

情報認識

「最尤推定法におけるモデル選択 (第7章)」

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

最大事後確率則に基づくパターン認識

- 事後確率:

$$p(y | x) \propto p(x | y)p(y)$$

- 最尤法による事前確率の推定:

$$\hat{p}(y) = \frac{n_y}{n}$$

- 最尤法による条件付き確率の推定:

$$\hat{p}(x | y) = q(x; \hat{\theta}_{ML})$$

- **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

- モデルはどうやって定めればよいか？
- 例えば、一口にガウスモデルといっても、分散共分散行列が
 - 任意の正値対称行列の場合 (自由度 $d(d+1)/2$)
 - 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る。
- どれを選べばよいか？

- d 次元の確率ベクトル: $x = (x^{(1)}, x^{(2)}, \dots, x^{(d)})^T$
- 2つのパラメータ:
 - d 次元ベクトル μ
 - d 次元正値行列 Σ

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right)$$

- 正規分布の期待値, 分散共分散行列

$$E[x] = \mu$$

$$V[x] = \Sigma$$

- 共分散がゼロ(即ち $\Sigma = \text{diag}(\sigma_i^2)$) のとき

$\text{diag}(\sigma_i^2)$: $\sigma_1^2, \sigma_2^2, \dots, \sigma_d^2$ を
対角成分に持つ対角行列

$$p(x; \mu, \{\sigma_i^2\}_{i=1}^d) = \frac{1}{(2\pi)^{d/2} \prod_{i=1}^d \sigma_i} \exp\left(-\sum_{i=1}^d \frac{(x^{(i)} - \mu^{(i)})^2}{2\sigma_i^2}\right)$$

- さらに分散が等しい(即ち $\Sigma = \sigma^2 I$) のとき

I : 単位行列

$$p(x; \mu, \sigma) = \frac{1}{(2\pi\sigma^2)^{d/2}} \exp\left(-\frac{(x - \mu)^T (x - \mu)}{2\sigma^2}\right)$$

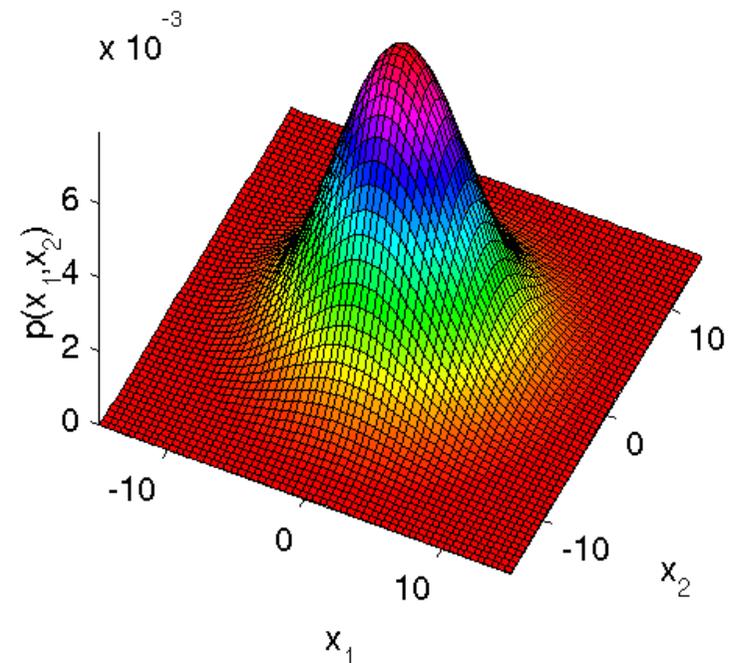
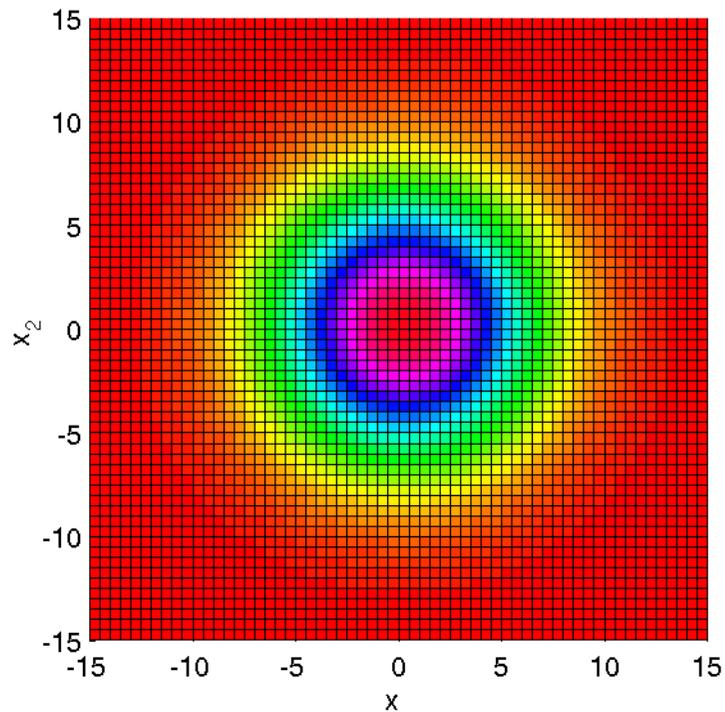
多次元正規分布の例(1)

