

統計解析

古谷知之

授業概要

- * 履修者の状況に応じて変更される場合がありますが、全体としては以下のような授業構成となります。
- * 講義の中でR演習を行うこともあります。

第1回	ガイダンス・単回帰分析	第8回	一般化線形回帰モデル(5)
第2回	重回帰分析(1)	第9回	一般化線形回帰モデル(6)
第3回	重回帰分析(2)	第10回	一般化線形混合モデル
第4回	一般化線形回帰モデル(1)	第11回	状態空間モデル
第5回	一般化線形回帰モデル(2)	第12回	R演習(1)
第6回	一般化線形回帰モデル(3)	第13回	R演習(2)
第7回	一般化線形回帰モデル(4)	第14回	R演習(3)

統計モデルの種類

	主な推定方法	データ分布	回帰係数
線形回帰モデル (単回帰・重回帰など)	最小二乗法	正規分布	一変数に一つ
一般化線形モデル	最尤推定法	正規分布以外 の分布も可能	一変数に一つ
一般化線形混合モデル			変数の個体差に 応じて推定可能
階層ベイズモデル	ベイズ推定		

本授業で扱う統計モデル

- 線形回帰モデル
 - 単回帰モデル、重回帰モデル
- 一般化線形回帰モデル
 - 離散：ポアソン回帰モデル、二項反応モデル（ロジスティック回帰モデル、プロビット回帰モデル、補対数対数モデル）、負の二項分布モデル、ゼロ過剰ポアソン回帰モデル、ゼロ過剰負の二項分布モデル
 - 連続：ガンマ回帰モデル、ベータ回帰モデル、指数-ガウス回帰モデル
 - スパース：Lasso回帰モデル、Ridge回帰モデル
- 一般化線形混合モデル
 - マルチレベルモデル
- 状態空間モデル

マルチレベルモデル

- グループ j に属する地域 i について、以下のような線形回帰モデルを考える

$$y_{ij} = \beta_0 + \beta_1 x_{ij1} + \dots + \beta_k x_{ijk} + \varepsilon_{ij}$$
$$\varepsilon_{ij} \sim N(0, \sigma_y^2)$$

- ここで、誤差項は未知である
- 定数項と回帰係数について、グループ毎の違いを認めるかどうかによって、柔軟なモデル推定ができる

マルチレベルモデル

- ① 定数項と回帰係数が変動しないモデル

$$y_{ij} = \beta_0 + \beta_1 x_{ij1} + \dots + \beta_k x_{ijk} + \varepsilon_{ij}$$

- ② 定数項のみが変動するモデル

$$y_{ij} = \beta_{j0} + \beta_1 x_{ij1} + \dots + \beta_k x_{ijk} + \varepsilon_{ij}$$

- ③ 回帰係数のみ変動するモデル

$$y_{ij} = \beta_0 + \beta_{j1} x_{ij1} + \dots + \beta_{jk} x_{ijk} + \varepsilon_{ij}$$

- ④ 定数項と回帰係数の両方が変動するモデル

$$y_{ij} = \beta_{j0} + \beta_{j1} x_{ij1} + \dots + \beta_{jk} x_{ijk} + \varepsilon_{ij}$$

ランダム効果と固定効果

- ランダム効果
 - 回帰係数や定数項について、個人や地域・グループでの変動を認めるモデル
- 固定効果
 - 上記のような変動を認めないモデル
- 混合効果
 - ランダム効果と固定効果が混在したモデル

マルチレベルモデル（混合効果モデル）

- レベル 1 = lower level

$$y_{ij} = \beta_0 + \beta_1 x_{ij1} + \varepsilon_{ij}$$
$$\varepsilon_{ij} \sim N(0, \sigma_y^2)$$
$$y_{ij} \sim N(\beta_0 + \beta_1 x_{ij1}, \sigma_y^2)$$

- レベル 2 = upper level

$$\beta_{j0} = \gamma_{00} + \gamma_{10} u_{j0} + \eta_{j0}$$
$$\beta_{j1} = \gamma_{01} + \gamma_{11} u_{j1} + \eta_{j1}$$
$$\eta_{jk} \sim N(0, \sigma_\eta^2)$$
$$\beta_{jk} \sim N(0, \sigma_{\beta_k}^2)$$
$$k = \{0, 1\}$$

