

Lineární statistické modely II

Domácí úkol

Vojtěch Matulík

505487

Obor Statistika a analýza dat

Přírodovědecká fakulta
Masarykova univerzita

26. května 2023

Příklad 1

Mějme lineární regresní model tvaru

$$Y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \varepsilon_i, \quad i = 1, 2, \dots, n,$$

kde pro náhodné chyby platí $\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$. Podle Scheffého věty platí

$$P\left(\{\mathbf{b}^T(A\hat{\boldsymbol{\beta}} - A\boldsymbol{\beta})\}^2 \leq mF_{1-\alpha}(m, n-p)\hat{\sigma}^2\mathbf{b}^T A(\mathbf{X}^T \mathbf{X})^{-1} A^T \mathbf{b}; \forall \mathbf{b} \in \mathbb{R}^m\right) = 1 - \alpha,$$

kde A je rozměrů $m \times p$ a b je rozměrů $m \times 1$. Pro případ kdy $b = (1, x, x^2)^T$ a $A = I_3$, kde I_3 je jednotková matice rozměrů 3×3 , pak platí

$$P\left(\{(1, x, x^2)(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})\}^2 \leq 3F_{1-\alpha}(3, n-3)\hat{\sigma}^2(1, x, x^2)(\mathbf{X}^T \mathbf{X})^{-1}(1, x, x^2)^T\right) = 1 - \alpha.$$

Pás spolehlivosti odvozený z Scheffého věty s pravděpodobností pokrytí $100(1-\alpha)\%$ je pro tento speciální případ tvaru

$$P\left(|(\hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 x^2) - (\beta_0 + \beta_1 x + \beta_2 x^2)| \leq \sqrt{3F_{1-\alpha}(3, n-3)\hat{\sigma}^2(1, x, x^2)(\mathbf{X}^T \mathbf{X})^{-1}(1, x, x^2)^T}\right) = 1 - \alpha.$$

Příklad 2

V datovém souboru `baseball_hit.Rdata` jsou zaznamenány údaje o baseballových odpalech. Máme k dispozici proměnnou vzdálenost (`distance`), která udává horizontální vzdálenost (v metrech) mezi pálkařem a dopadem míčku, a proměnnou úhel (`angle`), jež zachycuje velikost úhlu (ve stupních) mezi trajektorií míčku a zemí při odpalu.

(a) Lineární model

Na základě vizuálního posouzení, koeficientu determinace a hodnot Akaikova informačního kritéria (AIC) jsme vybrali kvadratický model, který nejlépe popisuje data `baseball_hit.Rdata`. Tento kvadratický model vypadá následujícím způsobem

```
1 mod <- lm(distance ~ angle + I(angle^2), data = data)
```

(b) Scheffého pásy spolehlivosti

Definujeme si funkci `ScheffePS1`, která slouží k výpočtu Scheffého pásu spolehlivosti lineárního regresního modelu s pravděpodobností pokrytí 95% v daném bodě.

```

2 ScheffePS1 <- function(b, A, mod, alpha){
3   X <- model.matrix(mod)
4   Y.hat <- predict(mod, newdata = data.frame(angle = b[2]))
5   n <- dim(X)[1]
6   m <- dim(A)[1]
7   p <- dim(A)[2]
8   M.alpha <- sqrt(m * qf(df1 = m, df2 = n - p, p = 1 - alpha))
9   s <- summary(mod)$sigma
10  result1 <- c(Y.hat - M.alpha * s *
11              sqrt(t(b) %*% A %*% solve(t(X) %*% X) %*% t(A) %*% b),
12              Y.hat + M.alpha * s *
13              sqrt(t(b) %*% A %*% solve(t(X) %*% X) %*% t(A) %*% b))
14  return(result1)
15 }

```

Tato funkce přijímá vektor regresních koeficientů b , maticovou reprezentaci kontrastů A , lineární regresní model mod a hladinu významnosti α . Výstupem funkce je dolní a horní mez Scheffého pásu spolehlivosti pro predikované hodnoty Y .

