

# MA012 Statistika II

## 8. Autokorelace a multikolinearita v lineárním modelu

Ondřej Pokora (pokora@math.muni.cz)

Ústav matematiky a statistiky, Přírodovědecká fakulta, Masarykova univerzita, Brno

(podzim 2015)



# Klasický lineární regresní model

$$Y \sim \mathcal{L}(X\beta, \sigma^2 I_n)$$

$$Y = X\beta + \varepsilon$$

- nesystematické chyby:

$$E\varepsilon_i = 0, \quad i = 1, \dots, n, \quad \text{tzn. } EY = X\beta,$$

- homogenita rozptylu chyb (měření):

$$D\varepsilon_i = \sigma^2 > 0, \quad i = 1, \dots, n,$$

- nekorelovanost chyb (měření):

$$C(\varepsilon_i, \varepsilon_j) = 0, \quad i, j = 1, \dots, n, \quad i \neq j, \quad \text{tzn. } D\varepsilon = DY = \sigma^2 I_n$$

V některých případech (často v časových řadách) hodnoty náhodné chyby  $\varepsilon_i$  závisí i na předchozích hodnotách  $\varepsilon_{i-k}$ ,  $k = 1, 2, \dots$ , což má za následek, že efekt náhodných chyb není jen okamžitý, ale je pocíťován i v budoucnosti. Tento případ se nazývá **autokorelace**.

Častým typem autokorelace je tzv. **autoregrese 1. řádu**, označovaná AR(1):

## Definice 1 (AR(1))

$$\varepsilon_i = \theta \varepsilon_{i-1} + u_i,$$

kde  $\theta$  je neznámý parametr,  $|\theta| < 1$ ,

$$Eu_i = 0, \quad Du_i = \sigma_u^2, \quad C(u_i, u_j) = 0 \text{ pro } i \neq j.$$

# AR(1)

Postupným aplikováním AR(1) formule na náhodné chyby obdržíme vyjádření

$$\begin{aligned}\varepsilon_i &= \theta\varepsilon_{i-1} + u_i = \theta(\theta\varepsilon_{i-2} + u_{i-1}) + u_i = \theta^2\varepsilon_{i-2} + \theta u_{i-1} + u_i = \\ &= \theta^2(\theta\varepsilon_{i-3} + u_{i-1}) + \theta u_{i-1} + u_i = \theta^3\varepsilon_{i-3} + \theta^2 u_{i-2} + \theta u_{i-1} + u_i = \\ &\dots = \sum_{l=0}^{\infty} \theta^l u_{i-l}.\end{aligned}$$

Počítejme střední hodnotu náhodných chyb:

$$E\varepsilon_i = E\left(\sum_{l=0}^{\infty} \theta^l u_{i-l}\right) = \sum_{l=0}^{\infty} \underbrace{\theta^l}_{\rightarrow 0} \underbrace{Eu_{i-l}}_0 = 0.$$

Při výpočtu rozptylu náhodných chyb využijeme nekorelovanosti veličin  $u$ :

$$D\varepsilon_i = D\left(\sum_{l=0}^{\infty} \theta^l u_{i-l}\right) = \sum_{l=0}^{\infty} \underbrace{\theta^{2l}}_{\rightarrow 0} \underbrace{Du_{i-l}}_{\sigma_u^2} = \sigma_u^2 \sum_{l=0}^{\infty} \underbrace{(\theta^2)^l}_{\frac{1}{1-\theta^2}} = \frac{\sigma_u^2}{1-\theta^2} = \sigma^2.$$

# AR(1)

Podobně spočítáme kovariance náhodných chyb (pro  $l = 1, 2, \dots$ ):

$$C(\varepsilon_i, \varepsilon_{i-l}) = \sum_{r=0}^{\infty} \sum_{s=0}^{\infty} \theta^r \theta^s C(u_{i-r}, u_{i-l-s}) = \theta^l \sigma_u^2 \sum_{r=0}^{\infty} \theta^{2r} = \frac{\theta^l \sigma_u^2}{1 - \theta^2} = \theta^l \sigma^2.$$

