

応用システム工学
第7回 指数型分布族と
一般化線形モデル
平成22年7月16日

統計的推測

統計的推測

- 信頼区間(区間推定)
 - 検定結果:信頼区間が値を含むかどうか
 - 検定精度:信頼区間の幅
- 仮説検定
 - 2つの関連するモデルのデータへの適合度の判定
 - 2つの一般化線形モデル(同一確率分布, 連結関数)
 - パラメータの少ないモデル→帰無仮説
 - パラメータの多いモデル→対立仮説
 - 簡単なモデルが複雑なモデル程度に当てはまるかを評価する
 - 要約統計量(適合度統計量)を用いた当てはめの適合性
 - 尤度関数の最大値
 - 対数尤度関数の最大値
 - 平方和基準の最小値
 - 残差に基づく複合型の統計量

仮説検定

- 作業仮説・実験仮説を対立仮説 H_1
- 作業仮説を否定する仮説を帰無仮説 H_0
- 帰無仮説を仮説検定の対象とする。
 - 帰無仮説が棄却されると対立仮説が支持される。
 - 帰無仮説が支持される場合
 - 真に対立仮説が誤っている
 - 対立仮説は正しいが、標本の大きさが十分でなく、帰無仮説を積極的に棄却できない

仮説検定のプロセス

- モデルを決める
 - 帰無仮説 H_0 に対応するモデル M_0
 - 対立仮説 H_1 に対応するモデル M_1
- モデルに対する適合統計量を求める
 - $M_0 \rightarrow G_0$, $M_1 \rightarrow G_1$
- 適合のよさの比較・評価を行う
 - $G_1 - G_0$ または G_1 / G_0
- 対立仮説 $G_1 \neq G_0$ に対する帰無仮説 $G_1 = G_0$ の検定
 - $G_1 - G_0$ の標本分布を評価
 - $G_1 = G_0$ が棄却されない $\rightarrow H_0$ は棄却されない $\rightarrow M_0$ がより良いモデル
 - $G_1 = G_0$ が棄却される $\rightarrow H_0$ は棄却される $\rightarrow M_1$ がより良いモデル

一般化線形モデルの標本分布で当てはまりを評価

スコア統計量の標本分布

- 一般化線形モデル

- 互いに独立な確率変数 Y_1, \dots, Y_N

- パラメータ β

- 期待値 $\mu_i = E[Y_i]$

- 連結関数 $g(\mu_i) = x_i^T \beta = \eta_i$

- スコア統計量 $U_j = \frac{\partial l}{\partial \beta_j} = \sum_{i=1}^N \left[\frac{(Y_i - \mu_i)}{\text{var}(Y_i)} x_{ij} \left(\frac{\partial \mu_i}{\partial \eta_j} \right) \right], j = 1, \dots, p$

- スコア統計量の期待値

- $E[Y_i] = \mu_i$ より $E[U_j] = 0, j = 1, \dots, p$

スコア統計量の標本分布

- 情報行列 \mathfrak{S}

→スコア統計量の分散共分散行列

– 要素

$$\mathfrak{S}_{jk} = E[U_j U_k]$$

- 単一のパラメータ β に対するスコア統計量の標本分布について考える

スコア統計量の標本分布

- 尤度関数

$$l = \sum l_i$$

$$l_i = \log f(Y_i, \beta)$$

- スコア統計量

$$U = \sum U_i$$

$$l_i = \frac{\partial \log f(Y_i, \beta)}{\partial \beta}$$

確率変数 Y_i
の関数

スコア統計量の標本分布

- 確率変数 $Y_i, i=1, \dots, n$ が互いに独立
 - $U_i, i=1, \dots, n$ も互いに独立
 - 確率分布も同一となる

$$E[U] = 0, \text{var}[U] = \mathfrak{S}$$

多変量正規分布に漸近する(漸近標本分布)

→ $U \sim N(0, \mathfrak{S})$

もとの標本では

→ $U^T \mathfrak{S}^{-1} U \sim \chi^2(p)$

対数尤度関数の近似

- 統計量の漸近的な標本分布を得る

– テイラー級数展開

$$f(x) = f(t) + (x-t) \left[\frac{df}{dx} \right]_{x=t} + \frac{1}{2} (x-t)^2 \left[\frac{d^2 f}{dx^2} \right]_{x=t} + \dots$$

