

VĚTA 3.5

Nechť $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ je LM (**), kde $h(\mathbf{X}) = m + 1$ a $\mathbf{e} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$. Potom

- 1) $\hat{\boldsymbol{\beta}}$ a s_n^2 jsou nezávislé náhodné veličiny,
- 2) $(n - m - 1) \frac{s_n^2}{\sigma^2} \sim \chi^2(n - m - 1)$,
- 3) jestliže $v_i = (\mathbf{X}^T \mathbf{X})_{ii}^{-1}$, potom $T_i = \frac{\hat{\beta}_i - \beta_i}{s_n \sqrt{v_i}} \sim t(n - m - 1)$.
- 4) Nechť $\mathbf{C} \in \mathbb{R}^{r, m+1}$ takové, že $h(\mathbf{C}) = r$. Potom kvadratická forma

$$\frac{q}{\sigma^2} = \frac{(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^T \mathbf{C}^T \left[\mathbf{C}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{C}^T \right]^{-1} \mathbf{C}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})}{\sigma^2} \sim \chi^2(r).$$

Důkaz.

Výsledky z LA:

- **Spektrální rozklad matice:**

$\mathbf{A}_{n \times n}$ symetrická matice \Rightarrow existuje ortogonální matice \mathbf{Q} a diagonální matice $\mathbf{\Lambda}$ tak, že $\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$, sloupce \mathbf{Q} jsou ON vlastní vektory matice \mathbf{A} a diagonální prvky matice $\mathbf{\Lambda}$ jsou jim odpovídající vlastní čísla.

- $\mathbf{A}_{n \times n}$ idempotentní matice \Rightarrow vlastní čísla jsou pouze 0 nebo 1 a $\text{h}(\mathbf{A}) = \text{tr}(\mathbf{A})$

Vlastnosti vektoru reziduí $\hat{\mathbf{e}}$

VĚTA 3.6

Uvažujme model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$, kde e_1, \dots, e_n jsou nekorelované a $e_i \sim (0, \sigma^2)$. Nechť $\hat{\boldsymbol{\beta}}$ je OLS $\boldsymbol{\beta}$ a $\hat{\mathbf{e}} = \mathbf{Y} - \hat{\mathbf{Y}}$ je vektor reziduí. Potom platí:

- 1) $E\hat{\mathbf{e}} = \mathbf{0}$, $\text{Cov}(\hat{\mathbf{e}}) = \sigma^2(\mathbf{I}_n - \mathbf{H})$,
- 2) pokud navíc $\mathbf{e} \sim N_n(0, \sigma^2\mathbf{I}_n)$, potom $\hat{\mathbf{e}} \sim N_n(\mathbf{0}, \sigma^2(\mathbf{I}_n - \mathbf{H}))$,
- 3) jestliže má model intercept, tj. $\beta_0 \neq 0$, potom $\sum_{i=1}^n \hat{e}_i = 0$,
- 4) $\sum_{i=1}^n \hat{e}_i \hat{y}_i = 0$.

DŮSLEDEK: Použitím bodů 3) a 4) dostaneme (stejně jako u jednorozměrné regrese)

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

tedy

$$SST = SSR + SSE$$

(v modelu s interceptem)

Důkaz.

3.2 Gauss-Markov theorem

- e_i i.i.d $N(0, \sigma^2) \Rightarrow$ OLS $\hat{\beta}$ je MLE, tzn. je eficientní (MVUE par. β)
- nenormální chyby:
 - ukážeme, že OLS $\hat{\beta}$ je **BLUE** (best linear unbiased estimator) par. β (za jistých podmínek)
 - mohou ale existovat lepší lineární vychýlené odhady nebo nelineární odhady

DEFINICE 3.1

Nechť β je vektor regresních parametrů v LM. Řekneme, že $\hat{\beta}$ je **lineární odhad** β , jestliže každé $\hat{\beta}_j$ je LK pozorování Y_i , $i = 1, \dots, n$, tedy

$$\hat{\beta}_j = \sum_{i=1}^n a_{ij} Y_i, \quad j = 0, \dots, m.$$

V maticovém zápisu $\hat{\beta} = \mathbf{A}\mathbf{Y}$, kde $\mathbf{A}^T = (a_{ij})$, $i = 1, \dots, n$, $j = 0, \dots, m$.

