

01RAD - přednáška 5, 8.10.2024

3. Vícerozměrná regrese

- předp., že kromě y_i máme pro každé $i \in \hat{n}$ k dispozici také m nezávislých prom. $x_{i1}, x_{i2}, \dots, x_{im}$
- model:

$$Y_i = \beta_0 + \sum_{j=1}^m x_{ij}\beta_j + e_i, \quad i = 1, \dots, n,$$

kde e_1, \dots, e_n jsou nezávislé (nekorelované) chyby a $e_i \sim (0, \sigma^2)$

- na základě pozorování $(x_{i1}, \dots, x_{im}, y_i), i = 1, \dots, n$, chceme odhadnout par. $\beta = (\beta_0, \beta_1, \dots, \beta_m)^T$
- předpoklad: $n > m + 1$
- maticový zápis: $\mathbf{Y} = (Y_1, \dots, Y_n)^T, \quad \mathbf{y} = (y_1, \dots, y_n)^T, \quad \mathbf{e} = (e_1, \dots, e_n)^T,$

$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1m} \\ 1 & x_{21} & \cdots & x_{2m} \\ \vdots & \vdots & \cdots & \vdots \\ 1 & x_{n1} & \cdots & x_{nm} \end{pmatrix} \quad \text{matice modelu (regresní matice)}$$

Model:

$$\mathbf{Y}_{n \times 1} = \mathbf{X}_{n \times (m+1)} \beta_{(m+1) \times 1} + \mathbf{e}_{n \times 1} \quad (**)$$

Nejdříve budeme předpokládat e_1, \dots, e_n nezávislé, $e_i \sim N(0, \sigma^2)$, tzn.

$$\mathbf{e} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I}_n) \quad \text{a} \quad \mathbf{Y} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n).$$

Věrohodnostní funkce:

Pro pevné σ^2 je

$$\max_{\boldsymbol{\beta}} L(\boldsymbol{\beta}, \sigma^2) \quad \Longleftrightarrow \quad \min_{\boldsymbol{\beta}} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \min_{\boldsymbol{\beta}} g(\boldsymbol{\beta}).$$

VĚTA 3.1

Uvažujme model $(**)$ a necht' $\mathbf{e} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$. Potom $\hat{\beta}$ je MLE parametru β právě tehdy, když $\hat{\beta}$ je řešením soustavy rovnic

$$\mathbf{X}^T \mathbf{X} \beta = \mathbf{X}^T \mathbf{y}. \quad (\text{soustava normálních rovnic}) \quad (3.1)$$

Je-li matice $\mathbf{X}^T \mathbf{X}$ nesingulární, má tato soustava jednoznačné řešení ve tvaru

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

Důkaz.

Lemma 3.1

Soustava lineárních rovnic $\mathbf{Ax} = \mathbf{y}$ má řešení $\iff \langle \mathbf{y}; \mathbf{z} \rangle = 0$ pro všechna \mathbf{z} splňující $\mathbf{A}^T \mathbf{z} = \mathbf{0}$.

VĚTA 3.2

Soustava normálních rovnic $\mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^T \mathbf{y}$ má vždy alespoň jedno řešení.

Důkaz.



POZNÁMKA 3.1

- z vět plyne, že MLE $\hat{\boldsymbol{\beta}}$ může být nalezen řešením $m + 1$ lineárních rovnic o $m + 1$ neznámých
- málokdy existuje analytické řešení, je třeba použít numerické metody
- matice $\mathbf{X}^T \mathbf{X}$ může být špatně podmíněná, což ovlivňuje numerickou přesnost $\hat{\boldsymbol{\beta}}$
- často se užívají metody jako Choleského rozklad, QR rozklad, singulární rozklad (SVD)

- odvodili jsme pro normální chyby, minimalizace $g(\beta)$ lze ale použít i pro jiné druhy chyb
- nalezené $\hat{\beta}$ se pak nazývá **ordinary least squares estimate (OLS)**

Jak poznat, že existuje jednoznačné řešení NR bez nutnosti výpočtu $\mathbf{X}^T \mathbf{X}$?

VĚTA 3.3

Matice $\mathbf{X}^T \mathbf{X}$ je nesingulární \iff jsou sloupce matice \mathbf{X} lineárně nezávislé.

Důkaz.



POZNÁMKA 3.2

- $n > m + 1$, $h(\mathbf{X}) = m + 1 \implies$ ex. jednoznačné řešení NR a to $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$,
- pokud navíc $\mathbf{e} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$, je to MLE

POZNÁMKA 3.3

- pokud jsou sloupce \mathbf{X} LZ, je $\mathbf{X}^T \mathbf{X}$ singulární (většinou detekováno num. metodou výpočtu $\hat{\beta}$)
- horší situace je, pokud jsou sloupce \mathbf{X} „téměř“ LZ, tzv. **multikolinearita**
- způsobuje problémy při výpočtu $\hat{\beta}$, protože $\mathbf{X}^T \mathbf{X}$ je „téměř“ singulární
- jak detekovat a řešit multikolinearitu probereme na konci přednášky

Odhad parametru σ^2

- pro normální chyby získáme MLE σ^2 derivací $\ln L(\beta, \sigma^2)$, tedy

$$\sigma^2 = \frac{1}{n} SSE = \frac{1}{n} (\mathbf{y} - \mathbf{X}\hat{\beta})^T (\mathbf{y} - \mathbf{X}\hat{\beta}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2,$$

kde $\hat{y}_i = (\mathbf{X}\hat{\beta})_i = \mathbf{x}_i^T \hat{\beta}$, $i = 1, \dots, n$, a \mathbf{x}_i^T značí i -tý řádek matice \mathbf{X}