– 単一パラメータを持つ対数尤度関数のパラメータ β の推定値 b 近傍での近似

$$l(\beta) = l(b) + (\beta - b)U(b) + \frac{1}{2} (\beta - b)^2 U'(b) + \dots$$

$$\cong l(b) + (\beta - b)U(b) + \frac{1}{2} (\beta - b)^2 U'(b)$$

対数尤度関数の近似

- スコア関数($\beta = b$)

$$U(b) = \frac{dl}{d\beta}$$

- スコア関数の一階微分の期待値による近似

$$U' = \frac{d^2l}{d\beta^2} \quad E(U') = -\mathfrak{I}$$

$$l(\beta) \cong l(b) + (\beta - b)U(b) - \frac{1}{2}(\beta - b)^2 \mathfrak{I}(b)$$

対数尤度関数の近似

- パラメータベクトル β の対数尤度関数の近似

$$l(\beta) \cong l(b) + (\beta - b)U'(b) - \frac{1}{2}(\beta - b)^T \mathfrak{I}(b)(\beta - b)$$

- スコア関数のテイラー級数近似

$$U(\beta) = U(b) + (\beta - b)U'(b) + \frac{1}{2}(\beta - b)^2 U''(b) + \dots$$

$$\cong U(b) + (\beta - b)U'(b)$$

U' を期待値 $E[U']$ で近似 $\rightarrow -\mathfrak{I}$

$$\cong U(b) - (\beta - b)\mathfrak{I}(b)$$

- ベクトルの場合 $U(\beta) \cong U(b) - \mathfrak{I}(b)(\beta - b)$

最尤推定量の標本分布

- 最尤推定量 $b = \hat{\beta}$

– 対数尤度関数を最大にする $U(b) = 0$

$$U(\beta) \cong U(b) - \mathfrak{J}(b)(\beta - b) = -\mathfrak{J}(b)(\beta - b)$$

- 情報統計量が正則なら

$$\mathfrak{J}(b)^{-1} U(\beta) = -(\beta - b)$$

- 情報統計量が一定ならスコア関数の期待値は0より

$$E[\mathfrak{J}(b)^{-1} U(\beta)] = \mathfrak{J}(b)^{-1} E[U(\beta)] = E[-(\beta - b)] = 0$$

b と β は一致する

最尤推定量の標本分布

- 情報行列 $\mathfrak{J} = E[UU^T]$
 - 対称性 $(\mathfrak{J}^{-1})^T = \mathfrak{J}^{-1}$
 - 最尤推定量の分散共分散行列

$$E[(b - \beta)(b - \beta)^T] = \mathfrak{J}^{-1} E[UU^T] \mathfrak{J}^{-1} = \mathfrak{J}^{-1}$$

- パラメータ数 p (自由度 p)のカイ二乗分布

$$(b - \beta)^T \mathfrak{J}^{-1} (b - \beta) \sim \chi^2(p)$$

- ワルド統計量という

単一パラメータの時 $b \sim N(\beta, \mathfrak{J}^{-1})$ と表せる

対数尤度比統計量

- モデルの適切さの評価
 - 飽和モデル(最大モデル, フルモデル) と比較
 - 推定されうる最大個数のパラメータを含む最も一般的なモデル
 - N 個の観測値 $Y_i, i=1, \dots, N$ があれば, N 個の線形成分 $X_iT \beta$ を対応させる N 個のパラメータでモデルを表せる
 - 推定できるパラメータの数は, 相異なる線形成分の数に等しい → 繰り返し
 - 同じ線形成分(同じ共変量)をもつ観測値では, 連続的な説明変数も同じになる

