

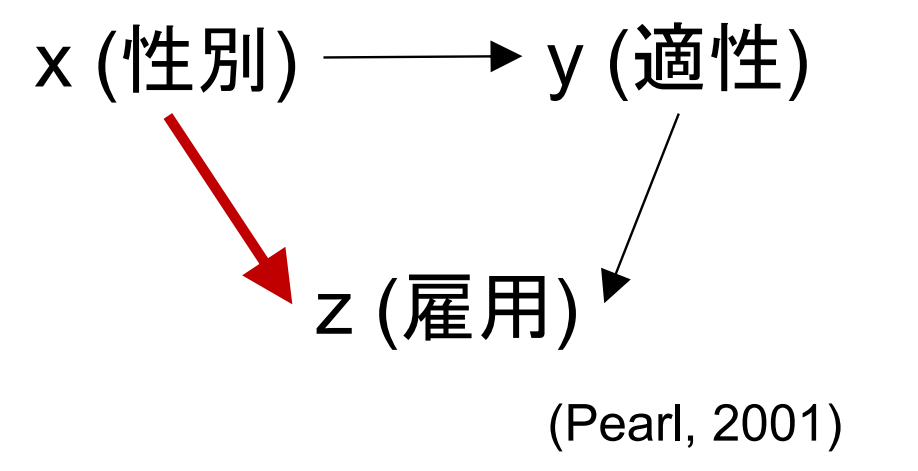
統計的因果探索とは

- データを用いて**因果グラフを推測**するための方法論



AIのための因果推論でも要(かなめ)

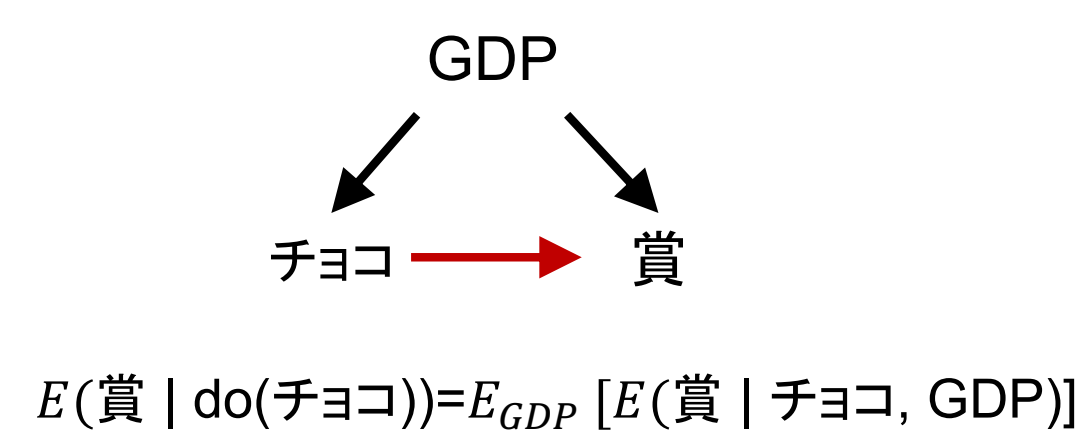
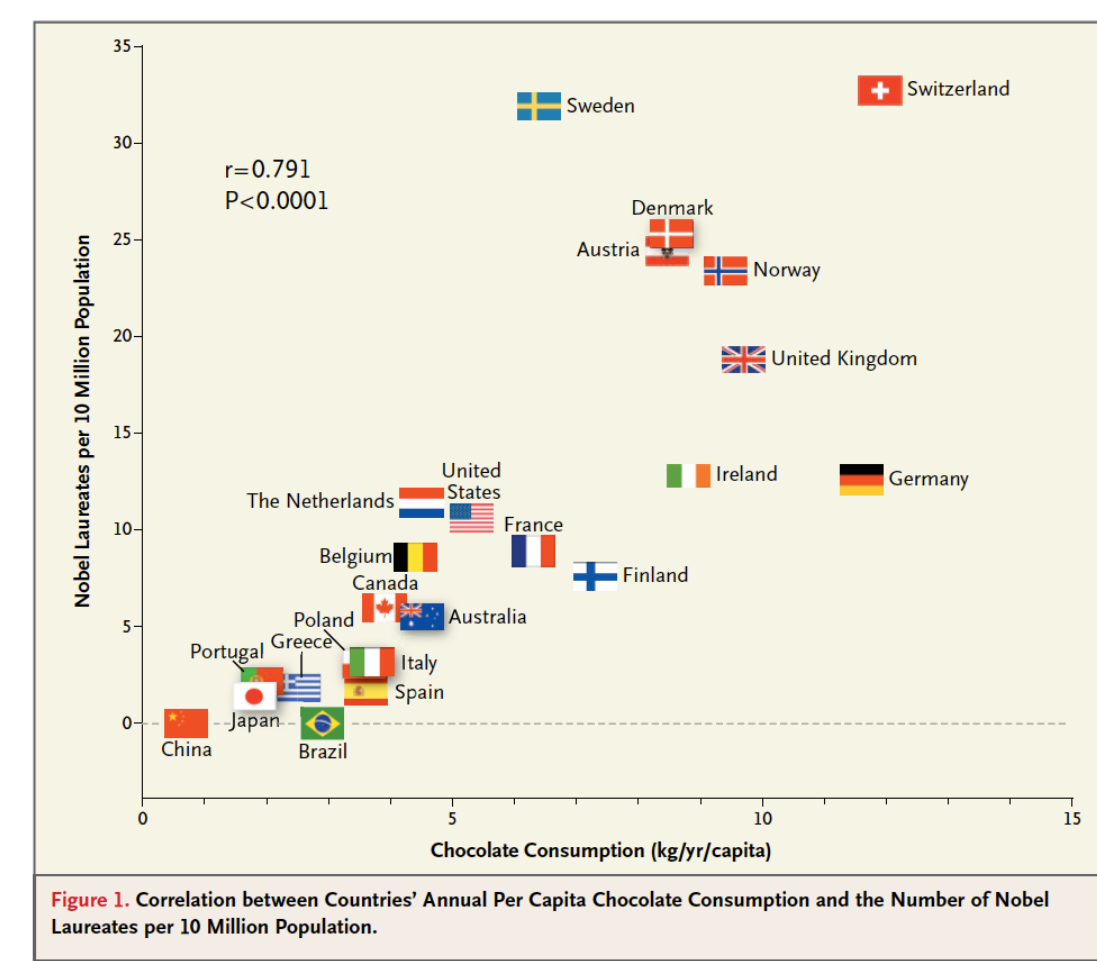
- 公平性 (Kusner et al., 2017)
- 説明性
 - 原因の確率 (Galhotra et al., 2021)
 - 予測メカニズム解析 (Bobaum et al., 2017; Sani et al., 2020)
 - 個体レベルの最適介入 (Kiritoshi et al., 2021)
- 転移学習 (Zhang et al., 2013; Zhang et al., 2020; Bareinboim et al., 2016)
- 科学的知識の取り込み (Teshima et al., 2021)
- さまざまな**因果に関するクエリ**(介入効果等)に答えられるかを判定するために**因果グラフが必要**



