

# Matematická statistika 2

Mirko Navara  
Centrum strojového vnímání  
katedra kybernetiky FEL ČVUT  
Karlovo náměstí, budova G, místnost 104a  
<http://cmp.felk.cvut.cz/~navara>

6. listopadu 2012

## Obsah

<b>1</b>	<b>Podmíněná rozdělení pravděpodobnosti</b>	<b>2</b>
1.1	Motivační úloha . . . . .	2
1.2	Opakování podmíněné pravděpodobnosti . . . . .	3
1.3	Rozdělení podmíněné <b>jevem</b> . . . . .	3
1.4	Rozdělení podmíněné <b>hodnotou diskrétní náhodné veličiny</b> . . . . .	4
1.5	Rozdělení podmíněné <b>hodnotou spojitě náhodné veličiny</b> . . . . .	4
1.6	Vlastnosti podmíněných rozdělení . . . . .	5
1.7	Příklady podmíněných rozdělení . . . . .	5
<b>2</b>	<b>Regrese</b>	<b>6</b>
2.1	Odhad konstanty . . . . .	7
2.2	Odhad konstanty metodou maximální věrohodnosti . . . . .	7
2.3	Řešení bez předpokladu normality . . . . .	8
2.4	Odhad funkce: obecný nelineární model . . . . .	8
2.5	Odhad náhodné veličiny: lineární model dimenze 1 . . . . .	9
2.6	Řešení metodou maximální věrohodnosti: lineární model dimenze 1 . . . . .	9
2.7	Regresní přímka 1 . . . . .	10
2.8	Regresní přímka 2 . . . . .	10
2.9	Interpretace regresních koeficientů . . . . .	11
2.10	Chyba lineární regrese . . . . .	11
2.11	Složky rozptylu regresního odhadu . . . . .	12
2.12	Odhad rozptylu . . . . .	12
2.13	Rozdělení odhadů a testy hypotéz o nich . . . . .	13
2.14	Testy korelace . . . . .	13
2.15	Příklad použití . . . . .	14
2.15.1	Odhady koeficientů 1. regresní přímky . . . . .	15
2.15.2	Odhady koeficientů 2. regresní přímky . . . . .	15
2.15.3	Odhady rozptylů, kovariance a korelace . . . . .	16
2.15.4	Chyba lineární regrese . . . . .	16
2.15.5	Složky rozptylu . . . . .	17

2.15.6	Odhady rozptylu původního rozdělení . . . . .	17
2.15.7	Testy korelace . . . . .	17
2.16	Odhad náhodné veličiny: lineární model dimenze $k$ . . . . .	18
2.17	Intervalové odhady regresních koeficientů při <b>známém</b> rozptylu . . . . .	20
2.18	Odhady rozptylu $\sigma^2$ původního rozdělení . . . . .	20
2.19	Intervalové odhady regresních koeficientů při <b>neznámém</b> rozptylu . . . . .	21

<b>3</b>	<b>Volba vysvětlujících proměnných</b>	<b>21</b>
3.1	Kritéria pro výběr modelu . . . . .	22
3.2	Volba vysvětlujících proměnných . . . . .	22
3.3	Rozdělení na trénovací a testovací data, cross-validation . . . . .	22

# 1 Podmíněná rozdělení pravděpodobnosti

## 1.1 Motivační úloha

Příklad (počet es):

- Hráč má 10 z balíčku 32 karet, v němž jsou 4 esa. Náhodná veličina  $Y$  je počet es v jeho ruce. Pravděpodobnost, že má  $y$  es, je

$$p_Y(y) = \frac{\binom{4}{y} \binom{32-4}{10-y}}{\binom{32}{10}}, \quad y \in \{0, 1, 2, 3, 4\},$$

(hypergeometrické rozdělení),

$y$	0	1	2	3	4
$p_Y(y)$	0.203	0.428	0.289	0.073	0.006

$$EY = \frac{10}{32} 4 = \frac{5}{4}.$$

- Podívám se na svých 10 karet (ze stejného balíčku) a vidím, že nastal jev  $B$ : mám 2 esa. Tím se pravděpodobnost počtu es prvního hráče mění na

$$p_{Y|B}(y) = \frac{\binom{2}{y} \binom{22-2}{10-y}}{\binom{22}{10}}, \quad y \in \{0, 1, 2\},$$

(vybíráme z 22 karet, mezi nimiž jsou 2 esa),

$y$	0	1	2
$p_{Y B}(y)$	$\frac{2}{7} \doteq 0.286$	$\frac{40}{77} \doteq 0.519$	$\frac{15}{77} \doteq 0.195$

$$E(Y|B) = \frac{10}{22} 2 = \frac{10}{11}.$$

- Obecněji: Náhodná veličina  $X$  je počet es v mé ruce a  $x$  její realizace („skutečná hodnota“), tj. nastal jev  $X = x$ .

Pravděpodobnost počtu es prvního hráče se mění na

$$p_{Y|X=x}(y) = p_{Y|X}(y|x) = \frac{\binom{4-x}{y} \binom{22-4+x}{10-y}}{\binom{22}{10}}, \quad x \in \{0, \dots, 4\}, \quad y \in \{0, \dots, 4-x\},$$

(vybíráme z  $22 - 4 + x$  karet, mezi nimiž je  $4 - x$  es),

$$E(Y|X = x) = \frac{10}{22} (4 - x) = \frac{5}{11} (4 - x), \quad x \in \{0, \dots, 4\}.$$

Tento výsledek popisuje rozdělení náhodné veličiny

$$E(Y|X) = \frac{5}{11} (4 - X),$$

která nabývá hodnot

$$E(Y|X = x) = \frac{5}{11} (4 - x) \in \left\{ \frac{20}{11}, \frac{15}{11}, \frac{10}{11}, \frac{5}{11}, 0 \right\}$$

s pravděpodobnostmi

$$p_X(x) = p_Y(x) = \frac{\binom{4}{x} \binom{32-4}{10-x}}{\binom{32}{10}}, \quad x \in \{0, 1, 2, 3, 4\}.$$

$x$	0	1	2	3	4
$E(Y X = x)$	$\frac{20}{11}$	$\frac{15}{11}$	$\frac{10}{11}$	$\frac{5}{11}$	0

Náhodná veličina  $E(Y|X)$  má střední hodnotu

$$\begin{aligned} E(E(Y|X)) &= \frac{20}{11} p_X(0) + \frac{15}{11} p_X(1) + \frac{10}{11} p_X(2) + \frac{5}{11} p_X(3) \\ &= E\left(\frac{5}{11} (4 - X)\right) = \left(\frac{5}{11} (4 - EX)\right) = \left(\frac{5}{11} \left(4 - \frac{5}{4}\right)\right) = \frac{5}{4} = EY. \end{aligned}$$

## 1.2 Opakování podmíněné pravděpodobnosti

Pravděpodobnost (=pravděpodobnostní míra)  $P$

Značení:  $P(\cdot)$  pravděpodobnostní míra,

$P[\cdot]$  pravděpodobnost jevu popsaného v závorce

(je-li argument jev, obě značení se shodují)

Jev  $B$ ,  $P(B) \neq 0 \neq P(\bar{B})$

$\forall$  jev  $A$ :  $P(A) = c P(A|B) + c' P(A|\bar{B})$ ,

kde  $c = P(B)$ ,  $c' = P(\bar{B})$ , tj.  $c + c' = 1$ .

$P$  je lineární (dokonce **konvexní**) kombinace podmíněných pravděpodobností  $P(\cdot|B)$ , resp.  $P(\cdot|\bar{B})$ .  
což jsou obyčejné pravděpodobnosti popisující jiné modely (kde  $B$  je jev jistý, resp. nemožný).

Nadále předpoklad  $P(B) \neq 0$  budeme potřebovat,

předpoklad  $P(\bar{B}) \neq 0$  nikoli.

## 1.3 Rozdělení podmíněné jevem

**Podmíněná pravděpodobnostní funkce**  $p_{Y|B}$  *diskrétní* náhodné veličiny  $Y$  za podmínky  $B$ :

$$p_{Y|B}(y) = P[Y = y|B] = \frac{P[Y = y, B]}{P[B]},$$

kde  $P[Y = y, B] = P(\{\omega \in \Omega \mid Y(\omega) = y, \omega \in B\}) = P(\{\omega \in B \mid Y(\omega) = y\})$ .

**Podmíněná distribuční funkce**  $F_{Y|B}$  náhodné veličiny  $Y$  za podmínky  $B$ :

$$F_{Y|B}(y) = P[Y \leq y|B] = \frac{P[Y \leq y, B]}{P[B]},$$

kde  $P[Y \leq y, B] = P(\{\omega \in B \mid Y(\omega) \leq y\})$ .

Pro daný jev  $B$  má všechny vlastnosti distribuční funkce.

