

4. 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 ネットワークのおおざっぱな説明
- 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

ビック・データ時代

- 90年代-00年代 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。

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

$P(\text{男})=0.5, P(\text{女})=0.5$

同時確率分布からの確率推論

- 「その人は髪が短い」ことがわかった

$$P(\text{男、髪短い、背高い})=0.2$$

$$P(\text{男、髪短い、背低い})=0.125$$

$$P(\text{女、髪短い、背高い})=0.05$$

$$P(\text{女、髪短い、背低い})=0.125$$

$$\text{男の確率}=0.325/0.5=0.65$$

同時確率分布からの確率推論

- さらに「その人は背が高い」ことがわかった

$$P(\text{男、髪短い、背高い})=0.2$$

$$P(\text{女、髪短い、背高い})=0.05$$

$$\text{男の確率}=0.2/0.25=0.8$$

確率推論の数学的定式化

データ x_d が得られたときの x_i の確率は

$$p(x_i|x_d) = \sum_{j \neq i} p(x_1, x_2, \dots, x_N | x_d)$$

世界中のすべての変数の同時確率分布を知ればなんでも推論できる！！

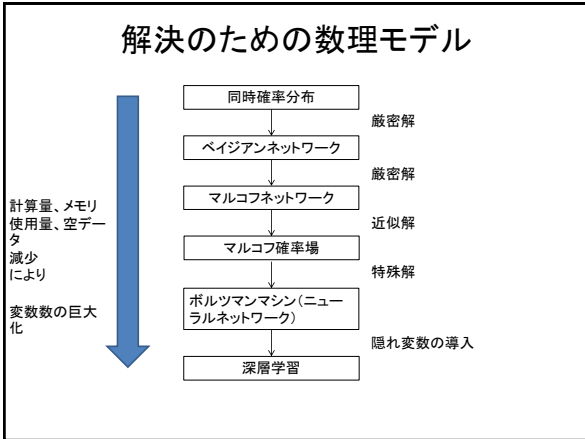
問題

- 変数が増えると状態数が指数的に増える
- すべて2値しかとらなくても 10変数で1024パターンの同時確率を推定しないといけない。
- 100変数で
1267650000000000000000000000000
パターンの同時確率を推定しなければならない。

二つの問題

- 計算量が指数的に爆発する。
- データ数よりパターン数のほうが多くなってしまおうと、各パターンを推定するためのデータが0になるものが大量発生。
(大量のデータがあっても空データだらけになる)

ビッグデータ問題の課題は、スパースデータ(空データの増加)と計算量(メモリ、計算速度)



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！！

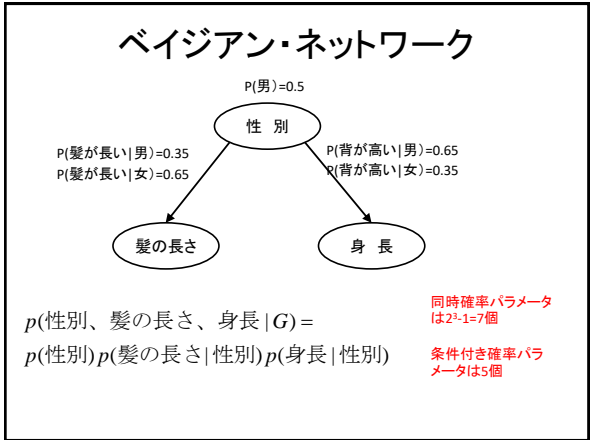
現在考えられる最もよい同時確率分布の推定値
⇒ 推論の予測精度が最高のはず！！

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

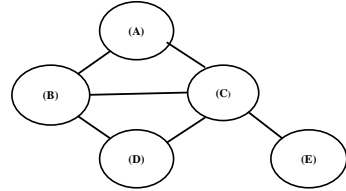
マルコフネットワークの同時確率分布 クリーク分解定理

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

Cはクリーク集合

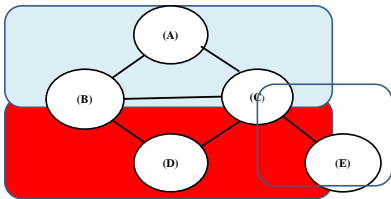
クリーク

- 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



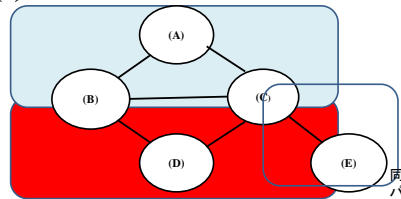
クリーク

- 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

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

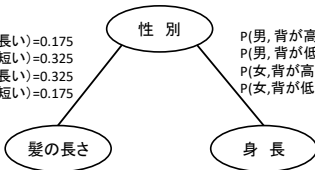


同時確率分布の
パラメータ数
31
⇒
7+7+3=17

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

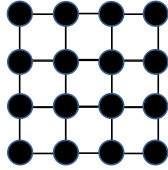
- パラメータ数はクリークの大きさに対して指数的に増え計算量が増えてしまうのでベイジアンネットワークより計算量が大い。
- 数学的厳密性を緩和して計算が簡易なモデルの必要性

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。

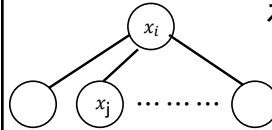
すなわち

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

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



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

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

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

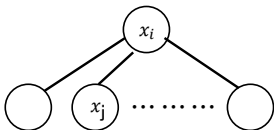
$E(x_i = 0) = b_i, E(x_i = 0, x_j) = w_{ij}x_j$ とおくと

ボルツマンマシン

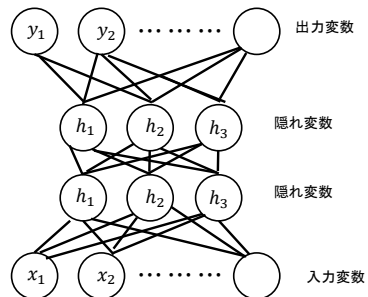
(ニューラルネットワーク)

$E(x_i = 0) = b_i, E(x_i = 0, x_j) = w_{ij}x_j$ とおくと

$$P(x_i = 1 | x_j \neq x_i) = \frac{1}{1 + \exp(-\sum_j w_{ij}x_j - b_i)}$$



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？

レポート課題

授業で学んだことで 面白かったこと、驚いたこと
や発見があれば それをレポートにまとめよ。

A4 1枚以上

送り先

sugahara@ai.lab.uec.ac.jp