivation



# Structural Equation Models



$$X_1 = f_1(E_1)$$

$$X_2 = f_2(E_2)$$

$$X_3 = f_3(X_1, X_2, E_3)$$

$$X_4 = f_4(X_3, E_4)$$

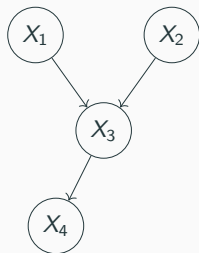
$$\mathbf{E} \sim \mathbb{P}_E$$

An SCM  $\mathcal{M}_X = (\mathcal{S}_X, \mathbb{P}_E)$  consists of

- $\mathcal{S}_X$  a set of structural equations  $X_i = f_i(X, E_i)$  ;
- $\mathbb{P}_E$  is a distribution over the exogenous variables  $E$ ;

We consider *perfect interventions*, e. g.  $\text{do}(X_3 = 0)$

## Desired problem setup



$$X_1 = f_1(E_1)$$

$$X_2 = f_2(E_2)$$

$$X_3 = f_3(X_1, X_2, E_3)$$

$$X_4 = f_4(X_3, E_4)$$

$$\mathbf{E} \sim \mathbb{P}_E$$

We assume the **causal graph is known**.

**Problem 1: Estimating  $\mathcal{M}$**  from data  $\mathcal{D}$  consisting of draws  $X^{(n)} \sim \mathbb{P}_X^{\text{do}(i_n)}$  from various interventions  $i_n$ . (= estimate  $f_i$  and  $\mathbb{P}_E$ )

**Problem 2: Active learning.** Choose informative intervention so that next estimate  $\widehat{\mathcal{M}}$  is 'closer'.

**How should we define distance between  $\mathcal{M}$  and  $\widehat{\mathcal{M}}$ ?**

## Open question

How should we define distance between  $\mathcal{M}$  and  $\widehat{\mathcal{M}}$ ?

$$L_{\text{KL}}^i(\widehat{\mathcal{M}}||\mathcal{M}) = \text{KL} \left[ \mathbb{P}_X^{\text{do}(i)} || \widehat{\mathbb{P}}_X^{\text{do}(i)} \right]$$

$$L_{\text{KL}}^{\mathcal{I}}(\widehat{\mathcal{M}}||\mathcal{M}) = \sup_{i \in \mathcal{I}} L_{\text{KL}}^i(\widehat{\mathcal{M}}||\mathcal{M})$$

$$L_{\text{MMD}_l}^i(\widehat{\mathcal{M}}||\mathcal{M}) = \text{MMD}_l \left[ \mathbb{P}_X^{\text{do}(i)} || \widehat{\mathbb{P}}_X^{\text{do}(i)} \right]$$

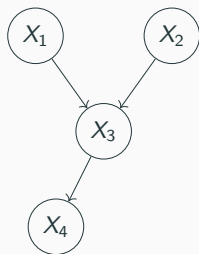
$$L_{\text{MMD}_l}^{\mathcal{I}}(\widehat{\mathcal{M}}||\mathcal{M}) = \sup_{i \in \mathcal{I}} L_{\text{MMD}_l}^i(\widehat{\mathcal{M}}||\mathcal{M})$$

**Goal:** make statements like

$$L(\widehat{\mathcal{M}}||\mathcal{M}) < \epsilon \implies ???$$

## Our problem setup

Assume **causal graph known** and **additive Gaussian noise with known diagonal covariance**



$$X_1 = f_1 + E_1$$

$$X_2 = f_2 + E_2$$

$$X_3 = f_3(X_1, X_2) + E_3$$

$$X_4 = f_4(X_3) + E_4$$

$$\mathbf{E} \sim \mathcal{N}(0, \Lambda)$$

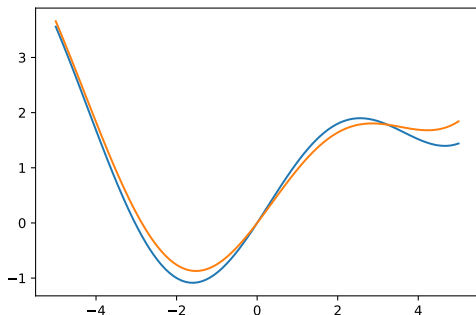
Only unknown quantities are functions  $f_n : \mathcal{X}_{pa(n)} \rightarrow \mathcal{X}_n$

Assume we are given measures  $\Pi_n$  on  $\mathcal{X}_{pa(n)}$

$$L_n(\hat{f}_n || f_n) = \int_{\mathcal{X}_{pa(n)}} (f_n(x) - \hat{f}_n(x))^2 d\Pi_n(x) \quad L(\hat{f} || f) = \sum_{n=1}^N \alpha_n L_n(\hat{f}_n || f_n)$$

## Our problem setup

$$L_n(\hat{f}_n || f_n) = \int_{\mathcal{X}_{pa(n)}} (f_n(x) - \hat{f}_n(x))^2 d\Pi_n(x) \quad L(\hat{f} || f) = \sum_{n=1}^N \alpha_n L_n(\hat{f}_n || f_n)$$