Definována, i když  $Y$  není diskrétní

⇒ **podmíněná kvantilová funkce**  $q_{Y|B}$

⇒ **podmíněná hustota**  $f_{Y|B}$  *spojité* náhodné veličiny

$$F_{Y|B}(y) = \int_0^y f_{Y|B}(t) dt$$

⇒ **podmíněná střední hodnota**  $E(Y|B)$

⇒ **podmíněný rozptyl**  $D(Y|B)$

## 1.4 Rozdělení podmíněné hodnotou diskrétní náhodné veličiny

*Diskrétní* náhodnou veličinu  $Y$  lze podmiňovat jevem  $X = x$ , pokud  $p_X(x) = P[X = x] \neq 0$ ; pak

$$p_{Y,X}(y, x) = P[Y = y, X = x] = P[Y = y \mid X = x] \cdot P[X = x] = p_{Y|X}(y|x) \cdot p_X(x),$$

$$\begin{aligned} F_{Y,X}(y, x) &= \sum_{t:t \leq x} \sum_{u:u \leq y} p_{Y,X}(u, t) = \sum_{t:t \leq x} \sum_{u:u \leq y} p_{Y|X}(u|t) \cdot p_X(t) \\ &= \sum_{t:t \leq x} F_{Y|X}(y|t) \cdot p_X(t). \end{aligned}$$

Obecněji (pro *libovolnou* náhodnou veličinu  $Y$ )

$$F_{Y,X}(y, x) = \sum_{t:t \leq x} F_{Y|X}(y|t) \cdot p_X(t).$$

## 1.5 Rozdělení podmíněné hodnotou spojité náhodné veličiny

Nelze podmiňovat jevem  $X = x$ , neboť  $P[X = x] = 0$ ; použijeme místo toho *hustotu*  $f_X(x)$  (je-li nenulová) a definujeme analogicky podmíněnou hustotu  $f_{Y|X}: \mathbb{R}^2 \rightarrow \langle 0, \infty \rangle$  jako libovolnou funkci splňující

$$f_{Y,X}(y, x) = f_{Y|X}(y|x) \cdot f_X(x),$$

přesněji a obecněji

$$F_{Y,X}(y, x) = \int_{-\infty}^x F_{Y|X}(y|t) \cdot f_X(t) dt = \int_{-\infty}^x \int_{-\infty}^y f_{Y|X}(u|t) \cdot f_X(t) du dt.$$

## 1.6 Vlastnosti podmíněných rozdělání

Můžee existovat

**podmíněná střední hodnota**  $E(Y|X = x)$

resp. **podmíněný rozptyl**  $D(Y|X = x)$

jako *funkce*  $E(Y|X)$ , resp.  $D(Y|X)$  proměnné  $x$ .

Pokud argument  $x$  neznáme, ale nahradíme náhodnou veličinou  $X$ , pak podmíněná střední hodnota, resp. podmíněný rozptyl je náhodná veličina a může mít obvyklé charakteristiky, např.  $E(E(Y|X)), D(E(Y|X)), E(D(Y|X))...$  Ty, pokud existují, splňují rovnosti

$$E(E(Y|X)) = EY,$$

$$DY = D(E(Y|X)) + E(D(Y|X)).$$

## 1.7 Příklady podmíněných rozdělání

**Příklad:** První hráč hodí čtyřmi stejnými mincemi; ty, na kterých padne líc, si ponechá jako výhru (náhodná veličina  $X$ ). Poté druhý hráč hodí zbylými mincemi; ty, na kterých padne líc, tvoří jeho výhru (náhodná veličina  $Y$ ). Popište sdružené rozdělání  $X, Y$ , marginální rozdělání  $X$  a  $Y$ , podmíněná rozdělání  $Y|X$  a  $X|Y$ . U všech rozdělání stanovte (podmíněnou) střední hodnotu a (podmíněný) rozptyl.

*Řešení:* (Čísła za tabulkami jsou pravděpodobnostní funkce marginálních rozdělání, což v případě sdruženého rozdělání jsou řádkové, resp. sloupcové součty.)

$p_{X,Y}$	$y$	0	1	2	3	4	$p_X(x)$
$x$	0	$\frac{1}{256}$	$\frac{1}{64}$	$\frac{3}{128}$	$\frac{1}{64}$	$\frac{1}{256}$	$\frac{1}{16}$
	1	$\frac{1}{32}$	$\frac{3}{32}$	$\frac{3}{32}$	$\frac{1}{32}$	0	$\frac{4}{32}$
	2	$\frac{3}{32}$	$\frac{3}{16}$	$\frac{3}{32}$	0	0	$\frac{8}{32}$
	3	$\frac{1}{8}$	$\frac{1}{8}$	0	0	0	$\frac{4}{8}$
	4	$\frac{1}{16}$	0	0	0	0	$\frac{4}{16}$
$p_Y(y)$		$\frac{8}{256}$	$\frac{27}{64}$	$\frac{27}{128}$	$\frac{3}{64}$	$\frac{1}{256}$	

$$EX = 2, \quad EY = \frac{27}{64} \cdot 1 + \frac{27}{128} \cdot 2 + \frac{3}{64} \cdot 3 + \frac{1}{256} \cdot 4 = 1$$

$p_{Y X}$	$y$	0	1	2	3	4	$p_X(x)$	$E(Y X = x)$
$x$	0	$\frac{1}{16}$	$\frac{1}{4}$	$\frac{3}{8}$	$\frac{1}{4}$	$\frac{1}{16}$	$\frac{1}{16}$	2
	1	$\frac{1}{8}$	$\frac{3}{8}$	$\frac{3}{8}$	$\frac{1}{8}$	0	$\frac{4}{8}$	$\frac{3}{2}$
	2	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$	0	0	$\frac{3}{4}$	1
	3	$\frac{1}{2}$	$\frac{1}{2}$	0	0	0	$\frac{4}{2}$	$\frac{1}{2}$
	4	1	0	0	0	0	$\frac{4}{1}$	0
$p_Y(y)$		$\frac{8}{256}$	$\frac{27}{64}$	$\frac{27}{128}$	$\frac{3}{64}$	$\frac{1}{256}$		

$$E(E(Y|X)) = \frac{1}{16} \cdot 2 + \frac{1}{4} \cdot \frac{3}{2} + \frac{3}{8} \cdot 1 + \frac{1}{4} \cdot \frac{1}{2} = 1 = EY$$

$p_{X Y}$	$y$	0	1	2	3	4	$p_X(x)$
$x$							
0		$\frac{1}{81}$	$\frac{1}{27}$	$\frac{1}{49}$	$\frac{1}{33}$	1	$\frac{1}{16}$
1		$\frac{8}{81}$	$\frac{12}{27}$	$\frac{14}{49}$	$\frac{13}{33}$	0	$\frac{1}{4}$
2		$\frac{8}{81}$	$\frac{14}{49}$	$\frac{9}{9}$	0	0	$\frac{3}{8}$
3		$\frac{27}{32}$	$\frac{8}{27}$	0	0	0	$\frac{1}{4}$
4		$\frac{81}{16}$	0	0	0	0	$\frac{1}{16}$
$p_Y(y)$		$\frac{81}{256}$	$\frac{27}{64}$	$\frac{27}{128}$	$\frac{3}{64}$	$\frac{1}{256}$	
$E(X Y=y)$		$\frac{8}{3}$	2	$\frac{4}{3}$	$\frac{2}{3}$	0	

$$E(E(X|Y)) = \frac{81}{256} \cdot \frac{8}{3} + \frac{27}{64} \cdot 2 + \frac{27}{128} \cdot \frac{4}{3} + \frac{3}{64} \cdot \frac{2}{3} = 2 = EX$$

**Příklad** [Wassermann 3.8.22, rozšířeno]: Necht'  $0 < a < b < 1$ . Náhodná veličina  $X$  má spojité rovnoměrné rozdělení v  $\langle 0, 1 \rangle$ , pomocí ní jsou definovány náhodné veličiny  $Y, Z$ :

$$Y = \begin{cases} 1 & \text{pro } X < b, \\ 0 & \text{jinak,} \end{cases}$$

$$Z = \begin{cases} 1 & \text{pro } X > a, \\ 0 & \text{jinak.} \end{cases}$$

Jsou veličiny  $Y, Z$  nezávislé? (Zdůvodněte.)

Najděte  $E(Y|Z), D(Y|Z), E(E(Y|Z)), E(D(Y|Z)), D(E(Y|Z))$  a ověřte rovnost  $DY = D(E(Y|Z)) + E(D(Y|Z))$ .

*Řešení:*

$p_{Y,Z}$	$y$	0	1	$p_Z(z)$
$z$				
0		0	$a$	$a$
1		$1-b$	$b-a$	$1-a$
$p_Y(y)$		$1-b$	$b$	

$p_{Y Z}$	$y$	0	1	$p_Z(z)$	$E(Y Z=z)$	$D(Y Z=z)$
$z$						
0		0	1	$a$	1	0
1		$\frac{1-b}{1-a}$	$\frac{b-a}{1-a}$	$1-a$	$\frac{b-a}{1-a}$	$\frac{(1-b)(b-a)}{(1-a)^2}$

$$E(E(Y|Z)) = 1 \cdot a + \frac{b-a}{1-a} \cdot (1-a) = b = EY$$

$$E(D(Y|Z)) = 0 \cdot a + \frac{(1-b)(b-a)}{(1-a)^2} \cdot (1-a) = \frac{(1-b)(b-a)}{1-a}$$

$$D(E(Y|Z)) = E((E(Y|Z))^2) - \underbrace{(E(E(Y|Z)))^2}_{(EY)^2} = 1^2 \cdot a + \left(\frac{b-a}{1-a}\right)^2 \cdot (1-a) - b^2 = \frac{a(1-b)^2}{1-a}$$

$$D(E(Y|Z)) + E(D(Y|Z)) = \frac{(1-b)(b-a)}{1-a} + \frac{a(1-b)^2}{1-a} = b(1-b) = DY$$

## 2 Regrese

neboli náhrada pozorované závislosti vhodnou funkcí.