マルチレベルモデル（混合効果モデル）

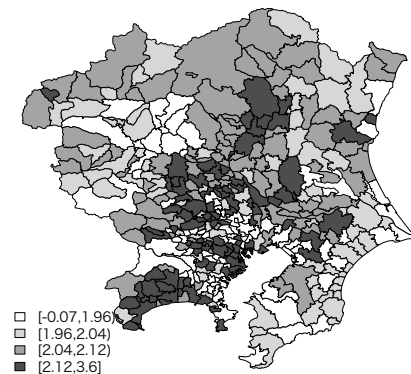
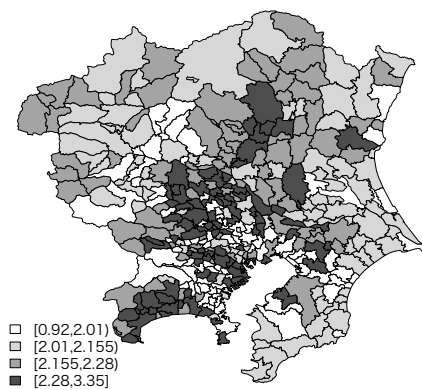
- レベル1とレベル2を合わせると、以下のように式変形される

$$\begin{aligned} y_{ij} &= \beta_0 + \beta_1 x_{ij1} + \varepsilon_{ij} \\ &= (\gamma_{00} + \gamma_{10} u_{j0}) + (\gamma_{01} + \gamma_{11} u_{j0}) x_{ij1} + \eta_j + \varepsilon_{ij} \\ &= \underbrace{(\gamma_{00} + \gamma_{01} x_{ij1})}_{\text{固定効果}} + \underbrace{(\gamma_{10} u_{j0} + \gamma_{11} u_{j0} x_{ij1})}_{\text{ランダム効果}} + \eta_j + \varepsilon_{ij} \end{aligned}$$

- この式は、次式のように整理できる

$$\begin{aligned} \mathbf{y}_i &= \mathbf{X}_i \mathbf{B} + \mathbf{Z}_i \mathbf{b}_i + \varepsilon_{ij} \\ \mathbf{b}_i &\sim N(0, \sigma_{b_i}^2) \\ \varepsilon_{ij} &\sim N(0, \sigma_y^2) \end{aligned}$$

マルチレベルモデルの推定結果の例



マルチレベルモデルの最尤推定

- 混合効果モデル

$$\begin{aligned} \mathbf{y}_i &= \mathbf{X}_i \mathbf{B} + \mathbf{Z}_i \mathbf{b}_i + \varepsilon_{ij} \\ \mathbf{b}_i &\sim N(0, \sigma_{b_i}^2) \\ \varepsilon_{ij} &\sim N(0, \sigma_y^2) \end{aligned}$$

- \mathbf{y}_i の共分散行列 \mathbf{V}_{y_i} は以下のようなになる

$$\mathbf{V}_{y_i} = \mathbf{Z}_i (\sigma_{b_i}^2 \mathbf{I}) \mathbf{Z}_i' + \sigma_y^2 \mathbf{I}$$

- $\mathbf{y}_i \sim N(\mathbf{X}_i \mathbf{B} + \mathbf{Z}_i \mathbf{b}_i, \mathbf{V}_{y_i})$ および $\mathbf{b}_i \sim N(0, \sigma_{b_i}^2)$ であることから、ランダム効果 $\mathbf{Z}_i \mathbf{b}_i$ に対する期待値 $E(\mathbf{Z}_i \mathbf{b}_i) = 0$ である
- したがって、 \mathbf{y}_i の周辺分布は $\mathbf{y}_i \sim N(\mathbf{X}_i \mathbf{B}, \mathbf{V}_{y_i})$ となる

マルチレベルモデルの最尤推定

- このとき、混合効果モデルの尤度関数 ℓ は次式で表せる

$$\ell = \frac{\exp\left[-\frac{1}{2}(\mathbf{y} - \mathbf{XB})' \mathbf{V}_y^{-1} (\mathbf{y} - \mathbf{XB})\right]}{(2\pi)^{(n/2)} \sqrt{|\mathbf{V}_y|}}$$

