

# ベイジアンネットワークの性質に関する研究

2016SS004 大道裕矢

指導教員：松田真一

## 1 はじめに

ベイジアンネットワークは、明確な仮説がないところから分析が可能であるという点から様々な分野で応用されている。しかし、吉見ら [5] にもあるように、ベイジアンネットワークで抽出されたグラフの因果関係の妥当性は常に保証されているとは限らない。本研究では、このベイジアンネットワークの構造パターンとその安定性について検討する。

## 2 ベイジアンネットワーク (BN)

ベイジアンネットワーク (BN: *Bayesian Network*) とは、ベイズの定理の考え方を基に因果関係をグラフィカルモデルとして可視化したものである。これは、非循環有向グラフと条件付き確率表 (CPT) または同時確率分布表等によって表される。BN は、不確実性を含む事象の予測や合理的な意思決定など、幅広い分野に応用されている。

(木村・岩崎 [2], 植野 [4] 参照)

## 3 BN を構築するソフトウェア

BN を構築するソフトウェアは、BayoLink や R 上のパッケージ等を含めて様々であるが、本研究では、BayoLink を中心に解析した。BayoLink は、NTT データ数理システムが提供する BN 構築支援システムである。GUI 上で、情報量基準に基づく最適なネットワーク構造を探索し、BN を構築することができる。また、構築したモデルについての確率推論や検証等を行うことができる。

(BayoLink 7.2, 操作マニュアル [1] 参照)

## 4 本研究で使用するデータについて

### 4.1 *birthwt*(低出生体重に関連する危険因子)

本データは、R の MASS ライブラリにある *birthwt* を使用して作成している。これは、1986 年にマサチューセッツ州のメディカルセンターで収集された低出生体重に関連する危険因子に関するデータである。このデータから、*low*(出生時体重が 2.5kg 未満であれば低い、それ以上であれば高いとする)、*race*(母親の人種：白人、黒人、その他)、*smoke*(母親の妊娠時の喫煙の有無)、*ht*(母親の高血圧の病歴の有無)、*ui*(母親の子宮過敏性の有無) の 5 つの変数を取り出して作成した。データ数は 189 である。

### 4.2 *birthrate*(出生率に関連する因子)

政府統計の総合窓口 [3] の人口動態調査より、都道府県別の合計特殊出生率に関連するデータを収集し、重回帰分析を行い、変数選択し、データを離散化することで、本データを作成した。*FMW* は女性の平均初婚年齢 (2017)、

*FMH* は男性の平均初婚年齢 (2017)、*OSR* は持家比率 (2013)、*BR* は合計特殊出生率 (2017) を指している。データ数は 47 である。

## 5 元データを用いた BN の構築

前章で示した *birthwt* データと *birthrate* データを用いて BN を BayoLink で構築する。ここで、各データの特徴から、作製上でそれぞれ二つの制約を設定した。*birthwt* データでは、*low* が子ノードを持たないように、*race* が親ノードを持たないように指定した。また、*birthrate* データでは、*BR* が子ノードを持たないように、*FMW* と *OSR* が親ノードを持たないように指定した。

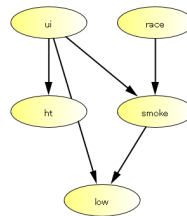


図1 *birthwt* の BN

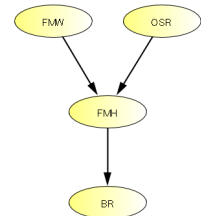


図2 *birthrate* の BN

## 6 シミュレーションによる実験

モンテカルロシミュレーション (MCS) や leave-one-out 法を用いて、以下の計 5 通りの実験を行った。また、それぞれの実験で BN を構築する際の制約や条件に関しては、前章で元データを用いた場合と同様とする。

### 6.1 実験 1 の概要と生成された BN

前章で構築した *birthwt* の BN を基に、*birthwt* の疑似データ生成プログラムを CPT 等に基づく MCS で作成する。このプログラムから 100 個のサンプルを抽出し、それぞれの BN を BayoLink で構築する。結果として 53 通りの BN が生成された。ここでは、生成数トップ 4 のみを示しており、BN1 が 8%、BN2, BN3 が 6%、BN4 が 5% とばらつきが顕著であった。

