

5. ベイジアンネットワークと他の機械学習モデルとの関係

植野真臣
電気通信大学
情報理工学研究科 情報数理プログラム

本日の目標

- ベイジアンネットワークと他の機械学習モデルとの関係を理解する

ビック・データ時代

- 90年代—00時代 インフラストラクチャー時代
いかにデータを蓄えるか、データを蓄えさせる時代
- 有り余るデータとその有効活用が課題:大量のデータから高精度に高度な処理ができる手法→次世代の企業コンピーテンシー
- 簡単にまねできない高精度な処理技法が目される時代に入
- 総合格闘技→数理統計、アルゴリズム、データベースの統合的技術

古典的人工知能(ルール)

- 古典的AI (アリストテレス)
- 論理推論、 IF then rule
- 人は死ぬ、ソクラテスは人である、ソクラテスは死ぬ
- 古典的AIの問題
- 当たり前のことしか推論できない
- 不確実な現象を推論できない
- 例外が多い
- 学習ができない

確率推論では

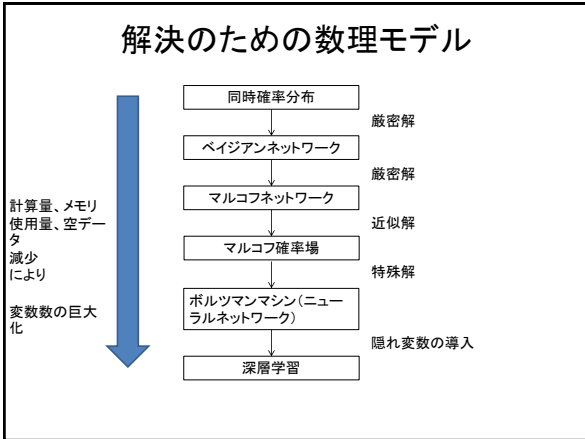
- より一般化した表現
- 同時確率分布

例: 性別、髪の長さ、背の高さの同時確率分布

データより性別、髪の長さ、背の高さの同時確率分布が以下であることが分かっているとします。

$P(\text{男、髪短い、背高い})=0.2$
 $P(\text{男、髪長い、背高い})=0.125$
 $P(\text{男、髪長い、背低い})=0.05$
 $P(\text{男、髪短い、背低い})=0.125$
 $P(\text{女、髪短い、背高い})=0.05$
 $P(\text{女、髪長い、背高い})=0.125$
 $P(\text{女、髪長い、背低い})=0.2$
 $P(\text{女、髪短い、背低い})=0.125$

$P(\text{男})=0.5, P(\text{女})=0.5$



ベイジアンネットワーク

- 確率構造が非循環有向グラフであれば、同時確率分布が条件付確率の積に因数分解できることが数学的に証明できる。
- 確率有向グラフが確率因果構造が対応し、ものごとの因果もわかる！！

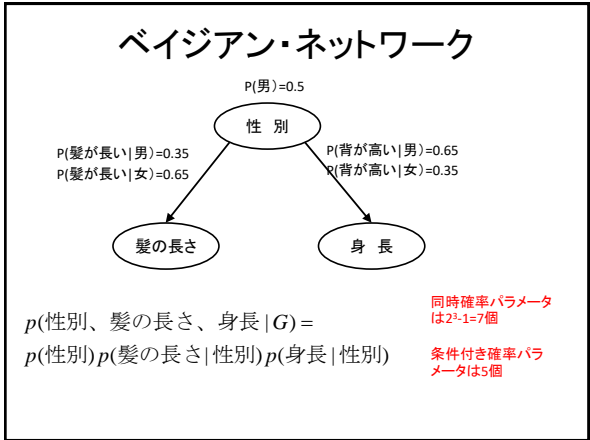
現在考えられる最もよい同時確率分布の推定値
⇒ 推論の予測精度が最高のはず！！

ベイジアンネットワークの学習

未知のデータへの予測を最大化する構造は

$$P(G|X) \propto P(G) \prod_{i=1}^N \prod_{j=1}^{q_i} \frac{\Gamma(\alpha_{ijk})}{\Gamma(\alpha_{ij} + n_{ij})} \prod_{k=0}^{r_i-1} \frac{\Gamma(\alpha_{ijk} + n_{ijk})}{\Gamma(\alpha_{ijk})}$$

