

Ordinal Causal Discovery

Yang Ni Bani Mallick

Department of Statistics Texas AM University

Presented by Wenbo Zhang

Table of Contents

1 Preliminaries: Causal Discovery

2 Motivation

3 Ordinary Causal Discovery

- Bivariate Ordinal Causal Discovery
- Identifiability
- Extension to Multivariate Ordinal Causal Discovery

4 Experiments

- Synthetic Data
- Real Data

Table of Contents

1 Preliminaries: Causal Discovery

2 Motivation

3 Ordinary Causal Discovery

- Bivariate Ordinal Causal Discovery
- Identifiability
- Extension to Multivariate Ordinal Causal Discovery

4 Experiments

- Synthetic Data
- Real Data

Causal Discovery

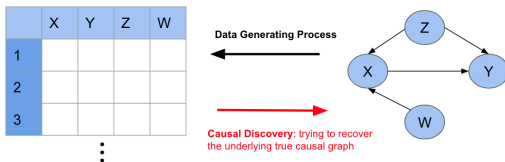
- Given a causal graph, we can do causal inference. What if we don't have the causal graph?

Causal Discovery

- Given a causal graph, we can do causal inference. What if we don't have the causal graph?
- Goal of Causal Discovery: identify the causal graph from data. (Here we only consider using observational data)

Causal Discovery

- Given a causal graph, we can do causal inference. What if we don't have the causal graph?
- Goal of Causal Discovery: identify the causal graph from data. (Here we only consider using observational data)
- Example:



X: Drug Dose

Y: Treatment Effects

Z: Age

W: Income

Definition

Consider a set of random variables $X = (X_1, \dots, X_p)$ with index set $V := \{1, \dots, p\}$. A causal graph of X is a directed acyclic graph (DAG). Each node corresponds to a single variable X_i , where $i \in \{1, \dots, p\}$ and each edge represents a direct causal relationship between variables.

Definition

Consider a set of random variables $X = (X_1, \dots, X_p)$ with index set $V := \{1, \dots, p\}$. A causal graph of X is a directed acyclic graph (DAG). Each node corresponds to a single variable X_i , where $i \in \{1, \dots, p\}$ and each edge represents a direct causal relationship between variables.

- Structural Equation Model (SEM)

$$X_j = f_j(\text{PA}_j, N_j), \quad j = 1, \dots, p$$

Definitions

Definition: d-separation

Two disjoint subsets of vertices A and B are d-separated by a third (also disjoint) subset C if every path between nodes in A and B is blocked by C .

Definition: Markov

The joint distribution $\mathcal{L}(X)$ is said to be Markov with respect to the DAG \mathcal{G} if

$$A, B \text{ d-sep. by } C \Rightarrow A \perp B \mid C$$

Notation

We denote by $\mathcal{M}(\mathcal{G})$ the set of distributions that are Markov with respect to \mathcal{G} :

$$\mathcal{M}(\mathcal{G}) := \{\mathcal{L}(X) : \mathcal{L}(X) \text{ is Markov w.r.t. } \mathcal{G}\}$$

Markov Equivalence Class (MEC)

Definition: Markov Equivalent

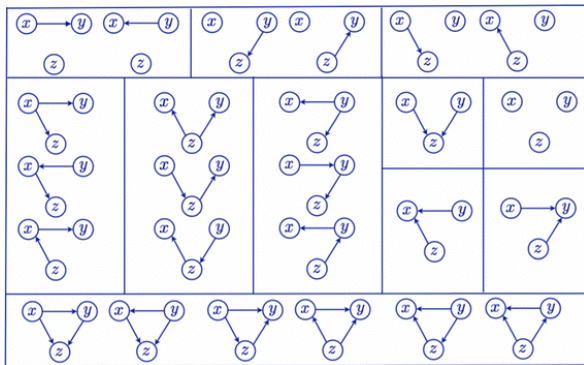
Two DAGs \mathcal{G}_1 and \mathcal{G}_2 are Markov equivalent if $\mathcal{M}(\mathcal{G}_1) = \mathcal{M}(\mathcal{G}_2)$.

