

7 Následky nesplnění předpokladů

V lineárním modelu jsme předpokládali, že známe prostor možných středních hodnot, že všechna pozorování mají stejný rozptyl, že jsou nekorelovaná (resp. nezávislá) a že mají normální rozdělení. Nyní se pokusíme popsat následky, které má nesplnění některého z uvedených předpokladů.

7.1 Prostor středních hodnot

Předpokládejme, že platí

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} + \mathbf{e}, \quad \mathbf{e} \sim (\mathbf{0}, \sigma^2\mathbf{I}), \quad (7.1)$$

přestože my předpokládáme platnost modelu $\mathbf{Y} \sim (\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$.

7.1.1 Vlastnosti odhadu $E\mathbf{Y}$

Označme $\mathbf{G} = (\mathbf{X}, \mathbf{Z})$ a $\boldsymbol{\delta} = (\boldsymbol{\beta}', \boldsymbol{\gamma}')'$ a veškeré statistiky vztažené k modelu $\mathbf{Y} \sim (\mathbf{G}\boldsymbol{\delta}, \sigma^2\mathbf{I})$ označíme dolním indexem g . Běžný odhad vektoru $E\mathbf{Y}$ je tedy

$$\hat{\mathbf{Y}}_g = \mathbf{G}(\mathbf{G}'\mathbf{G})^{-1}\mathbf{G}'\mathbf{Y}, \quad (7.2)$$

což je, jak víme např. z (3.13), průmět \mathbf{Y} do $\mathcal{M}(\mathbf{X}, \mathbf{Z}) = \mathcal{M}(\mathbf{X}, \mathbf{MZ})$. S použitím druhého vyjádření dostaneme

$$\begin{aligned} \hat{\mathbf{Y}}_g &= (\mathbf{X}, \mathbf{MZ}) \begin{pmatrix} \mathbf{X}'\mathbf{X} & \mathbf{0} \\ \mathbf{0} & \mathbf{Z}'\mathbf{MZ} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{X}' \\ \mathbf{Z}'\mathbf{M} \end{pmatrix} \mathbf{Y} \\ &= \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} + \mathbf{MZ}(\mathbf{Z}'\mathbf{MZ})^{-1}\mathbf{Z}'\mathbf{M}\mathbf{Y} \\ &= \hat{\mathbf{Y}} + \mathbf{MZ}(\mathbf{Z}'\mathbf{MZ})^{-1}\mathbf{Z}'\mathbf{u} & (7.3) \\ &= \mathbf{X}\mathbf{b}_g + \mathbf{Z}\mathbf{c}_g, & (7.4) \end{aligned}$$

kde \mathbf{b}_g a \mathbf{c}_g jsou obecně nějaká řešení příslušné normální rovnice.

Když přepíšeme (7.3) tak, aby bylo patrné jakou lineární kombinací sloupců matic \mathbf{X} , \mathbf{Z} je vektor $\hat{\mathbf{Y}}_g$ (co mohou být vektory \mathbf{b}_g , \mathbf{c}_g), dostaneme po úpravě (vyjádříme \mathbf{M} pomocí \mathbf{X})

$$\hat{\mathbf{Y}}_g = \mathbf{X}(\mathbf{b} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Z}\mathbf{c}_g) + \mathbf{Z}\mathbf{c}_g, \quad (7.5)$$

když jsme označili

$$\mathbf{c}_g = (\mathbf{Z}'\mathbf{M}\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{u}. \quad (7.6)$$

Můžeme tedy psát

$$\mathbf{b}_g = \mathbf{b} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Z}\mathbf{c}_g, \quad (7.7)$$

odkud je zřetelný zejména vztah mezi \mathbf{b} a \mathbf{b}_g .