- jedná se o vychýlený odhad, proto se obecně používá nestranný odhad

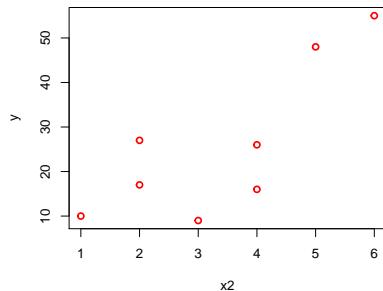
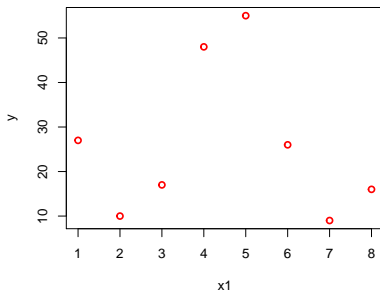
$$s_n^2 = \frac{1}{n - (m + 1)} SSE = \frac{1}{n - m - 1} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

a $s_n = \sqrt{s_n^2}$ jako odhad σ

- pro $e_i \sim (0, \sigma^2)$ se také používají statistiky s_n^2, s_n

PŘÍKLAD 3.1 (Umělá data)

x1	x2	y
2	1	10
3	2	17
4	5	48
1	2	27
5	6	55
6	4	26
7	3	9
8	4	16



Jednorozměrné modely:

$$y = \beta_0 + \beta_1 x_1 + e: \quad \hat{\beta}_0 = 28.5, \quad p_{val} = 0.093, \quad \hat{\beta}_1 = -0.57, \quad p_{val} = 0.847, \quad R^2 = 0.00668$$

$$y = \beta_0 + \beta_1 x_2 + e: \quad \hat{\beta}_0 = -1.34, \quad p_{val} = 0.891, \quad \hat{\beta}_1 = 8.10, \quad p_{val} = 0.018, \quad R^2 = 0.6356$$

Dvourozměrný model:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + e: \quad \hat{\beta}_0 = 8, \quad \hat{\beta}_1 = -5, \quad \hat{\beta}_2 = 12, \quad p_{val} = 0.000, \quad R^2 = 1$$

PŘÍKLAD 3.2 (Porodní váha a gestační stáří)

Sex	Age	Weight
boy	40	2968
boy	38	2795
girl	40	3163
⋮	⋮	⋮

Definujeme proměnné:

$$y = \text{Weight}, \quad x_1 = \text{Age}, \quad x_2 = \begin{cases} 0, & \text{Sex} = \text{boy}; \\ 1, & \text{Sex} = \text{girl}. \end{cases} \quad (\text{dummy variable})$$

Model: $Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + e$

Uvažujeme vlastně dvě rovnoběžné přímky

$$Y_b = \beta_0 + \beta_1 x_1 + e \quad \text{pro chlapce}, \quad Y_g = (\beta_0 + \beta_2) + \beta_1 x_1 + e \quad \text{pro dívky}$$

Pro různé směrnice, můžeme přidat **interakci** x_1 a x_2 , tj. $Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + e$

```
mod <- lm(Weight ~ Age + Sex + Age:Sex)
```

```

      Estimate Std. Error t value Pr(>|t|)
## (Intercept) -1268.67    1114.64  -1.138 0.268492
## Age          111.98      29.05   3.855 0.000986 ***
## Sexgirl      -872.99    1611.33  -0.542 0.593952
## Age:Sexgirl   18.42     41.76   0.441 0.663893
## ---
## Residual standard error: 180.6 on 20 degrees of freedom
## Multiple R-squared:  0.6435,    Adjusted R-squared:  0.59
## F-statistic: 12.03 on 3 and 20 DF,  p-value: 0.000101

```

3.1 Vlastnosti odhadů $\hat{\beta}$, s_n^2

VĚTA 3.4

Nechť $\hat{\beta}$ je OLS odhad parametru β v modelu (**), kde $h(\mathbf{X}) = m + 1$ a e_1, \dots, e_n jsou nekorelované a $e_i \sim (0, \sigma^2)$. Potom platí

- 1) $E(\hat{\beta}) = \beta$, ($\hat{\beta}$ je nestranný)
- 2) $\text{Cov}(\hat{\beta}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$,
- 3) $E(s_n^2) = \sigma^2$, (s_n^2 je nestranný)
- 4) pokud navíc $e_i \sim N(0, \sigma^2)$, $i = 1, \dots, n$, potom $\hat{\beta} \sim N_{m+1}(\beta, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1})$,
speciálně $\hat{\beta}_i \sim N(\beta_i, \sigma^2 v_i)$, kde $v_i = (\mathbf{X}^T \mathbf{X})_{ii}^{-1}$.

Důkaz.

POZNÁMKA 3.4 (Vlastnosti projekční matice \mathbf{H})

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T, \quad \hat{\mathbf{Y}} = \mathbf{H} \mathbf{Y}, \quad \mathbf{H}^T = \mathbf{H}, \quad (\mathbf{I}_n - \mathbf{H})^T = (\mathbf{I}_n - \mathbf{H}) \quad (\text{symetrie})$$

$$\mathbf{H}^2 = \mathbf{H}, \quad (\mathbf{I}_n - \mathbf{H})^2 = (\mathbf{I}_n - \mathbf{H}), \quad (\text{idempotentnost})$$

$$\mathbf{H} \mathbf{X} = \mathbf{X}, \quad \mathbf{H}(\mathbf{I}_n - \mathbf{H}) = (\mathbf{I}_n - \mathbf{H})\mathbf{H} = \mathbf{0}, \quad \text{tr}(\mathbf{H}) = \sum_{i=1}^n h_{ii} = m + 1.$$