ここで、
 $\alpha_{ijk} = 1/K$
 K : モデルのパラメータ数
 n_{ijk} : 変数 i が j 番目の親ノードパターンを条件として k をとる頻度



ベイジアンネットワークの問題

- 欠点として計算量の多さ
- 現在 厳密学習では、2000ノードのネットワーク(Natori, Uto, Ueno 2017))
- 厳密推論では200ノードのネットワーク(Li and Ueno 2017)
- 将来的にはこれを克服すれば最強ツールになる！！

マルコフネットワーク

- 無向グラフ構造
- ベイジアンネットワークの互換モデル

利点

- 非循環有向グラフ構造の仮定がいらない。

欠点

- ベイジアンネットワークから変換できるがその逆は不可
- 構造の学習ができない
- パラメータ数が多い

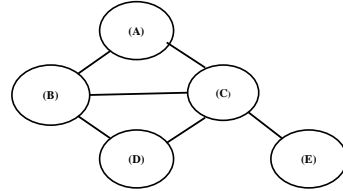
マルコフネットワークの同時確率分布 クリーク分解定理

- $P(x_1, x_2, \dots, x_N | G) = \frac{1}{Z(\theta)} \prod_{c \in C} \phi_c(x_c | \theta_c)$
- をギブス分布 (Gibbs Distribution) と呼ぶ。

Cはクリーク集合

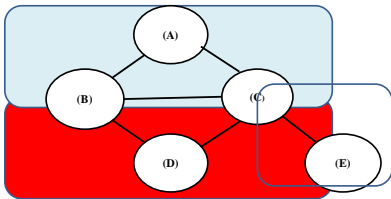
クリーク

- 頂点の部分集合が完全グラフ(すべての頂点間に辺がある)である場合



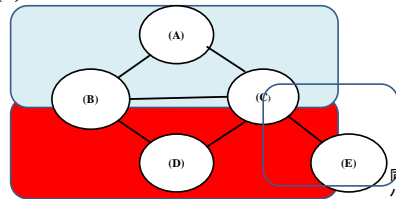
クリーク

- 頂点の部分集合が完全グラフ(すべての頂点間に辺がある)である場合



計算例 2値の場合

- $P(x_A, x_B, \dots, x_E | G) = \frac{1}{Z(\theta)} p(x_A, x_B, x_C) p(x_B, x_C, x_D) p(x_C, x_E)$

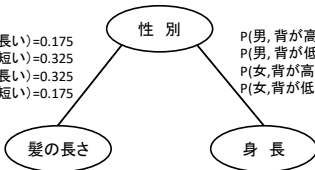


同時確率分布の
パラメータ数
31
⇒
7+7+3=17

マルコフ・ネットワーク

P(男, 髪が長い)=0.175
P(男, 髪が短い)=0.325
P(女, 髪が長い)=0.325
P(女, 髪が短い)=0.175

P(男, 背が高い)=0.325
P(男, 背が低い)=0.175
P(女, 背が高い)=0.175
P(女, 背が低い)=0.325



同時確率パラメータは
2³-1=7個
パラメータは6個

マルコフネットワークの問題

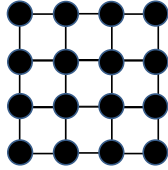
- パラメータ数はクリークの大きさに対して指数的に増え計算量が増えてしまうのでベイジアンネットワークより計算量が大い。
- 数学的厳密性を緩和して計算が簡易なモデルの必要性