- 対数尤度関数 $\log(\ell)$ は次式のようなになる

$$\log(\ell) = -\frac{\log(2\pi)}{2} - \frac{1}{2} [\log|\mathbf{V}_y| + (\mathbf{y} - \mathbf{XB})' \mathbf{V}_y^{-1} (\mathbf{y} - \mathbf{XB})]$$

マルチレベルモデルの最尤推定

- 対数尤度関数 $\log(\ell)$ が最大となるとき、
$$\mathbf{X}'\mathbf{V}_y^{-1}(\mathbf{y} - \mathbf{XB}) = 0$$

- 固定効果 \mathbf{B} の最尤推定量 $\hat{\mathbf{B}}$

$$\hat{\mathbf{B}} = (\mathbf{X}'\mathbf{V}_y^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}_y^{-1}\mathbf{y}$$

- $\hat{\mathbf{B}}$ の分散 $\mathbf{V}_{\hat{\mathbf{B}}}$

$$\mathbf{V}_{\hat{\mathbf{B}}} = (\mathbf{X}'\mathbf{V}_y^{-1}\mathbf{X})^{-1}$$

マルチレベルモデルの最尤推定

- ランダム効果対数 \mathbf{b} の尤度関数 $\log(\ell_{\mathbf{b}})$ は次式のようになる

$$\begin{aligned} \ell_{\mathbf{b}} &\propto |(\sigma_y^2 I)|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} (\mathbf{y} - \mathbf{XB} - \mathbf{Zb})' (\sigma_y^2 I)^{-1} (\mathbf{y} - \mathbf{XB} - \mathbf{Zb}) \right] \\ &\quad \times |(\sigma_b^2 I)|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} \mathbf{b}' (\sigma_b^2 I)^{-1} \mathbf{b} \right] \end{aligned}$$

- このとき対数尤度関数 $\log(\ell_{\mathbf{b}})$ は、

$$\begin{aligned} &\log(\ell_{\mathbf{b}}) \\ &= -\frac{\log(2\pi)}{2} \\ &\quad -\frac{1}{2} \left[\log|\sigma_y^2 I| + (\mathbf{y} - \mathbf{XB} - \mathbf{Zb})' (\sigma_y^2 I)^{-1} (\mathbf{y} - \mathbf{XB} - \mathbf{Zb}) + \log|\sigma_b^2 I| \right. \\ &\quad \left. + \mathbf{b}' (\sigma_b^2 I)^{-1} \mathbf{b} \right] \end{aligned}$$

マルチレベルモデルの最尤推定

- 対数尤度関数 $\log(\ell_{\mathbf{b}})$ が最大になるとき、 $\delta \log(\ell_{\mathbf{b}})/\delta \mathbf{b} = 0$

$$\begin{aligned}\delta \log(\ell_{\mathbf{b}})/\delta \mathbf{b} &= \mathbf{Z}'(\sigma_y^2 \mathbf{I})^{-1}(\mathbf{y} - \mathbf{XB} - \mathbf{Zb}) - (\sigma_b^2 \mathbf{I})^{-1} \mathbf{b} \\ &= \mathbf{Z}'(\sigma_y^2 \mathbf{I})^{-1}(\mathbf{y} - \mathbf{XB}) - \left(\mathbf{Z}'(\sigma_y^2 \mathbf{I})^{-1} \mathbf{Z} + (\sigma_b^2 \mathbf{I})^{-1}\right) \mathbf{b} = 0\end{aligned}$$

- ランダム効果 \mathbf{b} の最尤推定量 $\hat{\mathbf{b}}$ は、

$$\hat{\mathbf{b}} = \left(\mathbf{Z}'(\sigma_y^2 \mathbf{I})^{-1} \mathbf{Z} + (\sigma_b^2 \mathbf{I})^{-1}\right)^{-1} \mathbf{Z}'(\sigma_y^2 \mathbf{I})^{-1}(\mathbf{y} - \mathbf{XB})$$

