

Проверка гипотез в дисперсионном анализе повторных неполных наблюдений

Н.П. Алексеева, Е.В.Вербицкая

СПбГУ, математико-механический факультет
СПбГМУ, лаб.биомедицинской статистики

2018 г.

Пусть в i -й группе, $i = 1, 2, \dots, I$ имеется ν_i индивидов, $\nu_1 + \dots + \nu_I = n$, которые наблюдаются в моменты времени $t = 1, 2, \dots, T$. Выборку наблюдений обозначим через x_{ijt} .

Модель

$$x_{ijt} = \mu + \alpha_i + e_{ij}^1 + \beta_t + \gamma_{it} + e_{ijt}$$

- μ параметр генерального среднего,
- α_i дифф.эффекты групп,
- β_t дифф. эффекты фактора времени,
- γ_{it} дифф. эффекты взаимодействия факторов времени и разделения на группы,
- $e_{ij}^1 \sim \mathcal{N}(0, \sigma_1)$, $e_{ijt} \sim \mathcal{N}(0, \sigma)$ независимые случайные ошибки.

Для каждого индивида имеется множество наблюдений в точках $N_{ij} \subset (1, 2, \dots, T)$, $\text{card}(N_{ij}) = n_{ij}$.

В i -й группе моменты в времени t имеются данные об индивидах $M_{it} \subset (1, 2, \dots, \nu_i)$, $\text{card}(M_{it}) = m_{it}$,

$$\sum_{i=1}^I m_{it} = m_{i.}, \quad \sum_{t=1}^T m_{it} = m_{.t}, \quad \sum_{t=1}^T \sum_{i=1}^I m_{it} = m_{..}$$

Замечание

Необходимо предположение о наличии полных данных хотя бы в одной точке, и $\forall t, i \ m_{it} > 0$.

В случае полных данных устраняется индивидуальное влияние, измеряемое индивидуальным средним $x_{ij\cdot}$ и модель разделяется на две модели:

$$x_{ij\cdot} = \mu + \alpha_i + e_{ij\cdot}^1, \quad x_{ijt} - x_{ij\cdot} = \beta_t + \gamma_{it} + e_{ijt}$$

Лемма

В случае неполных данных имеет место смещение

$$E x_{ij\cdot} = \mu + \alpha_i + W_{ij}, \quad \text{где} \quad W_{ij} = \frac{1}{n_{ij}} \sum_{t \in N_{ij}} (\beta_t + \gamma_{it})$$

Определим групповую поправку как

$$G = \sum_{i=0}^{\infty} (MN)^i (ML - MNK), \text{ где вектора } L = \{x_{..t} - x_{...}\}_{t=1}^T, \\ K = \{x_{i..} - x_{...}\}_{i=1}^I, \text{ матрицы } M = \{m_{i1}, \dots, m_{iT}\}_{i=1}^I, \\ N = \{m_{1t}, \dots, m_{it}\}_{t=1}^T.$$

Лемма

Пусть $P_0 = MN$, $Q_0 = (I - P_0 + P_0^\infty)^{-1}$, $A = P_0 Q_0$, $B = Q_0 M$.
Тогда компонента групповой поправки имеет вид

$$G_i = \frac{1}{m_i} \sum_{t=1}^T m_{it} \beta_t + \varepsilon_i, \quad \varepsilon_i = \sum_{k=1}^I \sum_{t=1}^T m_{kt} \left(\frac{b_{it}}{m_t} - \frac{a_{ik}}{m_k} \right) e_{i \cdot t}.$$

Замечание

В случае однородности пропусков групповая поправка равна нулю.

Пусть J^i матрица инцидентности i -й группы, Λ_{iT} и Λ_{ν_i} диагональные матрицы с векторами (m_{it}) и (ν_{ij}) на главной диагонали.

Вектор индивидуальной поправки $H = \{H_{ij}\}_{i=1, j=1}^{I, \nu_i}$, имеет вид $H = (I - C)V + DU$, где $V_i = \{x_{ij}\}_{j=1}^{\nu_i}$, $U_i = \{x_{i \cdot t}\}_{t=1}^T$,

$R_i = \Lambda_{\nu_i} J^i$, $P_i = R_i \Lambda_{iT} (J^i)^T$, $P_i^\infty = \lim_{k \rightarrow \infty} P_i^k$,

$Q_i = (I - P_i + P_i^\infty)^{-1}$, $C_i = P_i^\infty + Q_i$, $D_i = Q_i R_i$.