This is the case if and only if \mathcal{G}_1 and \mathcal{G}_2 satisfy the same set of d-separations, which means the Markov condition entails the same set of (conditional) independence conditions.

Definition: Markov Equivalence Class (MEC)

If the set of DAGs are Markov equivalent to some DAG (e.g., encode the same set of conditional independence conditions), then they are in the Markov equivalence class (MEC).

Example of MEC with 3 Nodes



From <https://link.springer.com/article/10.1007/s41060-016-0038-6>

Lemma

Two DAGs are Markov equivalent if and only if they have the same skeleton and the same immoralities.

Important Facts

Assumption: Faithfulness and Causal Sufficiency

- Faithfulness: A, B d -sep. by $C \Leftarrow A \perp B \mid C$
- Causal Sufficiency: There are no unobserved confounders of any of the variables in the graph.

Facts

Without loss of generality, given the above assumptions and no further parametric or semi-parametric assumptions, we can only identify the DAG up to its Markov equivalence class (MEC) from observational data,

But not all further assumptions are helpful. There are a few examples

- Linear Gaussian: only identify MEC
- Multinomial: only identify MEC
- Nonlinear Additive Noise: can identify true DAG ...

Table of Contents

1 Preliminaries: Causal Discovery

2 Motivation

3 Ordinary Causal Discovery

- Bivariate Ordinal Causal Discovery
- Identifiability
- Extension to Multivariate Ordinal Causal Discovery

4 Experiments

- Synthetic Data
- Real Data

- While causal discovery for purely observational categorical data has been extensively studied, the vast majority of existing methods have exclusively focused on nominal (unordered) categorical variables.

Motivation

- While causal discovery for purely observational categorical data has been extensively studied, the vast majority of existing methods have exclusively focused on nominal (unordered) categorical variables.
- However, many categorical data contain **ordinal** information.

- While causal discovery for purely observational categorical data has been extensively studied, the vast majority of existing methods have exclusively focused on nominal (unordered) categorical variables.
- However, many categorical data contain **ordinal** information.
- Example:
psychologists often use questionnaires to measure latent traits such as personality and depression. The responses to those questionnaires are often categorical, say, with five levels (5-point Likert scale): “strongly disagree”, “disagree”, “neutral”, “agree”, and “strongly agree”.

- While causal discovery for purely observational categorical data has been extensively studied, the vast majority of existing methods have exclusively focused on nominal (unordered) categorical variables.
- However, many categorical data contain **ordinal** information.
- Example:
psychologists often use questionnaires to measure latent traits such as personality and depression. The responses to those questionnaires are often categorical, say, with five levels (5-point Likert scale): “strongly disagree”, “disagree”, “neutral”, “agree”, and “strongly agree”.
- Utilizing ordinal information of categorical data in causal discovery should be beneficial.

Table of Contents

1 Preliminaries: Causal Discovery

2 Motivation

3 Ordinary Causal Discovery

- Bivariate Ordinal Causal Discovery
- Identifiability
- Extension to Multivariate Ordinal Causal Discovery

4 Experiments

- Synthetic Data
- Real Data

- Firstly, an Ordinal Causal Discovery (OCD) method for bivariate data is introduced. Let $(X, Y) \in \{1, \dots, S\} \times \{1, \dots, L\}$ denote a pair of ordinal variables with S and L levels.

- Firstly, an Ordinal Causal Discovery (OCD) method for bivariate data is introduced. Let $(X, Y) \in \{1, \dots, S\} \times \{1, \dots, L\}$ denote a pair of ordinal variables with S and L levels.
- Model:
The bivariate OCD considers the following probability distribution for causal model $X \rightarrow Y$,

$$p_{X \rightarrow Y}(X, Y) = p(X)p(Y | X),$$

where $p(X)$ is a multinomial/categorical distribution with probabilities $\pi = (\pi_1, \dots, \pi_S)$ with $\sum_{s=1}^S \pi_s = 1$, and $p(Y | X)$ is defined by an ordinal regression model:

$$\Pr(Y \leq \ell | X = s) = F(\gamma_\ell - \beta_s), \quad \ell = 1, \dots, L, \quad s = 1, \dots, S$$

- Above equation implies the conditional probability distribution:

$$\Pr(Y = \ell \mid X = s) = F(\gamma_\ell - \beta_s) - F(\gamma_{\ell-1} - \beta_s),$$

where $\gamma_0 = -\infty$ and $\gamma_L = \infty$.