POZNÁMKA: Pokud v modelu (**) platí $h(\mathbf{X}) = m + 1$, potom je OLS $\hat{\beta}$ lineární, neboť $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$.

VĚTA 3.7 (Gauss-Markov)

Uvažujme model $\mathbf{Y} = \mathbf{X}\beta + \mathbf{e}$, kde matice \mathbf{X} má plnou hodnotu, e_1, \dots, e_n jsou nekorelované a $e_i \sim (0, \sigma^2)$. Potom OLS odhad $\hat{\beta}$ je BLUE parametru β .

Důkaz.

3.3 Testování modelu - tabulka ANOVA

Celkový F-test (overall F-test)

- Je model statisticky signifikantní? Tj. je alespoň jeden z koeficientů β_1, \dots, β_m nenulový?
- Mohli bychom testovat jednotlivé koeficienty $H_0 : \beta_j = 0$ pomocí alternativy t-testu.
- Celková chyba I. druhu by takto ale mohla být velká, pokud máme hodně proměnných. Museli bychom hodně snížit α pro jednotlivé testy, což zvýší pravděpodobnost chyby II. druhu
- Navíc je zde problém **multikolinearity**, jejíž jedním efektem jsou velké stand. chyby odhadů. To může vést k akceptování všech koeficientů jeho 0, i když je model celkově významný.

Bylo by dobré mít jednu statistiku pro test

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_m = 0 \quad \times \quad H_1 : \exists i \in \hat{m}, \beta_i \neq 0$$

ANOVA přístup pro jedn. regresi naznačuje, že statistika

$$F = \frac{\frac{SSR}{m}}{s_n^2}$$

by mohla být užitečná (vyplyne i z obecnějších přístupů k testování později)

Označení: $\bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{ij}$ – průměr j -tého sloupce matice \mathbf{X} ,

$$\bar{\mathbf{X}} = \begin{pmatrix} \bar{x}_0 & \bar{x}_1 & \cdots & \bar{x}_m \\ \vdots & \vdots & & \vdots \\ \bar{x}_0 & \bar{x}_1 & \cdots & \bar{x}_m \end{pmatrix}_{n \times m+1} \quad (\mathbf{X}_c)_{ij} = x_{ij} - \bar{x}_j, \quad i = 1, \dots, n, \quad j = 1, \dots, m$$

VĚTA 3.8

V modelu $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ tvaru (**), kde $h(\mathbf{X}) = m + 1$, e_i jsou nekorelované a $e_i \sim (0, \sigma^2)$, $i = 1, \dots, n$, platí

$$E \left[\frac{SSR}{m} \right] = \sigma^2 + \frac{\boldsymbol{\beta}^T (\mathbf{X} - \bar{\mathbf{X}})^T (\mathbf{X} - \bar{\mathbf{X}}) \boldsymbol{\beta}}{m} = \sigma^2 + \frac{\boldsymbol{\beta}_s^T \mathbf{X}_c^T \mathbf{X}_c \boldsymbol{\beta}_s}{m},$$

kde $\boldsymbol{\beta}_s = (\beta_1, \dots, \beta_m)^T$.

Věta z PRA: Nechť $Z = \mathbf{Y}^T \mathbf{A} \mathbf{Y}$ je kvadratická forma a nechť $E\mathbf{Y} = \boldsymbol{\mu}$ a $\text{Cov}\mathbf{Y} = \boldsymbol{\Sigma}$. Potom platí:

$$EZ = \text{tr}(\mathbf{A}\boldsymbol{\Sigma}) + \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu}.$$

Důkaz.

POZNÁMKA 3.5

- pokud $\beta_s = 0$, potom $E\left(\frac{SSR}{m}\right) = \sigma^2 = E s_n^2$, $\beta_s \neq 0$ implikuje, že $E\left(\frac{SSR}{m}\right) > \sigma^2$
- tedy velké hodnoty $F = \frac{SSR/m}{s_n^2}$ budou znamenat zamítnutí $H_0 : \beta_s = 0$