

# 一般化線形混合モデルにおけるモデル選択

吉田拓真<sup>1</sup> 勘場 大<sup>1</sup> 内藤 貫太<sup>2</sup>

<sup>1</sup> 島根大学大学院総合理工学研究科

<sup>2</sup> 島根大学総合理工学部 E-mail: naito@riko.shimane-u.ac.jp

## 1 はじめに

一般化線形スプライン回帰において、含まれるパラメータの推定に正則化法を用いるとき、あらかじめ、ノットの個数と平滑化パラメータを決定しておく必要があり、一般化情報量規準 GIC(小西, 北川(2004)) などが利用される。しかし、情報量規準による平滑化パラメータの決定はグリッドサーチによるため、多くの平滑化パラメータを含むモデルの場合、その計算量は膨大なものとなり、効率的であるとはいえない。

その一方で、一般化スプライン回帰は、スプラインの部分の係数をランダム効果として捉えることにより一般化線形混合モデルとなり、罰則付き擬似尤度 (PQL) から推定量、予測量を求めることができ、一般化スプライン回帰における平滑化パラメータが直接推定できる。本講演では、これらの議論を GIC に応用させることで、時間効率の良い新たな情報量規準を提案する。

## 2 問題設定

$n$  個のデータ  $\{(y_i, x_{i1}, \dots, x_{ip}) | i = 1, \dots, n\}$  が観測されたとき、未知の回帰関数  $f(x) = E[y|x]$ ,  $(x = [x_1 \dots x_p]^T)$  の推定を考える。そこで、 $x_i = [x_{i1} \dots x_{ip}]^T$  を与えたときの  $y_i$  が指数分布族の密度

$$f(y_i; \eta_i) = \exp\left(\frac{y_i \eta_i - b(\eta_i)}{\phi} + h(y_i, \phi)\right), \quad i = 1, \dots, n$$

を持つと仮定する。ここで、 $b'(\eta_i) = \mu_i = E[y_i|x_i]$  であり、 $\phi b''(\eta_i) = V[y_i|x_i]$  である。 $\eta(\cdot)$  を未知の回帰関数として、 $\eta_i = \eta(x_i)$  を扱い、

$$\eta(x_i) = \beta_0 + \sum_{j=1}^p \beta_j x_{ij} + \sum_{k=1}^{K_1} u_{1k} (x_{i1} - \kappa_{1k})_+ + \dots + \sum_{k=1}^{K_p} u_{pk} (x_{ip} - \kappa_{pk})_+, \quad i = 1, \dots, n$$

を想定する。ここで、 $\kappa_{11}, \dots, \kappa_{1K_1}, \dots, \kappa_{p1}, \dots, \kappa_{pK_p}$  はノット、 $(x)_+ = \max\{x, 0\}$  である。また、 $\eta = [\eta(x_1) \dots \eta(x_n)]^T$ ,  $\beta = [\beta_0 \beta_1 \dots \beta_p]^T$ ,  $\mathbf{u}_j = [u_{j1} \dots u_{jK_j}]^T$ ,  $\mathbf{u} = [\mathbf{u}_1^T \dots \mathbf{u}_p^T]^T$ ,  $\mathbf{z}_i = [(x_{i1} - \kappa_{11})_+ \dots (x_{i1} - \kappa_{1K_1})_+ \dots (x_{ip} - \kappa_{p1})_+ \dots (x_{ip} - \kappa_{pK_p})_+]^T$ ,  $\mathbf{Z} = [\mathbf{z}_1 \dots \mathbf{z}_n]^T$ ,  $\mathbf{X} = [x_1 \dots x_n]^T$  を導入すると、 $\eta = \mathbf{X}\beta + \mathbf{Z}\mathbf{u}$  と表せ、 $\mathbf{c}_i = [x_i^T \mathbf{z}_i^T]^T$ ,  $\mathbf{C} = [\mathbf{X} \mathbf{Z}] = [\mathbf{c}_1 \dots \mathbf{c}_n]^T$  とおき、 $\mathbf{w} = [\beta^T \mathbf{u}^T]^T$  とすると  $\eta = \mathbf{C}\mathbf{w}$  と表せる。

このとき、正則化法による  $\beta$ ,  $\mathbf{u}$  の推定量は、更新式

$$\begin{bmatrix} \beta^{new} \\ \mathbf{u}^{new} \end{bmatrix} \leftarrow (\mathbf{C}^T \mathbf{W} \mathbf{C} + n\phi\lambda)^{-1} \mathbf{C}^T \mathbf{W} \left\{ \mathbf{C} \begin{bmatrix} \beta^{old} \\ \mathbf{u}^{old} \end{bmatrix} + \mathbf{W}^{-1}(\mathbf{y} - \mu) \right\} \quad (1)$$