112

$$d = 2$$

$$\mu = (0,0)^T$$

$$\Sigma = \begin{pmatrix} 20 & 0 \\ 0 & 20 \end{pmatrix}$$



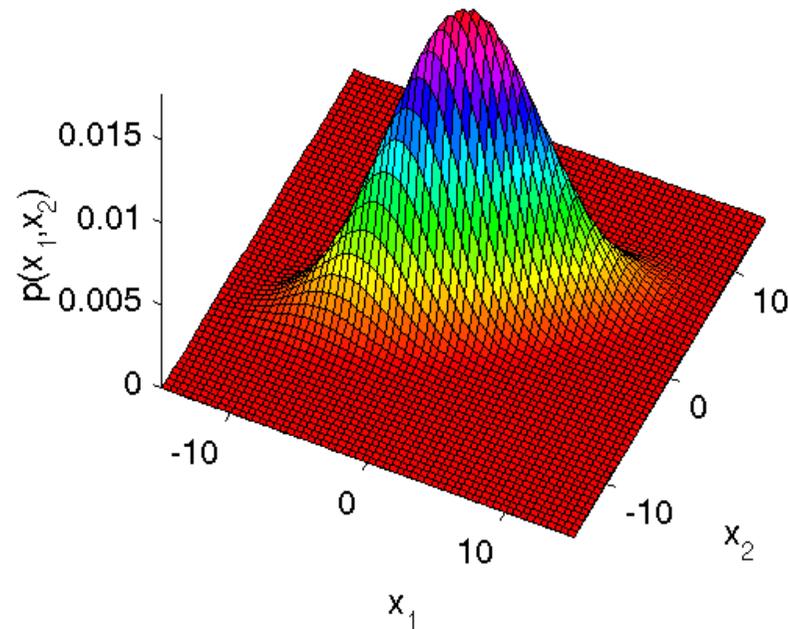
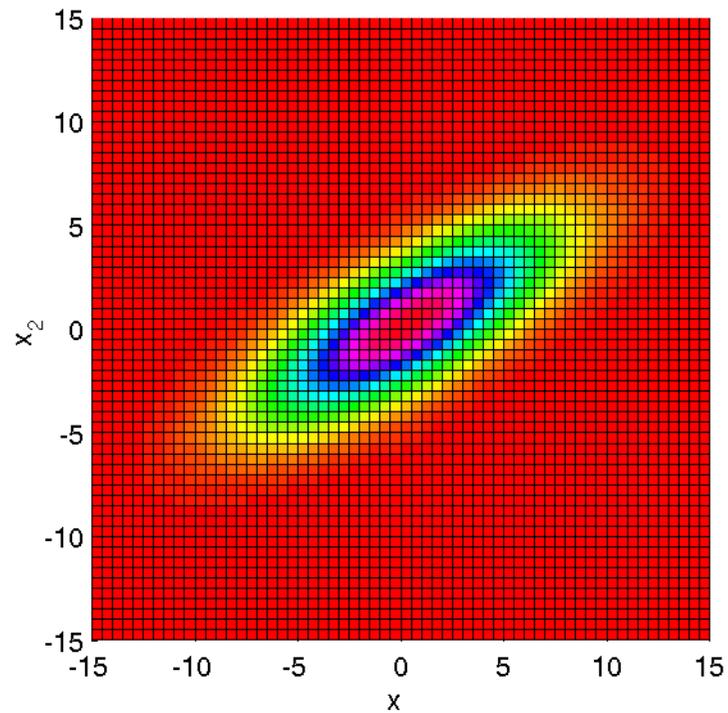
多次元正規分布の例(2)

