

# KOEFICIENT DETERMINACE V REGRESI S CHYBAMI V OBOU PROMĚNNÝCH

KAREL ZVÁRA

ABSTRACT. Koeficient determinace lze v lineární regresi s absolutním členem definovat řadou způsobů. Který z nich použít, když chceme koeficient determinace použít v modelu jednoduché lineární funkční závislosti?

## 0. ÚVOD

Kvalitu vysvětlení dat vyšetřovanou *regresní* závislostí měříme často pomocí koeficientu determinace. Když chceme tento pojem rozšířit i na jiný model závislosti (ne nutně regresní), můžeme se pokusit k tomu účelu použít některé z dostupných zobecnění koeficientu determinace. V tomto článku se budeme věnovat modelům funkční a strukturní závislosti.

První z těchto úloh lze formulovat následovně. Nechť

$$(0.1) \quad \begin{pmatrix} x_i \\ y_i \end{pmatrix} \sim N \left( \begin{pmatrix} \xi_i \\ \alpha + \beta \xi_i \end{pmatrix}, \sigma^2 \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right), \quad 1 \leq i \leq n,$$

jsou nezávislé náhodné vektory, v nichž  $\alpha, \beta, \xi_1, \dots, \xi_n, \sigma^2$  jsou neznámé parametry. V literatuře se v této souvislosti hovoří o *funkční* závislosti nebo o regresní přímce s chybami v obou proměnných. V případě, že bychom  $\xi_1, \dots, \xi_n$  chápali jako nezávislé realizace náhodné veličiny  $\xi$ , šlo by o *strukturní* model a v (0.1) by bylo rozdělení měřených veličin podmíněné neznámými hodnotami  $\xi_1, \dots, \xi_n$ . Standardní předpoklad  $\xi \sim N(\delta, \omega^2)$  by pak vedl k nepodmíněnému rozdělení

$$(0.2) \quad \begin{pmatrix} x_i \\ y_i \end{pmatrix} \sim N \left( \begin{pmatrix} \delta \\ \alpha + \beta \delta \end{pmatrix}, \begin{pmatrix} \omega^2 + \sigma^2 & \beta \omega^2 \\ \beta \omega^2 & \beta^2 \omega^2 + \sigma^2 \end{pmatrix} \right), \quad 1 \leq i \leq n,$$

V tomto případě mají všechny uvažované vektory stejné normální rozdělení s korelačním koeficientem

$$\rho = \frac{\beta \omega^2}{\sqrt{(\omega^2 + \sigma^2)(\beta^2 \omega^2 + \sigma^2)}}.$$

Při asymptotických úvahách musíme ve funkčním modelu předpokládat cosi o asymptotickém chování parametrů  $\xi_1, \dots, \xi_n$ . Předpoklady

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \xi_i = \delta,$$
$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n (\xi_i - \delta)^2 = \omega^2$$

způsobí, že mnohé asymptotické výrazy budou ve funkčním i strukturním modelu formálně stejné.

## 1. DEFINICE KOEFICIENTU DETERMINACE

Připomeňme několik možných definic koeficientu determinace. (Podrobný přehled je uveden např. v [2].) Všechny definice vycházejí z předpokladu, že

$$y_i \sim N(\alpha + x_i' \beta, \sigma^2), \quad 1 \leq i \leq n,$$

jsou nezávislé náhodné veličiny,  $\alpha, \beta$  a  $\sigma^2$  jsou neznámé parametry a  $x_1, \dots, x_n$  jsou dané vektory. Označíme-li jako  $a$  a  $b$  odhady parametrů  $\alpha$  a  $\beta$  metodou nejmenších čtverců, lze vyrovnané hodnoty počítat pomocí

$$\hat{y}_i = a + x_i' b, \quad 1 \leq i \leq n.$$

Rezidua  $u_i$  se pak definují jako  $u_i = y_i - \hat{y}_i$ . Součet čtverců reziduí označíme symbolem  $RSS$ .