## 2.1 Odhad konstanty

**Úloha 1:** Odhadujeme spojitou náhodnou veličinu  $Y$ , předpokládáme

$$Y = \vartheta + \mathcal{E},$$

kde  $\vartheta \in \mathbb{R}$  je neznámý parametr,

$\mathcal{E}$  je náhodná veličina (chyba, šum) s rozdělením  $N(0, \sigma^2)$ ,

$\sigma^2$  je (konstantní, známý nebo neznámý) rozptyl,

$$f_{\mathcal{E}}(t) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{t^2}{2\sigma^2}\right),$$
$$f_Y(t) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(t - \vartheta)^2}{2\sigma^2}\right).$$

**Vstupy:** posloupnost realizací vzniklých nezávislými pokusy (realizace náhodného výběru z rozdělení n. v.  $Y$ )

$$\mathbf{y} = (y_1, \dots, y_n).$$

Předpokládáme

$$y_j = \vartheta + e_j,$$

kde  $(e_1, \dots, e_n)$  je realizace náhodného výběru z rozdělení  $N(0, \sigma^2)$ .

## 2.2 Odhad konstanty metodou maximální věrohodnosti

$$\Lambda(\vartheta) = \prod_j f_Y(y_j) = \prod_j \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(y_j - \vartheta)^2}{2\sigma^2}\right),$$
$$\lambda(\vartheta) = \ln \Lambda(\vartheta) = \sum_j \ln f_Y(y_j) = \sum_j \left(-\ln \sigma \sqrt{2\pi} - \frac{(y_j - \vartheta)^2}{2\sigma^2}\right) =$$
$$= \underbrace{-n \ln \sigma \sqrt{2\pi}}_{\text{konst.}} - \underbrace{\frac{1}{2\sigma^2}}_{\text{konst.}} \underbrace{\sum_j (y_j - \vartheta)^2}_{\kappa(\vartheta)},$$

$$\hat{\vartheta} = \operatorname{argmax}_{\vartheta} \Lambda(\vartheta) = \operatorname{argmax}_{\vartheta} \lambda(\vartheta) = \operatorname{argmin}_{\vartheta} \kappa(\vartheta) =$$

$$= \operatorname{argmin}_{\vartheta} \sum_j (y_j - \vartheta)^2 = \operatorname{argmin}_{\vartheta} \sum_j e_j^2$$

$\Rightarrow$  metoda nejmenších čtverců

$$0 = \frac{\partial \kappa(\hat{\vartheta})}{\partial \hat{\vartheta}} = \frac{\partial}{\partial \hat{\vartheta}} \left( \sum_j y_j^2 - 2\hat{\vartheta} \sum_j y_j + n\hat{\vartheta}^2 \right) = -2 \sum_j y_j + 2n\hat{\vartheta},$$

$$\hat{\vartheta} = \frac{1}{n} \sum_j y_j = \bar{y}.$$

Realizace výběrového průměru  $\bar{y}$  = odhad metodou nejmenších čtverců = max. věrohodný odhad (za předpokladu normálního rozdělení chyb)

## 2.3 Řešení bez předpokladu normality

Obecněji (pokud  $\mathcal{E}$  má střední hodnotu)

$$EY = \vartheta + E\mathcal{E}.$$

Realizace výběrového průměru  $\bar{y}$  = odhad metodou nejmenších čtverců  $\neq$  max. věrohodný odhad (ten může být vychýlený)

## 2.4 Odhad funkce: obecný nelineární model

**Úloha 2:** Odhadujeme spojitou náhodnou veličinu  $Y$ , předpokládáme

$$Y = g(X) + \mathcal{E},$$

kde  $X$  je nezávislá **vysvětlující** náhodná veličina, jejíž hodnoty můžeme měřit (spojitá nebo diskrétní),

$g: \mathbb{R} \rightarrow \mathbb{R}$  je neznámá funkce (závislá na neznámých parametrech),

$\mathcal{E}$  je náhodná veličina s rozdělením  $N(0, \sigma^2)$ ,

$\sigma^2$  je (konstantní, známý nebo neznámý) rozptyl,

$$f_{Y|X}(y|x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y-g(x))^2}{2\sigma^2}\right).$$

**Vstupy:** realizace náhodného výběru ze sdruženého rozdělení n. vektoru  $(Y, X)$

$$(\mathbf{y}, \mathbf{x}) = ((y_1, x_1), \dots, (y_n, x_n)).$$

Předpokládáme  $y_j = g(x_j) + e_j$ ,

kde  $(e_1, \dots, e_n)$  je realizace náhodného výběru z rozdělení  $N(0, \sigma^2)$ .

**Řešení:** Pokud víme, že  $X = x$ , přejdeme k podmíněným pravděpodobnostem a jejich středním hodnotám:

$$E(Y|X = x) = E(g(X) + \mathcal{E}|X = x) = E(g(X)|X = x) = g(x),$$

Odhad (hodnot) funkce  $g(x)$  = odhad podmíněných středních hodnot  $E(Y|X = x)$ .

**Problém:** Jak je odhadnout z realizace  $((y_1, x_1), \dots, (y_n, x_n))$ ?

**Řešení:** Odhad  $\hat{g}(x)$  = aritmetický průměr těch  $y_j$ , pro která  $x_j = x$ , tj.

$$\hat{g}(x) = \frac{1}{|M(x)|} \sum_{j \in M(x)} y_j, \quad \text{kde } M(x) = \{j \in \{1, \dots, n\} : x_j = x\}.$$

**Problém 1:** Lze použít pouze pro velký rozsah výběru ve srovnání s počtem možných hodnot n. v.  $X$ , rozhodně pouze pro  $X$  diskrétní!

**Problém 2:** I pro rozsáhlá data není postup vhodný, nabývá-li  $X$  mnoha hodnot – přetrénování: do modelu zahrnujeme i statistické chyby jednotlivých výsledků.

**Řešení:** Omezíme se na modely (funkce  $g$ ), které mají mnohem méně parametrů, než je hodnot n. v.  $X$ .

Tím je umožněna též **statistická indukce**, tedy zobecnění výsledků i na ty hodnoty n. v.  $X$ , které dosud nebyly pozorovány.



## 2.5 Odhad náhodné veličiny: lineární model dimenze 1

**Úloha 3:** Odhadujeme spojitou náhodnou veličinu  $Y$ , předpokládáme

$$Y = \vartheta_0 + \vartheta_1 X + \mathcal{E},$$

kde  $\vartheta_0, \vartheta_1 \in \mathbb{R}$  jsou neznámé parametry (**regresní koeficienty**),

$\mathcal{E}$  je náhodná veličina s rozdělením  $N(0, \sigma^2)$ ,

$\sigma^2$  je (konstantní, známý nebo neznámý) rozptyl,

$X$  je nezávislá **vysvětlující** náhodná veličina,

$$f_{Y|X}(y|x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(y - \vartheta_0 - \vartheta_1 x)^2}{2\sigma^2}\right).$$

**Vstupy:** realizace náhodného výběru ze sdruženého rozdělení n. vektoru  $(Y, X)$

$$(\mathbf{y}, \mathbf{x}) = ((y_1, x_1), \dots, (y_n, x_n)).$$

Předpokládáme  $y_j = \vartheta_0 + \vartheta_1 x_j + e_j$ , kde  $(e_1, \dots, e_n)$  je realizace náhodného výběru z rozdělení  $N(0, \sigma^2)$ .