113

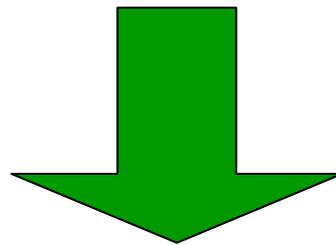
$$d = 2$$

$$\mu = (0,0)^T$$

$$\Sigma = \begin{pmatrix} 20 & 10 \\ 10 & 9 \end{pmatrix}$$



- 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。

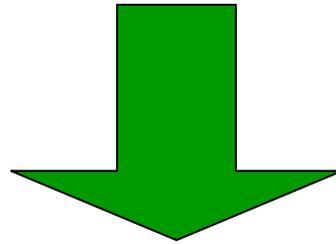


- 真の確率密度関数を含むよう、**パラメータ数の多い表現力の豊かなモデル**を選ぶべき。

単純なモデルがよい理由

115

- 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。



- 推定量の分散が十分小さくなるよう、**パラメータ数の少ないモデル**を選ぶべき。

- 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。
- 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. 各々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¹¹⁷

- 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.

■ 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。

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

- 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。**

- 負の対数尤度はKL情報量の近似:

$n \rightarrow \infty$ のとき, 一致する (大数の法則)

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

- 尤度を最大にするモデルを選べばよい?
- 複雑なモデルほど尤度は大きいので, 尤度最大のモデルを選ぶと, **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(Akaike's information criterion):

$$AIC = - \underbrace{\sum_{j=1}^n \log \hat{p}(x_j)}_{\text{負の最大対数尤度}} + \underbrace{\dim \theta}_{\text{パラメータ数}}$$

負の最大対数尤度 パラメータ数

- 訓練標本が十分に多いとき

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

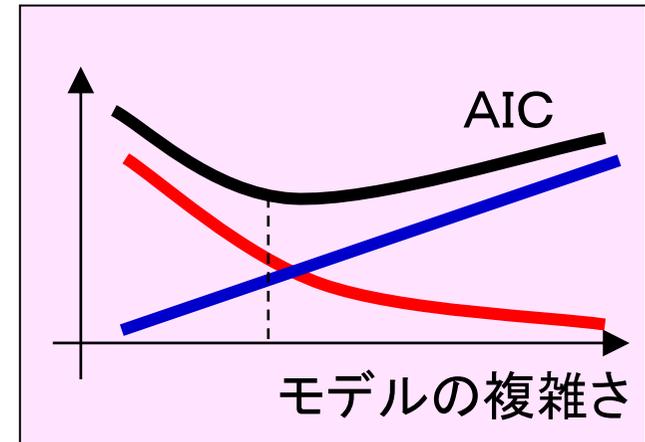
- AICを最小にするモデルを選ぶ.