Предложение

$$H_{ij} = W_{ij} - G_i + \mathcal{E}_{ij},$$

$$\mathcal{E}_{ij} = e_{ij} - \sum_{l=1}^{\nu_i} c_{jl}^i e_{il} + \sum_{\tau=1}^T d_{j\tau}^i e_{i \cdot \tau}$$

Теорема

Пусть $\Delta_{ij} = e_{ij} - \mathcal{E}_{ij} - \varepsilon_i$. Имеют место несмещенные модели

$$y_{ijt} = x_{ij} - x_{ijt} + (H_{ij} + G_i) = \beta_t + \gamma_{it} + e_{ijt} - \Delta_{ij},$$

$$z_{ij} = x_{ij} - (H_{ij} + G_i) = \alpha_i + e_{ij}^1 + \Delta_{ij},$$

$$\begin{aligned} \tilde{D}_{ij,kl} &= \frac{1}{\sigma^2} \mathbb{E} \Delta_{ij} \Delta_{kl} = F_0(i, k) + I_{[i=k]} \left(\sum_{\ell=1}^{\nu_i} \frac{c_{j\ell}^i c_{l\ell}^k}{n_{i\ell}} + \sum_{t=1}^T \frac{d_{jt}^i d_{lt}^k}{m_{it}} \right) - \\ &\quad - I_{[i=k]} \sum_{\ell=1}^{\nu_i} \sum_{\tau \in N_{i\ell}} \left(\frac{c_{j\ell}^i d_{l\tau}^k}{n_{i\ell} m_{k\tau}} + \frac{d_{j\tau}^i c_{l\ell}^k}{n_{k\ell} m_{i\tau}} \right), \\ F_0(i, k) &= \sum_{t=1}^T \frac{b_{it} b_{kt}}{m_{\cdot t}} - \sum_{\ell=1}^I \frac{a_{i\ell} a_{k\ell}}{m_{\cdot \ell}}. \end{aligned}$$

$$\mathbb{E}(e_{ij}^1 + \Delta_{ij})(e_{kl}^1 + \Delta_{kl}) = \sigma_1^2 I_{[i=k, j=l]} + \sigma^2 \tilde{D}_{ij,kl}.$$

Определим ковариационную матрицу ошибок

$$\Sigma = \mathbb{E}(e_{ijt} - \Delta_{ij})(e_{kl\tau} - \Delta_{kl}) = \sigma^2 \Lambda.$$

$$\Sigma = \sigma^2 (I_{[i=k, j=l, t=\tau]} - M_{ij,kl\tau} - M_{kl,ijt} + \tilde{D}_{ij,kl}),$$

$$M_{ij,kl\tau} = \frac{1}{\sigma^2} \mathbb{E} \Delta_{ij} e_{kl\tau} = I_{[i=k]} \left(\frac{c_{jl}^i}{n_{il}} - \frac{d_{j\tau}^i}{m_{i\tau}} \right) + \left(\frac{a_{ik}}{m_k} - \frac{b_{i\tau}}{m_{i\tau}} \right).$$

Матрица плана Н модели $\mathbb{E}y_{ijt} = \beta_t + \gamma_{it}$

$$q_{it} = -\frac{m_{it}}{m_{iT}}, \quad q_{\cdot t} = -\frac{m_{\cdot t}}{m_{\cdot T}}, \quad p_{it} = -\frac{m_{it}}{m_{It}}, \quad p_{i\cdot} = -\frac{m_{i\cdot}}{m_{I\cdot}}, \quad r_{it} = \frac{m_{it}}{m_{IT}}$$

	β_1	γ_{11}	γ_{21}	...	$\gamma_{\ell,1}$...	β_τ	$\gamma_{1,\tau}$	$\gamma_{2,\tau}$...	$\gamma_{\ell,\tau}$
Y_{11}	1	1	0	...	0	...	0	0	0	...	0
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
$Y_{1,\tau}$	0	0	0	...	0	...	1	1	0	...	0
$Y_{1,T}$	$q_{\cdot 1}$	q_{11}	0	...	0	...	$q_{\cdot \tau}$	$q_{1,\tau}$	0	...	0
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
$Y_{\ell,1}$	1	0	0	...	1	...	0	0	0	...	0
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
$Y_{\ell,\tau}$	0	0	0	...	0	...	1	0	0	...	1
$Y_{\ell,T}$	$q_{\cdot 1}$	0	0	...	$q_{\ell,1}$...	$q_{\cdot \tau}$	0	0	...	$q_{\ell,\tau}$
$Y_{I,1}$	1	p_{11}	p_{21}	...	$p_{\ell,1}$...	0	0	0	...	0
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
$Y_{I,\tau}$	0	0	0	...	0	...	1	$p_{1,\tau}$	$p_{2,\tau}$...	$p_{\ell,\tau}$
$Y_{I,T}$	$q_{\cdot 1}$	r_{11}	r_{21}	...	$r_{\ell,1}$...	$q_{\cdot \tau} v$	$r_{1,\tau}$	$r_{2,\tau}$...	$r_{\ell,\tau}$