## 2.6 Řešení metodou maximální věrohodnosti: lineární model dimenze 1

Při známém rozptylu  $\sigma^2$ :

$$\Lambda(\vartheta_0, \vartheta_1) = \prod_j f_{Y|X}(y_j|x_j) = \prod_j \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(y_j - \vartheta_0 - \vartheta_1 x_j)^2}{2\sigma^2}\right),$$

$$\lambda(\vartheta_0, \vartheta_1) = \ln \Lambda(\vartheta_0, \vartheta_1) = \underbrace{-n \ln \sigma \sqrt{2\pi}}_{\text{konst.}} - \underbrace{\frac{1}{2\sigma^2}}_{\text{konst.}} \underbrace{\sum_j (y_j - \vartheta_0 - \vartheta_1 x_j)^2}_{\kappa(\vartheta_0, \vartheta_1)},$$

$$\hat{\vartheta} = \operatorname{argmax}_{\vartheta_0, \vartheta_1} \Lambda(\vartheta_0, \vartheta_1) = \operatorname{argmax}_{\vartheta_0, \vartheta_1} \lambda(\vartheta_0, \vartheta_1) = \operatorname{argmin}_{\vartheta_0, \vartheta_1} \kappa(\vartheta_0, \vartheta_1)$$

$$= \operatorname{argmin}_{\vartheta_0, \vartheta_1} \sum_j (y_j - \vartheta_0 - \vartheta_1 x_j)^2 = \operatorname{argmin}_{\vartheta_0, \vartheta_1} \sum_j e_j^2$$

$\Rightarrow$  **metoda nejmenších čtverců**

odhad metodou nejmenších čtverců = max. věrohodný odhad (za předpokladu normálního rozdělení chyb)

$$0 = \frac{\partial \kappa(\hat{\vartheta}_0, \hat{\vartheta}_1)}{\partial \hat{\vartheta}_0} = 2n \hat{\vartheta}_0 - 2 \sum_j y_j + 2 \hat{\vartheta}_1 \sum_j x_j,$$

$$\hat{\vartheta}_0 = \frac{1}{n} \sum_j y_j - \hat{\vartheta}_1 \frac{1}{n} \sum_j x_j,$$

$$\hat{\vartheta}_0 = \bar{y} - \hat{\vartheta}_1 \bar{x}.$$

Dosadíme do

$$\begin{aligned}
 0 &= \frac{\partial \kappa(\hat{\vartheta}_0, \hat{\vartheta}_1)}{\partial \hat{\vartheta}_1} = 2 \hat{\vartheta}_1 \sum_j x_j^2 - 2 \sum_j x_j y_j + 2 \hat{\vartheta}_0 \sum_j x_j = \\
 &= 2 \hat{\vartheta}_1 \sum_j x_j^2 - 2 \sum_j x_j y_j + 2 n \bar{x} \underbrace{(\bar{y} - \hat{\vartheta}_1 \bar{x})}_{\hat{\vartheta}_0}, \\
 \hat{\vartheta}_1 &= \frac{\sum_j x_j y_j - n \bar{x} \bar{y}}{\sum_j x_j^2 - n \bar{x}^2} = \frac{\sum_j (x_j - \bar{x})(y_j - \bar{y})}{\sum_j (x_j - \bar{x})^2}
 \end{aligned}$$

a do odhadů chyb

## 2.7 Regresní přímka 1

Je tvořena body  $(x, y)$ , které splňují rovnici

$$y = \hat{\vartheta}_0 + \hat{\vartheta}_1 x.$$

**Věta:** Bod  $(\bar{x}, \bar{y})$  leží na regresní přímce,

$$\bar{y} = \hat{\vartheta}_0 + \hat{\vartheta}_1 \bar{x}.$$

**Důkaz:** Už jsme odvodili  $\hat{\vartheta}_0 = \bar{y} - \hat{\vartheta}_1 \bar{x}$ .

Odečtením dostáváme rovnici regresní přímky ve tvaru

$$y - \bar{y} = \hat{\vartheta}_1 (x - \bar{x}).$$

Sklon (směrnice) regresní přímky  $\hat{\vartheta}_1$  se nezmění, pokud přičteme konstantu ke všem hodnotám názávisle, resp. závisle proměnné. Mohli jsme od nich např. odečíst realizace výběrových průměrů  $\bar{x}$ , resp.  $\bar{y}$  a zjednodušit si výrazy.

## 2.8 Regresní přímka 2

**Poznámka:** Lze také odhadnout lineární závislost  $X$  na  $Y$  dle modelu

$$X = \vartheta_0^* + \vartheta_1^* Y + \mathcal{E}^*,$$

směrnice regresní přímky vyjde

$$\hat{\vartheta}_1^* = \frac{\sum_j (x_j - \bar{x})(y_j - \bar{y})}{\sum_j (y_j - \bar{y})^2}.$$

Výsledkem je obecně **jiná** regresní přímka, jejíž rovnici lze psát ve tvarech

$$\begin{aligned}
 x &= \hat{\vartheta}_0^* + \hat{\vartheta}_1^* y, \\
 x - \bar{x} &= \hat{\vartheta}_1^* (y - \bar{y}), \\
 y - \bar{y} &= \frac{1}{\hat{\vartheta}_1^*} (x - \bar{x}).
 \end{aligned}$$

## 2.9 Interpretace regresních koeficientů

Označme odhady rozptylů a kovariance:

$$\hat{\sigma}_{\mathbf{x}}^2 := \frac{1}{n} \sum_j (x_j - \bar{x})^2 = \frac{n-1}{n} s_{\mathbf{x}}^2 = \text{D Emp}(\mathbf{x}),$$

$$\hat{\sigma}_{\mathbf{y}}^2 := \frac{1}{n} \sum_j (y_j - \bar{y})^2 = \frac{n-1}{n} s_{\mathbf{y}}^2 = \text{D Emp}(\mathbf{y}),$$

$$c_{\mathbf{x},\mathbf{y}} := \frac{1}{n} \sum_j (x_j - \bar{x})(y_j - \bar{y}) = \text{cov}(\text{Emp}(\mathbf{x}, \mathbf{y})).$$

Odhad korelace = realizace výběrového koeficientu korelace =

$$r_{\mathbf{x},\mathbf{y}} := \frac{\sum_j (x_j - \bar{x})(y_j - \bar{y})}{\sqrt{\sum_j (x_j - \bar{x})^2 \sum_j (y_j - \bar{y})^2}} = \frac{c_{\mathbf{x},\mathbf{y}}}{\hat{\sigma}_{\mathbf{x}} \hat{\sigma}_{\mathbf{y}}} = \varrho(\text{Emp}(\mathbf{x}, \mathbf{y})).$$

Lze psát

$$\hat{\vartheta}_1 = \frac{c_{\mathbf{x},\mathbf{y}}}{\hat{\sigma}_{\mathbf{x}}^2}.$$

Rovnice regresní přímky pro závislost  $Y$  na  $X$ :

$$\begin{aligned} y - \bar{y} &= \hat{\vartheta}_1 (x - \bar{x}), \\ y - \bar{y} &= \frac{c_{\mathbf{x},\mathbf{y}}}{\hat{\sigma}_{\mathbf{x}}^2} (x - \bar{x}), \\ \frac{y - \bar{y}}{\hat{\sigma}_{\mathbf{y}}} &= r_{\mathbf{x},\mathbf{y}} \frac{x - \bar{x}}{\hat{\sigma}_{\mathbf{x}}}, \end{aligned}$$

rovnice regresní přímky pro závislost  $X$  na  $Y$  (to je **jiná** přímka!):

$$\begin{aligned} x - \bar{x} &= \hat{\vartheta}_1^* (y - \bar{y}), \\ x - \bar{x} &= \frac{c_{\mathbf{x},\mathbf{y}}}{\hat{\sigma}_{\mathbf{y}}^2} (y - \bar{y}), \\ \frac{x - \bar{x}}{\hat{\sigma}_{\mathbf{x}}} &= r_{\mathbf{x},\mathbf{y}} \frac{y - \bar{y}}{\hat{\sigma}_{\mathbf{y}}}. \end{aligned}$$

Směrnice obou regresních přímek mají stejná znaménka a součin  $\hat{\vartheta}_1 \hat{\vartheta}_1^* = \frac{c_{\mathbf{x},\mathbf{y}}^2}{\hat{\sigma}_{\mathbf{x}}^2 \hat{\sigma}_{\mathbf{y}}^2} = r_{\mathbf{x},\mathbf{y}}^2$ , takže

$$r_{\mathbf{x},\mathbf{y}} = \sqrt{\hat{\vartheta}_1 \hat{\vartheta}_1^*} \text{sign}(\hat{\vartheta}_1).$$

## 2.10 Chyba lineární regrese

Odhadli jsme lineární regresní funkci  $\hat{g}$ ,

$$\hat{g}(x) := \hat{\vartheta}_0 + \hat{\vartheta}_1 x,$$

pomocí ní hodnoty nezávisle proměnné v jednotlivých realizacích

$$\begin{aligned} \hat{y}_j &:= \hat{g}(x_j) = \hat{\vartheta}_0 + \hat{\vartheta}_1 x_j, \\ \hat{\mathbf{y}} &:= (\hat{y}_1, \dots, \hat{y}_n) \end{aligned}$$

a chyby (**rezidua**)

$$\begin{aligned}\hat{e}_j &:= y_j - \hat{y}_j = y_j - \hat{\vartheta}_0 - \hat{\vartheta}_1 x_j = y_j - \bar{y} - \hat{\vartheta}_1 (x_j - \bar{x}), \\ \hat{\mathbf{e}} &:= (\hat{e}_1, \dots, \hat{e}_n).\end{aligned}$$

**Věta:**  $\frac{1}{n} \sum_j \hat{y}_j = \bar{y}$ .

**Důkaz:**  $\frac{1}{n} \sum_j \hat{y}_j = \frac{1}{n} \sum_j (\hat{\vartheta}_0 + \hat{\vartheta}_1 x_j) = \hat{\vartheta}_0 + \hat{\vartheta}_1 \bar{x} = \bar{y}$ .

## 2.11 Složky rozptylu regresního odhadu

(Někdy se nedělí  $n$ .)