Kovarianční matice náhodných chyb je tedy tvaru

$$D\varepsilon = \underbrace{\frac{\sigma_u^2}{1 - \theta^2}}_{\sigma^2} \underbrace{\begin{pmatrix} 1 & \theta & \theta^2 & \dots & \theta^{n-1} \\ \theta & 1 & \theta & \dots & \theta^{n-2} \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \theta \\ \theta^{n-1} & \dots & \theta^2 & \theta & 1 \end{pmatrix}}_W = \sigma^2 W.$$

# Rozšířený lineární model

Máme tedy lineární regresní model  $Y \sim \mathcal{L}(X\beta, \sigma^2 W)$ :

## Definice 2 (Rozšířený lineární model)

$$Y = X\beta + \varepsilon, \quad E\varepsilon = \mathbf{0}, \quad D\varepsilon = \sigma^2 W,$$

kde kovariační matice již není diagonální.

## Věta 3 (Aitkenův odhad, zobecněná metoda nejmenších čtverců)

Mějme regresní model  $Y \sim \mathcal{L}(X\beta, \sigma^2 W)$  plné hodnosti, kde  $|W| > 0$ .  
Pak odhad vektoru parametrů zobecněnou metodou nejmenších čtverců je roven

$$\hat{\beta} = (X'W^{-1}X)^{-1}X'W^{-1}Y.$$

Z věty tedy plyne, že pokud známe parametr  $\theta$ , dokážeme najít odhady  $\hat{\beta}$  i v lineárním modelu s autokorelovanými náhodnými chybami.

# Detekce autokorelace

Graficky je možno autokorelaci detekovat tak, že vynášíme hodnoty reziduí  $r_i$  klasického lineárního modelu v závislosti na  $r_i$ . Je-li z grafu patrná přibližná lineární závislost, svědčí to o autoregresi 1. řádu (AR(1)) nebo o špatné volbě modelu.

## Věta 4 (Asymptotický test)

Statistika  $U$  má asymptoticky standardizované normální rozdělení pravděpodobnosti,

$$U = \frac{\hat{\theta} - \theta}{\sqrt{\frac{1-\theta^2}{n}}} \stackrel{as.}{\sim} N(0, 1).$$

Při testování hypotézy  $H_0 : \theta = 0$  proti  $H_1 : \theta \neq 0$  tedy pro dostatečně velká  $n$  ( $n \geq 30$ ) platí

$$U = \hat{\theta}\sqrt{n} \stackrel{as.}{\sim} N(0, 1),$$

a  $H_0$  zamítáme na hladině významnosti  $\alpha$ , pokud

$$|\hat{\theta}\sqrt{n}| > u_{1-\alpha/2}.$$

# Durbinův-Watsonův test

## Definice 5 (Durbinova-Watsonova statistika)

Durbinova-Watsonova statistika je definována pomocí reziduí v lineárním modelu:

$$D = \frac{\sum_{i=2}^n (r_i - r_{i-1})^2}{\sum_{i=1}^n r_i^2}.$$

## Věta 6 (Durbinův-Watsonův test)

Vždy platí  $0 \leq D \leq 4$ .

*Při nekorelovanosti náhodných chyb má statistika hodnotu  $D = 2$ .*

*Kladná korelace chyb se projeví hodnotou  $D \rightarrow 0$ ,*

*záporná korelace chyb se projeví hodnotou  $D \rightarrow 4$ .*

*Hodnoty kritických hodnot pro Durbinův-Watsonův test závisí na počtu pozorování, počtu parametrů a na hladině významnosti a jsou tabelovány.*



# Odhad parametru $\theta$

Parametr  $\theta$  v AR(1) modelu náhodných chyb lze odhadnout více způsoby:

- Odhadujeme jako regresní koeficient v modelu pro rezidua

$$r_i = \theta r_{i-1} + u_i, \quad i = 2, \dots, n$$

metodou nejmenších čtverců. Odtud pak

$$\hat{\theta} = \frac{\sum_{i=2}^n r_i r_{i-1}}{\sum_{i=2}^n r_{i-1}^2}.$$

- Pomocí Durbin – Watsonovy statistiky:

$$\hat{\theta} = 1 - \frac{D}{2}.$$

# Model bez autokorelace 1. řádu

- 1 Nalezneme odhad  $\hat{\theta}$ .
- 2 Vytvoříme nový model

$$Y_i^* = Y_{i+1} - \hat{\theta} Y_i, \quad x_{i,j}^* = x_{i+1,j} - \hat{\theta} x_{i,j}, \quad i = 1, \dots, n-1, j = 1, \dots, k,$$

tj. vznikne model

$$Y^* = X^* \beta^* + \varepsilon^*, \quad E\varepsilon^* = \mathbf{0}, \quad D\varepsilon^* = \sigma_{\varepsilon^*}^2 I_n$$

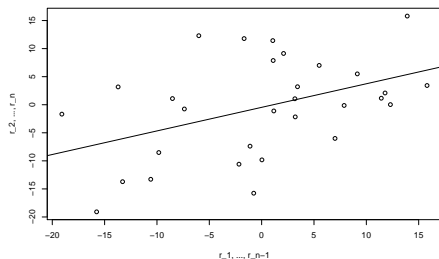
v němž počítáme odhady  $\hat{\beta}^*$  standardním způsobem.

## Příklad 1

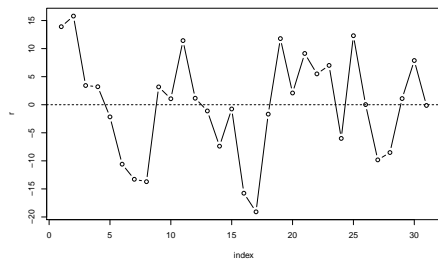
*V letech 1953 – 1983 byly měřeny ztráty vody při distribuci do domácností. Výsledky měření jsou uloženy v souboru [voda.csv](#). Proměnná  $x$  označuje množství vyrobené vody, proměnná  $Y$  ztrátu. Ověřte, zda se v datech vyskytuje autokorelace 1. řádu a případně ji odstraňte.*

Graficky: je patrná lineární závislost mezi rezidui

residual plot



residual plot



■ **Asymptotický test:**

$$U = |\hat{\theta} \sqrt{n}| = 2,339.$$

Nulovou hypotézu tedy **zamítáme**, neboť  $|U| > u_{1-\alpha/2} = 1,96$ .

■ **Durbinův-Watsonův test:**

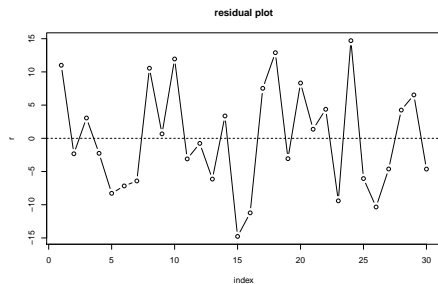
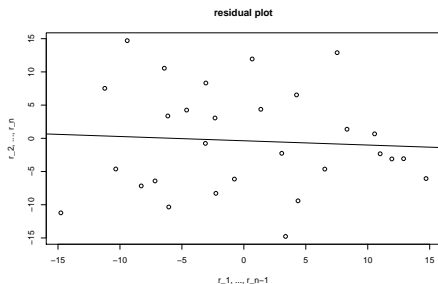
$$D = 1,082$$

a  $p$ -hodnota testu je 0,0016, takže také **zamítáme** nulovou hypotézu.