Kdykoliv budeme v dalším textu psát explicitní vzorce pro jednoduchou regresní přímku nebo pro model z Úvodu, bude užitečné označení

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2, \quad S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}), \quad S_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2,$$

$$s_{xx} = S_{xx}/n, \quad s_{xy} = S_{xy}/n, \quad s_{yy} = S_{yy}/n.$$

1.1. Víme, že lineární regresní funkce maximalizuje koeficient mnohonásobné korelace mezi náhodnou veličinou (závisle proměnnou)  $Y$  a nezávisle proměnnými  $X$ . V případě jednoduché lineární regrese jde o klasický Pearsonův korelační koeficient. Čtveřec maximální hodnoty výběrového koeficientu mnohonásobné korelace se pak nazývá koeficient determinace a značí  $R^2$ . V případě mnohorozměrného normálního rozdělení ( $Y$  i  $X$  chápeme jako náhodné) je  $R_1^2$  konzistentním odhadem čtverce koeficientu mnohonásobné korelace  $\rho_{Y,X}^2$ . Je známo, že tento koeficient mnohonásobné korelace je roven Pearsonovu koeficientu korelace mezi náhodnou veličinou  $Y$  a její nejlepší lineární predikcí pomocí  $X$ , kterou můžeme označit symbolem  $\hat{Y}(X)$ . Tento vztah platí i pro výběrové protějšky, takže je také

$$(1.1) \quad R_1^2 = r_{y, \hat{y}}^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sum_{i=1}^n (y_i - \bar{y})^2 \sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2}.$$

Speciálně pro jednoduchou regresní přímku dostaneme

$$(1.2) \quad R_1^2 = \frac{S_{xy}^2}{S_{xx} S_{yy}}.$$

1.2. Jiná interpretace koeficientu determinace má původ v analýze rozptylu. Když se snažíme vysvětlit variabilitu náhodné veličiny  $Y$  pomocí funkce konkrétní realizace  $x$  náhodného vektoru  $X$ , zjistíme, že nejlepším prediktorem  $\hat{Y}(x)$  je podmíněná střední hodnota  $E(Y|X = x)$ . V případě sdruženého normálního rozdělení je tato střední hodnota lineární funkcí vektoru  $x$ . Bez ohledu na předpoklad normality platí pro rozptyly vztah

$$(1.3) \quad \text{var}(Y) = \text{var}(Y \cdot x) + \text{var}(\hat{Y}(x)).$$

Veličina

$$1 - \frac{\text{var}(Y.x)}{\text{var}(Y)}$$

ukazuje, jak velký díl variability náhodné veličiny  $Y$  vysvětluje lineární závislost. Výběrovým protějškem posledního vztahu je alternativní definice koeficientu determinace

$$(1.4) \quad R_2^2 = 1 - \frac{RSS}{S_{yy}}.$$

Není těžké dokázat, že v doposud uvažovaném modelu platí  $R_1^2 = R_2^2$ .

1.3. Další dvě vyjádření koeficientu determinace odvodíme z testů hypotézy  $H_0 : \beta = 0$ . Inspiraci lze najít v [3]. Klasickým postupem je test poměrem věrohodnosti. Sdružená hustota náhodného vektoru  $\mathbf{y}$  je rovna

$$h(\mathbf{y}) = (2\pi\sigma^2)^{n/2} \exp\left(-\sum_{i=1}^n (y_i - \alpha - \mathbf{x}_i'\beta)^2 / (2\sigma^2)\right),$$

takže logaritmus věrohodnostní funkce je roven

$$l(\alpha, \beta, \sigma^2) = c - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \alpha - \mathbf{x}_i'\beta)^2.$$

Maximalizace dá pro  $\alpha, \beta$  odhady totožné s odhady metodou nejmenších čtverců a  $\widehat{\sigma^2} = RSS/n$ . Výsledná hodnota logaritmu věrohodnostní funkce je rovna  $\widehat{l} = c_1 - \frac{n}{2} \ln \widehat{\sigma^2}$ . Podobně za platnosti hypotézy  $H_0 : \beta = 0$  dostaneme odhady  $\widetilde{\alpha} = \bar{y}$  a  $\widetilde{\sigma^2} = S_{yy}/n$ . Proto je  $\widetilde{l} = c_1 - \frac{n}{2} \ln \widetilde{\sigma^2}$ .

