

応用システム工学
第4回 指数型分布族と
一般化線形モデル
平成22年6月25日

最尤推定1

- 確率変数ベクトル $y = [Y_1, Y_2, \dots, Y_n]^t$
 - 確率変数n個
- パラメータベクトル $\theta = [\theta_1, \theta_2, \dots, \theta_p]^t$
 - パラメータ数p個
- Y_i の同時確率密度関数 $f(y; \theta)$
 - パラメータ θ を固定した時の確率変数 y が対象
- 尤度関数 $L(\theta; y)$
 - 同時確率密度関数の数式表現と同一
 - θ と y の役割の入替
 - 確率変数 y を固定した時のパラメータ θ が対象
 - L は確率変数となる
 - 確率変数ベクトル y の関数

最尤推定2

- パラメータ空間 Ω $\theta \in \Omega$
 - パラメータベクトル θ がとる全ての値を含む集合
- パラメータベクトル θ の最尤推定量 $\hat{\theta}$
 - 尤度関数 L を最大にする θ
$$L(\hat{\theta}; y) \geq L(\theta; y), \theta \in \Omega$$
- 対数尤度関数 l
 - 尤度関数 L の対数をとったもの $l(\theta; y) = \log L(\theta; y)$
 - 対数関数は単調関数である
 - 対数尤度関数も単調関数となる
 - 最尤推定量は対数尤度関数を最大にする

$$l(\hat{\theta}; y) \geq l(\theta; y), \theta \in \Omega$$

最尤推定3

- 最尤推定量 $\hat{\theta}$ の求め方
 - 対数尤度関数の各パラメータに関する微分
 - 連立方程式

$$\frac{\partial l(\theta; y)}{\partial \theta_j} = 0, j = 1, 2, \dots, p \quad \text{極値の条件}$$

- 極大値との対応
 - 二階導関数が負定値行列

$$\frac{\partial^2 l(\theta; y)}{\partial \theta_j \partial \theta_k} \bigg|_{\theta = \hat{\theta}} < 0$$

- » 全ての局所最大(極大)値の中で、 l が最大値をとる θ が最尤推定量

最尤推定4

- 最尤推定量の不変性
 - パラメータ θ に対する任意の関数 $g(\theta)$
 - パラメータベクトル θ の最尤推定量 $\hat{\theta}$
 - $g(\theta)$ の最尤推定量も $g(\hat{\theta})$ となる
- 最尤推定量の求め方
 - 最尤推定を行うのに便利な関数を見つける
 - この関数に対して尤度関数を最大化する θ を見つける
 - 最尤推定量の不変性を用いる
 - 得られた θ は必要なパラメータの最尤推定量となる

残差とモデルの検証

- 反応変数 Y_i ,正規分布 $E(Y_i) = \mu_i; Y_i \sim N(\mu_i, \sigma^2)$
 - 期待値 μ_i ,分散 σ^2
- 残差 $y_i - \hat{\mu}_i$
 - 推定値 $\hat{\mu}_i$
- 標準化残差 $r_i = \frac{y_i - \hat{\mu}_i}{\hat{\sigma}}$
 - 未知パラメータ σ の推定値 $\hat{\sigma}$
- 残差平方和 $\sum (y_i - \hat{\mu}_i)^2$
 - モデルの適切さを計る統計量を与える

残差とモデルの検証

- 反応変数は互いに独立
 - 残差は近似的に平均0, 分散一定の正規分布に従うべき
- 標準化残差の正規分布との比較検討
 - 度数分布の正規分布との適合性の検討
 - 標準化残差は分布による異なる

推測と解釈

- 測定値(観測値)の構成
 - － 信号と雑音
- 節約の原則に基づくモデル化
 - － 効果を説明する以上の原因を仮定すべきでない
 - 説明されない変化を残さない複雑なモデルより, データを適度に説明する簡単なモデルが好ましい
 - － 節約的なモデルに用いるパラメータの仮説検定

統計モデル

- 正規分布に対するデータの分析
 - 独立な確率変数 Y_i
 - $Y_i \sim N(\mu_i, \sigma^2)$
 - 期待値の線形モデル
 - $E(Y_i) = \mu_i = X_i^T \beta$
 - 線形モデルの一般化
 - 正規分布以外の確率変数
 - 指数分布族(より広いクラスの分布) → 本授業の対象
 - 確率変数と説明変数の関係が非線形
 - 連結関数の適用 $g(\mu_i) = X_i^T \beta$
 - 一般化加法モデル
 - 大量の計算処理必要 → 本授業では扱わない
- 説明変数(デザイン行列): X
パラメータ: β

指数型分布族の確率変数

- 指数型分布族の条件

- 指数関数で表される確率分布

- 確率変数 Y , 単一パラメータ $\theta \rightarrow$ 最尤推定に便利

$$f(y; \theta) = s(y)t(\theta)\exp[a(y)b(\theta)]$$

- 但し, a, b, s, t は既知の関数

- $s(y) = \exp[d(y)], t(\theta) = \exp[c(\theta)]$ とおく

$$\begin{aligned} f(y; \theta) &= \exp[d(y)]\exp[c(\theta)]\exp[a(y)b(\theta)] \\ &= \exp[a(y)b(\theta) + c(\theta) + d(y)] \end{aligned}$$

指数型分布族の確率変数

$$f(y; \theta) = \exp[a(y)b(\theta) + c(\theta) + d(y)]$$

- $a(y)=y$ のとき「正準形」の分布と言い、このとき $b(\theta)$ を分布の自然パラメータと呼ぶ。
- 局外パラメータ
 - θ 以外のパラメータ
 - 既知として扱う必要あり
 - 単一パラメータを仮定しているため
- 指数型分布族の例
 - 正規分布, ポアソン分布, ワイブル分布, 二項分布