対数尤度比統計量

- 対象とするモデルは、飽和モデルと確率分布、連結関数が同じ一般化線形モデル
 - 飽和モデルのパラメータ数: m
 - 飽和モデルのパラメータベクトル: β_{\max}
 - パラメータベクトルの最尤推定量: b_{\max}
 - 飽和モデルの尤度関数は最も大きい $L(b_{\max}; y)$
 - 関心のあるモデルの尤度関数との比でモデルの適合度を評価: 尤度比

$$\lambda = \frac{L(b_{\max}; y)}{L(b; y)}$$

対数尤度比統計量

- 尤度比の対数→対数尤度関数の差

$$\log \lambda = l(b_{\max}; y) - l(b; y)$$

- $\log \lambda$ が大きいとモデルの適合が悪い
 - 標本分布を基に棄却域を決定
 - $2\log \lambda$: 逸脱度
 - カイ二乗分布

逸脱度の標本分布

- 逸脱度→対数尤度比統計量:D

$$D = 2[l(b_{\max}; y) - l(b; y)]$$

- パラメータ β の最尤推定量 b $U(b) = 0$

– パラメータベクトル β の対数尤度関数の近似
(テーラー展開)

$$l(\beta) \cong l(b) + (\beta - b)U(b) - \frac{1}{2}(\beta - b)^T \mathfrak{I}(b)(\beta - b)$$

$$l(\beta) - l(b) = -\frac{1}{2}(\beta - b)^T \mathfrak{I}(b)(\beta - b)$$

$$2[l(\beta) - l(b)] = -(\beta - b)^T \mathfrak{I}(b)(\beta - b) \sim \chi^2(p)$$

逸脱度の標本分布

- 逸脱度D

$$D = 2[l(b_{\max}; y) - l(b; y)]$$

$$= 2[l(b_{\max}; y) - l(\beta_{\max}; y)]$$

$$- 2[l(b; y) - l(\beta; y)]$$

$$+ 2[l(\beta_{\max}; y) - l(\beta; y)]$$

飽和モデル:パラメータ数m:自由度mのカイ二乗分布

対象モデル:パラメータ数p:自由度pのカイ二乗分布

対象モデルが飽和モデルに適合すると零に近い正の数

逸脱度の標本分布の近似 $D \sim x^2(m-p, \nu)$ 非心パラメータ ν

逸脱度の標本分布

- 確率変数 Y_i が正規分布
 - 逸脱度 D はカイ二乗分布
 - ただし分散 $\text{var}[Y_i] = \sigma^2$ に依存
- 確率変数 Y_i が他の分布
 - 逸脱度は D 近似的にカイ二乗分布
 - ちゃんと計算できるものも或る

仮説検定

- 検定の方法
 - ワールド統計量による検定
 - スコア統計量による検定
 - 逸脱度統計量の差を基にした検定
 - 単純化したモデル M_0 →帰無仮説 H_0
 - 一般的(複雑な)モデル M_1 →対立仮説 H_1
 - 検定を行うには, M_0 の線形成分が M_1 の線形成分の特別な場合となっている必要有

仮説検定

- モデルM0の帰無仮説H0

$$H_0 : \beta = \beta_0 = [\beta_1, \dots, \beta_q]^t$$

- モデルM1の対立仮説H1

$$H_1 : \beta = \beta_1 = [\beta_1, \dots, \beta_p]^t$$

– ただし $q < p < N$

- 対立仮説H1に対する帰無仮説H0の検定

仮説検定

- 対立仮説H1に対する帰無仮説H0の検定

- 逸脱度統計量の差

$$\begin{aligned}\Delta D &= D_0 - D_1 = 2[l(b_{\max}; y) - l(b_0; y)] - 2[l(b_{\max}; y) - l(b_1; y)] \\ &= 2[l(b_1; y) - l(b_0; y)]\end{aligned}$$

- 帰無仮説H0の支持

$$D_0 \sim \chi^2(N - q), D_1 \sim \chi^2(N - p) \Rightarrow \Delta D \sim \chi^2(p - q)$$

- どちらのモデルもよく一致する→より簡単なM0を支持

- 帰無仮説H0の棄却, 対立仮説H1支持

- ΔD が棄却域。カイ二乗分布上側 $100 \times \alpha\%$ より大

- M1の方がM0より有意によく記述している

- » M1がよくあてはまるとは限らない