- Above equation implies the conditional probability distribution:

$$\Pr(Y = \ell \mid X = s) = F(\gamma_\ell - \beta_s) - F(\gamma_{\ell-1} - \beta_s),$$

where $\gamma_0 = -\infty$ and $\gamma_L = \infty$.

- Let $\beta = (\beta_1, \dots, \beta_S)$ and $\gamma = (\gamma_2, \dots, \gamma_{L-1})$.

We can denote the model $p_{X \rightarrow Y}$ by $p_{X \rightarrow Y}(X, Y \mid \pi, \beta, \gamma)$

Similarly, we denote the model $p_{Y \rightarrow X}$ by $p_{Y \rightarrow X}(Y, X \mid \rho, \alpha, \eta)$.

- Above equation implies the conditional probability distribution:

$$\Pr(Y = \ell \mid X = s) = F(\gamma_\ell - \beta_s) - F(\gamma_{\ell-1} - \beta_s),$$

where $\gamma_0 = -\infty$ and $\gamma_L = \infty$.

- Let $\beta = (\beta_1, \dots, \beta_S)$ and $\gamma = (\gamma_2, \dots, \gamma_{L-1})$.

We can denote the model $p_{X \rightarrow Y}$ by $p_{X \rightarrow Y}(X, Y \mid \pi, \beta, \gamma)$

Similarly, we denote the model $p_{Y \rightarrow X}$ by $p_{Y \rightarrow X}(Y, X \mid \rho, \alpha, \eta)$.

- Method: given observations of (X, Y) , we can get two maximal likelihood estimates (MLE) $\hat{p}_{X \rightarrow Y}$ and $\hat{p}_{Y \rightarrow X}$.

If $\hat{p}_{X \rightarrow Y} > \hat{p}_{Y \rightarrow X}$, then $X \rightarrow Y$ is deemed a more likely data generating causal model.

Does the joint distribution $p(X, Y)$ tell us whether it has been induced by X to Y or from Y to X ? In other words, is the structure identifiable from the joint distribution?

Definition: Distribution Equivalence

$p_{X \rightarrow Y}(X, Y \mid \pi, \beta, \gamma)$ and $p_{Y \rightarrow X}(Y, X \mid \rho, \alpha, \eta)$ are distribution equivalent if for any values of (π, β, γ) there exist values of (ρ, α, η) such that $p_{X \rightarrow Y}(X, Y \mid \pi, \beta, \gamma) = p_{Y \rightarrow X}(Y, X \mid \rho, \alpha, \eta)$ for any X, Y , and vice versa.

- Distribution equivalent causal models are clearly not distinguishable from each other by examining their observational distributions.

- Example:
The multinomial causal Bayesian Networks are distribution equivalent.

$P(Y X)$	$X = 1$	$X = 2$	$X = 3$
$Y = 1$	0.16	0.84	0.16
$Y = 2$	0.34	0.14	0.34
$Y = 3$	0.50	0.02	0.50

(a)

X	$P(X)$
1	0.25
2	0.25
3	0.50

×

$P(X,Y)$	$X = 1$	$X = 2$	$X = 3$
$Y = 1$	0.040	0.210	0.080
$Y = 2$	0.085	0.035	0.170
$Y = 3$	0.125	0.005	0.250

(b)

$P(X Y)$	$Y = 1$	$Y = 2$	$Y = 3$
$X = 1$	0.12	0.29	0.33
$X = 2$	0.64	0.12	0.01
$X = 3$	0.24	0.59	0.66

(c)

Y	$P(Y)$
1	0.33
2	0.29
3	0.38

×

- Example:
The multinomial causal Bayesian Networks are distribution equivalent.