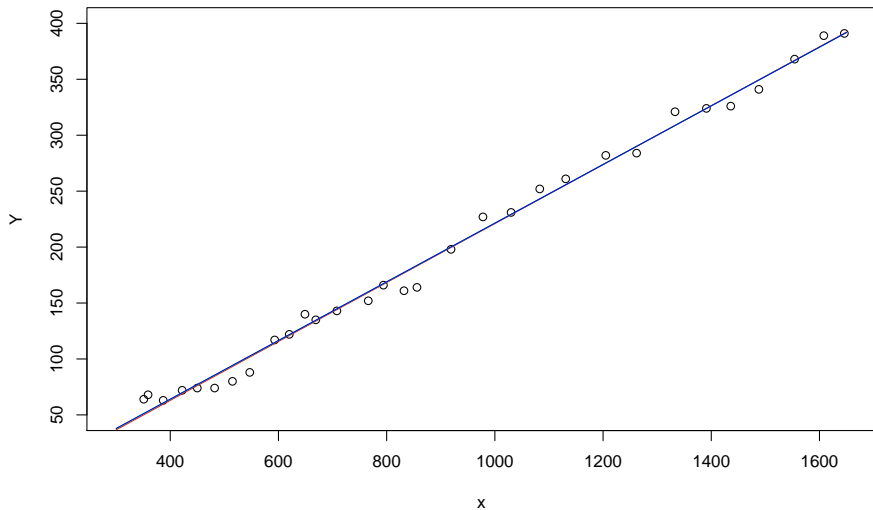
## Odstranění autokorelace:

Odhad metodou nejmenších čtverců:  $\hat{\theta} = 0,42$ , z D-W statistiky:  $\hat{\theta} = 0,459$ .

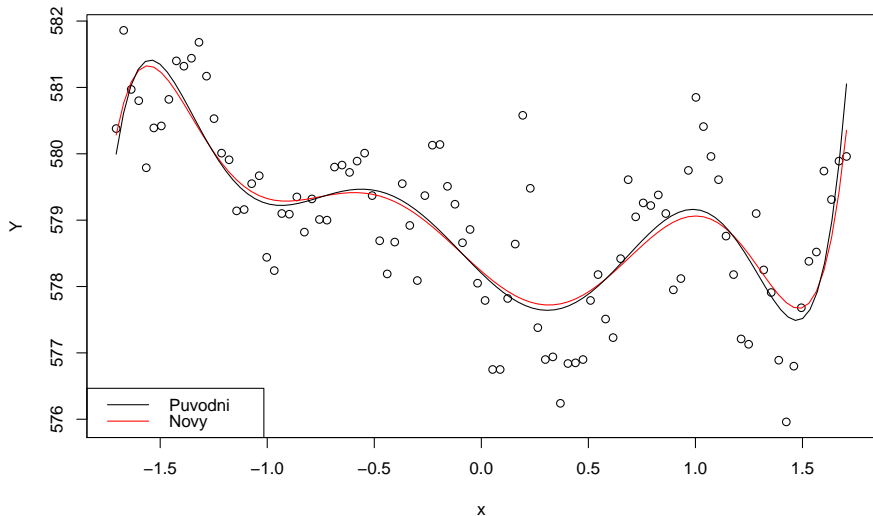
V nově vzniklém modelu vykreslíme rezidua:



Také D-W test již **nezamítá** nulovou hypotézu ( $p$ -hodnota je 0,4).



# Příklad LakeHuron: řešení





**Multikolinearitou** se rozumí vzájemná lineární závislost vysvětlujících proměnných. Přesnou multikolinearitou se rozumí případ, kdy jednotlivé sloupce  $x_j$ ,  $j = 1, \dots, p$  matice plánu  $X$  bez sloupce pro absolutní člen jsou lineárně závislé, takže pro aspoň jednu nenulovou konstantu  $c_j$  platí

$$c_1 x_1 + \dots + c_p x_p = \mathbf{0}.$$

V praxi bychom se s tímto případem neměli setkávat, neboť při rozumně sestaveném regresním modelu využijeme lineární kombinaci a zmenšíme počet vysvětlujících proměnných. Podobně nereálný je v praxi případ ortogonálních vysvětlujících proměnných, kdy matice  $X$  je ortogonální a platí, že  $X'X = I_p$ . V praxi se tedy **multikolinearitou** rozumí případ, kdy **přibližně** platí rovnice vyjadřující lineární kombinaci vysvětlujících proměnných. V případě silné multikolinearity je determinant informační matice  $X'X$  blízký nule, nejmenší vlastní číslo je rovněž blízké nule a matice  $X'X$  je „skoro singulární“. O multikolinearitě svědčí i vysoké hodnoty poměru největšího a nejmenšího vlastního čísla.

