

# スパースな共分散構造を持つ多変量線形回帰モデルにおける 高次元情報量規準

大阪大学大学院基礎工学研究科 片山 翔太

広島大学大学院理学研究科 伊森 晋平

次の多変量線形回帰モデルを考える。

$$\mathbf{y}_i = \mathbf{B}^{*T} \mathbf{x}_i + \boldsymbol{\varepsilon}_i = \sum_{j=1}^{K_n} x_{ij} \mathbf{b}_j^* + \boldsymbol{\varepsilon}_i, \quad i = 1, \dots, n.$$

ここで、 $\mathbf{y}_i$  は  $p_n$  次元結果変数ベクトル、 $\mathbf{x}_i = (x_{i1}, \dots, x_{iK_n})^T$  は  $K_n$  次元説明変数ベクトル、 $\mathbf{B}^* = (\mathbf{b}_1^*, \dots, \mathbf{b}_{K_n}^*)^T$  は未知の  $K_n \times p_n$  係数行列をそれぞれ表し、 $\boldsymbol{\varepsilon}_i$  はそれぞれ独立に  $p_n$  次元多変量正規分布  $N_{p_n}(\mathbf{0}, \boldsymbol{\Sigma}^*)$  に従う。いま、真のモデル  $M_{k_n}^* = \{i | \mathbf{b}_i^* \neq \mathbf{0}\}$  をある候補モデルの集合  $\mathcal{M}$  から情報量規準によって選択することを考える。なお下添字  $k_n^*$  はその集合の要素数を表す。この問題に対して、Nishii(1984) で提案された一般化情報量規準 (GIC) は、ある候補モデル  $M_{k_n} \in \mathcal{M}$  に対して次で与えられる。

$$\text{GIC}_{a_n}(M_{k_n}) = \log \det(\mathbf{S}_{k_n}) + a_n p_n k_n.$$

ここで、 $\mathbf{S}_{k_n} = n^{-1} \mathbf{Y}^T (\mathbf{I}_n - \mathbf{P}_{k_n}) \mathbf{Y}$ 、 $\mathbf{P}_{k_n} = \mathbf{X}_{k_n} (\mathbf{X}_{k_n}^T \mathbf{X}_{k_n})^{-1} \mathbf{X}_{k_n}^T$ 、 $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)^T$ 、 $\mathbf{X}_{k_n} = (x_{ij} : 1 \leq i \leq n, j \in M_{k_n})$  であり、 $a_n$  はある正の調節パラメータである。

近年の情報技術の発展に伴い、 $p_n$  および  $K_n$  が標本サイズ  $n$  に比して大きなデータ (高次元データ) が、マーケティングリサーチ、ウェブ解析、生物情報など多種多様な分野で収集されている。しかしながら、このような高次元データに対して上記の GIC は、 $\mathbf{X}_{k_n}^T \mathbf{X}_{k_n}$  の特異性および  $\mathbf{S}_{k_n}$  の特異性といった2つの問題により定義出来ない。そこで本報告では、高次元データにおけるこれらの問題点を解決した新しい GIC を提案する。最初の問題点に対しては、Chen and Chen (2008) で導入された制限候補モデルを考える。すなわち候補モデルを、 $k_n \leq \min(n, K_n)$  を満たすものに制限する。残りの問題点に対しては、Yuan and Lin (2007) で提案された精度行列のスパース推定量を用いる。具体的には、 $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$  として、係数行列  $\mathbf{B}^*$  と精度行列  $\boldsymbol{\Omega}^* = \boldsymbol{\Sigma}^{*-1}$  に対する  $\ell_1$  制約付き対数尤度関数

$$Q(\mathbf{B}, \boldsymbol{\Omega}) = \frac{1}{n} \text{tr} \boldsymbol{\Omega} (\mathbf{Y} - \mathbf{X} \mathbf{B})^T (\mathbf{Y} - \mathbf{X} \mathbf{B}) - \log \det(\boldsymbol{\Omega}) + \tau_n \|\boldsymbol{\Omega}^-\|_1$$

を考え、次のように高次元 GIC (HGIC) を構成する。

$$\text{HGIC}_{a_n}(M_{k_n}) = Q(\hat{\mathbf{B}}_{k_n}, \hat{\boldsymbol{\Omega}}_{k_n}) + a_n p_n k_n$$

なお、 $\|\boldsymbol{\Omega}\|_1$  は  $\boldsymbol{\Omega}$  の要素ごとの  $\ell_1$  ノルムを意味し、 $\tau_n > 0$ 、 $\boldsymbol{\Omega}^- = \boldsymbol{\Omega} - \text{diag}(\boldsymbol{\Omega})$  であり、

$$(\hat{\mathbf{B}}_{k_n}, \hat{\boldsymbol{\Omega}}_{k_n}) = \underset{(\mathbf{B}, \boldsymbol{\Omega}) \in \Theta(M_{k_n}) \times \Xi}{\text{argmin}} Q(\mathbf{B}, \boldsymbol{\Omega})$$

である。ここで、 $\Theta(M_{k_n}) = \{\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_{K_n}) | \mathbf{b}_i = \mathbf{0}, i \in M_{k_n}^c\}$ 、 $\Xi = \{\boldsymbol{\Omega} | \boldsymbol{\Omega} = \boldsymbol{\Omega}^T\}$  である。このようにして定義された HGIC は、調節パラメータ  $a_n$  を適切なレートで選べば  $(n, p_n, K_n) \rightarrow \infty$  の下で一致性を持つことが示される。詳細は当日報告する。