マルコフ確率場

$$P(x_1, x_2, \dots, x_N | G) = \frac{1}{Z(\theta)} \prod_i \phi(x_i)$$

$$\prod_{(i,j)} \phi(x_i, x_j)$$

- クリークをまともに計算せず、グラフの辺ごとに分離して同時確率分布を近似する。
- 画像処理などで用いられる。



マルコフネットワークに戻ろう Log-Linear モデル

マルコフネットワークのファクターを

$$\phi_c(x_c | \theta_c) = \exp(-E(x_c | \theta_c))$$

と定義する。

ここで、 $E(x_c | \theta_c) = -\log(\phi_c(x_c | \theta_c)) > 0$ はクリーク c のエネルギー関数と呼ばれる。

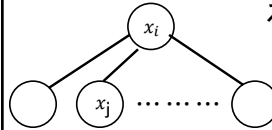
すなわち

$$P(x_1, x_2, \dots, x_N | G) = \frac{1}{Z(\theta)} \exp(-\sum_c E(x_c | \theta_c))$$

マルコフ確率場のLog-Linear モデル表現

$$P(x_1, x_2, \dots, x_N | G) = \frac{1}{Z(\theta)} \exp(-\sum_i E(x_i) - \sum_{(i,j)} E(x_i, x_j))$$

x は二値しかとらずに以下の構造を考える



$$P(x_i = 1 | x_1, x_2, \dots, x_N, G) = \frac{1}{Z(\theta)} \exp(-\sum_i E(x_i = 1) - \sum_{(i,j)} E(x_i = 1, x_j))$$

$$= \frac{\exp(-\sum_i E(x_i = 1) - \sum_{(i,j)} E(x_i = 1, x_j))}{\exp(-\sum_i E(x_i = 1) - \sum_{(i,j)} E(x_i = 1, x_j)) + \exp(-\sum_i E(x_i = 0) - \sum_{(i,j)} E(x_i = 0, x_j))}$$

$$= \frac{1}{1 + \exp(\sum_i E(x_i = 1) - \sum_i E(x_i = 0) + \sum_{(i,j)} E(x_i = 1, x_j) - \sum_{(i,j)} E(x_i = 0, x_j))}$$

$-\sum_i E(x_i = 1) + \sum_i E(x_i = 0) = b_i, \sum_{(i,j)} E(x_i = 1, x_j) - \sum_{(i,j)} E(x_i = 0, x_j) = w_{ij} x_j$ とおくと

ボルツマンマシン

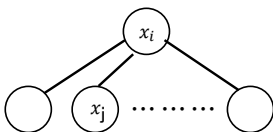
(ニューラルネットワーク)

$$-\sum_i E(x_i = 1) + \sum_i E(x_i = 0) = b_i,$$

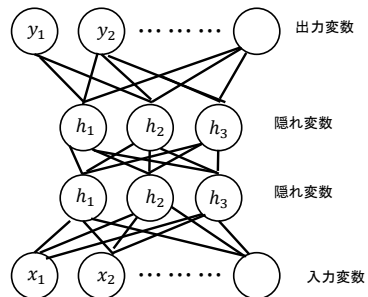
$$\sum_{(i,j)} E(x_i = 1, x_j) - \sum_{(i,j)} E(x_i = 0, x_j) = w_{ij} x_j$$

とおくと

$$P(x_i = 1 | x_1, x_2, \dots, x_N, G) = \frac{1}{1 + \exp(-\sum_j w_{ij} x_j - b_i)}$$



深層学習モデル (ディープラーニング)



隠れ変数の役割

- 隠れ変数を積分消去すると
- 全変数間に辺が引かれた完全グラフ構造となる。
- 完全グラフ構造において、各辺の重みを最適化することにより、マルコフグラフの構造も同時に推定できる。
- 計算不可能な複雑な構造を 隠れ変数を導入することにより、単純で計算可能な階層構造に変換している。
- 真の確率構造が複雑な場合、隠れ変数層を増やさなければならぬはず。
- ベイジアンネットワークで学習されるエッジ数が隠れ変数の数に関係している可能性が高い。