**Celkový rozptyl** (angl. *total variation*)  $:= \hat{\sigma}_{\mathbf{y}}^2 = \frac{1}{n} \sum_j (y_j - \bar{y})^2$ .

**Rozptyl modelu** (angl. *explained variation*)  $:= \hat{\sigma}_{\hat{\mathbf{y}}}^2 = \frac{1}{n} \sum_j (\hat{y}_j - \bar{y})^2$ .

**Reziduální rozptyl** (angl. *unexplained variation*)  $:=$

$$\hat{\sigma}_{\hat{\mathbf{e}}}^2 = \frac{1}{n} \sum_j \hat{e}_j^2 = \frac{1}{n} \sum_j (y_j - \hat{y}_j)^2.$$

**Věta:**  $\hat{\sigma}_{\mathbf{y}}^2 = \hat{\sigma}_{\hat{\mathbf{y}}}^2 + \hat{\sigma}_{\hat{\mathbf{e}}}^2$ .

**Věta:**  $r_{\mathbf{x}, \mathbf{y}}^2 = \frac{\hat{\sigma}_{\hat{\mathbf{y}}}^2}{\hat{\sigma}_{\mathbf{y}}^2} = 1 - \frac{\hat{\sigma}_{\hat{\mathbf{e}}}^2}{\hat{\sigma}_{\mathbf{y}}^2}$ .

## 2.12 Odhad rozptylu

Max. věrohodný odhad rozptylu *původního* rozdělení:

$$\begin{aligned}\lambda(\vartheta_0, \vartheta_1, \sigma) &= -n \ln \sigma \sqrt{2\pi} - \frac{1}{2\sigma^2} \sum_j (y_j - \vartheta_0 - \vartheta_1 x_j)^2 = \\ &= \underbrace{-n \ln \sqrt{2\pi}}_{\text{konst.}} - n \ln \sigma - \frac{1}{2\sigma^2} \sum_j e_j^2, \\ 0 &= \frac{\partial \lambda(\vartheta_0, \vartheta_1, \sigma)}{\partial \sigma} = -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_j e_j^2,\end{aligned}$$

řešením je

$$\hat{\sigma}_{\hat{\mathbf{e}}}^2 = \frac{1}{n} \sum_j \hat{e}_j^2 = \text{D Emp}(\hat{\mathbf{e}}).$$

Tento odhad je vychýlený; nestranný odhad je (*odvození viz [Likeš, Machek]*)

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_j \hat{e}_j^2 = \frac{1}{n-2} \sum_j (y_j - \hat{\vartheta}_0 - \hat{\vartheta}_1 x_j)^2.$$

## 2.13 Rozdělení odhadů a testy hypotéz o nich

**Věta:** Odhady  $\hat{\vartheta}_0, \hat{\vartheta}_1, \hat{\sigma}^2$  skutečných parametrů  $\vartheta_0, \vartheta_1, \sigma^2$  jsou nestranné, konzistentní, asymptoticky normální.

Odhady rozptylů regresních koeficientů

$$\hat{\sigma}_{\hat{\vartheta}_0}^2 = \frac{\hat{\sigma}^2}{n^2 \hat{\sigma}_{\mathbf{x}}^2} \sum_j x_j^2 = \frac{\hat{\sigma}_{\hat{\mathbf{e}}}^2}{n(n-2) \hat{\sigma}_{\mathbf{x}}^2} \sum_j x_j^2,$$
$$\hat{\sigma}_{\hat{\vartheta}_1}^2 = \frac{\hat{\sigma}^2}{n \hat{\sigma}_{\mathbf{x}}^2} = \frac{\hat{\sigma}_{\hat{\mathbf{e}}}^2}{(n-2) \hat{\sigma}_{\mathbf{x}}^2}$$

**nejsou nezávislé.**

Test absolutního členu, tj. nulové hypotézy  $H_0: \vartheta_0 = c$ :

$$\frac{\hat{\vartheta}_0 - c}{\frac{\hat{\sigma}_{\hat{\mathbf{e}}}}{\hat{\sigma}_{\mathbf{x}}} \sqrt{\frac{1}{n} \sum_j x_j^2}} \sqrt{n-2}$$

testujeme na rozdělení  $t(n-2)$ .

Test směrnice, tj. nulové hypotézy  $H_0: \vartheta_1 = c$ :

$$\frac{\hat{\vartheta}_1 - c}{\frac{\hat{\sigma}_{\hat{\mathbf{e}}}}{\hat{\sigma}_{\mathbf{x}}}} \sqrt{n-2}$$

testujeme na rozdělení  $t(n-2)$ .

Test chyby regrese  $\hat{\vartheta}_0 + \hat{\vartheta}_1 x$  pro dané  $x$ , tj. nulové hypotézy  $H_0: \vartheta_0 + \vartheta_1 x = c$ :

$$\frac{\hat{\vartheta}_0 + \hat{\vartheta}_1 x - c}{\hat{\sigma}_{\hat{\mathbf{e}}} \sqrt{n+1 + \frac{n(x-\bar{x})^2}{\hat{\sigma}_{\mathbf{x}}^2}}} \sqrt{n-2}$$

testujeme na rozdělení  $t(n-2)$ .

## 2.14 Testy korelace

**Test nekorelovanosti (nulovosti korelace):** Za předpokladu  $H_0: \varrho(X, Y) = 0$  testujeme

$$\frac{r_{\mathbf{x}, \mathbf{y}}}{\sqrt{1 - r_{\mathbf{x}, \mathbf{y}}^2}} \sqrt{n-2} = \frac{\hat{\sigma}_{\hat{\mathbf{y}}}}{\hat{\sigma}_{\hat{\mathbf{e}}}} \sqrt{n-2}$$

na rozdělení  $t(n-2)$ . (Zde  $\frac{\hat{\sigma}_{\hat{\mathbf{y}}}}{\hat{\sigma}_{\hat{\mathbf{e}}}}$  je poměr směrodatných odchylek odpovídajících vysvětlené a nevysvětlené části celkového rozptylu.)

**Test (nenulové) hodnoty korelace:** Označme funkci

$$h(t) := \frac{1}{2} \ln \frac{1+t}{1-t}.$$

Za předpokladu  $H_0: \varrho(X, Y) = c$ , kde  $c \neq 0$ , pochází

$$z_{\mathbf{x}, \mathbf{y}} := h(r_{\mathbf{x}, \mathbf{y}}) = \frac{1}{2} \ln \frac{1 + r_{\mathbf{x}, \mathbf{y}}}{1 - r_{\mathbf{x}, \mathbf{y}}}$$

z rozdělení přibližně  $N(\mu, \sigma^2)$ , kde

$$\begin{aligned}\mu &:= h(c) = \frac{1}{2} \ln \frac{1+c}{1-c}, \\ \sigma^2 &:= \frac{1}{n-3}, \\ \sigma &= \sqrt{\frac{1}{n-3}}.\end{aligned}$$

Normované kritérium

$$\frac{z_{\mathbf{x},\mathbf{y}} - \mu}{\sigma} = (h(r_{\mathbf{x},\mathbf{y}}) - h(c)) \sqrt{n-3}$$

testujeme na rozdělení  $N(0, 1)$ .

**Test rovnosti dvou korelací:** Z nezávislých výběrů rozsahu  $n$ , resp.  $m$ , vypočítáme výběrové koeficienty korelace  $r_{\mathbf{x},\mathbf{y}}$ , resp.  $r_{\mathbf{u},\mathbf{v}}$ , z nich

$$\begin{aligned}z_{\mathbf{x},\mathbf{y}} &:= h(r_{\mathbf{x},\mathbf{y}}) = \frac{1}{2} \ln \frac{1+r_{\mathbf{x},\mathbf{y}}}{1-r_{\mathbf{x},\mathbf{y}}}, \\ z_{\mathbf{u},\mathbf{v}} &:= h(r_{\mathbf{u},\mathbf{v}}) = \frac{1}{2} \ln \frac{1+r_{\mathbf{u},\mathbf{v}}}{1-r_{\mathbf{u},\mathbf{v}}}.\end{aligned}$$

Za předpokladu  $H_0: \varrho(X, Y) = \varrho(U, V)$  pochází  $z := z_{\mathbf{x},\mathbf{y}} - z_{\mathbf{u},\mathbf{v}} = h(r_{\mathbf{x},\mathbf{y}}) - h(r_{\mathbf{u},\mathbf{v}})$  z rozdělení přibližně  $N(0, \sigma^2)$ , kde

$$\begin{aligned}\sigma^2 &:= \frac{1}{n-3} + \frac{1}{m-3}, \\ \sigma &= \sqrt{\frac{1}{n-3} + \frac{1}{m-3}}.\end{aligned}$$

Normované kritérium

$$\frac{z_{\mathbf{x},\mathbf{y}} - z_{\mathbf{u},\mathbf{v}}}{\sigma} = \frac{h(r_{\mathbf{x},\mathbf{y}}) - h(r_{\mathbf{u},\mathbf{v}})}{\sqrt{\frac{1}{n-3} + \frac{1}{m-3}}}$$

testujeme na rozdělení  $N(0, 1)$ .