Dále si vytváříme vektor xx , který obsahuje 200 hodnot, které se rovnoměrně rozkládají mezi minimální a maximální hodnotou proměnné $angle$. Aplikujeme funkci `ScheffePS1` na každou hodnotu vektoru xx a výsledek ukládáme do proměnné `MatScheffe`.


```

16 xx <- seq(from = min(data$angle), to = max(data$angle),
17           length.out = 200)
18
19 MatScheffe <- sapply(1:length(xx), function(i){
20   result2 <- ScheffePS1(mod = mod, A = diag(3),
21                         b = c(1, xx[i], xx[i]^2),
22                         alpha = 0.05)
23   return(result2)
24 })

```

Tím jsme tedy získali hodnoty Scheffého pásu spolehlivosti pro náš model `mod` s pravděpodobností pokrytí 95 % (tyto hodnoty jsou uloženy v proměnné `MatScheffe`).

(c) Bodová a intervalová predikce

Predikujme (bodově i intervalově) vzdálenost dopadu míčku od pálkaře pokud odpálí míček pod úhlem 55° v  příkazem:

```

25 prediction <- predict(mod, newdata = data.frame(angle = 55),
26                      interval = "prediction", level = 0.95)

```

Intervalová predikce nám vychází (76.5906, 82.47312), bodová pak 70.70808.

(d) Simultánní intervalový odhad

Pomocí následujících řádků kódu vypočteme simultánní intervalový odhad střední hodnoty délky odpalu pro úhly odpalu 20° , 25° a 30° . Přičemž volíme hladinu významnosti $\alpha = 0.05$ a využíváme Bonferroniho adjustace.

```

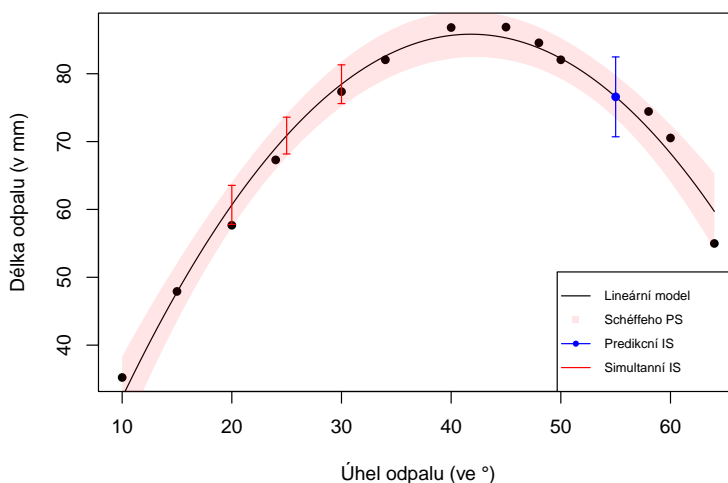
27 angles <- c(20, 25, 30)
28 alpha <- 0.05
29 bonferroni.alpha <- alpha / length(angles)
30 simultaneous.interval <- predict(mod, newdata =
31                                   data.frame(angle = angles),
32                                   interval = "confidence",
33                                   level = 1 - bonferroni.alpha)
34 lower.bound.simultaneous <- simultaneous.interval[, "lwr"]
35 upper.bound.simultaneous <- simultaneous.interval[, "upr"]

```

Výsledné intervalové odhady nám vycházejí (1587.934, 1628.435) pro úhel odpalu 20°, (1688.657, 1711.751) pro úhel odpalu 25° a (1775.560, 1808.885) pro úhel odpalu 30°.

(e) Vykreslení výsledků

Nyní zakreslíme výsledky z předchozích bodů do jednoho obrázku. Tj. vykreslíme bodový graf pozorování, přidáme křivku našeho modelu, zvýrazníme pás spolehlivosti a požadované intervalové odhady i predikce a nakonec vykreslíme legendu.