<table border="1" style="border-collapse: collapse; text-align: center;"> <thead> <tr> <th>$P(Y X)$</th> <th>$X = 1$</th> <th>$X = 2$</th> <th>$X = 3$</th> </tr> </thead> <tbody> <tr> <td>$Y = 1$</td> <td>0.16</td> <td>0.84</td> <td>0.16</td> </tr> <tr> <td>$Y = 2$</td> <td>0.34</td> <td>0.14</td> <td>0.34</td> </tr> <tr> <td>$Y = 3$</td> <td>0.50</td> <td>0.02</td> <td>0.50</td> </tr> </tbody> </table>	$P(Y X)$	$X = 1$	$X = 2$	$X = 3$	$Y = 1$	0.16	0.84	0.16	$Y = 2$	0.34	0.14	0.34	$Y = 3$	0.50	0.02	0.50	×	<table border="1" style="border-collapse: collapse; text-align: center;"> <thead> <tr> <th>X</th> <th>$P(X)$</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>0.25</td> </tr> <tr> <td>2</td> <td>0.25</td> </tr> <tr> <td>3</td> <td>0.50</td> </tr> </tbody> </table>	X	$P(X)$	1	0.25	2	0.25	3	0.50	≡	<table border="1" style="border-collapse: collapse; text-align: center;"> <thead> <tr> <th>$P(X, Y)$</th> <th>$X = 1$</th> <th>$X = 2$</th> <th>$X = 3$</th> </tr> </thead> <tbody> <tr> <td>$Y = 1$</td> <td>0.040</td> <td>0.210</td> <td>0.080</td> </tr> <tr> <td>$Y = 2$</td> <td>0.085</td> <td>0.035</td> <td>0.170</td> </tr> <tr> <td>$Y = 3$</td> <td>0.125</td> <td>0.005</td> <td>0.250</td> </tr> </tbody> </table>	$P(X, Y)$	$X = 1$	$X = 2$	$X = 3$	$Y = 1$	0.040	0.210	0.080	$Y = 2$	0.085	0.035	0.170	$Y = 3$	0.125	0.005	0.250
$P(Y X)$	$X = 1$	$X = 2$	$X = 3$																																									
$Y = 1$	0.16	0.84	0.16																																									
$Y = 2$	0.34	0.14	0.34																																									
$Y = 3$	0.50	0.02	0.50																																									
X	$P(X)$																																											
1	0.25																																											
2	0.25																																											
3	0.50																																											
$P(X, Y)$	$X = 1$	$X = 2$	$X = 3$																																									
$Y = 1$	0.040	0.210	0.080																																									
$Y = 2$	0.085	0.035	0.170																																									
$Y = 3$	0.125	0.005	0.250																																									
(a)				(b)																																								

<table border="1" style="border-collapse: collapse; text-align: center;"> <thead> <tr> <th>$P(X Y)$</th> <th>$Y = 1$</th> <th>$Y = 2$</th> <th>$Y = 3$</th> </tr> </thead> <tbody> <tr> <td>$X = 1$</td> <td>0.12</td> <td>0.29</td> <td>0.33</td> </tr> <tr> <td>$X = 2$</td> <td>0.64</td> <td>0.12</td> <td>0.01</td> </tr> <tr> <td>$X = 3$</td> <td>0.24</td> <td>0.59</td> <td>0.66</td> </tr> </tbody> </table>	$P(X Y)$	$Y = 1$	$Y = 2$	$Y = 3$	$X = 1$	0.12	0.29	0.33	$X = 2$	0.64	0.12	0.01	$X = 3$	0.24	0.59	0.66	×	<table border="1" style="border-collapse: collapse; text-align: center;"> <thead> <tr> <th>Y</th> <th>$P(Y)$</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>0.33</td> </tr> <tr> <td>2</td> <td>0.29</td> </tr> <tr> <td>3</td> <td>0.38</td> </tr> </tbody> </table>	Y	$P(Y)$	1	0.33	2	0.29	3	0.38	≡	
$P(X Y)$	$Y = 1$	$Y = 2$	$Y = 3$																									
$X = 1$	0.12	0.29	0.33																									
$X = 2$	0.64	0.12	0.01																									
$X = 3$	0.24	0.59	0.66																									
Y	$P(Y)$																											
1	0.33																											
2	0.29																											
3	0.38																											
(c)																												

- As the example shows we can find a set of parameters, i.e., the conditional $p(X | Y)$ and marginal $p(Y)$ probabilities of the reverse causal model $Y \rightarrow X$, which leads to the same joint distribution as the causal model. Therefore, not identifiable.