## 2.15 Příklad použití

**Příklad 1** [Markechová, Tirpáková, Stehlíková, př. 9.3, 9, 13]: Měřili jsme výšku  $X$  [cm] a hmotnost  $Y$  [kg]  $n = 10$  žáků, výsledky jsou v tabulce:

$x_j$ [cm]	151	154	165	146	158	172	142	169	176	138
$y_j$ [kg]	43	48	56	42	55	54	44	49	55	38

Odhadneme parametry jednotlivých rozdělení (která považujeme přibližně za normální):

$$\begin{aligned}\bar{x} &= \frac{1}{n} \sum_j x_j = 157.1 \text{ cm} \\ \hat{\sigma}_{\mathbf{x}}^2 &= \frac{1}{n} \sum_j (x_j - \bar{x})^2 \doteq 154.7 \text{ cm}^2 \\ \hat{\sigma}_{\mathbf{x}} &= \sqrt{\hat{\sigma}_{\mathbf{x}}^2} \doteq 12.4 \text{ cm} \\ \bar{y} &= \frac{1}{n} \sum_j y_j = 48.4 \text{ kg} \\ \hat{\sigma}_{\mathbf{y}}^2 &= \frac{1}{n} \sum_j (y_j - \bar{y})^2 \doteq 37.44 \text{ kg}^2 \\ \hat{\sigma}_{\mathbf{y}} &= \sqrt{\hat{\sigma}_{\mathbf{y}}^2} \doteq 6.12 \text{ kg}\end{aligned}$$

### 2.15.1 Odhady koeficientů 1. regresní přímky

$$\hat{\vartheta}_1 = \frac{\sum_j (x_j - \bar{x})(y_j - \bar{y})}{\sum_j (x_j - \bar{x})^2} = \frac{\sum_j x_j y_j - n \bar{x} \bar{y}}{\sum_j x_j^2 - n \bar{x}^2} \doteq 0.421 \text{ kg / cm}$$

$$\hat{\vartheta}_0 = \bar{y} - \hat{\vartheta}_1 \bar{x} \doteq -17.78 \text{ kg}$$

Rovnice 1. regresní přímky (ve stejných jednotkách jako výše):

$$y = \hat{\vartheta}_0 + \hat{\vartheta}_1 x \doteq -17.78 + 0.421 x,$$

nebo

$$\begin{aligned} y - \bar{y} &= \hat{\vartheta}_1 (x - \bar{x}), \\ y - 48.4 &= 0.421 (x - 157.1), \end{aligned}$$

např. pro výšku  $x = 160$  cm dostáváme odhad hmotnosti  $y = \hat{\vartheta}_0 + \hat{\vartheta}_1 x \doteq -17.78 + 0.421 \cdot 160 = 49.6$  kg.

Kontrola: Bod  $(\bar{x}, \bar{y})$  leží na 1. regresní přímce,  $\hat{\vartheta}_0 + \hat{\vartheta}_1 \bar{x} \doteq -17.78 + 0.421 \cdot 157.1 = 48.4$  kg.

### 2.15.2 Odhady koeficientů 2. regresní přímky

$$\hat{\vartheta}_1^* = \frac{\sum_j (x_j - \bar{x})(y_j - \bar{y})}{\sum_j (y_j - \bar{y})^2} \doteq 1.74 \text{ cm / kg}$$

$$\hat{\vartheta}_0^* = \bar{x} - \hat{\vartheta}_1^* \bar{y} \doteq 72.9 \text{ cm}$$

Rovnice 2. regresní přímky:

$$x = \hat{\vartheta}_0^* + \hat{\vartheta}_1^* y \doteq 72.9 + 1.74 y,$$

nebo

$$\begin{aligned} x - \bar{x} &= \hat{\vartheta}_1^* (y - \bar{y}), \\ x - 157.1 &= 1.74 (y - 48.4), \end{aligned}$$

$$y - \bar{y} = \frac{1}{\hat{\vartheta}_1^*} (x - \bar{x}),$$

$$y - 48.4 = 0.575 (x - 157.1),$$

např. pro hmotnost  $y = 50$  kg dostáváme odhad výšky  $x = \hat{\vartheta}_0 + \hat{\vartheta}_1 x \doteq 72.9 + 1.74 \cdot 50 = 160$  cm. (*To není totéž jako u 1. regresní přímky!*)

Kontrola: Bod  $(\bar{x}, \bar{y})$  leží na 2. regresní přímce,  $\hat{\vartheta}_0^* + \hat{\vartheta}_1^* \bar{y} \doteq 72.9 + 1.74 \cdot 48.4 = 157$  kg. (*Jen přibližná shoda kvůli zaokrouhlovacím chybám.*)

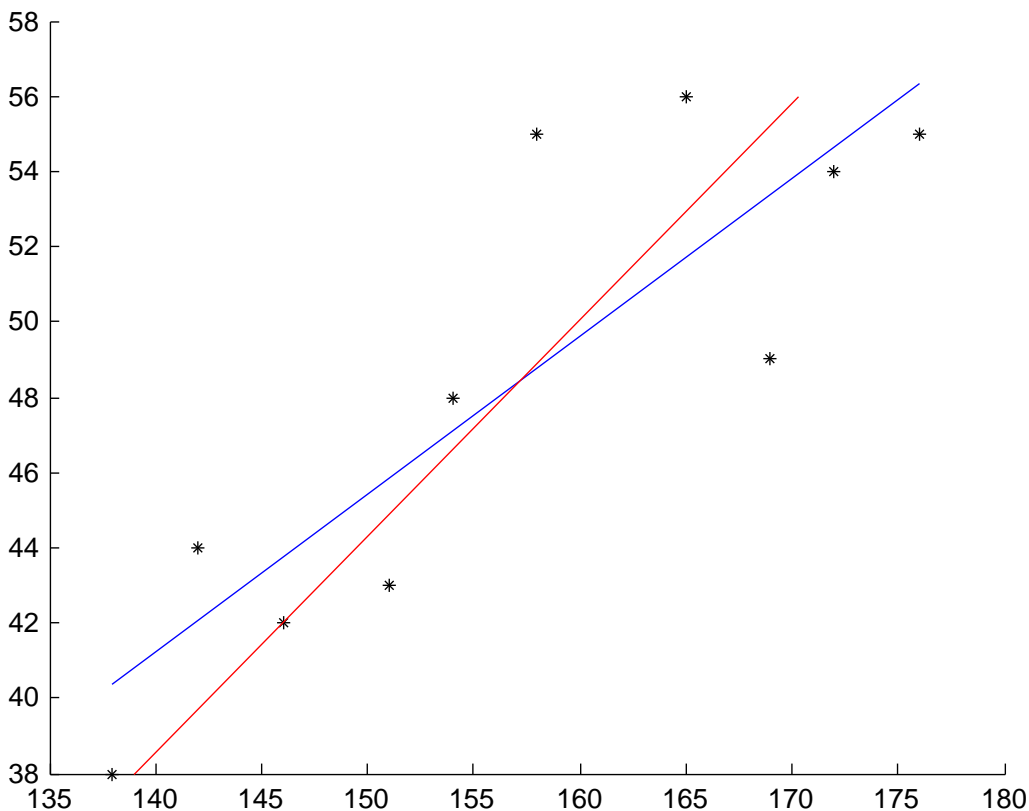
Součin směrnic obou regresních přímek:

$$\hat{\vartheta}_1 \hat{\vartheta}_1^* = 0.421 \cdot 1.74 = 0.733 = r_{x,y}^2,$$

$$r_{x,y} = \sqrt{0.733} = 0.856.$$

(Kladné znaménko vyjadřuje, že větším hodnotám výšky odpovídá v průměru větší hmotnost.)

<http://lib.stat.cmu.edu/DASL/Stories/USCrime.html>



### 2.15.3 Odhady rozptylů, kovariance a korelace

$$\hat{\sigma}_x^2 := \frac{1}{n} \sum_j (x_j - \bar{x})^2 \doteq 154.7 \text{ cm}^2,$$

$$\hat{\sigma}_y^2 := \frac{1}{n} \sum_j (y_j - \bar{y})^2 \doteq 37.44 \text{ kg}^2,$$

$$c_{x,y} := \frac{1}{n} \sum_j (x_j - \bar{x})(y_j - \bar{y}) \doteq 65.2 \text{ kg cm},$$

$$r_{x,y} := \frac{\sum_j (x_j - \bar{x})(y_j - \bar{y})}{\sqrt{\sum_j (x_j - \bar{x})^2 \sum_j (y_j - \bar{y})^2}} = \frac{c_{x,y}}{\hat{\sigma}_x \hat{\sigma}_y} \doteq \frac{65.2}{\sqrt{154.7 \cdot 37.44}} \doteq 0.856.$$

### 2.15.4 Chyba lineární regrese

Odhady hodnot nezávisle proměnné a chyby (**rezidua**)

$$\hat{y}_j = \hat{\vartheta}_0 + \hat{\vartheta}_1 x_j,$$

$$\hat{e}_j = y_j - \hat{\vartheta}_0 - \hat{\vartheta}_1 x_j,$$



$\hat{y}_j$ [kg]	45.8	47.1	51.7	43.7	48.8	54.7	42.0	53.4	56.4	40.4
$\hat{e}_j$ [kg]	-2.8	0.9	4.3	-1.7	6.2	-0.7	2.0	-4.4	-1.4	-2.4

Kontrola:

$$\frac{1}{n} \sum_j \hat{y}_j \doteq 48.4 \doteq \bar{y},$$

$$\frac{1}{n} \sum_j \hat{e}_j = 0.$$

### 2.15.5 Složky rozptylu

**Rozptyl modelu:**  $\hat{\sigma}_{\hat{\mathbf{y}}}^2 = \frac{1}{n} \sum_j (\hat{y}_j - \bar{y})^2 \doteq 27.44 \text{ kg}^2$ .