より得られ、収束値  $[\hat{\beta} \hat{\mathbf{u}}]$  を用いて  $\hat{\eta} = \mathbf{X}\hat{\beta} + \mathbf{Z}\hat{\mathbf{u}}$  となる。ただし、 $\mu = [\mu_1 \dots \mu_n]^T$ ,  $b'(\eta_i) = \mu_i = E[y_i|x_i]$ ,  $\mathbf{W} = \text{diag}[b''(\mathbf{X}\beta + \mathbf{Z}\mathbf{u})]$  であり、 $\lambda = \text{diag}[0 \dots 0 \lambda_1 \dots \lambda_p]$  は平滑化パラメータである。

### 3 統計的モデル選択

正則化法によって推定されたスプラインモデル  $f(y|\hat{w}, \hat{\phi})$  のノットの個数  $K = \{K_1, \dots, K_p\}$  , 平滑化パラメータ  $\lambda$  の選択は GIC を用いて ,  $GIC = GIC(K, \lambda)$  を最小とする  $K, \lambda$  として得られる . しかし , GIC を最小にする  $\lambda_1, \dots, \lambda_p$  は , グリッドサーチによって決定するものであるため , 時間コストの観点からは効率的であるとは言えない . そこで我々は GIC によるモデル選択において効率の良い平滑化パラメータの決定を以下に提案する .

### 4 一般化線形混合モデル

$y = [y_1 \dots y_n]^T$  の密度における  $u$  を  $u \sim N(0, G)$  とすると , モデルは一般化線形混合モデルとなる . ただし ,  $G = \text{diag}[\sigma_1^2 \mathbf{1}_{K_1} \dots \sigma_p^2 \mathbf{1}_{K_p}]$  である . そこで , PQL を用いて ,  $\beta$  の推定量と  $u$  の予測量を求めると , PQL による  $[\beta^T u^T]^T$  の更新式は

$$\begin{bmatrix} \beta^{new} \\ u^{new} \end{bmatrix} \leftarrow (C^T W C + \phi B)^{-1} C^T W \left\{ C \begin{bmatrix} \beta^{old} \\ u^{old} \end{bmatrix} + W^{-1}(y - \mu) \right\} \quad (2)$$

で与えられる . ここで ,  $B = \text{diag}[0 \dots 0 (1/\sigma_1^2) \mathbf{1}_{K_1} \dots (1/\sigma_p^2) \mathbf{1}_{K_p}]$  ,  $\mu = b'(X\beta + Zu)$  ,  $W = \text{diag}[b''(X\beta + Zu)]$  である . (2) は , (1) において ,  $B = n\lambda$  とした式と一致している . したがって ,  $B$  の中に含まれるパラメータを推定することができれば , 平滑化パラメータ  $\lambda_1, \dots, \lambda_p$  の推定値は  $\hat{\lambda}_i = 1/(n\hat{\sigma}_i^2)$  , ( $i = 1, \dots, p$ ) として得られることになる .  $B$  の推定量は , 擬似データ  $y_p = X\beta + Zu + W^{-1}(y - \mu) = X\beta + Zu + \varepsilon_p$  を考えたとき ,  $y_p \sim N(X\beta, V)$  ,  $V = ZGZ^T + \phi W$  を仮定し , 最尤法により推定できる (Ruppert et al(2003)) . そこで , GIC において ,  $\hat{\lambda} = \hat{B}/n$  としたとき , ノットの個数  $K$  を  $GIC_P = GIC(K, \hat{B}/n)$  を最小とするものとして選択する .

$GIC_P$  では平滑化パラメータを推定するので , 第 3 節で扱った手法のようにグリッドサーチにより探索する必要が無い分 , 計算コストの観点から効率的であると言える .

### 5 比較

線形モデルや線形加法モデル , 変化係数モデルに対して , 実データへの当てはめ , 予測精度 , シミュレーションにおける GIC と  $GIC_P$  の比較を行った結果 , 予測精度や曲線の形状はほとんど同じであったが , 計算時間は  $GIC_P$  の方がはるかに早いことがわかった .

### 6 まとめと展望

5 節より , GIC と  $GIC_P$  は , 共に良いモデル選択規準となることがわかったが , 計算時間が早い分 ,  $GIC_P$  は GIC の簡便版として有用である . また , セミパラメトリック回帰の枠組みの中で , 空間モデル (Ruppert et al(2003)) にも  $GIC_P$  を適用できるだろう .

### 参考文献

- [1] Ruppert,D., Wand,M.P. and Carroll,R.J.(2003). *Semiparametric Regression*, Cambridge University Press.
- [2] 小西貞則 , 北川源四郎 (2004) . 情報量規準 . 朝倉書店.