

情報認識

「最尤推定法におけるモデル選択」

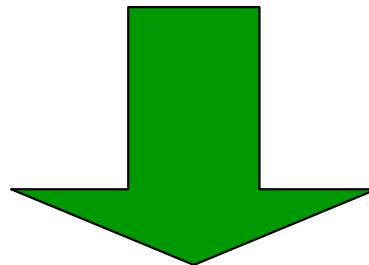
- 講師： 杉山 将（計算工学専攻）
- 居室： W8E-505
- 電子メール： sugi@cs.titech.ac.jp

最尤推定法とモデルの選択

- 最尤推定法: あらかじめ定めたパラメトリックモデルの中から最も尤もらしい確率密度関数を選ぶ方法
- これまで、ガウスモデルを扱ってきたが、一口にガウスモデルといっても、分散共分散行列が
 - 任意の正值行列の場合(自由度 d^2)
 - 対角行列で対角成分が異なる場合(自由度 d)
 - 対角行列で対角成分が等しい場合(自由度 1)などいろいろなものがあった。
- 実際にはどれを選べばよいのだろうか？

複雑なモデルがよい理由

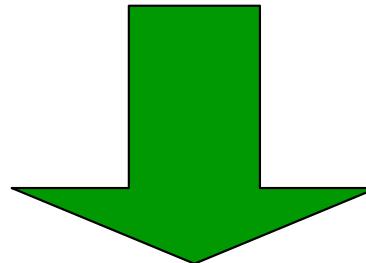
- パラメトリック法では、モデルの中に真の確率密度関数を良く近似するものが含まれていなければ、そもそもよい結果は得られない。



- 真の確率密度関数を含むよう、多くのパラメータを含む表現力の高い複雑なモデルを選ぶ。

単純なモデルがよい理由

- 訓練標本数がパラメータ数と比べてそれほど多くない場合、最尤推定法のよさは、理論的には保証されない。



- パラメータ数の少ないモデルを選ぶ。

モデル選択

- このように、モデルの選択には二つの相反する要請があり、実際には程よい複雑さのモデルを用いなければならない。
- 訓練標本を用いて適切なモデルを選ぶことを、**モデル選択(model selection)**という。

モデル選択の流れ

1. いくつかのパラメトリックモデルを用意する.

$$\{q_i(x; \theta)\}_i$$

2. それぞれのモデルに対して、最尤推定量 $\hat{\theta}_{ML_i}$ を求める.

3. それぞれのモデルから得られた確率密度関数の推定量を次のように定める.

$$\hat{p}_i(x) = q_i(x; \hat{\theta}_{ML_i})$$

4. $\{\hat{p}_i(x)\}_i$ から真の確率密度関数 $p(x)$ に最も「近い」ものを選ぶ.

確率密度関数の近さを測る規準

- カルバック・ライブラー情報量(Kullback-Leibler information):

$$KL(p \parallel \hat{p}) = \int_D p(x) \log \frac{p(x)}{\hat{p}(x)} dx$$

- KL情報量は常に非負で, $\hat{p}(x) = p(x)$ のときだけゼロになる.
- 従って, KL情報量が小さければ, $\hat{p}(x)$ は「よい」といえる.
- 最尤推定量は, KL情報量の近似を最小にする

$$\arg \max_{\theta} \left(\sum_{i=1}^n \log q(x_i; \theta) \right) = \arg \min_{\theta} \left(\sum_{i=1}^n \log \frac{p(x_i)}{q(x_i; \theta)} \right)$$

距離

■ 数学的な距離(distance)の定義

1. $d(x, y) \geq 0$
2. $d(x, y) = d(y, x)$
3. $d(x, y) = 0 \Leftrightarrow x = y$
4. $d(x, y) + d(y, z) \geq d(x, z)$

■ KL情報量は2と4を満たさないため, 厳密には距離ではないことに注意.

$$KL(p \parallel \hat{p}) \neq KL(\hat{p} \parallel p)$$

KL情報量の推定

- KL情報量には未知の確率密度関数 $p(x)$ が含まれているため、直接計算できない。
- 訓練標本からKL情報量を推定する。

$$KL(p \parallel \hat{p}) = \underbrace{\int_D p(x) \log p(x) dx - \int_D p(x) \log \hat{p}(x) dx}_{\text{エントロピー(entropy)}}$$

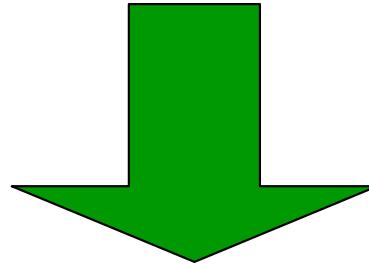
- エントロピーは定数なので、第二項目のみを推定すればよい。

KL情報量の単純な推定

- 負の尤度はKL情報量の二項目に収束する

$$-\frac{1}{n} \sum_{i=1}^n \log \hat{p}(x_i) \rightarrow -\int_D p(x) \log \hat{p}(x) dx$$