**Reziduální rozptyl:**  $\hat{\sigma}_{\hat{\mathbf{e}}}^2 = \frac{1}{n} \sum_j \hat{e}_j^2 = \frac{1}{n} \sum_j (y_j - \hat{y}_j)^2 \doteq 10 \text{ kg}^2$ .

**Celkový rozptyl:**  $\hat{\sigma}_{\mathbf{y}}^2 = \frac{1}{n} \sum_j (y_j - \bar{y})^2 \doteq 37.44 \text{ kg}^2 \doteq \hat{\sigma}_{\hat{\mathbf{y}}}^2 + \hat{\sigma}_{\hat{\mathbf{e}}}^2$ .

Kontrola:

$$r_{\mathbf{x}, \mathbf{y}}^2 = \frac{\hat{\sigma}_{\hat{\mathbf{y}}}^2}{\hat{\sigma}_{\mathbf{y}}^2} = 1 - \frac{\hat{\sigma}_{\hat{\mathbf{e}}}^2}{\hat{\sigma}_{\mathbf{y}}^2},$$

$$0.733 \doteq \frac{27.44}{37.44} \doteq 1 - \frac{10}{37.44}.$$

### 2.15.6 Odhady rozptylu původního rozdělení

Max. věrohodný:

$$\hat{\sigma}_{\hat{\mathbf{e}}}^2 = \frac{1}{n} \sum_j \hat{e}_j^2 \doteq 10 \text{ kg}^2,$$

nestranný:

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_j \hat{e}_j^2 \doteq 12.5 \text{ kg}^2.$$

### 2.15.7 Testy korelace

**Test nekorelovanosti (nulovosti korelace):** Nulovou hypotézu  $H_0 : \rho(X, Y) = 0$  posoudíme podle (realizace) statistiky

$$\frac{r_{\mathbf{x}, \mathbf{y}}}{\sqrt{1 - r_{\mathbf{x}, \mathbf{y}}^2}} \sqrt{n-2} \doteq \frac{0.856}{\sqrt{1 - 0.856^2}} \sqrt{10-2} \doteq 4.69,$$

kteřou lze určit i z poměru rozptylu modelu a reziduálního modelu,

$$\sqrt{\frac{\hat{\sigma}_{\hat{\mathbf{y}}}^2}{\hat{\sigma}_{\hat{\mathbf{e}}}^2}} (n-2) \doteq \sqrt{\frac{27.44}{10}} (10-2) \doteq 4.69.$$

Testujeme ji na Studentovo rozdělení  $t(8)$ . Pro oboustranný test (alternativní hypotéza  $H_1 : \varrho(X, Y) \neq 0$ ) vychází dosažená významnost 0.00078, pro jednostranný test (alternativní hypotéza  $H_1 : \varrho(X, Y) > 0$ ) 0.0016, což dovoluje obě hypotézy zamítnout na vysoké hladině významnosti.

(Samozřejmě nelze zamítnout nulovou hypotézu  $H_0 : \varrho(X, Y) \geq 0$  proti alternativní hypotéze  $H_1 : \varrho(X, Y) < 0$ .)

**Test (nenulové) hodnoty korelace:** Testujeme nulovou hypotézu  $H_0 : \varrho(X, Y) = c$  proti alternativní hypotéze  $H_1 : \varrho(X, Y) \neq c$  např. pro  $c = 0.6$ . Provedeme nelineární transformaci obou porovnávaných hodnot:

$$z_{\mathbf{x}, \mathbf{y}} = \frac{1}{2} \ln \frac{1 + r_{\mathbf{x}, \mathbf{y}}}{1 - r_{\mathbf{x}, \mathbf{y}}} = \frac{1}{2} \ln \frac{1 + 0.856}{1 - 0.856} \doteq 1.28,$$

$$\frac{1}{2} \ln \frac{1 + c}{1 - c} \doteq 0.693,$$

jejich rozdíl poměříme směrodatnou odchylkou

$$\sigma_{Z_{\mathbf{x}, \mathbf{y}}} = \sqrt{\frac{1}{n - 3}} \doteq 0.378,$$

testujeme

$$\frac{1.28 - 0.693}{0.378} \doteq 1.55$$

na rozdělení  $N(0, 1)$ , což nestačí k zamítnutí nulové hypotézy na hladině významnosti 5% (ani pro jednostranné hypotézy  $H_0 : \varrho(X, Y) \leq c$ ,  $H_1 : \varrho(X, Y) > c$ ).

**Test rovnosti dvou korelací:** Pokud v jiném průzkumu rozsahu  $m = 20$  vyšel odhad korelace stejných náhodných veličin např.  $d = 0.4$ , transformujeme i tuto hodnotu:

$$\frac{1}{2} \ln \frac{1 + d}{1 - d} = \frac{1}{2} \ln \frac{1 + 0.4}{1 - 0.4} \doteq 0.424$$

rozdíl poměříme směrodatnou odchylkou

$$\sqrt{\frac{1}{n - 3} + \frac{1}{m - 3}} \doteq 0.449,$$

testujeme

$$\frac{1.28 - 0.424}{0.449} \doteq 1.91$$

na rozdělení  $N(0, 1)$ , což nestačí k zamítnutí nulové hypotézy o rovnosti obou korelací na hladině významnosti 5% (ale můžeme na této hladině zamítnout jednostrannou hypotézu  $H_0 : \varrho(X, Y) \leq d$  proti  $H_1 : \varrho(X, Y) > d$ ).

## 2.16 Odhad náhodné veličiny: lineární model dimenze $k$

**Úloha 4:** Odhadujeme spojitou náhodnou veličinu  $Y$  pomocí  $k$  **vysvětlujících** náhodných veličin  $X_1, \dots, X_k$  na základě realizace náhodného výběru

$$((y_1, x_{11}, \dots, x_{1k}), \dots, (y_n, x_{n1}, \dots, x_{nk})),$$

kde  $n > k$ , obvykle  $n \gg k$ . Předpokládáme lineární model

$$Y = \sum_{i=1}^k \vartheta_i X_i + \mathcal{E},$$

kde  $\boldsymbol{\vartheta} = (\vartheta_1, \dots, \vartheta_k)^T \in \mathbb{R}^k$  je (sloupcový) vektor neznámých parametrů (**regresních koeficientů**),

$\mathcal{E}$  je náhodná veličina s rozdělením  $N(0, \sigma^2)$ ,

$\sigma^2$  je (konstantní, známý nebo neznámý) rozptyl,

**Otázka:** Kam se poděl absolutní člen  $\vartheta_0$  z modelu dimenze 1?

**Odpověď 1:** Je-li potřebný, lze jej zahrnout tak, že jedna z vysvětlujících proměnných je konstanta, BÚNO 1.

**Odpověď 2:** Od všech veličin  $Y, X_1, \dots, X_k$  lze odečíst jejich výběrové průměry.

Pro realizace dostáváme

$$y_j = \sum_{i=1}^k \vartheta_i x_{ji} + e_j,$$

maticový zápis:

$$\mathbf{y} = \mathbf{X} \boldsymbol{\vartheta} + \mathbf{e},$$

kde  $\mathbf{y} = (y_1, \dots, y_n)^T$ ,  $\mathbf{e} = (e_1, \dots, e_n)^T \in \mathbb{R}^n$  jsou (sloupcové) vektory realizací (*nezávislých*) náhodných veličin,

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1k} \\ x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nk} \end{pmatrix}.$$

Max. věrohodný odhad (za předpokladu normálního rozdělení chyb) = odhad metodou nejmenších čtverců  $\hat{\boldsymbol{\vartheta}}$ ; který minimalizuje **reziduální součet čtverců** (angl. **residual sum of squares**, **RSS**)

$$R_{SS} = \sum_{j=1}^n \hat{e}_j^2 = \sum_{j=1}^n \left( y_j - \sum_{i=1}^k \hat{\vartheta}_i x_{ji} \right)^2 = \|\mathbf{e}\|^2 = \|\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\vartheta}}\|^2.$$

Geometrické řešení:  $\mathbf{X} \hat{\boldsymbol{\vartheta}}$  je lineární kombinace sloupců matice  $\mathbf{X}$ , která má nejmenší euklidovskou vzdálenost od  $\mathbf{y}$ , tj. kolmý průmět  $\mathbf{y}$  do lineárního podprostoru generovaného sloupci matice  $\mathbf{X}$ .