統計的因果推論では因果グラフが要(かなめ)

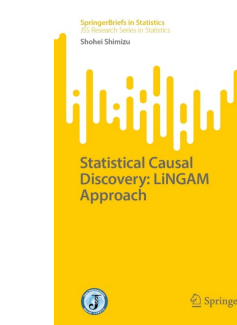
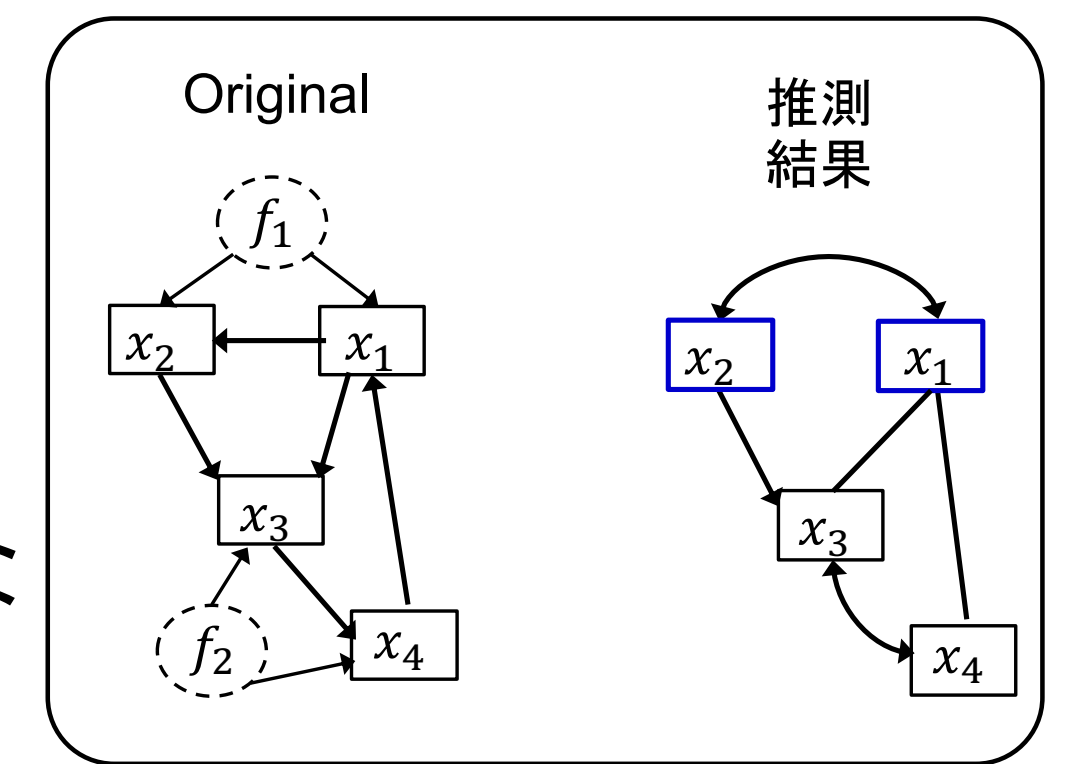
- データから**介入効果**を推定
 - チョコ消費量を変えるとノーベル賞受賞者の数はどのくらい増えるのか(減るのか)
 - 機械学習のする予測
 - チョコ消費がこのくらいならノーベル賞数このくらい?
 - ノーベル賞数がこのくらいならチョコ消費このくらい?
