

A variable selection criterion for GEE with large cluster sizes

広島大・理 佐藤 倫治

医療, 経済などのあらゆる分野において, 相関のあるデータに関する解析が議論されている. その中でも経時データは, 観測対象者を時間経過に伴って観測した繰り返し測定データのことで, 同一個体間のデータは相関を持ち, 異個体間のデータは独立であるという特徴を持つ. この相関のある経時データを解析する手法の一つに Liang & Zeger (1986) で提案された一般化推定方程式 (GEE) がある. 周辺モデルは以下のように表される.

目的変数ベクトル $\mathbf{y}_i = (y_{i1}, \dots, y_{im})'$ と説明変数行列 $\mathbf{X}_i = (\mathbf{x}_{i1}, \dots, \mathbf{x}_{im})'$ の組 $(\mathbf{y}_i, \mathbf{X}_i)$ が与えられているとする. ただし, $i = 1, \dots, n$ である. また, \mathbf{x}_{ij} は $(p \times 1)$ 説明変数ベクトルであり, $i \neq k$ ならば \mathbf{y}_i と \mathbf{y}_k は独立である. ただし, \mathbf{y}_i 内には相関があるとする. また, y_{ij} の密度関数は,

$$f(y_{ij}; \theta_{ij}, \phi) = \exp \{ \{y_{ij}\theta_{ij} - a(\theta_{ij})\} / \phi + b(y_{ij}, \phi) \},$$

であると仮定する. 平均は $E[Y_{ij}] = \dot{a}(\theta_{ij})$, 分散は $\text{Var}[Y_{ij}] = \ddot{a}(\theta_{ij})\phi$, と表せる. このとき, 次のモデル, $\eta_{ij} = \mathbf{x}'_{ij}\boldsymbol{\beta}$, $\theta_{ij} = u(\eta_{ij})$ を考える. $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)'$ は未知パラメータ, ϕ は尺度パラメータである. $\boldsymbol{\beta}$ は以下で定義される GEE により推定される.

$$\sum_{i=1}^n \mathbf{D}'_i(\boldsymbol{\beta}) \mathbf{V}_i^{-1}(\boldsymbol{\beta}) \{ \mathbf{y}_i - \boldsymbol{\mu}_i(\boldsymbol{\beta}) \} = \mathbf{0}_p.$$

ここで, $\boldsymbol{\mu}_i = (\dot{a}(\theta_{i1}), \dots, \dot{a}(\theta_{in}))'$, $\mathbf{D}_i(\boldsymbol{\beta}) = \mathbf{A}_i(\boldsymbol{\beta}) \boldsymbol{\Delta}_i(\boldsymbol{\beta}) \mathbf{X}_i$, $\mathbf{V}_i(\boldsymbol{\beta}) = \mathbf{A}_i^{1/2}(\boldsymbol{\beta}) \mathbf{R}(\boldsymbol{\alpha}) \mathbf{A}_i^{1/2}(\boldsymbol{\beta}) \phi$, $\mathbf{A}_i = \text{diag}\{\ddot{a}(\theta_{i1}), \dots, \ddot{a}(\theta_{in})\}$, $\boldsymbol{\Delta}_i = \text{diag}\{\partial\theta_{i1}/\partial\eta_{i1}, \dots, \partial\theta_{in}/\partial\eta_{in}\}$ である. また, $\mathbf{R}(\boldsymbol{\alpha})$ は作業用相関行列であり, $\boldsymbol{\alpha}$ は相関パラメータである.

Inatsu & Imori (2013) では, GEE を用いたモデルに対する変数選択規準を提案した. \mathbf{z}_i を \mathbf{y}_i と独立に同一の分布に従うデータとしたとき, モデルの良さを測る尺度として以下のリスクを用いた.

$$E_y \left[E_z \left[\sum_{i=1}^N (\mathbf{z}_i - \hat{\boldsymbol{\mu}}_i)' \text{Cov}[\mathbf{z}_i]^{-1} (\mathbf{z}_i - \hat{\boldsymbol{\mu}}_i) \right] \right]. \quad (1)$$

そして (1) の有効な推定量を求めることで, 相関構造の影響を反映させた変数選択規準として,

$$\text{PMSEG} = \sum_{i=1}^n \{ \mathbf{y}_i - \boldsymbol{\mu}_i(\hat{\boldsymbol{\beta}}) \}' \mathbf{A}_i^{-1/2}(\hat{\boldsymbol{\beta}}_f) \mathbf{R}_0^{-1}(\hat{\boldsymbol{\beta}}_f) \mathbf{A}_i^{-1/2}(\hat{\boldsymbol{\beta}}_f) \{ \mathbf{y}_i - \boldsymbol{\mu}_i(\hat{\boldsymbol{\beta}}) \} + 2p,$$

を提案した. ここで, $\hat{\boldsymbol{\beta}}$ と $\hat{\boldsymbol{\beta}}_f$ はそれぞれ候補モデルとフルモデルにおける GEE 推定量であり, \mathbf{R}_0 は真の相関行列の一致推定量である.

Inatsu & Imori (2013) では, 標本数 n だけが大きくなる大標本での変数選択規準を導出していた. 本発表では, 標本数 n だけでなく, 観測時点数 m も共に大きくなる大標本高次元での変数選択規準の導出について話す.

参考文献

- [1] Liang, K. Y. & Zeger, S. L. (1986). Longitudinal data analysis using generalized linear models. *Biometrika*, **73**, 13-22.
- [2] Inatsu, Y. & Imori, S. (2013). Model selection criterion based on the prediction mean squared error in generalized estimating equations. *TR 13-10*, *Statistical Research Group*, Hiroshima University, Hiroshima.