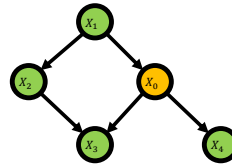
やはり脳モデルはすごい！！

- ビッグデータにおける同時確率分布の問題は、変数の値のパターンがコンピュータや人間のメモリに入らないこと、計算速度が遅すぎる、パターンが多すぎて空データが増えてしまうことである！！
- 脳モデルは、メモリに乗らないほどの変数パターンは計算せず、すべて独立変数のように扱い、隠れ変数が仲介する階層モデルにより、結果として変数間の依存性を補完する。
- 計算速度、メモリ使用量、欠損データ、近似精度のトレードオフをすべて解決する！！

隠れ変数の数と構造の最適化が今後のビッグチャレンジ

- **問題:** 周辺尤度や従来の情報量規準は隠れ変数の数と構造を決めることはできない。ただ、モデルが正則性を満たさないということだけではない。
- 隠れ変数の数と構造を最適にする規準(スコア)は何か？
- ベイジアンネットワークで得られる構造のパラメータ数とどのような関係にあるのか？
- 数学的に解明できるのか？
- その構造を学習できるのか？

ベイジアンネットワーク分類器 (Bayesian Network Classifier: BNC)



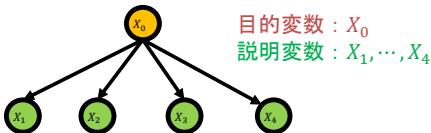
目的変数: X_0
 説明変数: X_1, \dots, X_4

$P(X_0, X_1, \dots, X_n | G)$ から $P(X_0 | X_1, \dots, X_n)$ を計算して、分類を行う。

Friedmanら (1997) による批判

周辺尤度で学習した構造の分類精度が、単純な構造をとるNaive Bayesより劣ることが多々ある。

Naive Bayesの例



N. Friedman, D. Geiger, and M. Goldszmidt, "Bayesian Network Classifiers," Machine Learning, vol.29, no.2, pp.131-163, 1997.

Friedmanら (1997) による批判

BDeu: 同時確率分布 $P(X_0, X_1, \dots, X_n | G)$ を表現する生成モデルをモデル化

分類が目的ならば、条件付確率分布 $P(X_0 | X_1, \dots, X_n, G)$ を表現する識別モデルを学習すべき

条件付き対数尤度スコア (Conditional Log Likelihood: CLL) の提案

$$\sum_{d=1}^N \log P(x_0^d | x_1^d, \dots, x_n^d, G, \theta)$$

N. Friedman, D. Geiger, and M. Goldszmidt, "Bayesian Network Classifiers," Machine Learning, vol.29, no.2, pp.131-163, 1997.

条件付き対数尤度の問題点

条件付き対数尤度スコアは周辺尤度とは異なり閉形式で表せないため、構造探索に効率的なアルゴリズムが適用できず、厳密学習に膨大な時間がかかる。

条件付き対数尤度の近似手法

- 構造探索に対してHill-Climbing法の適用した近似学習 (Grossman et al., 2004)
- 構造探索に効率的なアルゴリズムを用いることができるように修正した近似条件付き対数尤度スコア approximate GLL (aGLL) (Carvalho et al., 2013)

上記の近似学習手法の方が、周辺尤度を用いた学習法より分類精度の高い分類器が得られることが経験的に示されている。

A. M. Carvalho, P. Ado, and P. Mateus. Efficient Approximation of the Conditional Relative Entropy with Applications to Discriminative Learning of Bayesian Network Classifiers. *Entropy*, 15(7): 2716–2735, 2013.