- 介入効果を「正しく」推定するには**因果グラフが必要** (e.g., バックドア基準)

Messerli, (2012), New England Journal of Medicine



研究内容

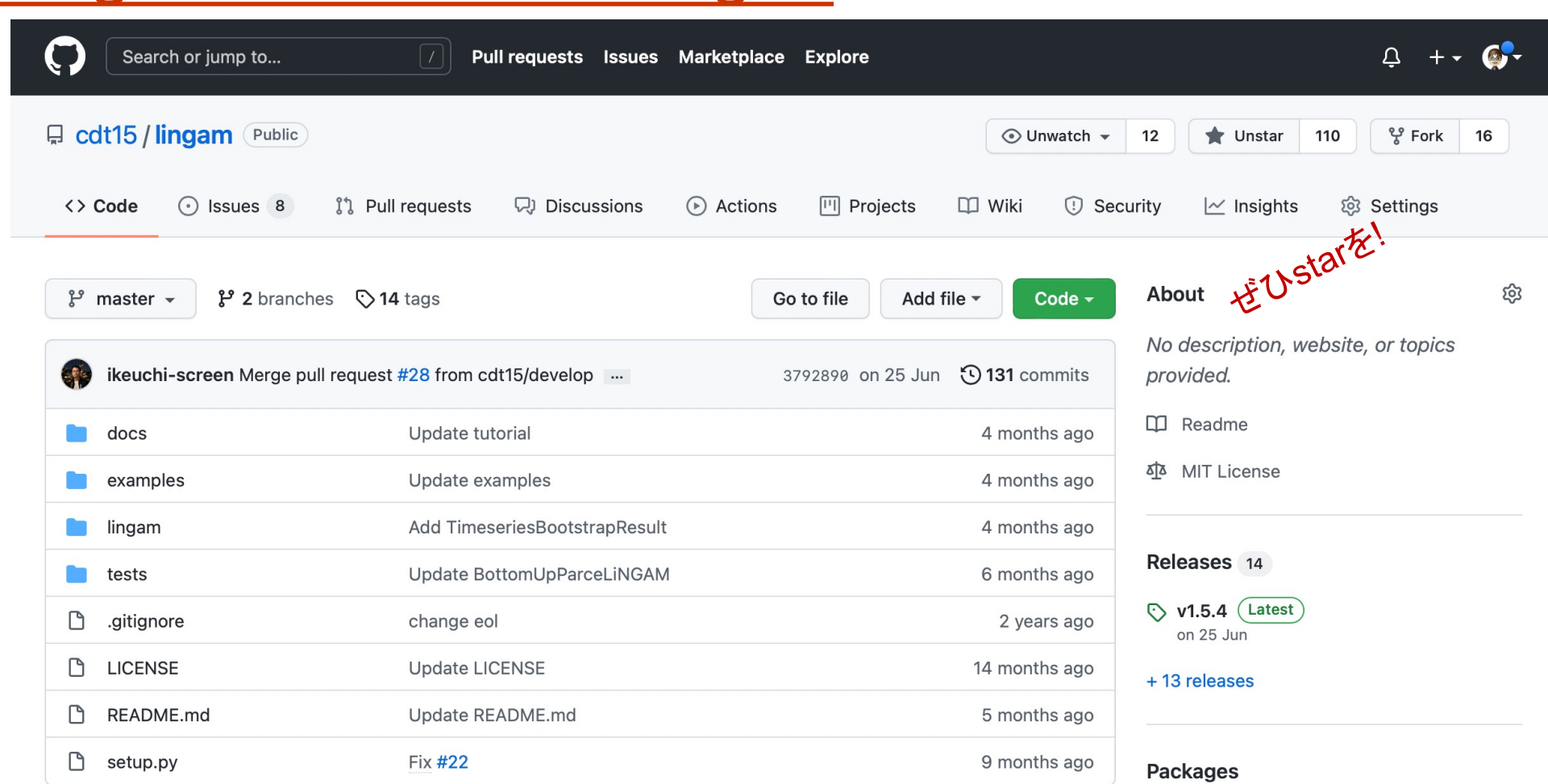
- 因果探索の研究開発
 - 未観測共通原因を許す
 - 離散と連続が混在
 - 非線形
 - 巡回
- データからわかる限界の理論的説明
- 領域知識とデータから、その時点で可能な最大限の因果グラフを分析者に提供