AICの直感的解釈

122

$$AIC = - \underbrace{\sum_{j=1}^n \log \hat{p}(x_j)}_{\text{負の最大対数尤度}} + \underbrace{\dim \theta}_{\text{パラメータ数}}$$

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(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)**」ともよばれる。

オッカムのかみそり(続き)

124

$$AIC = - \underbrace{\sum_{j=1}^n \log \hat{p}(x_j)}_{\text{負の最大対数尤度}} + \underbrace{\dim \theta}_{\text{パラメータ数}}$$

負の最大対数尤度 パラメータ数

- 「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!

- 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.

$$TIC = -\sum_{j=1}^n \log q(x_j; \hat{\theta}_{ML}) + \text{trace}(JH^{-1})$$

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

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

■ TICの近似性能

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

- $f(n) = o(g(n)) : n \rightarrow \infty$ 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 $1/n$ のオーダーまで不偏

■ 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.



まとめ

127

- 確率密度関数をうまく推定するためには、よいパラメトリックモデルを選ぶ必要がある。
- モデルのよさはKullback-Leibler情報量で測る。
- KL情報量は直接計算できないが、AICがそのよい近似を与える。
- AICはオッカムのかみそりの原理に沿っている。

$$q(x; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right)$$

1. 共分散がない, 即ち, 分散共分散行列が $\Sigma = \text{diag}(\sigma_i^2)$ で与えられるとき, σ_i^2 の最尤推定量は次式で与えられることを証明せよ.

$$\hat{\sigma}_i^2 = \frac{1}{n} \sum_{j=1}^n (x_j^{(i)} - \hat{\mu}^{(i)})^2$$

$$\hat{\mu} = \frac{1}{n} \sum_{j=1}^n x_j$$

但し, $x^{(i)}, \mu^{(i)}$ はベクトル x, μ の第 i 要素である.

2. 更に分散が等しいとき, 即ち, 分散共分散行列が $\Sigma = \sigma^2 I$ (I は単位行列) で与えられるとき, σ^2 の最尤推定量は次式で与えられることを証明せよ.

$$\hat{\sigma}^2 = \frac{1}{nd} \sum_{j=1}^n (x_j - \hat{\mu})^T (x_j - \hat{\mu}) = \frac{1}{d} \sum_{i=1}^d \hat{\sigma}_i^2$$

3. 以下を証明せよ.

- x : 期待値が0で共分散行列が単位行列である
d次元正規分布に従う確率変数
- μ : d次元ベクトル
- Σ : d次元正値対称行列
- $\Sigma^{1/2}$: $(\Sigma^{1/2})^2 = \Sigma$ を満たすd次元行列(常に存在)

このとき,

$$y = \Sigma^{1/2} x + \mu$$

は, 期待値が μ で共分散行列が Σ の正規分布に従う.

4. 計算機実験: 次の3つガウスモデルを考える.

- Σ : 半正値対称行列 (パラメータ数 $d(d+1)/2$)
- $\Sigma = \text{diag}(\sigma_i^2)$ (パラメータ数 d)
- $\Sigma = \sigma^2 I$ (パラメータ数 1)

2次元 ($d = 2$) の適当な確率分布から適当な個数だけ標本を生成し, AICを使ってモデル選択を行え. そして, 真の分布や標本の個数を変化させ, 結果がどのように変化するか考察せよ.

- **確率分布の例**: 二次元の一様分布, 様々な共分散行列を持つ二次元の正規分布

- $[0,1]^d$ 上の一様分布に従う n 個の確率変数の生成:
> `rand(d,n)`
- 期待値が0で共分散行列が単位行列の d 次元の正規分布に従う n 個の確率変数の生成:
> `randn(d,n)`
- 期待値が μ で共分散行列が Σ の d 次元の正規分布に従う n 個の確率変数の生成:
> `sqrtm(Sigma)*randn(d,n)+Mu`
(SigmaとMuは各自定義せよ)
- 適当な d 次元正値行列 Σ の生成例:
> `A=randn(d,d); Sigma=A'*A;`
(Aは逆を持てば何でも良い)

- 計算機室のoctaveでグラフが表示されない場合は, aqua termという端末を起動してからoctaveを起動すればよいようである.