Z (7.3) plyne, že rozdíl reziduálních součtů čtverců mezi uvažovaným modelem $\mathbf{Y} \sim (\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$ a skutečně platným modelem $\mathbf{Y} \sim (\mathbf{G}\boldsymbol{\delta}, \sigma^2\mathbf{I})$ je

$$\begin{aligned} RSS - RSS_g &= \|\mathbf{M}\mathbf{Z}(\mathbf{Z}'\mathbf{M}\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{u}\|^2 \\ &= \|\mathbf{M}\mathbf{Z}\mathbf{c}_g\|^2. \end{aligned} \quad (7.8)$$

Porovnejme ještě střední hodnoty obou reziduálních součtů čtverců. Protože platí model (7.1), je zřejmě $E\,RSS_g = (n - h(\mathbf{X}, \mathbf{Z}))\sigma^2$. Jinak to dopadne u reziduálního součtu čtverců RSS z (nesprávně) předpokládaného modelu. Postupnými úpravami dostaneme

$$E\,RSS = E\|\mathbf{M}\mathbf{Y}\|^2 = E\|\mathbf{M}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} + \mathbf{e})\|^2 = E\|\mathbf{M}\mathbf{Z}\boldsymbol{\gamma} + \mathbf{M}\mathbf{e}\|^2,$$

tedy (s ohledem na $E\mathbf{e} = \mathbf{0}$)

$$\begin{aligned} E\,RSS &= \|\mathbf{M}\mathbf{Z}\boldsymbol{\gamma}\|^2 + E\|\mathbf{M}\mathbf{e}\|^2 \\ &= \|\mathbf{M}\mathbf{Z}\boldsymbol{\gamma}\|^2 + (n - h(\mathbf{X}))\sigma^2. \end{aligned} \quad (7.9)$$

Vraťme se k odhadu $\hat{\mathbf{Y}}$. Jeho střední hodnota je rovna

$$E\hat{\mathbf{Y}} = \mathbf{H}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma}) = \mathbf{X}\boldsymbol{\beta} + \mathbf{H}\mathbf{Z}\boldsymbol{\gamma}.$$

Obecně tedy není nestranným odhadem pro $E\mathbf{Y}$, má vychýlení

$$\text{bias } \hat{\mathbf{Y}} = E\hat{\mathbf{Y}} - E\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{H}\mathbf{Z}\boldsymbol{\gamma} - (\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma}) = -\mathbf{M}\mathbf{Z}\boldsymbol{\gamma}. \quad (7.10)$$

Shrňme vlastnosti odhadů klasického modelu.

Věta 7.1. (Vychýlení odhadů, platí-li širší model) Nechť platí $\mathbf{Y} \sim (\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma}, \sigma^2\mathbf{I})$. Pro statistiky odvozené z modelu $\mathbf{Y} \sim (\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$ platí

$$\text{bias } \hat{\mathbf{Y}} = -\mathbf{M}\mathbf{Z}\boldsymbol{\gamma}, \quad (7.11)$$

$$\text{bias } S^2 = \frac{\|\mathbf{M}\mathbf{Z}\boldsymbol{\gamma}\|^2}{n - h(\mathbf{X})}, \quad (7.12)$$

Porovnejme nyní varianční matice odhadů $\hat{\mathbf{Y}}$ a $\hat{\mathbf{Y}}_g$. Snadno dostaneme

$$\begin{aligned}\text{var } \hat{\mathbf{Y}}_g &= \sigma^2(\mathbf{X}, \mathbf{MZ}) \begin{pmatrix} \mathbf{X}'\mathbf{X} & \mathbf{O} \\ \mathbf{O} & \mathbf{Z}'\mathbf{MZ} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{X}' \\ \mathbf{Z}'\mathbf{M} \end{pmatrix} \\ &= \sigma^2(\mathbf{H} + \mathbf{MZ}(\mathbf{Z}'\mathbf{MZ})^{-1}\mathbf{Z}'\mathbf{M}).\end{aligned}\quad (7.13)$$

Je-li matice \mathbf{MZ} nenulová, bude matice $\text{var } \hat{\mathbf{Y}}_g - \text{var } \hat{\mathbf{Y}}$ pozitivně definitní, takže vychýlený odhad je co do rozptylu lepší. Vychýlené odhady však neporovnáváme pomocí jejich rozptylu či varianční matice, ale pomocí střední čtvercové chyby. *Střední čtvercová chyba* odhadu \mathbf{T} parametru $\boldsymbol{\theta}$ je definována jako

$$\begin{aligned}\text{MSE}(\mathbf{T}) &= \mathbf{E}(\mathbf{T} - \boldsymbol{\theta})(\mathbf{T} - \boldsymbol{\theta})' \\ &= \text{var}(\mathbf{T}) + \text{bias}(\mathbf{T})\text{bias}(\mathbf{T})'.\end{aligned}$$