S. Shimizu. Statistical Causal Discovery: LiNGAM approach. SpringerBriefs in Statistics, Springer Japan, 2022.

LiNGAM Python package

- <https://github.com/cdt15/lingam>



T. Ikeuchi, M. Ide, Y. Zeng, T. N. Maeda, and S. Shimizu.
Python package for causal discovery based on LiNGAM.
Journal of Machine Learning Research, 24(14): 1-8, 2023.

- 有向道や有向辺のブートストラップ確率
- 例えば、閾値0.05を越えるものを解釈



モデル仮定の評価

- 分析前
 - Gaussianity test
 - ヒストグラム
 - 連続変数?
 - 多重共線性
 - 領域知識
- 分析後
 - 誤差(残差)の独立性評価
 - 例えば、HSIC (Gretton et al., 2005)
 - マルコフバウンダリーによる予測の良さで評価 (Biza et al., 2020) (実装準備中)
 - 複数のデータセットでの結果を比較
 - 領域知識による評価

連続変数と離散変数が混在する場合

Definition 1 (Linear Mixed causal model, LiM) If a causal model for mixed data satisfies assumptions A1-A3, then this SCM is called a Linear Mixed causal model, abbreviated as LiM.

因果グラフが識別可能

$$\text{Let } \mathcal{F} = \{f^{con}, f^{dis} | f^{con}(x, e) = bx + e + c, f^{dis}(x, e) = \begin{cases} 1, & bx + e + c > 0 \\ 0, & \text{otherwise} \end{cases}\} \text{ be a set}$$

of two functions which work on continuous and discrete variables, respectively, where x, e, b and c are a random variable of X , error term, coefficient and intercept, respectively. $\mathcal{P} = \{P^{con}, P^{dis}\}$ denotes the set of probabilistic distributions for continuous and discrete variables. Using these notations, our model can be rewritten as:

$$x_i = f_i(x_{pa(i)}, e_i), \quad e_i \sim P(e_i), \quad (3)$$

where $f_i \in \mathcal{F}$, and $P(e_i) \in \mathcal{P}$.

- A1. Observed variables x_i ($i = 1, \dots, p$) form a Directed Acyclic Graph (DAG).
- A2. The value assigned to each continuous variable x_i is a linear function of its parent variables denoted by $x_{pa(i)}$ plus a non-Gaussian error term e_i , that is,