Obrázek 1: Vykreslení výsledků v jednom grafu

Příklad 3

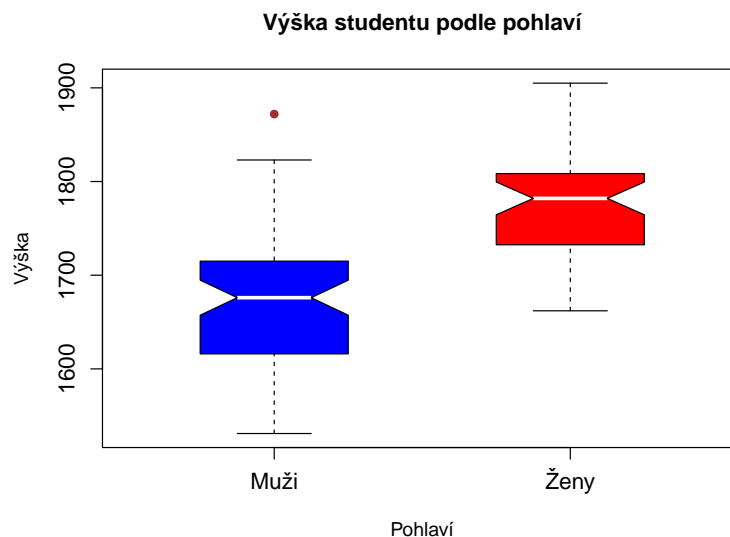
V tomto příkladě budeme pracovat s antropometrickými údaji studentů vysokých škol uloženými v souboru lrmfoot.txt. Data obsahují proměnné: pohlaví (sex), délku chodidla v milimetrech (foot.L) a tělesnou výšku v milimetrech (body.H). Bude nás zajímat efekt pohlaví na tělesnou výšku adjustovaný na délku chodidla.

(a) Čištění dat a vykreslení krabicových grafů

Provedeme čištění dat.


```
36 any(is.na(data)) # Chybející hodnoty?
37 unique(data$sex) # Prvky mimo obor hodnot?
38 sort(unique(data$foot.L)) # Prvky mimo obor hodnot?
39 sort(unique(data$body.H)) # Prvky mimo obor hodnot?
```

Zjišťujeme, že data jsou čistá. Pokračujeme vykreslením krabicových diagramů popisujících výšku studentů v závislosti na pohlaví.




Obrázek 2: Vykreslení krabicových diagramů


(b) Výběr modelu

Budeme modelovat závislost střední hodnoty tělesné výšky na pohlaví a délce chodidla. Vyzkoušíme si různé varianty složitosti modelu: model se vzájemnou interakcí pohlaví a délky chodidla (1 – všeobecný, různé sklony přímek), model bez interakce (2 – ANCOVA, stejné sklony přímek) a model bez vlivu proměnné pohlaví (3 – jedna přímka). Modely jsme si v  definovali následovně


```
40 # Model s interakcí
41 model1 <- lm(body.H ~ sex * foot.L, data = data)
42 summary(model1)
43
44 # Model bez interakce (ANCOVA)
45 model2 <- lm(body.H ~ sex + foot.L, data = data)
46 summary(model2)
47
48 # Model bez vlivu promenne foot.L
49 model3 <- lm(body.H ~ foot.L, data = data)
50 summary(model3)
```

Výstup z  posledních dvou řádků funkce summary aplikované na model1:

```
51 Multiple R-squared:  0.7252,      Adjusted R-squared:  0.7179
52 F-statistic: 99.42 on 3 and 113 DF,  p-value: < 2.2e-16
```

Výstup z  posledních dvou řádků funkce summary aplikované na model2:

```
53 Multiple R-squared:  0.7199,      Adjusted R-squared:  0.715
54 F-statistic: 146.5 on 2 and 114 DF,  p-value: < 2.2e-16
```

Výstup z  posledních dvou řádků funkce summary aplikované na model3:

```
55 Multiple R-squared:  0.715,      Adjusted R-squared:  0.7126
56 F-statistic: 288.6 on 1 and 115 DF,  p-value: < 2.2e-16
```