- Example:
The multinomial causal Bayesian Networks are distribution equivalent.

<table border="1" style="border-collapse: collapse; text-align: center;"> <thead> <tr> <th>$P(Y X)$</th> <th>$X = 1$</th> <th>$X = 2$</th> <th>$X = 3$</th> </tr> </thead> <tbody> <tr> <td>$Y = 1$</td> <td>0.16</td> <td>0.84</td> <td>0.16</td> </tr> <tr> <td>$Y = 2$</td> <td>0.34</td> <td>0.14</td> <td>0.34</td> </tr> <tr> <td>$Y = 3$</td> <td>0.50</td> <td>0.02</td> <td>0.50</td> </tr> </tbody> </table>	$P(Y X)$	$X = 1$	$X = 2$	$X = 3$	$Y = 1$	0.16	0.84	0.16	$Y = 2$	0.34	0.14	0.34	$Y = 3$	0.50	0.02	0.50	×	<table border="1" style="border-collapse: collapse; text-align: center;"> <thead> <tr> <th>X</th> <th>$P(X)$</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>0.25</td> </tr> <tr> <td>2</td> <td>0.25</td> </tr> <tr> <td>3</td> <td>0.50</td> </tr> </tbody> </table>	X	$P(X)$	1	0.25	2	0.25	3	0.50	≡	<table border="1" style="border-collapse: collapse; text-align: center;"> <thead> <tr> <th>$P(X, Y)$</th> <th>$X = 1$</th> <th>$X = 2$</th> <th>$X = 3$</th> </tr> </thead> <tbody> <tr> <td>$Y = 1$</td> <td>0.040</td> <td>0.210</td> <td>0.080</td> </tr> <tr> <td>$Y = 2$</td> <td>0.085</td> <td>0.035</td> <td>0.170</td> </tr> <tr> <td>$Y = 3$</td> <td>0.125</td> <td>0.005</td> <td>0.250</td> </tr> </tbody> </table>	$P(X, Y)$	$X = 1$	$X = 2$	$X = 3$	$Y = 1$	0.040	0.210	0.080	$Y = 2$	0.085	0.035	0.170	$Y = 3$	0.125	0.005	0.250
$P(Y X)$	$X = 1$	$X = 2$	$X = 3$																																									
$Y = 1$	0.16	0.84	0.16																																									
$Y = 2$	0.34	0.14	0.34																																									
$Y = 3$	0.50	0.02	0.50																																									
X	$P(X)$																																											
1	0.25																																											
2	0.25																																											
3	0.50																																											
$P(X, Y)$	$X = 1$	$X = 2$	$X = 3$																																									
$Y = 1$	0.040	0.210	0.080																																									
$Y = 2$	0.085	0.035	0.170																																									
$Y = 3$	0.125	0.005	0.250																																									
(a)				(b)																																								

<table border="1" style="border-collapse: collapse; text-align: center;"> <thead> <tr> <th>$P(X Y)$</th> <th>$Y = 1$</th> <th>$Y = 2$</th> <th>$Y = 3$</th> </tr> </thead> <tbody> <tr> <td>$X = 1$</td> <td>0.12</td> <td>0.29</td> <td>0.33</td> </tr> <tr> <td>$X = 2$</td> <td>0.64</td> <td>0.12</td> <td>0.01</td> </tr> <tr> <td>$X = 3$</td> <td>0.24</td> <td>0.59</td> <td>0.66</td> </tr> </tbody> </table>	$P(X Y)$	$Y = 1$	$Y = 2$	$Y = 3$	$X = 1$	0.12	0.29	0.33	$X = 2$	0.64	0.12	0.01	$X = 3$	0.24	0.59	0.66	×	<table border="1" style="border-collapse: collapse; text-align: center;"> <thead> <tr> <th>Y</th> <th>$P(Y)$</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>0.33</td> </tr> <tr> <td>2</td> <td>0.29</td> </tr> <tr> <td>3</td> <td>0.38</td> </tr> </tbody> </table>	Y	$P(Y)$	1	0.33	2	0.29	3	0.38	≡	
$P(X Y)$	$Y = 1$	$Y = 2$	$Y = 3$																									
$X = 1$	0.12	0.29	0.33																									
$X = 2$	0.64	0.12	0.01																									
$X = 3$	0.24	0.59	0.66																									
Y	$P(Y)$																											
1	0.33																											
2	0.29																											
3	0.38																											
(c)																												

