

# 線形混合モデルと小地域推定

- 特に変数選択規準の導入を巡って -

久保川 達也

( 東京大学 大学院経済学研究科, E-mail: tatsuya@e.u-tokyo.ac.jp )

## 1 線形混合モデルを利用した小地域推定

$X_i$  を  $n_i \times p$  行列,  $\beta = (\beta_0, \dots, \beta_{p-1})'$ ,  $Z_i$  を  $n_i \times m_i$  フルランク行列,  $G_i$  を  $m_i \times m_i$  正定値行列とする一般的な線形混合モデル

$$y_i = X_i \beta + Z_i v_i + \epsilon_i, \quad i = 1, \dots, k,$$

について考察した。ここで, 変量効果  $v_i$  と誤差項  $\epsilon_i$  は互いに独立に分布し,

$$v_i \sim \mathcal{N}_{m_i}(\mathbf{0}, \sigma^2 G_i), \quad \epsilon_i \sim \mathcal{N}_{n_i}(\mathbf{0}, \sigma^2 I_{n_i})$$

とする。  $N = \sum_{i=1}^k n_i$ ,  $M = \sum_{i=1}^k m_i$ ,  $\mathbf{y} = (y'_1, \dots, y'_k)'$ ,  $\mathbf{X} = (X'_1, \dots, X'_k)'$  とし,  $\mathbf{v}$ ,  $\epsilon$  も同様に定義し,  $\mathbf{Z} = \text{block diagonal}(\mathbf{Z}_1, \dots, \mathbf{Z}_k)$  とおくと,  $\mathbf{R} = \mathbf{I}_N$  に対して

$$\mathbf{y} = \mathbf{X}\beta + \mathbf{Z}\mathbf{v} + \epsilon, \quad \mathbf{v} \sim \mathcal{N}_M(\mathbf{0}, \sigma^2 \mathbf{G}), \quad \epsilon \sim \mathcal{N}_N(\mathbf{0}, \sigma^2 \mathbf{R}),$$

と表される。ただし,  $\mathbf{G} = \text{block diag}(\mathbf{G}_1, \dots, \mathbf{G}_k)$  である。枝分かれ誤差回帰モデル (NERM), Fay-Herriot モデル, 経時測定 線形混合モデルはこのモデルに含まれる。

Henderson (1950) の混合モデル方程式を解くことにより,  $\beta$ ,  $\mathbf{v}$  の推定量

$$\hat{\beta} = (\mathbf{X}'\Sigma^{-1}\mathbf{X})^{-1}\mathbf{X}'\Sigma^{-1}\mathbf{y}, \quad \hat{\mathbf{v}} = \mathbf{G}\mathbf{Z}'\Sigma^{-1}(\mathbf{y} - \mathbf{X}\hat{\beta})$$

が得られる。共分散行列  $\mathbf{G}$ ,  $\mathbf{R}$  が既知の場合に, 既知のベクトル  $\mathbf{a} \in R^p$ ,  $\mathbf{b} \in R^q$  に対して  $\mu = \mathbf{a}'\beta + \mathbf{b}'\mathbf{v}$  を推定したいときには,  $\mu$  の最良線形不偏予測量 (BLUP) は,

$$\hat{\mu} = \mathbf{a}'\hat{\beta} + \mathbf{b}'\mathbf{G}\mathbf{Z}'\Sigma^{-1}(\mathbf{y} - \mathbf{X}\hat{\beta})$$

で与えられる。 $\mathbf{G}$ ,  $\mathbf{R}$  が未知の場合には, ML, REML などにより推定量  $\hat{\mathbf{G}}$ ,  $\hat{\mathbf{R}}$  が求められ, これらを代入することにより, 経験最良線形不偏予測量 (EBLUP) が得られる。

集計データに基づいた地域レベルモデルとして用いられる Fay-Herriot モデル

$$\bar{y}_i = \bar{\mathbf{x}}'_i \beta + v_i + \epsilon_i, \quad i = 1, \dots, k$$

については,  $\mu_i = \bar{\mathbf{x}}'_i \beta + v_i$  の EBLUP は,

$$\hat{\mu}_i(\hat{\psi}) = \bar{\mathbf{x}}'_i \hat{\beta}(\hat{\psi}) + (1 - \hat{\gamma}_i) \left( \bar{y}_i - \bar{\mathbf{x}}'_i \hat{\beta}(\hat{\psi}) \right), \quad \hat{\gamma}_i = (1 + n_i \hat{\psi})^{-1},$$

