

Causal Discovery in the Presence of Missing Data

Ruibo Tu*1, Cheng Zhang*2, Paul Ackermann3, Karthika Mohan4, Hedvig Kjellström¹, and Kun Zhang^{*5} ¹KTH Royal Institute of Technology, ²Microsoft Research, Cambridge, ³Karolinska Institute, ⁴University of California, Berkeley, ⁵Carnegie Mellon University

What is the relationship

between IQ and gender?

Gender

female

female

female

male

male

male

College

not in

not in

in

in

in

in

IQ

high

low

high

low

medium

medium

Carnegie Mellon University Microsoft Research Karolinska Berkeley Institutet

INTRODUCTION

- Missing Data Problem: Analyze the influence of the missing data on the constraintbased causal discovery algorithms -- Propositions
- Solution: Propose a framework MVPC and two correction *methods* for correcting the results of the deletion-based PC algorithm.

DELETION-BASED PC

- List-wise deletion and Test-wise deletion
- Deletion might produce errors in the skeleton search step. Observed data distribution vs Full data distribution
- PC Recap



MISSINGNESS GRAPH

Missingness Graphs: missingness indicators and proxy variables



o Faithful observability

 $X \perp\!\!\!\perp Y \mid \{\mathbf{Z}, \mathbf{R}_{\mathbf{K}} = \mathbf{0}\} \Longleftrightarrow X \perp\!\!\!\perp Y \mid \{\mathbf{Z}, \mathbf{R}_{\mathbf{K}} = \mathbf{1}\}$

PROPOSITIONS

- Proposition 1: The erroneous edges are extraneous edges $\begin{array}{l} X \perp \!\!\!\!\perp Y \mid \{ \mathbf{Z}, R_x = 0, R_y = 0, \mathbf{R_z} = \mathbf{0} \} \Longrightarrow X \perp \!\!\!\!\!\perp Y \mid \mathbf{Z} \\ X \not \perp Y \mid \{ \mathbf{Z}, R_x = 0, R_y = 0, \mathbf{R_z} = \mathbf{0} \} \Longrightarrow X \not \perp Y \mid \mathbf{Z} \end{array}$
- **Proposition 2:** Determine circumstances where the extraneous edges happen. If there is an extraneous edge between X and Y, under the four assumptions there is at least a missingness indicator of X, Y, and Z is either the direct common effect or a descendant of the direct common effect of X and Y. (Z: any variable set which satisfies that given Z, X is independent of Y)



MVPC

- Skeleton Search (PC)
 - 1. Detecting direct causes of missingness indicators
 - 2. Detecting potential extraneous edges
 - 3. Recovering the true causal graph skeleton
- Determining the orientation (PC)

Recovering the true causal graph skeleton

- Get access to the data from the joint distribution P(X, Y, Z)
 - (A) Permutation-based correction:

$$\begin{split} P(X,Y,Z) &= \int_W P(X,Y,Z \mid W) P(W) dW \\ &= \int_W P(X,Y^*,Z \mid W,R_y = 0) P(W) dW \end{split}$$

(B) Density ratio weighted corection

$$P(X, Y, Z) = \frac{P(X, Y, Z \mid R_z = 0)P(R_z = 0)}{P(R_z = 0 \mid X, Y)}$$
$$= P(X, Y, Z \mid R_z = 0) \frac{P(X, Y)}{P(X, Y \mid R_z = 0)}$$

Experiments

Synthetic data

ideal 🗱 target 🗱 MVPC 🧮 TD-PC 💹 LD-PC Missing data in MAR Missing data in MNAR











Test-wise deletion PC



REFERENCES

- Mohan, Karthika, Judea Pearl, and Jin Tian. "Graphical models for inference with missing data." Advances in neural information processing systems. 2013.
- Strobl, Eric V., Shyam Visweswaran, and Peter L. Spirtes. "Fast causal inference with non-random missingness by test-wise deletion." International Journal of Data Science and Analytics (2017): 1-16.