本当に条件付き対数尤度は良いのか？

- 周辺尤度最大化より条件付き対数尤度最大化の方がなぜ良いのかという理由については未だ明らかになっていない。
- 周辺尤度を最大化する構造を厳密に学習できるにもかかわらず、既存研究では近似学習を行っている。
このため、探索精度の悪さが結果に影響したのかもしれない。

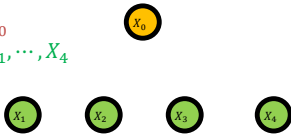
周辺尤度による厳密学習と条件付き対数尤度による近似学習の比較

- Sugahara, Uto, and Ueno(2018)は、サンプルサイズが大きいときは周辺尤度を用いた厳密学習の方が条件付き対数尤度の近似学習よりも高い分類精度を示すことを経験的に示した。
- 一方で、サンプルサイズの小さい場合、周辺尤度は目的変数の親変数が多い構造を学習する傾向があり、パラメータ数が指数的に増えてしまい分類精度が著しく悪くなることもわかった。

Sugahara, S.; Uto, M.; and Ueno, M. 2018. Exact learning augmented naive Bayes classifier. In Proceedings of the 9th International Conference on Probabilistic Graphical Models, volume 72 of Proceedings of Machine Learning Research, 439–450. PMLR.

提案手法

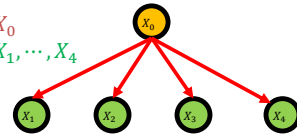
目的変数: X_0
説明変数: X_1, \dots, X_4



提案手法

まず強制的に目的変数から全説明変数へエッジを引く。

目的変数: X_0
説明変数: X_1, \dots, X_4

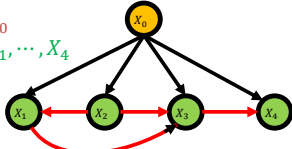


提案手法

まず強制的に目的変数から全説明変数へエッジを引く。

目的変数: X_0

説明変数: X_1, \dots, X_4



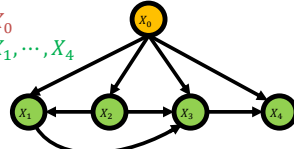
説明変数間の構造を学習。

提案手法

このように目的変数から全説明変数にエッジが引かれて
いる構造を, Augmented Naïve Bayes (ANB) と呼ぶ。

目的変数: X_0

説明変数: X_1, \dots, X_4



目的変数の親変数が0で, 全ての説明変数の子変数とするため,
先述した分類精度低下の問題を回避できる。

ANB厳密学習の漸近的性質

- 定理 (Sugahara and Ueno, 2021)

厳密学習されたANBは真の構造と分類確率が等価
な構造に漸的に一致する。

S. Sugahara and M. Ueno., Exact Learning Augmented Naive Bayes Classifier, arXiv:2107.03018.

Deep Learningとの比較

- Deep Learning では分類確率までは一貫性が
ないため, この性質は確率推論を含む意思
決定問題などで大きな意味を持つ。
- Deep Learningでは説明性・解釈性が乏しいこ
とが問題となっているが 因果をグラフィカル
モデルで表現できるベイジアンネットワークは
利点が多い

おわり

半年間 ご清聴ありがとうございました！！

復習問題1.

ベイズの定理を書け。

復習問題2.

被害者Xはある日狙撃された。この事象をEとしよう。

命中率8割のスナイパーAと2割のスナイパーBのどちらかが犯人であることが分かっている。今、どちらが犯人かは全くわからない。

Aが犯人である確率を求めよ。

例題つづき

- そのあとさらに2発Xに銃弾が打たれたが2発とも外れた。この事象をEとしてそれぞれが犯人である確率を求めよ。

復習問題3 (3 囚人問題)

ある監獄にアラン、バーナード、チャールズという3人の囚人がいて、それぞれ独房に入れられている。3人は近く処刑される予定になっていたが、恩赦が出て3人のうち1人だけ釈放されることになったという。誰が恩赦になるかは明かされておらず、それぞれの囚人が「私は釈放されるのか?」と聞いても看守は答えない。囚人アランは一計を案じ、看守に向かって「私以外の2人のうち少なくとも1人は死刑になるはずだ。その者の名前を知りたい。私のことじゃないんだから教えてくれてもよいだろう?」と頼んだ。すると看守は「バーナードは死刑になる」と教えてくれた。それを聞いたアランは「これで釈放される確率が $1/3$ から $1/2$ に上がった」とひそかに喜んだ。アランの釈放される確率を求めよ。