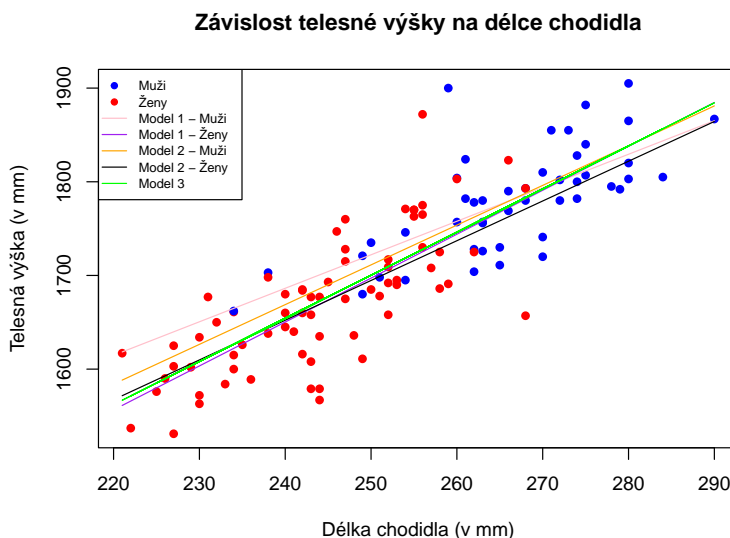
Pokud bychom měli vybrat nejlepší model pouze na základě těchto metrik, model 1 by mohl být mírně preferován, protože má nejvyšší hodnotu koeficientu determinace R^2 . Nicméně je důležité zvážit také interpretovatelnost a relevantnost jednotlivých koeficientů, stejně jako další diagnostické testy.

Celkově lze říci, že model 1, model 2 a model 3 jsou v tomto případě velmi podobné a nelze jednoznačně rozhodnout, který z nich je nejlepší. Avšak kvůli interpretovatelnosti volíme model 3. Odhad vektoru parametrů β pro model 3 je:

```
57 > coefficients(model3)
58 (Intercept)      foot.L
59 549.966107      4.600951
```

(c) Vykreslení modelů

Vykreslíme modely 1–3.



Obrázek 3: Vykreslení modelů v jednom grafu


(d) Vykreslení grafu a simultánní odhad

Vypočteme z dat minimální a maximální délku chodidla a v rámci tohoto rozsahu zkonstruujeme 95% Scheffého pás spolehlivosti pro model 3. V  provedeno následovně

```

60 xx <- seq(from = min.foot, to = max.foot, length.out = 200)
61
62 ScheffePS2 <- function(b, A, mod, alpha) {
63   X <- model.matrix(mod)
64   Y.hat <- predict(mod, newdata = data.frame(foot.L = b[2]))
65   n <- dim(X)[1]
66   m <- dim(A)[1]
67   p <- dim(A)[2]
68   M.alpha <- sqrt(m * qf(df1 = m, df2 = n - p, p = 1 - alpha))
69   s~<- summary(mod)$sigma
70   result1 <- c(Y.hat - M.alpha * s *
71               sqrt(t(b) %*% A %*% solve(t(X) %*% X) %*% t(A) %*% b),
72               Y.hat + M.alpha * s *
73               sqrt(t(b) %*% A %*% solve(t(X) %*% X) %*% t(A) %*% b))
74   return(result1)
75 }
76
77 MatScheffe <- sapply(1:length(xx), function(i) {
78   result2 <- ScheffePS2(mod = model3, A = diag(2), b = c(1, xx[i]),
79                       alpha = 0.05)
80   return(result2)
81 })