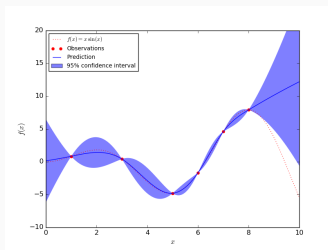
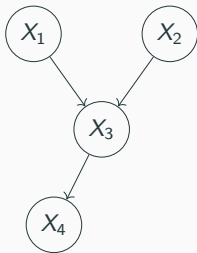


# Problem 1: Estimating $\mathcal{M}$

Estimating  $\mathcal{M}$  reduces to estimating each  $f_n$ .

Assume  $f_n \sim \mathcal{GP}(0, k_n)$

Let  $\mathcal{D}_n \subseteq \mathcal{D}$  be observations of  $(X_{pa(n)}, X_n)$  for  $i$  in which  $X_n$  is not intervened.



Posterior:

$$f_n | \mathcal{D}_n \sim \mathcal{GP}(\mu_{f_n | \mathcal{D}_n}, k_{f_n | \mathcal{D}_n})$$

## Problem 1: Estimating $\mathcal{M}$

Posterior:  $f_n | \mathcal{D}_n \sim \mathcal{GP}(\mu_{f_n | \mathcal{D}_n}, k_{f_n | \mathcal{D}_n})$

Which  $\hat{f}$  should be chosen, given the posterior over  $f$ ?

**Lemma 1:**

$$\mathbb{E}_{f | \mathcal{D}} [L(\hat{f} || f)] = \sum_{n=1}^N \alpha_n \int_{\mathcal{X}_{pa(n)}} \left( \hat{f}_n(x) - \mu_{f_n | \mathcal{D}_n}(x) \right)^2 + k_{n | \mathcal{D}_n}(x, x) d\Pi_n(x)$$

**Lemma 2:** Let  $\mu_{f | \mathcal{D}}$  be the tuple of functions  $(\mu_{f_n | \mathcal{D}_n})_{n=1, \dots, N}$ . Then

$$\mu_{f | \mathcal{D}} = \arg \min_{\hat{f}} \mathbb{E}_{f | \mathcal{D}} [L(\hat{f} || f)]$$

**Expected total risk:**  $\mathcal{R}(\mathcal{D}) = \mathbb{E}_{f | \mathcal{D}} [L(\mu_{f | \mathcal{D}} || f)]$

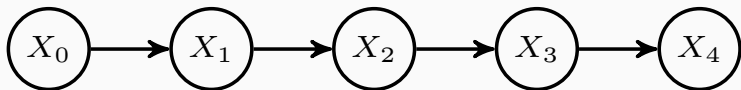
## Problem 2: Active Learning

**Goal:** choose intervention  $i$  so that  $\mathcal{R}(\mathcal{D} \cup (i, x))$  is minimised, taking into account cost of intervention.

$$V(i|\mathcal{D}) = \frac{\mathcal{R}(\mathcal{D}) - \mathbb{E}_{x \sim \mathbb{P}_X^{\text{do}(i)}} \mathcal{R}(\mathcal{D} \cup \{(i, x)\})}{c(i)}$$

**Problem:**  $\mathbb{P}_X^{\text{do}(i)}$  is unknown! Solution: replace with *estimate*  $\tilde{\mathbb{P}}_X^{\text{do}(i)}$ , from which we can sample.

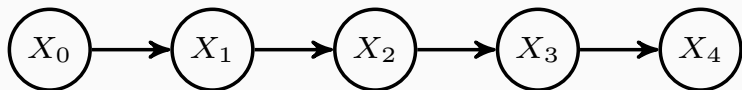
## Problem 2: Active Learning



**Algorithm 1:** Estimate  $\mathbb{E}_{x \sim \mathbb{P}_X^{\text{do}(i)}} \mathcal{R}(\mathcal{D} \cup \{(i, x)\})$  as

$$\frac{1}{T} \sum_{t=1}^T \mathcal{R}(\mathcal{D} \cup \{(i, x_t)\}), \quad x_t \sim \tilde{\mathbb{P}}_X^{\text{do}(i)}$$

## Problem 2: Active Learning



**Algorithm 2:** Dynamic programming (need 'nice' graph and interventions).

$$U_n(x_{n-1}) = \mathcal{R}_n(\mathcal{D}_n \cup \{(x_{n-1}, x_n)\}) \quad \text{for any value } x_n$$

$$\mathbb{E}_{x \sim \tilde{\mathbb{P}}_X^{\text{do}(i)}} [\mathcal{R}(\mathcal{D} \cup \{(i, x)\})] = \sum_{n=1}^m U_n^{\text{curr}} + \mathbb{E}_{x \sim \tilde{\mathbb{P}}_X^{\text{do}(i)}} \left[ \sum_{n=m+1}^N U_n(x_{n-1}) \right].$$