- As the example shows we can find a set of parameters, i.e., the conditional $p(X | Y)$ and marginal $p(Y)$ probabilities of the reverse causal model $Y \rightarrow X$, which leads to the same joint distribution as the causal model. Therefore, not identifiable.
- However, if incorporating the ordinal information, we will show that $p_{X \rightarrow Y}(X, Y | \pi, \beta, \gamma)$ and $p_{Y \rightarrow X}(Y, X | \rho, \alpha, \eta)$ are generally not distribution equivalent and are, therefore, identifiable.

Theorem (Identifiability of OCD)

Let $X \in \{1, \dots, S\}$ and $Y \in \{1, \dots, L\}$ where $S, L \geq 2$ and $\max\{S, L\} \geq 3$. Suppose $X \rightarrow Y$ is the data generating causal model and the observational probability distribution of (X, Y) is given by:

$$p(X, Y) = p_{X \rightarrow Y}(X, Y \mid \pi, \beta, \gamma)$$

For almost all (π, β, γ) with respect to the Lebesgue measure, the distribution cannot be equivalently represented by the reverse causal model, i.e., there does not exist (ρ, α, η) such that,

$$p(X, Y) = p_{Y \rightarrow X}(Y, X \mid \rho, \alpha, \eta), \forall X, Y$$

Ordinal Bayesian Networks

- Given a large enough dataset, in the causal direction, $\hat{p}_{X \rightarrow Y}(X, Y)$ can be arbitrarily close to the true $p(X, Y)$. However, there does not exist any set of parameter values in the reverse causal model $p_{Y \rightarrow X}(Y, X \mid \rho, \alpha, \eta)$ that produces the conditional $p(X \mid Y)$ and marginal $p(Y)$ probability such that $\hat{p}_{Y \rightarrow X}(X, Y)$ close $p(X, Y)$.

Ordinal Bayesian Networks

- Given a large enough dataset, in the causal direction, $\hat{p}_{X \rightarrow Y}(X, Y)$ can be arbitrarily close to the true $p(X, Y)$. However, there does not exist any set of parameter values in the reverse causal model $p_{Y \rightarrow X}(Y, X | \rho, \alpha, \eta)$ that produces the conditional $p(X | Y)$ and marginal $p(Y)$ probability such that $\hat{p}_{Y \rightarrow X}(X, Y)$ close $p(X, Y)$.
- Example: The ordinal causal Bayesian Networks

$P(Y X)$	$X = 1$	$X = 2$	$X = 3$
$Y = 1$	0.16	0.84	0.16
$Y = 2$	0.34	0.14	0.34
$Y = 3$	0.50	0.02	0.50

×

X	$P(X)$
1	0.25
2	0.25
3	0.50

(a)

$P(X, Y)$	$X = 1$	$X = 2$	$X = 3$
$Y = 1$	0.040	0.210	0.080
$Y = 2$	0.085	0.035	0.170
$Y = 3$	0.125	0.005	0.250

⊥

(b)

$P(X Y)$	$Y = 1$	$Y = 2$	$Y = 3$
$X = 1$	0.12	0.29	0.33
$X = 2$	0.64	0.12	0.01
$X = 3$	0.24	0.59	0.66

×

Y	$P(Y)$
1	0.33
2	0.29
3	0.38

(c)

$\hat{P}(X, Y)$	$X = 1$	$X = 2$	$X = 3$
$Y = 1$	0.104	0.087	0.141
$Y = 2$	0.065	0.071	0.154
$Y = 3$	0.078	0.089	0.211

(d)