```

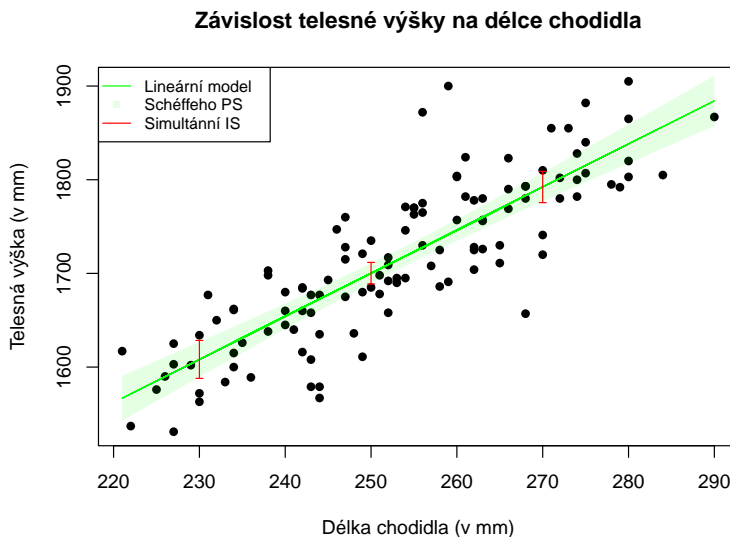
Dále odhadneme střední hodnotu tělesné výšky pro jedince s délkou chodidla 230 mm, 250 mm a 270 mm pomocí 95% simultánních oboustranných intervalů spolehlivosti se Šidákovou adjustací. V  provedeno následovně

```

82 foot.lengths <- c(230, 250, 270)
83 alpha <- 0.05
84 alpha.adjusted <- (1 - (1 - alpha) ^ (1 / length(foot.lengths))) / 2
85 simultaneous.interval <- predict(model3,
86                                newdata = data.frame(foot.L = foot.
87                                                        lengths),
88                                interval = "confidence",
89                                level = 1 - alpha.adjusted)
89 lower.bound.simultaneous <- simultaneous.interval[, "lwr"]
90 upper.bound.simultaneous <- simultaneous.interval[, "upr"]


```

Nakonec vytvoříme graf, který bude obsahovat všechna pozorování spolu s regresní přímkou odpovídající modelu 3. Do tohoto grafu zakreslíme námi vypočtený 95% Scheffého pás spolehlivosti a vypočtené 95% simultánní oboustranné intervaly spolehlivosti se Šidákovou adjustací.



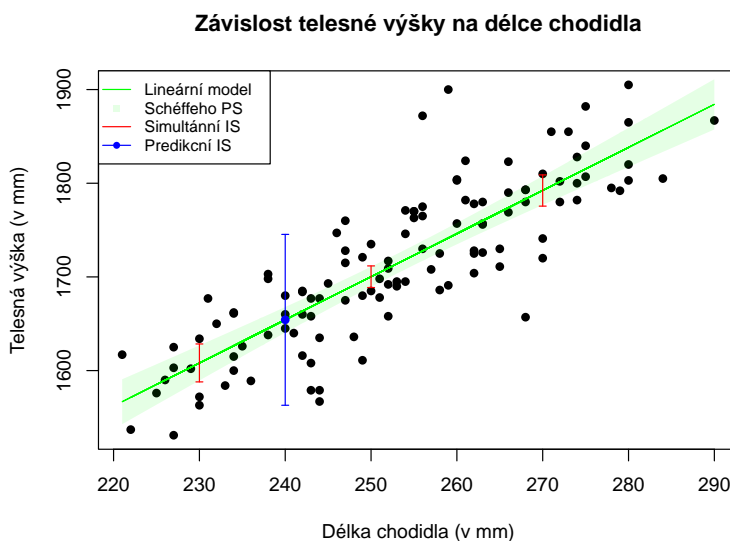
Obrázek 4: Vykreslení modelu, dat, Scheffého pásu a intervalů spolehlivosti

(e) Predikce výšky

Pomocí 95% intervalu spolehlivosti predikujeme výšku jedince, jenž má délku chodidla 240 mm. To uděláme v  příkazem

```
91 prediction <- predict(model3, newdata = data.frame(foot.L = 240),
92                      interval = "prediction", level = 0.95)
```

Tento interval je: (1654.194, 1745.393). Nakonec znázorníme tento interval do obrázku 4.



Obrázek 5: Vykreslení predikce