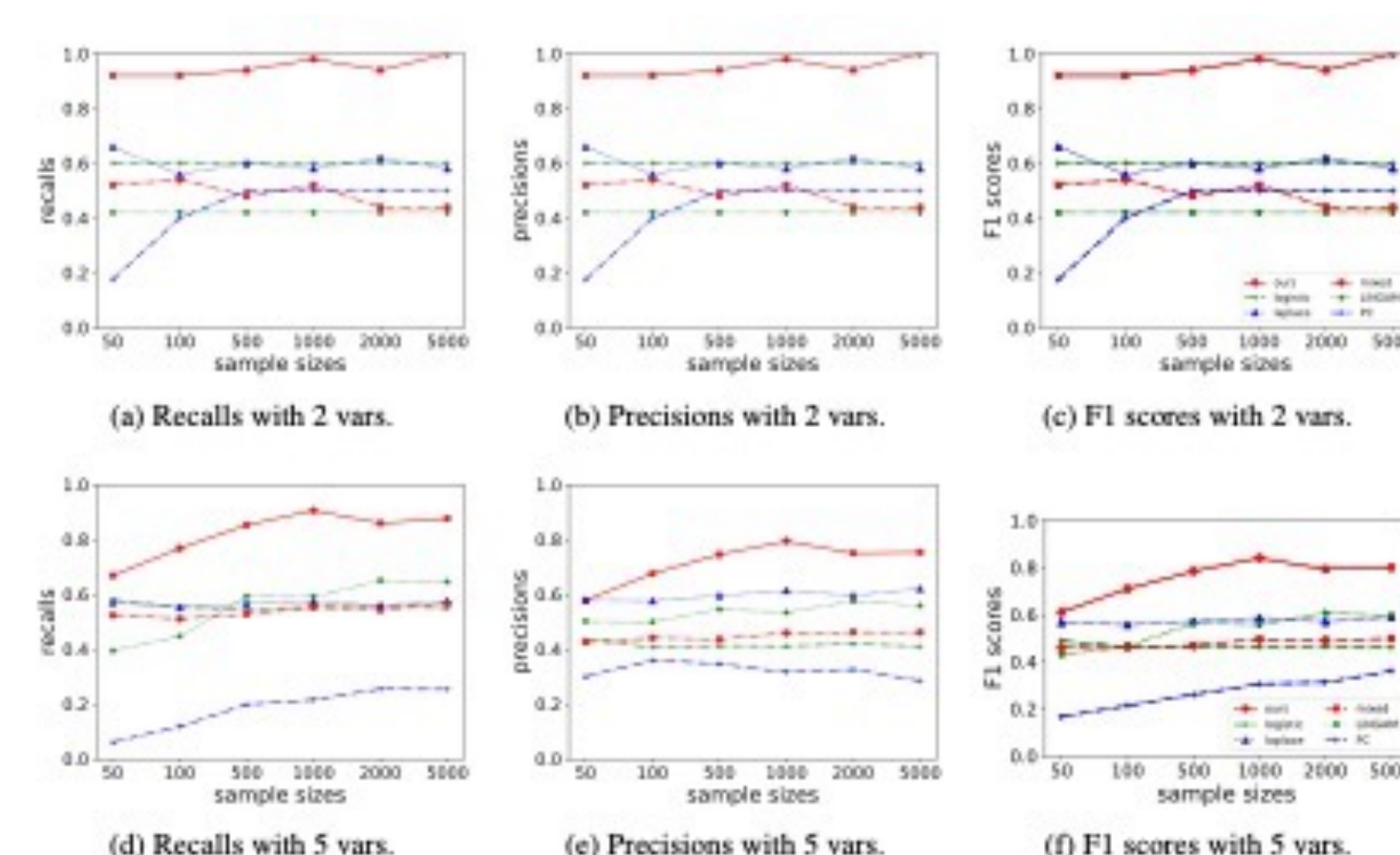
$$x_i = e_i + c_i + \sum_{j \in pa(i)} b_{ij} x_j, \quad e_i \sim \text{Non-Gaussian}(\cdot), \quad (1)$$

where the error terms e_i are continuous random variables with non-Gaussian densities, and the error variables e_i are independent of each other. The coefficients b_{ij} and intercepts c_i are constants.

- A3. For each discrete variable x_i , its value equals 1 if the linear function of its parent variables $x_{pa(i)}$ plus a Logistic error term e_i is larger than 0, otherwise, its value equals 0. That is,

$$x_i = \begin{cases} 1, & e_i + c_i + \sum_{j \in pa(i)} b_{ij} x_j > 0 \\ 0, & \text{otherwise} \end{cases}, \quad e_i \sim \text{Logistic}(0, 1), \quad (2)$$

where the error terms e_i are identical to those in Eq.(1), but follow the Logistic distribution.



Y. Zeng, S. Shimizu, H. Matsui, F. Sun.
Causal Discovery for Linear Mixed Data.
In Proc. The First Conference on Causal Learning and Reasoning (CLear2022), pages 994-1009, Eureka, CA,

今後: 仮定を緩めていく