Známy asymptotický test poměrem věrohodnosti je založen na statistice  $2(\widehat{l} - \widetilde{l})$ , která má v regulárním případě za platnosti testované hypotézy asymptoticky rozdělení  $\chi^2$ . V našem případě dostaneme

$$\begin{aligned} \widehat{l} - \widetilde{l} &= -\frac{n}{2} (\ln \widehat{\sigma^2} - \ln \widetilde{\sigma^2}) \\ &= -\frac{n}{2} \ln \frac{RSS}{S_{yy}} \\ &= -\frac{n}{2} \ln(1 - R_2^2). \end{aligned}$$

Odtud snadno dostaneme novou definici koeficientu determinace

$$(1.5) \quad R_3^2 = 1 - \exp\left(-\frac{2}{n}(\widehat{l} - \widetilde{l})\right).$$

1.4. Waldův test hypotézy, která tvrdí, že  $g(\theta) = 0$  (kde  $g(\theta)$  je vhodná hladká funkce a  $\theta$  je vektor neznámých parametrů), spočívá v tom, že do této funkce zkusíme dosadit maximálně věrohodný odhad vektoru parametrů a vezmeme v úvahu jeho asymptotickou varianční matici:

$$W = g(\widehat{\theta})' \left( \widehat{\text{var}(g(\widehat{\theta}))} \right)^{-1} g(\widehat{\theta}).$$

Za standardních předpokladů regularity je asymptotická varianční matice odhadu  $\theta$  dána inverzní maticí k Fisherově informační matici. V našem konkrétním případě lineární regrese je maximálně věrohodným odhadem vektoru  $\beta$  odpovídající odhad  $b$  metodou nejmenších čtverců s varianční maticí

$$\begin{aligned}\text{var}(b) &= \sigma^2 \left( \sum_{i=1}^n x_i x_i' \right)^{-1} \\ &= \sigma^2 (X'X)^{-1}.\end{aligned}$$

Použijme opět odhad  $\widehat{\sigma^2} = RSS/n$ . Protože v uvažovaném případě platí  $RSS = S_{yy} - b'(X'X)b$ , po malých úpravách ověříme, že i následující definice koeficientu determinace je ekvivalentní v případě lineární regrese s definicí (1.4) :

$$(1.6) \quad R_1^2 = \frac{W}{W+n}.$$

1.5. S posledními dvěma testy (poměrem věrohodnosti a Waldův test) se v učebnicích uvádí často ještě další, s nimi asymptoticky ekvivalentní test, založený na ověření, zda je příslušný Lagrangeův multiplikátor dostatečně blízko nule (viz např. [3] nebo (6e.3.6) v [4]). Označme jako  $v(\theta)$  vektor parciálních derivací logaritmu věrohodnostní funkce a jako  $F(\theta)$  odpovídající Fisherovu informační matici příslušnou  $n$  pozorováním. Označíme-li jako  $\tilde{\theta}$  maximálně věrohodný odhad vektoru parametrů za hypotézy, pak má Raova statistika tvar

$$S = v(\tilde{\theta})' \left( F(\tilde{\theta}) \right)^{-1} v(\tilde{\theta}).$$

Po určité námaze dostaneme

$$v(\alpha, \beta, \sigma^2) = \frac{1}{\sigma^2} \begin{pmatrix} 1'(y - \alpha 1 - X\beta) \\ X'(y - \alpha 1 - X\beta) \\ -(1/n) + (y - \alpha 1 - X\beta)'(y - \alpha 1 - X\beta)/(2\sigma^2) \end{pmatrix},$$

$$F(\alpha, \beta, \sigma^2) = \frac{1}{\sigma^2} \begin{pmatrix} 1'1 & 1'X & 0 \\ X'1 & X'X & 0 \\ 0 & 0' & n/(2\sigma^2) \end{pmatrix}.$$

Použijeme-li toho, že za platnosti hypotézy  $H_0 : \beta = 0$  je

$$\tilde{\sigma^2} = s_{yy},$$

dostaneme po úpravách

$$S = n \left( 1 - \frac{RSS}{S_{yy}} \right).$$

Zřejmě je tedy souvislost s koeficientem determinace dána vztahem

$$(1.7) \quad R_3^2 = \frac{1}{n} S.$$

## 2. FUNKČNÍ MODEL

Vraťme se k funkčnímu modelu z úvodu. Zlogaritmujeme-li sdruženou hustotu náhodného vektoru  $x, y$ , dostaneme logaritmus věrohodnostní funkce

$$l(\alpha, \beta, \sigma^2, \xi) = -n \ln(2\pi) - n \ln(\sigma^2) - \frac{1}{2\sigma^2} \left[ \sum_{i=1}^n (x_i - \xi_i)^2 + \sum_{i=1}^n (y_i - \alpha - \beta \xi_i)^2 \right].$$