- 複雑なモデルほど負の尤度は小さい



- 負の尤度を最小にするモデルを選ぶと、常に最も複雑なモデルが選ばれてしまう。
- もう少し精密なKL情報量の近似が必要！

赤池の情報量規準

- 赤池の情報量規準(Akaike's information criterion):

$$AIC = -2 \sum_{i=1}^n \log q(x_i; \hat{\theta}_{ML}) + 2 \dim \theta$$

- 訓練標本が十分に多いとき

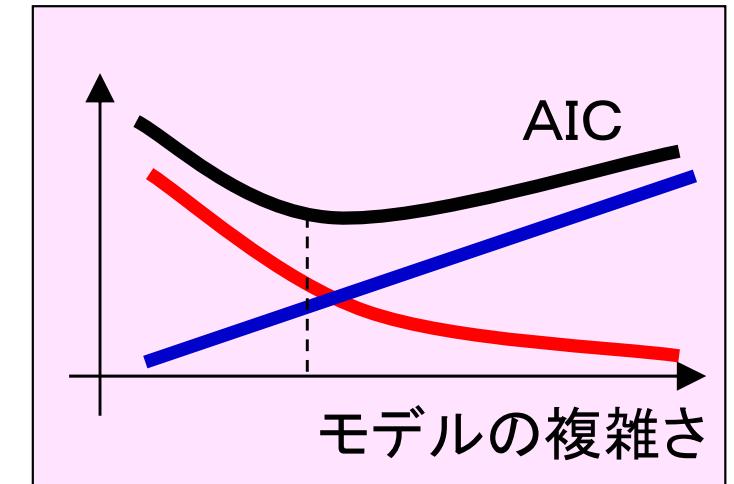
$$\frac{1}{2n} AIC \approx -\int_D p(x) \log \hat{p}(x) dx$$

- AICを最小にするモデルを選べばよい.

AICの直感的解釈

$$AIC = 2 \left(- \sum_{i=1}^n \log q(x_i; \hat{\theta}_{ML}) + \dim \theta \right)$$

負の最大対数尤度 パラメータ数



- モデルが複雑な場合、負の最大対数尤度は小さいがパラメータ数が大きいためAICは大きい。
- モデルが単純な場合、パラメータ数は小さいが負の最大対数尤度が大きいためAICは大きい。
- モデルが程よく複雑な場合、二つの項がバランスよく小さくなり、AICは小さい。

オッカムのかみそり

- オッカムのかみそり(Occam's Razor): 14世紀の哲学者オッカムによる「不必要に実体の数を増やしてはならない」という提言.
- 現代でも科学理論を構築する上での基本的な指針としてよく用いられる.
- 「現象を同程度うまく説明する仮説があるなら、単純な方を選べ」
- 「けちの原理(principle of parsimony)」とも呼ばれる.

オッカムのかみそり(続き)

$$AIC = 2 \left(- \underbrace{\sum_{i=1}^n \log q(x_i; \hat{\theta}_{ML})}_{\text{負の最大対数尤度}} + \underbrace{\dim \theta}_{\text{パラメータ数}} \right)$$

- 「現象を同程度うまく説明する仮説」: 尤度が等しい二つのモデル
- 「単純な方」: パラメータ数が少ない方
- AICはオッカムのかみそりの妥当性を理論的に裏付けている。

AICの精密化

- 理論的には、AICよりも次の**竹内の情報量規準(Takeuchi's information criterion)**の方がより精密。

$$TIC = 2 \left(- \sum_{i=1}^n \log q(x_i; \hat{\theta}_{ML}) + \text{trace}(JH^{-1}) \right)$$

$$J_{j,k} = \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta^{(j)}} \log q(x_i; \theta) \Bigg|_{\theta=\hat{\theta}_{ML}} \frac{\partial}{\partial \theta^{(k)}} \log q(x_i; \theta) \Bigg|_{\theta=\hat{\theta}_{ML}}$$

$$H_{j,k} = - \frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta^{(j)} \partial \theta^{(k)}} \log q(x_i; \theta) \Bigg|_{\theta=\hat{\theta}_{ML}}$$

AICの精密化(続き)

■ TICの近似性能

$$E[TIC] = E\left[-\int_D p(x) \log \hat{p}(x) dx\right] + o\left(\frac{1}{n}\right)$$

- $f(n) = o(g(n))$: $n \rightarrow \infty$ のとき $|f(n)| < Cg(n)$
- $f(n) = O(g(n))$: $n \rightarrow \infty$ のとき $f(n)/g(n) \rightarrow 0$

- TICは $1/n$ のオーダまで不偏推定である.
- モデルが真の確率密度関数を含むとき, 即ち $p(x) = q(x; \theta_{true})$ のとき, $TIC = AIC$.