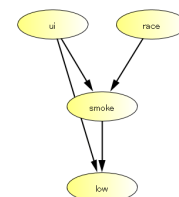


図3 BN1

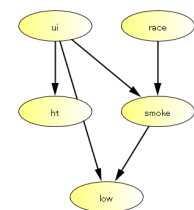


図4 BN2

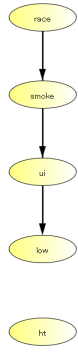


図5 BN3

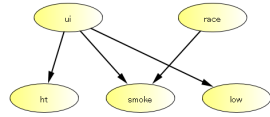


図6 BN4

### 6.2 実験2の概要と生成されたBN

前章で構築した *birthrate* のBNを基に、*birthrate* の疑似データ生成プログラムをCPT等に基づくMCSで作成する。このプログラムから100個のサンプルを抽出し、それぞれのBNをBayoLinkで構築する。結果として5通りのBNが生成された。ここで、主要なBNを図7~9に示す。以下の3通りで全体の91%を占めている。

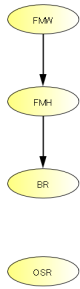


図7 BN5

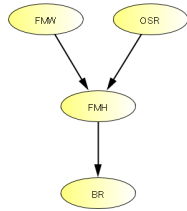


図8 BN6

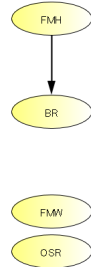


図9 BN7

### 6.3 実験3の概要と生成されたBN

実験2で作成した *birthrate* の疑似データ生成プログラムの生成データ数を200に変更して、同様に100個のサンプルを抽出し、それぞれのBNをBayoLinkで構築する。結果として、BN5, BN6, BN8, BN9の4通りのBNが生成され、BN6の構造のみで全体の93%を占めている。

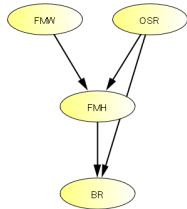


図10 BN8

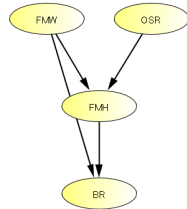


図11 BN9

### 6.4 実験4の概要

leave-one-out法によりモデル構造自身の安定性を確認する。本実験では、*birthrate* データを用いて、47通りの

データを作成し、それぞれのBNをBayoLinkで構築する。結果として、47通りのデータ全てで元データと同じ構造が再現された。

### 6.5 実験5の概要

*birthrate* データからランダムにデータを5個ずつ除外して、実験4と同様にBNをBayoLinkで構築する。結果として、BN5, BN6の2通りのBNが生成され、BN6の構造で全体の51%, BN5で49%を占めている。

## 7 考察

まず、実験4で元データの構造が100%再現されたことから、現実のデータには十分な説明能力が備わっていることが読み取れる。しかし、MCS(実験1, 2)やデータに一定数の欠損(実験5)がある場合に、BNの構造の不安定になることが確認できた。したがって、データの適切な補完や推論などの対策がやはり重要となる。

ここで、実験1, 2を比較することにより、BNの構造を安定させるためには、一定の制約とアイテム数の制限が有効であることが分かる。そのため、BNを用いてデータ分析を行う場合、吉見ら[5]にも示されているが、ドメイン知識を持った専門家の知見を導入できると良い。更に、実験2, 3を比較することにより、データ数が多い方が構造は安定することも分かった。また、実験3において、更にデータを増やすことで構造が淘汰され、元データの構造に帰着されることが予想される。

また、部分的なネットワークに注目することで、結果に対して直接的な関連性が薄くても、間接的に強い関連性を持つデータがあることが分かった。よって、変数選別する場合には、この点にも注意して多角的に考察する必要がある。

BNは構造が異なると推論結果に大きな影響が生じる。したがって、BNを用いた解析の際は、以上の点を念頭においておくことが重要であるといえる。

## 8 おわりに

BNの構造パターンと安定性について、実際のデータを用いたシミュレーションを通して検討できた。本研究で学んだことを今後活かしていきたい。

## 参考文献

- [1] BayoLink 7.2, 操作マニュアル, 2019.
- [2] 木村陽一・岩崎弘利:「ベイジアンネットワーク技術」, 東京電機大学出版局, 2006.
- [3] 政府統計の総合窓口, e-Stat, <https://www.e-stat.go.jp/>, (2019/11 閲覧).
- [4] 植野真臣:「ベイジアンネットワーク」, コロナ社, 2013.
- [5] 吉見将太・黒川悦子・橋本和夫:「ベイジアンネットワークにおけるインタラクティブモデル構築手法の検討」, 2011.