Při hledání odhadů metodou maximální věrohodnosti dostaneme dostaneme soustavu rovnic

$$(2.1) \quad \sum_{i=1}^n (y_i - \alpha - \beta \xi_i) = 0,$$

$$(2.2) \quad \sum_{i=1}^n \xi_i (y_i - \alpha - \beta \xi_i) = 0,$$

$$(2.3) \quad (x_i - \xi_i) + (y_i - \alpha - \beta \xi_i) = 0, \quad 1 \leq i \leq n,$$

$$(2.4) \quad -\frac{n}{\sigma^2} + \frac{1}{2\sigma^4} \left[ \sum_{i=1}^n (x_i - \xi_i)^2 + \sum_{i=1}^n (y_i - \alpha - \beta \xi_i)^2 \right] = 0.$$

Ze vztahů (2.1) a (2.3) dostaneme po úpravách

$$(2.5) \quad \frac{1}{n} \sum_{i=1}^n \hat{\xi}_i = \bar{x},$$

$$\hat{\alpha} = \bar{y} - \beta \bar{x},$$

$$(2.6) \quad \hat{\xi}_i = x_i + \frac{\beta}{1 + \beta^2} ((y_i - \bar{y}) - (\beta(x_i - \bar{x}))), \quad 1 \leq i \leq n,$$

$$(2.7) \quad \hat{\xi}_i = \bar{x} + \frac{1}{1 + \beta^2} ((x_i - \bar{x}) + (\beta(y_i - \bar{y}))), \quad 1 \leq i \leq n.$$

Když odtud dosadíme do (2.2), pak s využitím (2.1) dojdeme ke kvadratické rovnici pro  $\beta$ :

$$(2.8) \quad S_{xy}\beta^2 - (S_{yy} - S_{xx})\beta - S_{xy} = 0.$$

Lze dokázat, že věrohodnostní funkci minimalizuje pouze jeden kořen této kvadratické rovnice (předpokládáme  $S_{xy} \neq 0$ ):

$$(2.9) \quad \hat{\beta} = \frac{(S_{yy} - S_{xx}) + \sqrt{(S_{yy} - S_{xx})^2 + 4S_{xy}^2}}{2S_{xy}}.$$

Při dalších úpravách je užitečné vědět, že z (2.8) plyne vztah

$$(2.10) \quad S_{yy} - 2\hat{\beta}S_{xy} + \hat{\beta}^2 S_{xx} = (1 + \hat{\beta}^2)(S_{yy} - \hat{\beta}S_{xx}).$$

Odtud je mimo jiné

$$\begin{aligned}
 (2.11) \quad RSS_y &= \sum_{i=1}^n (y_i - \hat{\alpha} - \hat{\beta}\hat{\xi}_i)^2 \\
 &= \frac{1}{1 + \hat{\beta}^2} (S_{yy} - \hat{\beta}S_{xy}),
 \end{aligned}$$

$$\begin{aligned}
 (2.12) \quad RSS_x &= \sum_{i=1}^n (x_i - \hat{\xi}_i)^2 \\
 &= \frac{\hat{\beta}^2}{1 + \hat{\beta}^2} (S_{yy} - \hat{\beta}S_{xy}),
 \end{aligned}$$

takže je dohromady

$$\begin{aligned}
 (2.13) \quad RSS_{x+y} &= RSS_x + RSS_y \\
 &= S_{yy} - \hat{\beta}S_{xy}.
 \end{aligned}$$

V této souvislosti je zajímavé si připomenout, že u klasické regresní přímky platí

$$RSS = S_{yy} - bS_{xy}.$$

Dosadíme-li z (2.13) do (2.4), dostaneme odhad parametru  $\sigma^2$ :

$$\begin{aligned}
 s^2 &= \frac{RSS_{x+y}}{2n} \\
 &= \frac{s_{yy} - \hat{\beta}s_{xy}}{2}.
 \end{aligned}$$

Uvážíme-li, že za předpokladů uvedených v úvodu konvergují veličiny  $s_{xx}, s_{xy}, s_{yy}$  podle pravděpodobnosti k prvkům varianční matice, totiž po řadě k výrazům

$$\omega^2 + \sigma^2, \quad \beta\omega^2, \quad \beta^2\omega^2 + \sigma^2,$$

zjistíme konzistenci odhadů  $\hat{\alpha}, \hat{\beta}$ . Ovšem dále zjistíme, že je

$$\begin{aligned}
 \text{plim } s^2 &= \text{plim } \frac{s_{yy} - \hat{\beta}s_{xy}}{2} \\
 &= \frac{\beta^2\omega^2 + \sigma^2 - \beta\beta\omega^2}{2} \\
 &= \frac{\sigma^2}{2}.
 \end{aligned}$$