# Důvody multikolinearity

- Multikolinearitu způsobuje regresní rovnice obsahující **nadbytečné** vysvětlující proměnné. Statistickými technikami můžeme přebytečné proměnné identifikovat a vyloučit z regresní rovnice.
- Multikolinearitu jen ztěží odstraníme v úlohách, kdy vzájemná spříženost hodnot vysvětlujících proměnných je způsobena neuvažovanými veličinami nebo **formou statistického zjišťování**. Jde-li např. o údaje z časových řad, je podobný vývoj sledovaných veličin dostatečným důvodem vzniku multikolinearity. Vzhledem k tomu, že multikolinearitu hodnotíme výhradně na základě určitého souboru pozorování, stačí nesprávný výběr kombinací hodnot vysvětlujících proměnných, nerepresentujících obor možných hodnot, k existenci významné multikolinearity.
- Závažným důvodem multikolinearity je **skutečný vztah** vysvětlujících proměnných v rámci sledovaného jevu, procesu nebo systému. V tomto případě je třeba využít všechny informace nevýběrového charakteru k zlepšení kvality regresních odhadů.

# Důsledky multikolinearity

V případě přesné multikolinearity je matice  $X'X$  singulární a běžnou inverzí nepořídíme odhad neznámých parametrů  $\beta$  metodou nejmenších čtverců. Pro přibližnou multikolinearitu jsme sice schopni matici  $X'X$  invertovat, ale kvalita pořízených odhadů je poměrně nízká.

Snížení kvality se projevuje

- v kovarianční matici  $D\hat{\beta} = \sigma^2(X'X)^{-1}$
- v přesnosti prováděných výpočtů

neboť důsledkem vysokých rozptylů odhadů jsou příliš široké intervaly spolehlivosti, a tedy malá přesnost odhadu.

Logickým důsledkem multikolinearity je obtížné vyjádření individuálního vlivu jednotlivých vysvětlujících proměnných. Projevuje se to nízkými hodnotami testových kritérií v T-testech nedovolujícími potvrdit závažnost jednotlivých regresorů v regresní funkci. Závažným důsledkem je značná výpočetní nespolehlivost a nestabilní hodnoty regresních odhadů. Stačí malý zásah do statistických údajů a výsledné odhady jsou odlišné.

# VIF – Variance Inflation Factors

## Definice 7 (VIF = variance inflation factors)

Diagonální prvky matice  $(\mathbf{X}'\mathbf{X})^{-1}$ , tj.

$$(a_1, \dots, a_p) = \text{diag}(\mathbf{X}'\mathbf{X})^{-1},$$

označované jako **VIF – variance inflation factors**, úzce souvisí s **mnohonásobnými korelačními koeficienty**, vyjadřující vztah  $j$ -té vysvětlující proměnné a lineární funkce ostatních vysvětlujících proměnných. Lze je zapsat jako

$$a_j = \frac{1}{(1 - r_j^2) \mathbf{x}'_j \mathbf{x}_j},$$

kde  $r_j = r_{x_j \cdot x_1 x_2 \dots x_{j-1} x_{j+1} \dots x_p}$  je koeficient mnohonásobné korelace mezi  $j$ -tou proměnnou a všemi ostatními proměnnými (kromě  $j$ ).

Vysoký stupeň multikolinearity se projevuje vysokými hodnotami korelačních koeficientů  $r_j$ , ale i vysokými hodnotami některých jednoduchých korelačních koeficientů.

# Detekce multikolinarity

Test multikolinarity je založen na testování, že korelační matice nezávisle proměnných je jednotková,

$$H_0 : \mathbf{R} = \mathbf{I}_p, \quad H_1 : \mathbf{R} \neq \mathbf{I}_p,$$

kde  $\mathbf{R}$  je korelační matice nezávisle proměnných rozměru  $p \times p$ .