と表される。ここで、 $\psi = \sigma_v^2/\sigma^2$ ,  $V(\psi) = \psi I_k + D$  に対して、 $\hat{\beta}(\psi) = (X'V(\psi)^{-1}X)^{-1}X'V(\psi)^{-1}y$  である。

$\hat{\psi}$  の推定方法について紹介し、小地域の推定精度を高めるために EBLUP  $\hat{\mu}_i(\hat{\psi})$  が有用であること、特に線形混合モデルのどのような仕組みが推定精度を高めるのに役立つのかについて説明した。また、EBLUP の平均 2 乗誤差の推定、EBLUP に基づいた信頼区間の構成と補正方法、 $\beta$  に関する線形仮説に関する GLS 検定の Bartlett 型補正について説明した。

## 2 線形混合モデルにおけるベイズ情報量規準 (BIC)

次に、 $X$  と  $v$  についての変数選択規準として BIC, 修正 BIC を導出した。 $V = ZGZ' + I_N$ ,  $\|u\|_A^2 = u'Au$  に対して、 $\hat{\beta}(V) = (X'V^{-1}X)^{-1}X'V^{-1}y$ ,  $\hat{\sigma}_0^2(V) = \|y - X\hat{\beta}(V)\|_{V^{-1}}^2/N$  とおくと、BIC は、

$$BIC(\hat{\beta}(V), \hat{\sigma}_0^2(V)) = N[\log(2\pi\hat{\sigma}_0^2(V)) + 1] + \log|V| + r_{(X)} \log(N)$$

で与えられる。 $V$  もしくは  $G$  が未知母数  $\psi = (\psi_1, \dots, \psi_d)$  の関数で表されるときには、BIC は

$$BIC(\hat{\beta}(\hat{V}), \hat{\sigma}_0^2(\hat{V}), \hat{V}) = N[\log(2\pi\hat{\sigma}_0^2(\hat{V})) + 1] + \log|\hat{V}| + (r_{(X)} + d) \log(N)$$

と書かれる。ここで、 $\hat{\psi} = (\hat{\psi}_1, \dots, \hat{\psi}_d)$  は、 $\hat{\psi} - \psi = O_p(N^{-1/2})$  なる  $\psi$  の推定量で、 $G$ ,  $V$  は  $\hat{G} = G(\hat{\psi})$ ,  $\hat{V} = V(\hat{\psi})$  で推定されるものとする。

修正 BIC は、尤度関数を群内分散と群間分散の項に分割し、それぞれに BIC を当てはめることにより導かれ、

$$\begin{aligned} mBIC(\hat{G}) = & N[\log(2\pi\hat{\sigma}_1^2) + 1] + \log|\hat{V}| + \|\tilde{y}_2 - \tilde{X}_2\hat{\beta}(\hat{V})\|_{\tilde{W}_2^{-1}}^2/\hat{\sigma}_1^2 \\ & - M + r_{(\tilde{X}_1)} \log(N - M) + (r_{(\tilde{X}_2)} + d) \log(M) \end{aligned}$$

で与えられる。記号の具体的な意味については、論文を参照してほしい。

NERM と経時測定 線形混合モデルについて、BIC, 修正 BIC を導出した。数値実験により、真の変数の組を選択する意味において、BIC, 修正 BIC は AIC, CAIC より優れていることがわかった。特に、 $k$  が小さいときには、修正 BIC は BIC よりよいようである。真のモデルが変量効果を含む場合と含まない場合を調べたが、いずれの場合も修正 BIC の挙動がよいことが示された。

最後に、京浜急行電鉄本線及び久里浜線沿いの宅地物件について 1997 年から 2001 年までの 5 年間に公表された  $1m^2$  当たりの地価公示価格のデータを利用して、修正 BIC が選択する説明変数を調べてみた。各駅を 1 つの小地域と考え、また  $i$  番目の駅を最寄り駅とする物件のデータをその小地域からとられたデータと考えて、NERM と経時測定 線形混合モデルに修正 BIC を適用してみたところ、合理的な線形混合モデルが選択された。