Konzistentním odhadem parametru  $\sigma^2$  tedy bude nikoliv  $s^2$ , ale

$$\begin{aligned}
 (2.14) \quad \hat{\sigma}^2 &= \frac{RSS_{x+y}}{n} \\
 &= s_{yy} - \hat{\beta}s_{xy}
 \end{aligned}$$

Nyní nalezneme výrazy pro  $R^2$  podle jednotlivých vztahů z 1. kapitoly.

2.1. Protože je podle (1.5) průměrná hodnota odhadů  $\hat{\xi}_i$  totožná s  $\bar{x}$ , platí nutně  $\bar{\hat{y}} = \bar{y}$ .  
Dále je

$$\begin{aligned} \sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{y}) &= \sum_{i=1}^n (y_i - \bar{y})(\hat{\alpha} + \hat{\beta}\hat{\xi}_i - \bar{y}) \\ &= \sum_{i=1}^n (y_i - \bar{y})\hat{\beta}(\hat{\xi}_i - \bar{x}) \\ &= \frac{\hat{\beta}}{1 + \hat{\beta}^2} \sum_{i=1}^n (y_i - \bar{y}) \left( (x_i - \bar{x}) + \hat{\beta}(y_i - \bar{y}) \right) \\ &= \frac{\hat{\beta}}{1 + \hat{\beta}^2} (S_{xy} + \hat{\beta}S_{yy}), \end{aligned}$$

takže po dosazení do definice (1.1) dostaneme po malé úpravě

$$R_{1y}^2 = \frac{(S_{xy} + \hat{\beta}S_{yy})^2}{S_{yy}(S_{xx} + 2\hat{\beta}S_{xy} + \hat{\beta}^2S_{yy})}.$$

S využitím skutečnosti, že  $\hat{\beta}$  je řešením kvadratické rovnice (2.8), lze poslední vztah zjednodušit na

$$(2.15) \quad R_{1y}^2 = \frac{\hat{\beta}}{1 + \hat{\beta}^2} \frac{\hat{\beta}s_{yy} + s_{xy}}{s_{yy}}.$$

2.2. Do vztahu (1.4) snadno dosadíme místo reziduálního součtu čtverců výraz  $RSS_y$ , což vede k

$$\begin{aligned} R_{2y}^2 &= 1 - \frac{RSS_y}{S_{yy}}, \\ &= 1 - \frac{s_{yy} - \hat{\beta}s_{xy}}{(1 + \hat{\beta}^2)s_{yy}}, \\ &= R_{1y}^2. \end{aligned}$$

Symbole  $y$  v indexech již naznačily následující "kacířskou" myšlenku: místo vysvětlování pouze hodnot veličiny  $Y$  se pokusit vysvětlit variabilitu obou měřených veličin, protože obě měříme s náhodnou chybou. To vede k výrazu

$$(2.16) \quad \begin{aligned} R_{2x+y}^2 &= 1 - \frac{RSS_{x+y}}{S_{xx} + S_{yy}}, \\ &= 1 - \frac{S_{yy} - \hat{\beta}S_{xy}}{S_{xx} + S_{yy}}, \\ &= \frac{s_{xx} + \hat{\beta}s_{xy}}{s_{xx} + s_{yy}}. \end{aligned}$$

Jako cvičení nechávám čtenáři ověření tvrzení, že stejný výraz dostaneme, když vyjdem z definice (1.1).

Je zajímavé porovnat uvedené dvě verze koeficientu determinace  $R_2^2$ . Po snadné úpravě, s využitím skutečnosti, že s pravděpodobností 1 je  $RSS_{x+y} > 0$ , zjistíme, že s pravděpodobností 1 platí ekvivalence

$$(2.17) \quad R_{2y}^2 \leq R_{2x+y}^2 \Leftrightarrow \hat{\beta}^2 \leq \frac{s_{xz}}{s_{yy}}.$$