## Věta 8 (Test multikolinarity)

*Za platnosti  $H_0$  má statistika*

$$K = - \left[ n - 1 - \frac{1}{6}(2p + 7) \right] \ln |\mathbf{R}| \sim \chi^2 \left( \frac{1}{2}p(p - 1) \right)$$

*$\chi^2$  rozdělení pravděpodobnosti.*

*Hypotézu  $H_0$  tedy na hladině významnosti  $\alpha$  zamítáme, pokud*

$$K > \chi_{1-\alpha}^2 \left( \frac{1}{2}p(p - 1) \right).$$

# Multikolarita: identifikace proměnných

Pro identifikaci proměnných způsobujících multikolaritu se doporučují statistiky

$$F_j = \frac{n-p}{p-1} (d_j - 1),$$

kde  $d_k$  jsou diagonální prvky inverzní matice ke korelační matici nezávisle proměnných,  $(d_1, \dots, d_p) = \text{diag } \mathbf{R}^{-1}$ .

## Věta 9

*V případě, že proměnná  $X_j$  nezpůsobuje multikolaritu, má veličina  $F_j$  Fisherovo-Snedecorovo rozdělení,*

$$F_j \sim F(p-1, n-p).$$

*Pokud tedy  $F_j > F_{1-\alpha}(p-1, n-p)$ , přispívá  $j$ -tá proměnná k multikolaritě.*

# Metoda postupné regrese

Následující algoritmus stavby lineárního modelu bez multikolinearity je založen na tom, že do modelu postupně zařazujeme jen ty regresory (nezávisle proměnné), které významně přispívají ke zlepšení kvality odhadu  $\hat{\beta}$ . K výběru nejlepší podmnožiny regresorů používá tzv. **metodu postupné regrese**:

- 1 Spočteme korelační matici  $\mathbf{R}$  regresorů a provedeme test hypotézy  $H_0 : \mathbf{R} = \mathbf{I}_p$ . Pokud test  $H_0$  nezamítne, pracujeme s klasickým lineárním modelem, do něhož můžeme zařadit všechny regresory. Je-li korelace mezi regresory prokázána, pokračujeme dalším krokem.
- 2 Spočteme korelační koeficienty  $r_{Y, X_1}, \dots, r_{Y, X_k}$  a pro další krok vybereme regresor  $X_j$  s největší absolutní hodnotou korelačního koeficientu  $|r_{Y, X_j}|$ .
- 3 Sestavíme model  $Y = \beta_0 + \beta_1 X_j$  a odhadneme jeho parametry metodou nejmenších čtverců. Vypočteme hodnotu statistiky F-testu modelu,

$$F = (n - 2) \frac{R_{Y \cdot X_j}^2}{1 - R_{Y \cdot X_j}^2}, \text{ kde } R_{Y \cdot X_j}^2 \text{ je koeficient determinace tohoto modelu.}$$

Pokud  $F > F_{1-\alpha}(1, n - 2)$ , ponecháme regresor  $X_j$  v modelu.

# Metoda postupné regrese

- 4 Spočteme parciální korelační koeficienty

$r_{Y X_1 \cdot X_j}, \dots, r_{Y X_{j-1} \cdot X_j}, r_{Y X_{j+1} \cdot X_j}, \dots, r_{Y X_p \cdot X_j}$  a pro další krok vybereme ten regresor  $X_i$ , jehož  $|r_{Y X_i \cdot X_j}|$  je největší.

- 5 Sestavíme model  $Y = \beta_0 + \beta_1 X_j + \beta_2 X_i$  a odhadneme jeho parametry metodou nejmenších čtverců. Vypočteme hodnotu statistiky

$$F = (n - 3) \frac{R_{Y \cdot X_j X_i}^2 - R_{Y \cdot X_j}^2}{1 - R_{Y \cdot X_j X_i}^2},$$
 kde ve jmenovateli je koeficient determinace

nového modelu a v čitateli rozdíl koeficientů determinace nového a předchozího modelu. Připomeňme, že předchozí model je podmodelem nového modelu. Pokud  $F > F_{1-\alpha}(2, n - 3)$ , ponecháme regresor  $X_i$  v modelu.