Проверка значимости эффектов взаимодействия

$$\Theta = (\beta_1, \gamma_{11}, \dots, \gamma_{i,1}, \dots, \beta_{T-1}, \gamma_{1,T-1}, \dots, \gamma_{i,T-1})$$

Оценка $\hat{\Theta} = (H^T \Lambda^{-1} H)^{-1} H^T \Lambda^{-1} Y$ при минимизации

$$Q_{2e}(\Theta) = (Y - H\Theta^T)^T \Lambda^{-1} (Y - H\Theta^T).$$

Число степеней свободы $\mu_{2e} = m.. - n - I(T - 1)$.

Вектор усеченной модели параметров $\Theta_* = (\beta_1, \dots, \beta_{T-1})^T$

Оценка $\hat{\Theta}_* = (H_*^T \Lambda^{-1} H_*)^{-1} H_*^T \Lambda^{-1} Y$ при минимизации

$$Q_{2e*}(\Theta_*) = (Y - H_*\Theta_*^T)^T \Lambda^{-1} (Y - H_*\Theta_*^T),$$

число степеней свободы равно

$$\mu_{2e*} - \mu_{2e} = m.. - n - T + 1 - (m.. - n - I(T - 1)) = IT - T - I + 1.$$

$$F = \frac{(Q_{2e*}(\hat{\Theta}_*) - Q_{2e}(\hat{\Theta}))(m.. - n - I(T - 1))}{Q_{2e}(\hat{\Theta})(IT - T - I + 1)} \sim \mathcal{F}(\mu_{2e*} - \mu_{2e}, \mu_{2e})$$

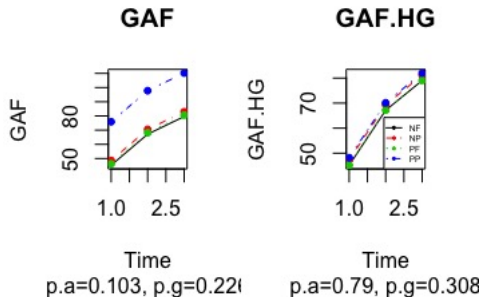


Figure: Динамика индекса глобального функционирования (ИГФ) при разной терапии при удалении наблюдений с пропусками ($n = 63$) и с учетом всей информации ($n = 332$) .

Различия по ИГФ не было, но у больных с изначально высоким ИГФ было больше шансов остаться в программе даже при отсутствии терапии.

Алгоритм построения линейного предиктора по неполным данным $F(\mathcal{Y}, \mathcal{X})$

$\Omega = (\omega_1, \dots, \omega_n)$ наблюдения по неполным данным:
 зависимой переменной $\mathcal{Y}(\Omega) = (y(\omega_1), \dots, y(\omega_n))$,
 независимых переменных $\mathcal{X} = (X_1(\Omega), \dots, X_p(\Omega))$.

Модель линейной регрессии $Y = X\beta + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma^2 \mathbb{I})$

- $Y = \mathcal{Y}(\Omega_1)$, $\Omega_1 \subseteq \Omega = (\omega_1, \dots, \omega_n)$;
- $X = (X_{t_1}(\Omega_1), \dots, X_{t_k}(\Omega_1))$, $(t_1, \dots, t_k) \subseteq (1, \dots, p)$
- $\hat{\beta}$, $\hat{\sigma}^2$ оценки МНК

Предиктор $\hat{Y}(\Omega) = X_* \hat{\beta}$, где $X_* = (X_{t_1}(\Omega), \dots, X_{t_k}(\Omega))$

Алгоритм замены пропущенных значений наблюдаемыми, наиболее близкими к предсказанным

Актуальность темы

для статистического анализа многомерных данных с повторностями нужна редукция размерности
для этого используется факторный анализ