Střední čtvercovou chybu $\hat{\mathbf{Y}}$ jako odhadu pro $\mathbf{E}\mathbf{Y}$ lze tedy psát

$$\text{MSE } \hat{\mathbf{Y}} = \text{var } \hat{\mathbf{Y}} + (\text{bias } \hat{\mathbf{Y}})(\text{bias } \hat{\mathbf{Y}})' = \sigma^2\mathbf{H} + \mathbf{MZ}\boldsymbol{\gamma}\boldsymbol{\gamma}'\mathbf{Z}'\mathbf{M}.\quad (7.14)$$

Protože $\hat{\mathbf{Y}}_g$ je nestranným odhadem $\mathbf{E}\mathbf{Y}$, platí $\text{MSE } \hat{\mathbf{Y}}_g = \text{var } \hat{\mathbf{Y}}_g$.

Porovnejme střední čtvercové chyby $\hat{\mathbf{Y}}_g$ a $\hat{\mathbf{Y}}$ jako odhadů vektoru $\mathbf{E}\mathbf{Y}$:

$$\text{MSE } \hat{\mathbf{Y}}_g - \text{MSE } \hat{\mathbf{Y}} = \sigma^2(\mathbf{MZ}(\mathbf{Z}'\mathbf{MZ})^{-1}\mathbf{Z}'\mathbf{M} - \mathbf{MZ}\boldsymbol{\gamma}\boldsymbol{\gamma}'\mathbf{Z}'\mathbf{M}/\sigma^2).$$

Nyní stačí použít tvrzení věty A.9 pro $\mathbf{A} = \mathbf{MZ}$ a $\mathbf{c} = \boldsymbol{\gamma}/\sigma$, abychom zjistili, že rozdíl středních čtvercových chyb dá pozitivně semidefinitní matici, právě když je $\|\mathbf{A}\mathbf{c}\|^2 = \|\mathbf{MZ}\boldsymbol{\gamma}/\sigma\|^2 \leq 1$. Došli jsme tak k tvrzení následující věty.

Věta 7.2. (Když je vychýlení malé) Nechť platí $\mathbf{Y} \sim (\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma}, \sigma^2\mathbf{I})$. Pro $\hat{\mathbf{Y}}_g$ z tohoto modelu a pro $\hat{\mathbf{Y}}$ z modelu $\mathbf{Y} \sim (\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$ platí ekvivalence

$$\text{MSE } \hat{\mathbf{Y}}_g \geq \text{MSE } \hat{\mathbf{Y}} \iff \|\text{bias } \hat{\mathbf{Y}}\|^2 \leq \sigma^2.\quad (7.15)$$

Při předpovědi budoucího pozorování tedy je výhodnější použít menší model, když je vychýlení způsobené touto volbou dostatečně malé.

Věta 7.3. (Důsledek) Nechť platí $\mathbf{Y} \sim (\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma}, \sigma^2\mathbf{I})$, nechť $\boldsymbol{\theta} = \mathbf{p}'\boldsymbol{\beta} + \mathbf{s}'\boldsymbol{\gamma}$ je odhadnutelný parametr v tomto modelu. Nechť \mathbf{b} je libovolné řešení normální rovnice $\mathbf{X}'\mathbf{X}\mathbf{b} = \mathbf{X}'\mathbf{Y}$. Potom je parametr $\tau = \mathbf{p}'\boldsymbol{\beta}$ odhadnutelný také v modelu $\mathbf{Y} \sim (\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$ a platí

$$\text{MSE } \hat{\boldsymbol{\theta}} \geq \text{MSE } \hat{\tau} \iff \|\mathbf{MZ}\boldsymbol{\gamma}\|^2 \leq \sigma^2.$$

Důkaz: Především je třeba dokázat, že τ je odhadnutelný parametr. Odhadnutelnost $\boldsymbol{\theta}$ je podle věty 2.4 ekvivalentní s existencí vektoru $\mathbf{q} \in \mathbb{R}^n$, pro který platí $\mathbf{q}'(\mathbf{X}, \mathbf{Z}) = (\mathbf{p}', \mathbf{s}')$. Speciálně to tedy znamená existenci \mathbf{q} , pro který platí $\mathbf{q}'\mathbf{X} = \mathbf{p}'$,