**2.3.** Odhady parametrů metodou maximální věrohodnosti jsou uvedeny v úvodu této kapitoly. Dosadíme-li tyto odhady do logaritmu věrohodnostní funkce, dostaneme

$$\hat{l} = -n \ln(s_{yy} - \hat{\beta}s_{zy}).$$

Za platnosti hypotézy  $H_0 : \beta = 0$  vyjde zřejmě

$$\tilde{\alpha} = \bar{y}, \quad \tilde{\xi}_i = x_i, \quad 1 \leq i \leq n.$$

Podobně jako dříve dostaneme, že odhad

$$\tilde{\sigma}^2 = s_{yy}$$

je konzistentním odhadem rozptylu  $\sigma^2$ . Je tedy (s uvážením, že tentokrát máme k dispozici  $2n$  pozorování)

$$(2.18) \quad \begin{aligned} R_3^2 &= 1 - \exp\left(-\frac{1}{n}(\hat{l} - \tilde{l})\right), \\ &= 1 - \frac{s_{yy} - \hat{\beta}s_{zy}}{s_{yy}}, \\ &= \hat{\beta} \frac{s_{zy}}{s_{yy}}. \end{aligned}$$

V této souvislosti je zajímavé, že v klasické situaci lze koeficient determinace (1.2) upravit do podobného tvaru

$$R^2 = b \frac{s_{zy}}{s_{yy}}.$$

Protože je  $R^2 \geq 0$ , platí ekvivalence

$$R_3^2 \geq r_{z,y}^2 \Leftrightarrow |\hat{\beta}| \geq |b|.$$

Lze ovšem snadno dokázat, že nerovnost na pravé straně platí vždy, kdy je  $s_{xz} > 0$ . Můžeme tedy s pravděpodobností 1 tvrdit, že je

$$(2.19) \quad R_3^2 \geq r_{z,y}^2.$$

S využitím skutečnosti, že s pravděpodobností 1 musí být  $s_{yy} - \hat{\beta}s_{zy} > 0$ , snadno zjistíme, že je

$$(2.20) \quad R_3^2 \leq R_{2y}^2,$$

$$(2.21) \quad R_3^2 \leq R_{2x+y}^2.$$



2.4. Abychom spočítali Fisherovu informační matici  $F$ , vyjádříme nejprve jednotlivé derivace logaritmu věrohodnostní funkce pomocí náhodných veličin

$$\begin{pmatrix} e_i \\ f_i \end{pmatrix} = \begin{pmatrix} y_i - \alpha - \beta \xi_i \\ x_i - \xi_i \end{pmatrix} \sim N \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \sigma^2 \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right), \quad 1 \leq i \leq n.$$

Snadno tak dostaneme

$$\begin{aligned} \frac{\partial l}{\partial \alpha} &= \frac{1}{\sigma^2} \sum_{i=1}^n e_i, \\ \frac{\partial l}{\partial \beta} &= \frac{1}{\sigma^2} \sum_{i=1}^n \xi_i e_i, \\ \frac{\partial l}{\partial \xi_i} &= \frac{1}{\sigma^2} (f_i + e_i), \quad 1 \leq i \leq n, \\ \frac{\partial l}{\partial \sigma^2} &= -\frac{1}{\sigma^2} + \frac{1}{2\sigma^2} \left( \sum_{i=1}^n (e_i^2 + f_i^2) \right). \end{aligned}$$

Odtud postupně vyjde

$$\begin{aligned} E \left( \frac{\partial l}{\partial \alpha} \right)^2 &= \frac{n}{\sigma^2}, \\ E \left( \frac{\partial l}{\partial \alpha} \right) \left( \frac{\partial l}{\partial \beta} \right) &= \frac{\sum_{i=1}^n \xi_i}{\sigma^2}, \\ E \left( \frac{\partial l}{\partial \alpha} \right) \left( \frac{\partial l}{\partial \xi_i} \right) &= \frac{\beta}{\sigma^2}, \quad 1 \leq i \leq n, \\ E \left( \frac{\partial l}{\partial \alpha} \right) \left( \frac{\partial l}{\partial \sigma^2} \right) &= 0, \\ E \left( \frac{\partial l}{\partial \beta} \right)^2 &= \frac{\sum_{i=1}^n \xi_i^2}{\sigma^2}, \\ E \left( \frac{\partial l}{\partial \beta} \right) \left( \frac{\partial l}{\partial \xi_i} \right) &= \frac{\beta \xi_i}{\sigma^2}, \quad 1 \leq i \leq n, \\ E \left( \frac{\partial l}{\partial \beta} \right) \left( \frac{\partial l}{\partial \sigma^2} \right) &= 0, \\ E \left( \frac{\partial l}{\partial \xi_i} \right) \left( \frac{\partial l}{\partial \xi_j} \right) &= \frac{1 + \beta^2}{\sigma^2} \delta_{ij}, \quad 1 \leq i, j \leq n, \\ E \left( \frac{\partial l}{\partial \xi_i} \right) \left( \frac{\partial l}{\partial \sigma^2} \right) &= 0, \quad 1 \leq i \leq n, \\ E \left( \frac{\partial l}{\partial \sigma^2} \right)^2 &= \frac{n}{\sigma^4}. \end{aligned}$$