指数型分布族の対数尤度関数 最尤推定の準備

- 指数型分布族の確率密度関数

$$f(y; \theta) = \exp[a(y)b(\theta) + c(\theta) + d(y)]$$

- 指数型分布族の尤度関数

$$L(\theta; y) = \exp[a(y)b(\theta) + c(\theta) + d(y)]$$

- 指数型分布族の対数尤度関数

$$l(\theta; y) = \log L(\theta, y) = a(y)b(\theta) + c(\theta) + d(y)$$

指数型分布族としての正規分布

- 正規分布 $Y \sim N(\mu, \sigma^2)$
 - 対称な分布を持つ連続データに対するモデル
 - 確率密度関数

$$f(y; \mu) = \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} \exp\left[-\frac{1}{2\sigma^2}(y - \mu)^2\right]$$

$$= \exp\left[\log \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} - \frac{1}{2\sigma^2}(y - \mu)^2\right]$$

$$= \exp\left[-\frac{1}{2}\log(2\pi\sigma^2) - \frac{1}{2\sigma^2}(y^2 - 2y\mu + \mu^2)\right]$$

次頁に
つづく

指数型分布族としての正規分布

- 確率密度関数

$$f(y; \mu) = \exp\left[\frac{y\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2} - \frac{y^2}{2\sigma^2} - \frac{1}{2}\log(2\pi\sigma^2)\right]$$
$$= \exp[a(y)b(\mu) + c(\mu) + d(y)] \quad \Rightarrow a(y)b(\mu) = \frac{y\mu}{\sigma^2}$$

- パラメータ μ

- 正準形 $a(y) = y$

- 自然パラメータ $b(\mu) = \frac{\mu}{\sigma^2}$

- 局外パラメータ σ^2

- その他の項 $c(\mu) = -\frac{\mu^2}{2\sigma^2} - \frac{1}{2}\log(2\pi\sigma^2), d(y) = -\frac{y^2}{2\sigma^2}$

指数型分布族としてのポアソン分布

- ポアソン分布 $Y \sim \text{Poisson}(\theta)$
 - 計数データのモデルとして用いる
 - 確率変数: $Y \rightarrow$ 離散値, $y=1,2,3,\dots$ 等
 - 確率密度関数
$$f(y; \theta) = \frac{\theta^y e^{-\theta}}{y!}$$
$$= \exp(y \log \theta - \theta - \log y!)$$
$$= \exp[a(y)b(\theta) + c(\theta) + d(y)]$$
- 正準形: $a(y)=y$ である
- 自然パラメータ $b(\theta) = \log \theta$
- $c(\theta) = -\theta$, $d(y) = -\log y!$

指数型分布族としてのワイブル分布

- ワイブル分布

- 故障時間(生存期間)の分布を表すモデル

- 確率変数: y (時間)

- パラメータ: λ (分布の形状), θ (尺度)

- 確率密度関数

$$f(y; \theta, \lambda) = \frac{\lambda y^{\lambda-1}}{\theta^\lambda} \exp\left[-\left(\frac{y}{\theta}\right)^\lambda\right]$$
$$= \exp\left[\log \lambda + (\lambda - 1)\log y - \lambda \log \theta - \left(\frac{y}{\theta}\right)^\lambda\right]$$

指数型分布族としてのワイブル分布

$$f(y; \theta, \lambda) = \exp \left[\log \lambda + (\lambda - 1) \log y - \lambda \log \theta - \left(\frac{y}{\theta} \right)^\lambda \right]$$
$$= \exp [a(y)b(\theta) + c(\theta) + d(y)]$$

- パラメータ: θ
- $a(y) = y^\lambda$ $\lambda = 1$ のとき正準形となる
- 自然パラメータ $b(\theta) = -\theta^{-\lambda}$
– 局外パラメータ: λ
- その他の項 $c(\theta) = -\lambda \log \theta$
 $d(y) = \log \lambda + (\lambda - 1) \log y$

指数型分布族としての二項分布

- 二項分布 $Y \sim \text{Binominal}(n, \pi)$
 - 二値(0,1)をとるn回の独立試行に対する分布
 - 確率変数: $Y \rightarrow 1$ の回数。離散値, $y=1, 2, 3, \dots, n$
 - 二項係数 \rightarrow n個からy個を選ぶ組み合わせの数

$$\binom{n}{y} = {}_n C_y = C(n, y) = \frac{n!}{y!(n-y)!}$$

- 確率密度関数

$$f(y; \pi) = \binom{n}{y} \pi^y (1 - \pi)^{(n-y)}$$

指数型分布族としての二項分布

- 確率密度関数

$$\begin{aligned} f(y; \pi) &= \exp \left[\log \binom{n}{y} + y \log \pi + (n - y) \log(1 - \pi) \right] \\ &= \exp \left[y \log \pi - y \log(1 - \pi) + n \log(1 - \pi) + \log \binom{n}{y} \right] \\ &= \exp \left[y \log \frac{\pi}{1 - \pi} + n \log(1 - \pi) + \log \binom{n}{y} \right] \\ &= \exp [a(y)b(\pi) + c(\pi) + d(y)] \end{aligned}$$

指数型分布族としての二項分布

- パラメータ π

- 正準形

$$a(y) = y$$

- 自然パラメータ

$$b(\pi) = \log \frac{\pi}{1 - \pi}$$

– 局外パラメータ n

- その他の項 $c(\pi) = n \log(1 - \pi)$

$$d(y) = \log \binom{n}{y}$$