# Model selection over partially ordered sets

Armeen Taeb\*, Peter Bühlmann°, Venkat Chandrasekaran†

Department of Statistics, University of Washington \*; Seminar for Statistics, ETH Zürich  $^{\circ}$ ; Departments of Computing and Mathematical Sciences and of Electrical Engineering, Caltech  $^{\dagger}$ 

#### Motivation

Model selection with Boolean-logical structure:

- formulate and test hypothesises, e.g. is this variable present?
- easy to define model complexity and false positives

What about for problems that lack Boolean-logical structure?

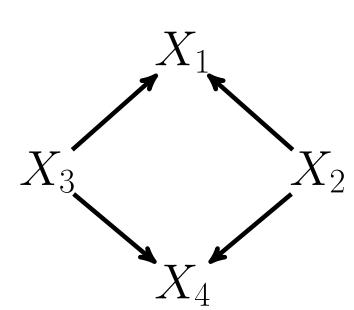
- ranking: global structure of transitivity
- clustering: global structure of set-partitions
- causal inference: global structure of acyclicity
- continuous problems, e.g. blind-source separation

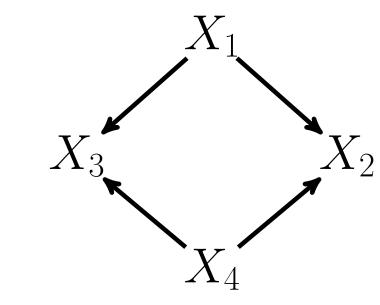
## Shortcomings of the standard perspective

Example I: clustering

true clusters =  $\{a,b\}$ ,  $\{c\}$  estimated clusters =  $\{a,b,c\}$ Boolean-logical perspective: FD = 2

Example II: causal structure learning





(a) true CPDAG

(b) estimated CPDAG

Boolean-logical perspective: FD = 4

# Model organization via posets

Models organized according to a poset  $\mathcal{L}$  with relations  $\preceq$ :

Attribute	Meaning
$\preceq$	containment between simpler & more complex models
least element	the "null" model representing no discoveries
$\operatorname{rank}(\cdot)$	measures complexity of a model

#### False discovery framework

**Similarity valuation**: A symmetric function  $\rho: \mathcal{L} \times \mathcal{L} \to \mathbb{R}$  with:

- $0 \le \rho(x, y) \le \min\{\operatorname{rank}(x), \operatorname{rank}(y)\}\$ for all  $x, y \in \mathcal{L}$ ,
- $\rho(x,y) \le \rho(z,y)$  for all  $x \le z$ ,
- $\bullet \rho(x,y) = \operatorname{rank}(x)$  if and only if  $x \leq y$ .

#### Definitions

Letting  $x^* \in \mathcal{L}$  be a true model and  $\hat{x} \in \mathcal{L}$  be an estimate.

$$\mathrm{TD}(\hat{x}, x^{\star}) \triangleq \rho(\hat{x}, x^{\star}),$$

$$FD(\hat{x}, x^*) \triangleq rank(\hat{x}) - \rho(\hat{x}, x^*) = rank(\hat{x}) - TD(\hat{x}, x^*),$$

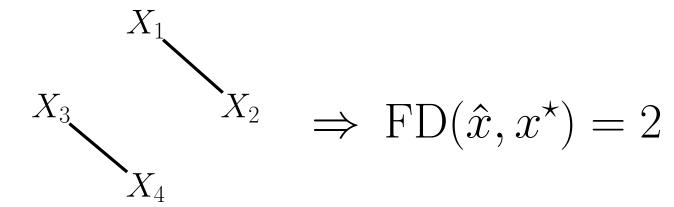
$$\mathrm{FDP}(\hat{x}, x^\star) \triangleq \frac{\mathrm{rank}(\hat{x}) - \rho(\hat{x}, x^\star)}{\mathrm{rank}(\hat{x})} = \frac{\mathrm{FD}(\hat{x}, x^\star)}{\mathrm{rank}(\hat{x})}.$$

Goal: maximize rank subject to false discovery control

Suitable similarity valuation:  $\rho_{\text{meet}}(\hat{x}, x^*) \triangleq \max_{z \leq \hat{x}, z \leq x^*} \text{rank}(z).$ 

- FD in clustering:
- # groups in the coarsest common refinement minus # groups in  $\hat{x}$ Example I: common refinement =  $\{a,b\}, \{c\} \Rightarrow \mathrm{FD}(\hat{x},x^*) = 1$
- FD in causal:

# edges in  $\hat{x}$  minus #edges in a densest CPDAG that contains conditional dependencies encoded in both  $\hat{x}, x^*$ 



Other appropriate similarity valuations in e.g. total ranking, subspace selection and blind source separation

## Greedy approaches to model selection

Starting from least model, greedily grow model complexity
Key ingredients:

- data-driven function  $\Psi$ : measures statistical significance for moving between neighboring models
- ullet minimal set of neighboring models  ${\mathcal S}$ : accounting for invariances

**Theorem:**  $\Psi_{\text{stable}}$ : based on subsampling and stability of a base procedure, and  $\Psi_{\text{test}}$ : based on testing; used-specified  $\alpha \in (0, 1)$ 

$$\Psi_{\text{stable}}: \quad \mathbb{E}[\text{FD}(\hat{x}, x^*)] \leq \sum_{k} \frac{q_k^2}{|\mathcal{S}_k|(1 - 2\alpha)},$$

$$\Psi_{\text{test}}: \quad \mathbb{P}\left(\text{FD}(\hat{x}, x^{\star}) > 0\right) \leq \alpha |\mathcal{S}|.$$

- $S_k$  = restriction of S to a specific rank
- $q_k = \text{avg.}$  discoveries by base procedure w.r.t. specific rank

# Experiments

Ranking educational systems: improving ranking of countries based on new PISA test scores: base ranking from 2015 scores

• new ranking from 2018 test scores using our algorithm with  $\Psi_{test}$  with family-wise-error control at level 0.05

Causal discovery from biological data: identifying causal relationships among proteins from Sachs dataset

• CPDAG estimated using our algorithm with desired FD level = 2