Informační matici vztahenou k vektoru parametrů  $(\alpha, \beta, \xi_1, \dots, \xi_n, \sigma^2)$  tedy lze zapsat jako

$$(2.22) \quad F = \frac{1}{\sigma^2} \begin{pmatrix} n & \sum \xi_i & \beta \mathbf{1}' & 0 \\ \sum \xi_i & \sum \xi_i^2 & \beta \xi' & 0 \\ \mathbf{1} & \beta \xi & (1 + \beta^2) \mathbf{I} & 0 \\ 0 & 0 & \mathbf{0}' & n/\sigma^2 \end{pmatrix}.$$

Protože se hypotéza vztahuje k jedinému parametru  $\beta$ , stačí, když nalezneme prvek na místě (2, 2) matice  $F^{-1}$ . Označme levý horní roh této inverzní matice rozměru  $2 \times 2$  symbolem  $Q$ . Vzhledem k speciálnímu tvaru invertované matice platí podle často používaného vzorce pro inverzi matice rozdělené na 4 podmatice (viz např. [5], str. 33)

$$\begin{aligned} \sigma^2 Q^{-1} &= \begin{pmatrix} n & \sum \xi_i \\ \sum \xi_i & \sum \xi_i^2 \end{pmatrix} - \frac{\beta^2}{1 + \beta^2} \begin{pmatrix} n & \sum \xi_i \\ \sum \xi_i & \sum \xi_i^2 \end{pmatrix}, \\ &= \frac{1}{1 + \beta^2} \begin{pmatrix} n & \sum \xi_i \\ \sum \xi_i & \sum \xi_i^2 \end{pmatrix}. \end{aligned}$$

Odtud je hledaný prvek inverze Fisherovy informační matice roven

$$\begin{aligned} q_{22} &= \sigma^2 (1 + \beta^2) \frac{n}{\det Q}, \\ &= \sigma^2 \frac{1 + \beta^2}{\sum (\xi_i - \bar{\xi})^2}. \end{aligned}$$

Opět si připomeňme klasickou regresní přímku, kde platí

$$\text{var}(b) = \sigma^2 \frac{1}{\sum (x_i - \bar{x})^2}.$$

Z (2.7) plyne, že je

$$\sum_{i=1}^n (\hat{\xi}_i - \bar{\xi})^2 = \frac{S_{xx} + 2\hat{\beta}S_{xy} + \hat{\beta}^2 S_{yy}}{(1 + \hat{\beta}^2)^2}.$$

Waldova statistika  $W$  je tedy rovna

$$(2.23) \quad \begin{aligned} W &= \frac{\hat{\beta}^2}{\text{var}(\hat{\beta})}, \\ &= \frac{\hat{\beta}^2}{\sigma^2 (1 + \hat{\beta}^2)^3} (S_{xx} + 2\hat{\beta}S_{xy} + \hat{\beta}^2 S_{yy}). \end{aligned}$$

Zbývá zapsat koeficient determinace podle (1.6), kam dosadíme skutečně  $n$  (jmenovatel z použitého konzistentního odhadu rozptylu  $\sigma^2$ ) a nikoliv počet pozorování  $2n$ . Po úpravách vyjde

$$(2.24) \quad \begin{aligned} R_4^2 &= \frac{W}{W + n} \\ &= \frac{\hat{\beta}^2}{1 + \hat{\beta}^2 (1 + \hat{\beta}^2 + \hat{\beta}^4) (s_{yy} - \hat{\beta} s_{xy}) + \hat{\beta}^2 (s_{xx} + s_{yy})}. \end{aligned}$$

