

二値データに対する 因果推論手法の提案

稲積 孝紀^{*1}, 鷺尾 隆^{*1}, 清水 昌平^{*1}, 鈴木 譲^{*2},
山本章博^{*3}, 河原 吉伸^{*1}

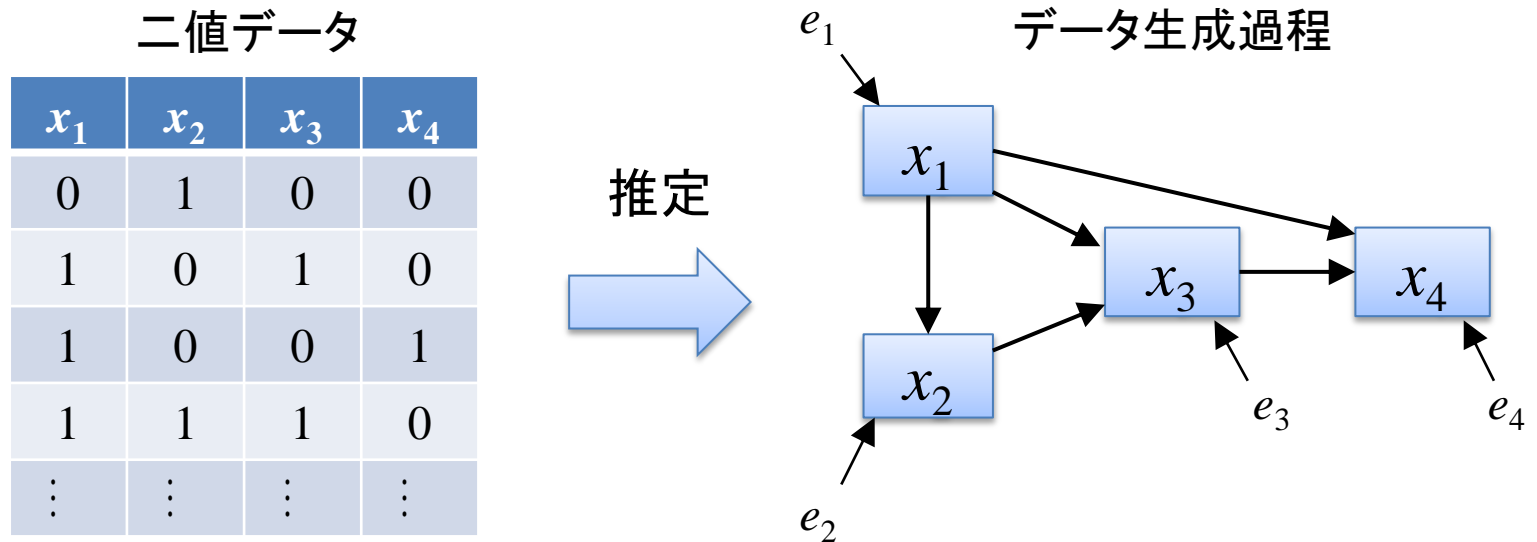
^{*1} 大阪大学 産業科学研究所

^{*2} 大阪大学 大学院理学研究科

^{*3} 京都大学 大学院情報学研究科

研究背景と目的

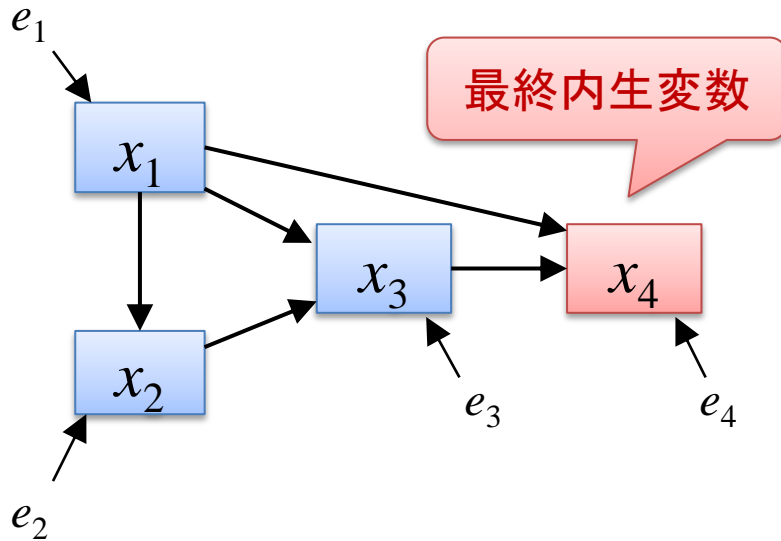
- 与えられた**二値データ** $X = \{x_1, x_2, \dots, x_d\}$ がどのようなデータ生成過程によって生成されているかを推定する.



- 従来手法 ... edge の向きを一意に同定できない場合がある。
→ データ生成過程を**一意に同定できる**
新しい因果推論手法 BExSAM を提案する (UAI2011 採
択).

提案手法: BE_xSAM モデル

- 因果構造として **DAG 構造** を考え, 対応する式を定義する.
 - 外乱 e_i は **排他的論理和** によって加わるものとする.