Ordinal Bayesian Networks

- Given a large enough dataset, in the causal direction, $\hat{p}_{X \rightarrow Y}(X, Y)$ can be arbitrarily close to the true $p(X, Y)$. However, there does not exist any set of parameter values in the reverse causal model $p_{Y \rightarrow X}(Y, X \mid \rho, \alpha, \eta)$ that produces the conditional $p(X \mid Y)$ and marginal $p(Y)$ probability such that $\hat{p}_{Y \rightarrow X}(X, Y)$ close $p(X, Y)$.
- Example: The ordinal causal Bayesian Networks

$P(Y X)$	$X = 1$	$X = 2$	$X = 3$
$Y = 1$	0.16	0.84	0.16
$Y = 2$	0.34	0.14	0.34
$Y = 3$	0.50	0.02	0.50

×

X	$P(X)$
1	0.25
2	0.25
3	0.50

+

$P(X, Y)$	$X = 1$	$X = 2$	$X = 3$
$Y = 1$	0.040	0.210	0.080
$Y = 2$	0.085	0.035	0.170
$Y = 3$	0.125	0.005	0.250

‡

$P(X Y)$	$Y = 1$	$Y = 2$	$Y = 3$
$X = 1$	0.12	0.29	0.33
$X = 2$	0.64	0.12	0.01
$X = 3$	0.24	0.59	0.66

×

Y	$P(Y)$
1	0.33
2	0.29
3	0.38

+

$\hat{P}(X, Y)$	$X = 1$	$X = 2$	$X = 3$
$Y = 1$	0.104	0.087	0.141
$Y = 2$	0.065	0.071	0.154
$Y = 3$	0.078	0.089	0.211

- Therefore, $p_{X \rightarrow Y}(X, Y \mid \pi, \beta, \gamma)$ can be distinguished from $p_{Y \rightarrow X}(Y, X \mid \rho, \alpha, \eta)$.

Multivariate Ordinal Causal Discovery

Let $\mathbf{X} = (X_1, \dots, X_p) \in \{1, \dots, L_1\} \times \dots \times \{1, \dots, L_p\}$ denote p ordinal variables. Let $G = (V, E)$ denote a causal BN with a set of nodes $V = \{1, \dots, p\}$ directed edges E . Let $pa(j) = \{k \mid k \rightarrow j\} \subseteq V$ denote the set of parents of node j in G and let $\mathbf{X}_{pa(j)} = \{X_k \mid k \in pa(j)\}$. Given G , the joint distribution of \mathbf{X} factorizes,

$$p(\mathbf{X} \mid G) = \prod_{j=1}^p p(X_j \mid \mathbf{X}_{pa(j)}),$$

where each conditional distribution $p(X_j \mid \mathbf{X}_{pa(j)})$ is an ordinal regression model of which the cumulative distribution is given by, for $\ell = 1, \dots, L_j$

$$\Pr(X_j \leq \ell \mid \mathbf{X}_{pa(j)}) = F \left(\gamma_{j\ell} - \sum_{k \in pa(j)} \beta_{jk} X_k - \alpha_j \right)$$

The implied conditional probability distribution

$$\begin{aligned} & \Pr(X_j = \ell \mid \mathbf{X}_{pa(j)} = \mathbf{s}) \\ &= F\left(\gamma_{j\ell} - \sum_{k \in pa(j)} \beta_{jkh_k} - \alpha_j\right) - F\left(\gamma_{j,\ell-1} - \sum_{k \in pa(j)} \beta_{jkh_k} - \alpha_j\right) \end{aligned}$$

where $\ell = 1, \dots, L_j$ and $\mathbf{s} \in \prod_{k \in pa(j)} \{1, \dots, L_k\}$.

Causal Graph Structure Learning

A score-and-search learning algorithm is proposed to estimate the structure of causal graphs.

- Score. We score causal graphs by the Bayesian information criterion (BIC). Let $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ denote n realizations of \mathbf{X} . The score of G (smaller is better) is given by

$$\text{BIC}(G | \mathbf{x}) = -2 \sum_{i=1}^n \log \hat{p}(\mathbf{x}_i | G) + K \log(n)$$

Causal Graph Structure Learning

A score-and-search learning algorithm is proposed to estimate the structure of causal graphs.