2.5. Za platnosti hypotézy dostáváme konzistentní odhady ve tvaru

$$\tilde{\alpha} = \bar{y}, \quad \tilde{\beta} = 0, \quad \tilde{\xi} = x, \quad \tilde{\sigma}^2 = s_{yy},$$

což vede k

$$v(\tilde{\theta}) = \frac{1}{s_{yy}} \begin{pmatrix} 0 \\ S_{xy} \\ 0 \\ 0 \end{pmatrix},$$

$$F(\tilde{\theta})^{-1} = \frac{1}{s_{yy}} \begin{pmatrix} * & * & 0' & 0 \\ * & S_{xx}^{-1} & 0' & 0 \\ 0 & 0 & I & 0 \\ 0 & 0 & 0' & s_{yy}/n \end{pmatrix},$$

kde jsem označili hvězdičkou ty výrazy, které v dalším výpočtu nepotřebujeme. Po úpravě tedy dostaneme

$$S = n \frac{S_{xy}^2}{S_{xx} S_{yy}},$$

což po dosazení do (1.7) dá možná poněkud překvapivý výsledek

$$(2.25) \quad R_3^2 = \frac{S_{xy}^2}{S_{xx} S_{yy}} = r_{xy}^2.$$

### 3. PŘÍKLAD

Numerickou ilustraci provedeme na známých datech o koncentraci kyseliny mléčné v krvi u matek ( $x$ ) a u jejich novorozenců ( $y$ ) (viz [1]). Na obrázku je vyznačeno výchozích 6 párů pozorování, souvislou čarou je vyznačena standardní regresní přímka pro závislost  $y$  na  $x$ , čárkovaně je vyznačena přímka spočítaná podle 2. kapitoly i s odhadnutými body o souřadnicích  $(\hat{\xi}_i, \hat{y}_i)$  na této přímce. Na silnou závislost ukazují hodnoty koeficientů spočítaných podle (2.15), (2.16), (2.18), (2.24) a (2.25):

$$R_{1y}^2 = 0,961,$$

$$R_{2x+y}^2 = 0,968,$$

$$R_3^2 = 0,929,$$

$$R_4^2 = 0,881,$$

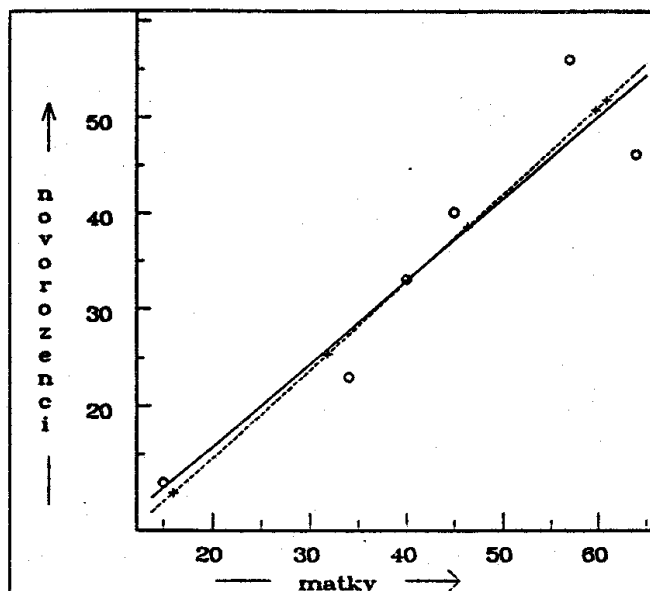
$$R_5^2 = 0,874.$$

Pro zajímavost, statistiky  $W$ ,  $S$  a  $2(\hat{l} - \bar{l})$ , které asymptoticky mají za platnosti hypotézy stejné  $\xi^2$  rozdělení s jedním stupněm volnosti dají

$$W = 44,5,$$

$$S = 5,2$$

$$2(\hat{l} - \bar{l}) = 15,9.$$



OBR. 1. Koncentrace kyseliny mléčné [mg/ml]

Článek vznikl za finanční podpory grantu číslo 313/93/0616 Grantové agentury České republiky.

#### REFERENCES

1. Anděl J., *Matematická statistika*, SNTL, Praha, 1978.
2. Kvålseth T.O., *Cautionary note about  $R^2$* , *The American Statistician* 39 (1985), 279–285.
3. Magge L.,  *$R^2$  measures based on Wald and likelihood ratio joint significance tests*, *The American Statistician* 44 (1990), 251–253.
4. Rao, C.R., *Lineární metody statistické indukce a jejich aplikace*, Academia, Praha, 1978.
5. Zvára K., *Regresní analýza*, Academia, Praha, 1989.