- 6 Spočteme parciální korelační koeficienty  $r_{Y X_1 \cdot X_j X_i}, \dots, r_{Y X_p \cdot X_j X_i}$  mezi  $Y$  a dosud nezařazenými regresory za vyloučení vlivu regresorů do modelu již zařazených. Ze skupiny dosud nezařazených regresorů vybereme regresor s největší absolutní hodnotou parciálního korelačního koeficientu, sestavíme s ním nový model a analogicky opakujeme kroky 5–6, dokud lze zařazovat další regresory.



# Model standardizovaných proměnných

Místo původních proměnných  $y_i$  a  $x_{ij}$  pracujeme s proměnnými ve tvaru

$$q_i = \frac{y_i - \bar{y}}{s_y}, \quad z_{ij} = \frac{x_{ij} - \bar{x}_j}{s_{x_j}},$$

kde  $s_y$  a  $s_{x_j}$  jsou směrodatné odchylky jednotlivých proměnných. Standardizací vysvětlujících proměnných dostáváme při použití metody nejmenších čtverců místo matice  $\mathbf{X}'\mathbf{X}$  korelační matici  $\mathbf{R} = \mathbf{Z}'\mathbf{Z}/n$ . Vektor  $\mathbf{Z}'\mathbf{q}/n$  obsahuje jednoduché korelační koeficienty  $r_{YX_j}$ . Standardizací proměnných se zmenšují zaokrouhlovací chyby a zlepšují se možnosti hodnocení individuálního vlivu proměnných pomocí regresních parametrů.

# Model v kanonickém tvaru

Místo modelu ve tvaru

$$Y = X'\beta + \varepsilon$$

pracujeme s modelem

$$Y = U'\gamma + \varepsilon,$$

kde matice  $U = XV$ , vektor  $\gamma = V'\beta$  a  $V$  je matice standardizovaných vlastních vektorů odpovídajících vlastním číslům matice  $X'X$ . Odhady parametrů v kanonickém tvaru:

$$\hat{\gamma} = L^{-1}U'Y,$$

kde  $L$  je diagonální matice s vlastními čísly matice  $X'X$ . Kovarianční matice odhadů  $D(\hat{\gamma}) = \sigma^2 L^{-1}$  ukazuje, že i v tomto případě jsou odhady nezávislé. Reziduální součet čtverců se transformací nemění.

# Hřebenová regrese (ridge regression)

Autoři tzv. **hřebenové regrese** vyšli z faktu, že vliv multikolinearity se silně projevuje při výpočtu inverzní matice  $(\mathbf{X}'\mathbf{X})^{-1}$ , a to konkrétně tím, že diagonální prvky této inverzní matice jsou příliš velké.

Navrhli proto jiný odhad

$$\hat{\beta} = (\mathbf{X}'\mathbf{X} + m \mathbf{I}_k)^{-1} \mathbf{X}'\mathbf{Y},$$

závisející na parametru  $m > 0$ . Výhodou takového odhadu je především to, že vhodnou volbou  $m$  lze numerické vlastnosti potřebné pro výpočet inverzní matice výrazně zlepšit a získat tak relativně přesné hodnoty. Nevýhodou je to, že tento odhad není nestranný a že neexistuje jednotný algoritmický postup pro určení vhodné hodnoty parametru  $m$ . Obvykle se vychází z modelu v kanonickém tvaru a postupně se zkouší různé volby hodnot  $m$ , dokud se nedosáhne „stabilizace“, kdy model dostává charakter ortogonálního systému.

## Příklad 2

V souboru `vydaje.csv` jsou uložena data o 20 náhodně vybraných domácnostech. Sloupce proměnné „domácnosti“ obsahují postupně tyto údaje: výdaje za potraviny a nápoje ( $Y$ ), počet členů domácnosti ( $X_1$ ), počet dětí ( $X_2$ ), průměrný věk výdělečně činných ( $X_3$ ) a příjem domácnosti ( $X_4$ ). Metodou postupné regrese zkonstruuje model s nejlepší podmíněností regresorů.