- Score. We score causal graphs by the Bayesian information criterion (BIC). Let $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ denote n realizations of \mathbf{X} . The score of G (smaller is better) is given by

$$\text{BIC}(G | \mathbf{x}) = -2 \sum_{i=1}^n \log \hat{p}(\mathbf{x}_i | G) + K \log(n)$$

- Search. Exhaustive Search are utilized for small networks and Greedy Search is used for moderate-sized networks.

Algorithm 1 Greedy Search

Input: data \mathbf{x} , initial graph G
Compute $\text{BIC}(G|\mathbf{x})$ and set $\text{BIC}_* = \text{BIC}(G|\mathbf{x})$.
repeat
 Initialize $\text{Improvement} = \text{false}$.
 for all graphs G' reachable from G **do**
 Compute $\text{BIC}(G'|\mathbf{x})$.
 if $\text{BIC}(G'|\mathbf{x}) < \text{BIC}_*$ **then**
 Set $G = G'$ and $\text{BIC}_* = \text{BIC}(G'|\mathbf{x})$
 Set $\text{Improvement} = \text{true}$.
 end if
 end for
until Improvement is false
Output: graph G

Table of Contents

1 Preliminaries: Causal Discovery

2 Motivation

3 Ordinary Causal Discovery

- Bivariate Ordinal Causal Discovery
- Identifiability
- Extension to Multivariate Ordinal Causal Discovery

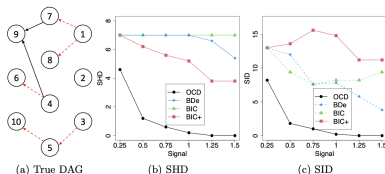
4 Experiments

- Synthetic Data
- Real Data

Synthetic Data

Let $n = 500$ and the number of categories $L = 5$ for each node.

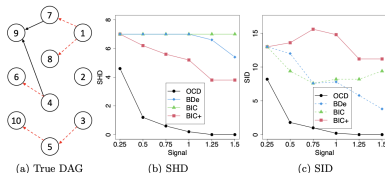
- Low-Dimensional Multivariate Ordinal Data



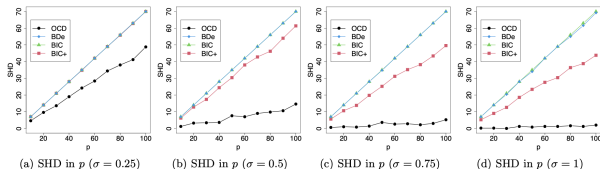
Synthetic Data

Let $n = 500$ and the number of categories $L = 5$ for each node.

- Low-Dimensional Multivariate Ordinal Data



- High-Dimensional Multivariate Ordinal Data



- Sachs's Single-Cell Flow Cytometry Data
853 cells and 11 phosphorylated proteins

	OCD	bQCD*	IGCI*	GR-AN*
SHD	14	15	15	16
SID	62	69	82	80
	HCR*	SLOPE*	ANM	LiNGAM
SHD	16	17	17	17
SID	76	86	78	86
	PC	CPC	GES	IAMB
SHD	18	18	18	20
SID	50-83	50-80	50-80	79-70
	BIC	BDe	MXM	RESIT
SHD	20	21	21	40
SID	53-77	49-104	49-104	45

- Sachs's Single-Cell Flow Cytometry Data
853 cells and 11 phosphorylated proteins

	OCD	bQCD*	IGCI*	GR-AN*
SHD	14	15	15	16
SID	62	69	82	80
	HCR*	SLOPE*	ANM	LiNGAM
SHD	16	17	17	17
SID	76	86	78	86
	PC	CPC	GES	IAMB
SHD	18	18	18	20
SID	50-83	50-80	50-80	79-70
	BIC	BDe	MXM	RESIT
SHD	20	21	21	40
SID	53-77	49-104	49-104	45

- Single-Cell RNA-Sequencing Data
6701 pairs and 2 variables (transcription factor X and its target Y)

	OCD	HCR	bQCD	SLOPE
ACC	0.61	0.36	0.45	0.50
CPU	19m	22m	3.4h	2h

1. Explain why identifiability doesn't hold for Multinomial Bayesian Networks (Hint: the example from the paper I show).

Thank you!