対応



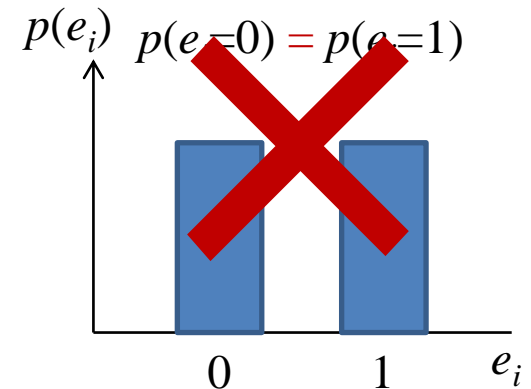
$$x_1 = e_1$$

$$x_2 = x_1 \oplus e_2$$

$$x_3 = x_1 x_2 \oplus e_3$$

$$x_4 = (x_1 + x_3) \oplus e_4$$

- さらに, 外乱 e_i の分布に偏りを仮定する:
 $0 < p(e_i=1) < 1$ かつ $p(e_i=1) \neq 0.5$



提案手法: モデル推定アルゴリズム

観測データの度数表

x_1	x_2	出現回数
0	0	20
0	1	10
1	0	15
1	1	30

x_2 の比
 } 2 : 1
 } 1 : 2

- 最終内生変数 x_k の同定

x_2 の比が $x_1 = 0, 1$ について
 項の交換を許した上で等しい.



x_2 は最終内生変数である.

- 最終内生変数 x_k の親となる変数集合の同定

観測データの度数表

x_2	$x_1 = 0$ のときの 出現回数	$x_1 = 1$ のときの 出現回数
0	20	15
1	10	30

x_2 の比が x_1 の値によって異なる.

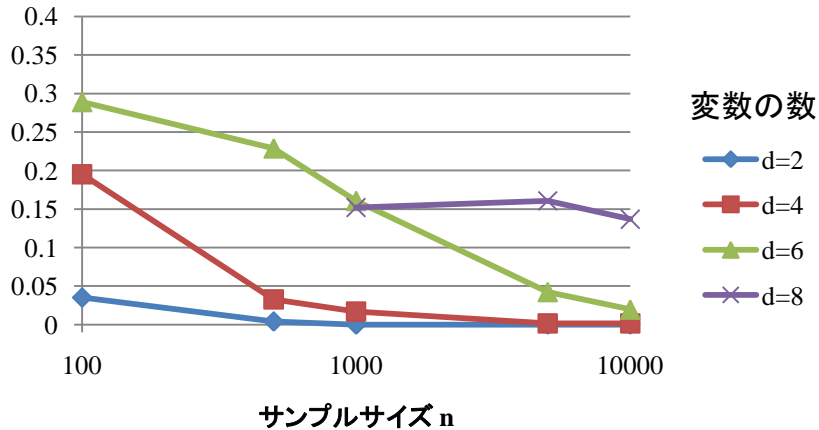


x_1 は x_2 の親である.

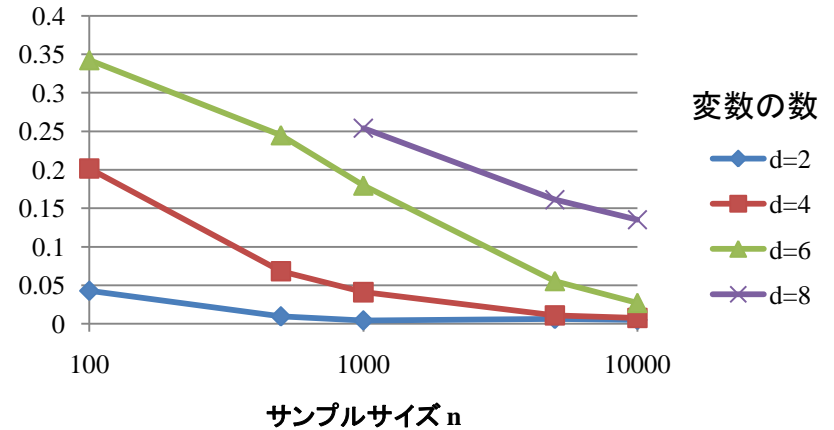
- x_k を除いたモデルを新たな入力モデルとして繰り返す.

評価実験: 人工データ

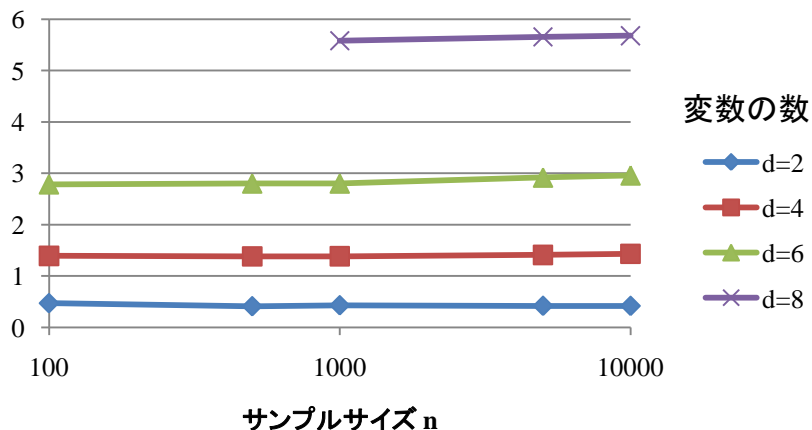
因果的順序の推定エラー率 ER_0



DAG 構造全体の推定エラー率 ER_s



計算時間 CT (ミリ秒)



- エラー率 (左上, 右上):
 - サンプルサイズの増加に従い減少する。
- 計算時間 (左下):
 - 主に変数の数の増加に従い増加する。

評価実験: 実世界データ

- ウィスコンシン州の高校生の
大学進学意思, 親の働きかけ, 社会的ステータスの関係
[Sewell and Shah 1968, Spirtes et al. 2000]
 - 男子高校生のみを抽出.
 - IQ(4値変数)が最も低い学生を除く.
 - サンプルサイズ 3756
- BExSAM を適用した結果(検定の有意水準 0.05)