Uvažujme nejdřív model se všemi regresory. Spočtěme nejprve pro ilustraci  $|(\mathbf{X}'\mathbf{X})^{-1}| = 4,65 \times 10^{-16}$ . Také hodnoty VIF jsou pro první dva regresory vysoké:

| clenu | deti  | vek  | prijem |
|-------|-------|------|--------|
| 21,23 | 16,18 | 1,31 | 3,4    |

Testujeme-li hypotézu  $H_0 : \mathbf{R} = \mathbf{I}_4$ , hodnota testové statistiky

$$K = - \left[ 19 - \frac{15}{6} \right] \ln |\mathbf{R}| = 64,94$$

výrazně převyšuje kritickou hodnotu  $\chi_{0,95}^2(6) = 12,59$ . Hypotézu  $H_0$  tedy na hladině významnosti 0,05 zamítáme.

Pro identifikaci proměnných způsobujících multikolinearitu můžeme spočítat dílčí statistiky  $F_j$

| clenu  | deti  | vek  | prijem |
|--------|-------|------|--------|
| 107,88 | 80,94 | 1,66 | 12,83  |

kteřé porovnáme s kritickou hodnotou  $F_{0,95}(3, 16) = 3,24$ .

Postupná regrese

- 1 Spočteme korelační koeficienty  $r_{Y,X_1}, \dots, r_{Y,X_4} = (0,77; 0,67; 0,18; 0,73)$ . Vybereme regresor  $X_1$ , neboť jeho korelace je v absolutní hodnotě největší.

- 2 Sestavíme model  $Y = \beta_0 + \beta_1 X_1$ . Vypočteme hodnotu statistiky

$$F = \frac{(n-2)R_{Y \cdot X_1}^2}{1-R_{Y \cdot X_1}^2} = \frac{18 \cdot 0,5897}{1-0,5897} = 25,87. \text{ Tato hodnota je větší než}$$

$F_{0,95}(1, 18) = 4,41$ , takže regresor  $X_1$  ponecháme v modelu.

- 3 Spočteme parciální korelační koeficienty

$r_{Y,X_2 \cdot X_1}, r_{Y,X_3 \cdot X_1}, r_{Y,X_4 \cdot X_1} = (0,36; 0,499; 0,32)$ . Vybereme regresor  $X_3$ , jehož parciální korelační koeficient je v absolutní hodnotě největší.

- 4 Sestavíme model  $Y = \beta_0 + \beta_1 X_1 + \beta_3 X_3$ . Vypočteme hodnotu statistiky

$$F = \frac{(n-3)(R_{Y \cdot X_1}^2 - R_{Y \cdot X_1 X_3}^2)}{1 - R_{Y \cdot X_1 X_3}^2} = \frac{17 \cdot 0,102}{1 - 0,69} = 5,64. \text{ Tato hodnota je větší než}$$

$F_{0,95}(2, 17) = 3,59$ , tedy ponecháme regresor  $X_3$  v modelu.

- 5 Spočteme parciální korelační koeficienty

$r_{Y, X_2 \cdot (X_1, X_3)}, r_{Y, X_4 \cdot (X_1, X_3)} = (0,19; 0,17)$ . Vybereme regresor  $X_2$ , jehož parciální korelační koeficient je v absolutní hodnotě největší.

- 6 Sestavíme model  $Y = \beta_0 + \beta_1 X_1 + \beta_3 X_3 + \beta_2 X_2$ . Vypočteme hodnotu

$$\text{statistiky } F = \frac{(n-3)(R_{Y \cdot X_1 X_3}^2 - R_{Y \cdot X_1 X_3 X_2}^2)}{1 - R_{Y \cdot X_1 X_3 X_2}^2} = \frac{16 \cdot 0,012}{1 - 0,704} = 0,63. \text{ Tato hodnota je}$$

menší než  $F_{0,95}(3, 16) = 3,24$ , a tedy regresor  $X_2$  již nezahrneme do modelu.

Výsledný model je tedy tvaru

$$Y = \beta_0 + \beta_1 X_1 + \beta_3 X_3.$$