Ten je charakterizován tím, že každý jeho vektor  $\mathbf{z}$  je kolmý na sloupce matice  $\mathbf{X}$ , tj.

$$\mathbf{X}^T \mathbf{z} = \mathbf{0}.$$

⇒ **soustava normálních rovnic**

$$\mathbf{X}^T (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\vartheta}}) = \mathbf{0},$$

$$\mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\vartheta}} = \mathbf{X}^T \mathbf{y}.$$

**Předp.:**  $\mathbf{X}$  má plnou hodnot,  $k$ .

**Pozn.:** To je obvyklý případ, ale pro reálná data nelze zaručit. Pokud to není splněno, lze nějakou vysvětlující proměnnou vyjádřit jako lineární kombinaci ostatních a v modelu ji vynechat.

**Důsledek:** (Symetrická) matice soustavy normálních rovnic  $\mathbf{X}^T \mathbf{X}$  je regulární, řešení je

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

**Pozn.:** Matice  $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$  je **pseudoinverze** obdélníkové matice  $\mathbf{X}$ . Vynásobíme-li jí matici  $\mathbf{X}$  zleva, dostaneme jednotkovou matici.

**Věta:**  $\forall i$ :  $\hat{\boldsymbol{\vartheta}}_i$  (přesněji  $\hat{\Theta}_i$ ) je nejlepší nestranný odhad  $\boldsymbol{\vartheta}_i$ .

Kovarianční matice vektoru odhadů  $\hat{\Theta}$  je

$$\boldsymbol{\Sigma}_{\hat{\Theta}} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}.$$

**Věta:** Hodnota regresní funkce v bodě  $\mathbf{x} = (x_1, \dots, x_k)^T \in \mathbb{R}^k$ ,

$$\mathbf{x}^T \hat{\boldsymbol{\vartheta}} = \sum_{i=1}^k \hat{\boldsymbol{\vartheta}}_i x_i,$$

je nejlepší nestranný odhad vysvětlované náhodné veličiny  $Y$  v bodě  $\mathbf{x}$ .

## 2.17 Intervalové odhady regresních koeficientů při známém rozptylu

**Věta:** (Symetrický)  $(1 - \alpha)$ -konfidenční interval pro odhad parametru  $\boldsymbol{\vartheta}_i$  je

$$\hat{\boldsymbol{\vartheta}}_i \pm \Phi^{-1} \left(1 - \frac{\alpha}{2}\right) \sqrt{(\boldsymbol{\Sigma}_{\hat{\Theta}})_{ii}} = \hat{\boldsymbol{\vartheta}}_i \pm \Phi^{-1} \left(1 - \frac{\alpha}{2}\right) \sigma \sqrt{c_{ii}},$$

kde  $(\boldsymbol{\Sigma}_{\hat{\Theta}})_{ii}$ , resp.  $c_{ii}$ , je  $i$ -tý prvek na diagonále matice  $\boldsymbol{\Sigma}_{\hat{\Theta}}$ , resp.  $\mathbf{C} := (\mathbf{X}^T \mathbf{X})^{-1}$ .

**Pozn.:** I pro intervalový odhad jednoho regresního koeficientu potřebujeme vypočítat matici  $(\mathbf{X}^T \mathbf{X})^{-1}$ .

Pokud rozptyl neznáme, je potřeba jej nahradit odhadem (*viz dále*).

## 2.18 Odhady rozptylu $\sigma^2$ původního rozdělení

Maximálně věrohodný:

$$\hat{\sigma}_{\hat{\boldsymbol{\epsilon}}}^2 = \frac{1}{n} R_{SS} = \frac{1}{n} \sum_{j=1}^n \hat{\epsilon}_j^2 = \frac{1}{n} \sum_{j=1}^n \left( y_j - \sum_{i=1}^k \hat{\boldsymbol{\vartheta}}_i x_{ji} \right)^2,$$

nestranný:

$$\hat{\sigma}^2 = \frac{1}{n-k} R_{SS} = \frac{1}{n-k} \sum_{j=1}^n \hat{\epsilon}_j^2 = \frac{1}{n-k} \sum_{j=1}^n \left( y_j - \sum_{i=1}^k \hat{\boldsymbol{\vartheta}}_i x_{ji} \right)^2.$$

$\frac{R_{SS}}{\sigma^2}$  pochází z rozdělení  $\chi^2(n-k)$ .

Odhady  $R_{SS}$  a  $\hat{\boldsymbol{\vartheta}}$  jsou nezávislé. (Nikoli však jednotlivé složky vektoru  $\hat{\boldsymbol{\vartheta}}$  mezi sebou!)

## 2.19 Intervalové odhady regresních koeficientů při **neznámém rozptylu**

**Věta:** (Symetrický)  $(1 - \alpha)$ -konfidenční interval pro odhad parametru  $\vartheta_i$  je

$$\hat{\vartheta}_i \pm q_{t(n-k)} \left(1 - \frac{\alpha}{2}\right) \hat{\sigma} \sqrt{c_{ii}} = \hat{\vartheta}_i \pm q_{t(n-k)} \left(1 - \frac{\alpha}{2}\right) \cdot \sqrt{\frac{R_{SS} c_{ii}}{n - k}},$$

kde  $c_{ii}$  je  $i$ -tý prvek na diagonále matice  $\mathbf{C} = (\mathbf{X}^T \mathbf{X})^{-1}$ .

## 3 Volba vysvětlujících proměnných

Můžeme vyhodnotit, zda se jednotlivé regresní koeficienty statisticky významně liší od 0, např. pro závislost kriminality na různých sociálních faktorech ve státech USA v r. 1960 (převzato z [Wasserman]) a

<http://lib.stat.cmu.edu/DASL/Stories/USCrime.html>:

Popis proměnných v originále (Intercept=absolutní člen):

Crime rate: # of offenses reported to police per million population

The number of males of age 14-24 per 1000 population

Indicator variable for Southern states (0 = No, 1 = Yes)

Mean # of years of schooling  $\times 10$  for persons of age 25 or older

1960 per capita expenditure on police by state and local government

1959 per capita expenditure on police by state and local government

Labor force participation rate per 1000 civilian urban males age 14-24

The number of males per 1000 females

State population size in hundred thousands

Unemployment rate of urban males per 1000 of age 14-24

Unemployment rate of urban males per 1000 of age 35-39

Median value of transferable goods and assets or family income in tens of \$

Covariate	$\hat{\beta}_j$	$\hat{\text{se}}(\hat{\beta}_j)$	t value	p-value
(Intercept)	-589.39	167.59	-3.51	0.001 **
Age	1.04	0.45	2.33	0.025 *
Southern State	11.29	13.24	0.85	0.399
Education	1.18	0.68	1.7	0.093
Expenditures	0.96	0.25	3.86	0.000 ***
Labor	0.11	0.15	0.69	0.493
Number of Males	0.30	0.22	1.36	0.181
Population	0.09	0.14	0.65	0.518
Unemployment (14-24)	-0.68	0.48	-1.4	0.165
Unemployment (25-39)	2.15	0.95	2.26	0.030 *
Wealth	-0.08	0.09	-0.91	0.367

Může se zdát, že např. vyšší výdaje na vzdělání a prevenci kriminality souvisí s vyšší kriminalitou. To ovšem může být z mnoha jiných důvodů než je kauzální závislost (v kterémkoli směru).

Naopak některé vysvětlující proměnné, např. „jižní státy“, jsou „zbytečné“. Jejich vynecháním se model téměř nezhorší a přitom zjednoduší.

### 3.1 Kritéria pro výběr modelu

**underfitting:** příliš jednoduchý model, nepřesný

**overfitting:** zbytečně složitý model

Kritéria vyvažují počet vysvětlujících proměnných a chybu modelu (obvykle  $R_{SS}$ ), např.

$$R_{SS}(\mathcal{S}) + 2 \hat{\sigma}^2 |\mathcal{S}| ,$$

kde  $\mathcal{S}$  je množina vybraných proměnných,  $|\mathcal{S}|$  jejich počet a  $R_{SS}(\mathcal{S})$  je reziduální součet čtverců modelu používajícího pouze proměnné z  $\mathcal{S}$ .

$$\lambda_{\mathcal{S}} - c |\mathcal{S}| ,$$

kde  $\lambda_{\mathcal{S}}$  je logaritmus věrohodnosti modelu používajícího pouze proměnné z  $\mathcal{S}$ . Doporučené hodnoty  $c$  např. 1 nebo  $\frac{1}{n} \ln n$ .

Zde nezáleží na použitých jednotkách, neboť porovnáváme pouze relativní přesnost jednotlivých modelů.

### 3.2 Volba vysvětlujících proměnných

A. Postupně přidáváme vždy tu proměnnou, která nejvíc vylepší kritérium.

B. Použijeme všechny proměnné a postupně ubíráme vždy tu, která nejvíc vylepší kritérium.

C. Začneme od jakékoli podmnožiny kritérií a náhodně ji upravujeme (přidání/ubrání/výměna), větší pravděpodobnost dáme změnám, které víc vylepšují kritérium.

### 3.3 Rozdělení na trénovací a testovací data, cross-validation

Místo skutečné chyby používáme její odhad na základě trénovacích dat.

Testovací data by měla být s trénovacími disjunktní.

Úspornější postupy:

**cross-validation:** rozdělíme data na  $m$  disjunktních množin, vždy použijeme jednu jako testovací a ostatní jako trénovací,

**leave-one-out cross-validation:** extrémní případ předchozího pro  $m = n$ .