マルチレベル・モデルのベイズ推定

$$y_i \sim N(X_i B, \sigma_y^2)$$

$$\begin{pmatrix} \beta_{j0} \\ \beta_{j1} \end{pmatrix} \sim N \left(\begin{pmatrix} \gamma_{00} + \gamma_{10} u_{j0} \\ \gamma_{01} + \gamma_{11} u_{j0} \end{pmatrix}, \begin{pmatrix} \sigma_{\beta_0}^2 & \rho \sigma_{\beta_0} \sigma_{\beta_1} \\ \rho \sigma_{\beta_0} \sigma_{\beta_1} & \sigma_{\beta_1}^2 \end{pmatrix} \right)$$

$$\begin{aligned}B &\sim N(0, V_B) \\ \varepsilon_{ij} &\sim N(0, \sigma_y^2) \\ \eta_{jk} &\sim N(0, \sigma_\eta^2)\end{aligned}$$

マルチレベル・モデルの階層ベイズ推定

- 未知パラメータは B を事前情報とする階層ベイズ法を用いて推定できる
- このとき σ_y^2 と V_B の事前情報を与える必要がある
- σ_y^2 の事前情報は、スケールパラメータ $s_{y_i}^2 (= \sigma_{y_i}^2)$ を用いて、 χ^2 分布の逆分布として与えることができる

$$\sigma_y^2 | s_{y_i}^2 \sim \frac{\eta_j s_{y_i}^2}{\chi_{\eta_j}^2}$$

- V_B の事前分布は、逆ウィシャート分布（後述）により与えることができる

$$V_B | \nu, V \sim IW(\nu, V)$$

- ここで、 ν, V は V_B の自由度と共分散のパラメータである

マルチレベル・モデルの階層ベイズ推定

- したがって、マルチレベルモデルを階層ベイズ推定する場合は、以下のような条件付き分布を次々に生成させれば良い

$$\sigma_y^2 | s_{y_i}^2 \sim \frac{\eta_j s_{y_i}^2}{\chi_{\eta_j}^2}$$

$$V_B | \nu, V \sim IW(\nu, V)$$

$$y_i | X_i, B, \sigma_y^2 \sim N(X_i B, \sigma_y^2)$$

$$B | Z_i, \sigma_{b_i}^2, V_B \sim N(0, V_B)$$

$$b_i | \sigma_{b_i}^2 \sim N(0, \sigma_{b_i}^2)$$

$$\sigma_{b_i}^2 | \eta_j \sim N(0, \sigma_{\eta}^2)$$

ウィシャート分布

- X_i が n 個の p 次元多変量正規分布に従うとき

$$X_i \sim N(0, \Sigma)$$

- 標本分散共分散行列 W は逆ウィシャート分布に従う

$$X = \sum_{i=1}^n X_i \cdot X_i^T$$

$$X \sim W(X; \Sigma, n, p)$$

- ベイズ統計では多変量正規分布の共分散行列の自然共役事前分布として用いられる

ウィシャート分布

- ウィシャート行列はカイ二乗分布を多次元化したものである

$$W(X; \Sigma, n, p) = \frac{|\Sigma|^{-\frac{n}{2}} |X|^{-\frac{n-p-1}{2}}}{2^{\frac{np}{2}} \cdot \Gamma_p\left(\frac{n}{2}\right)} \cdot \exp\left(-\frac{1}{2} \text{tr} \Sigma^{-1} X\right)$$

- $\Gamma_p(a)$ は p 次元多変量 Γ 関数

$$\Gamma_p(a) = \pi^{\frac{p(p-1)}{4}} \prod_{i=1}^p \left[a - \frac{1}{2}(j-1) \right]$$

逆ウィシャート分布

- $W(X; \Sigma, n, p)$ に従う X を $V = X^{-1}$ とすると、 V の分布を逆ウィシャート分布という

$$W(V; \Sigma, n, p) = \frac{|\Sigma|^{-\frac{n}{2}} |V|^{-\frac{n-p-1}{2}}}{2^{\frac{np}{2}} \cdot \Gamma_p\left(\frac{n}{2}\right)} \cdot \exp\left(-\frac{1}{2} \text{tr} \Sigma^{-1} V^{-1}\right)$$