

3 Linearna regresija

3.1 Linearna regresija

Zanima nas povezanost števila ur učenja na teden z rezultatom na izpitu iz statistike. Vzemimo, da vemo, da se rezultat na izpitu v populaciji porazdeljuje pogojno normalno: $Y|X \sim N(\beta_0 + \beta_1 X, \sigma^2)$.

- Iz spodnjega izpisa preberite ocene populacijskih parametrov. Interpretirajte rezultate, katere ničelne domneve so testirane in kako?

```
Call:
lm(formula = y ~ x)

Residuals:
    Min       1Q   Median       3Q      Max
-20.683  -4.746   2.844   4.512  14.693

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  19.2049     7.5172   2.555 0.033921 *
x              3.6850     0.6217   5.927 0.000351 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 11.44 on 8 degrees of freedom
Multiple R-squared:  0.8145,    Adjusted R-squared:  0.7913
F-statistic: 35.13 on 1 and 8 DF,  p-value: 0.0003508
```

Ocene parametrov so $\hat{\beta}_0 = 19,2$, $\hat{\beta}_1 = 3,7$, $\hat{\sigma} = 11,4$. Testirani sta dve ničelni domnevi: $H_{0int} : \beta_0 = 0$ in $H_0 : \beta_1 = 0$. Pri linearni regresiji nas ponavadi zanima le druga - saj ta govori o povezanosti med spremenljivkama v populaciji. Pri iskanju porazdelitve cenilke $\hat{\beta}_1$ bi se lahko oprli na teorijo metode največjega verjetja, vendar pa v tem primeru aproksimacija ni potrebna. Cenilka je namreč linearna kombinacija vrednosti Y (to smo izpeljali v nalogi 3.2), zato je normalno porazdeljena. Njena varianca (standardna napaka) je ocenjena iz podatkov, zato je standardizirana vrednost cenilke porazdeljena kot t . Testna statistika

$$T = \frac{\hat{\beta}_1}{\widehat{SE}_{\beta_1}} = \frac{3,7}{0,6} = 5,9$$

je torej porazdeljena kot t z 8 stopinjami prostosti (pri ocenjevanju SE porabimo dve stopinji prostosti). Ta test se imenuje Waldov test.

- Kako bi ničelno domnevo $H_0 : \beta_1 = 0$ preverili s posplošenim testom razmerja verjetij?

Namig: Kjer je le mogoče, uporabite rezultate iz prejšnje naloge

Začnimo z ocenami pod ničelno domnevo. Pod ničelno domnevo, je povprečje za vse posameznike enako, neposredno torej lahko uporabimo rezultate iz prejšnje naloge, le da namesto μ pišemo β_0 , zato je maksimum funkcije verjetja pod ničelno domnevo enak

$$\begin{aligned} L_0(y, x, \hat{\beta}_0, \hat{\sigma}) &= \frac{1}{(\sqrt{2\pi}\hat{\sigma})^n} \exp\left\{-\frac{\sum_{i=1}^n (y_i - \hat{\beta}_0)^2}{2\hat{\sigma}^2}\right\} \\ &= \frac{1}{(\sqrt{2\pi}\hat{\sigma})^n} \exp\left\{-\frac{n \sum_{i=1}^n (y_i - \hat{\beta}_0)^2}{2 \sum_{i=1}^n (y_i - \hat{\beta}_0)^2}\right\} \\ &= \frac{1}{(\sqrt{2\pi}\hat{\sigma})^n} \exp\left\{-\frac{n}{2}\right\} \end{aligned}$$

Pod alternativno domnevo na enak način uporabimo rezultat, da je ocena $\hat{\sigma}$ enaka

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - \hat{\beta}_0 - x_i \hat{\beta}_1)^2}{n}$$

Funkcija verjetja je enaka:

$$L(y, x, \beta_0, \beta_1, \sigma) = \frac{1}{(\sqrt{2\pi}\sigma)^n} \exp\left\{-\frac{\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2}{2\sigma^2}\right\}$$

In tako je maksimum funkcije verjetja enak

$$\begin{aligned} L_A(y, x, \hat{\beta}_0, \hat{\beta}_1, \hat{\sigma}) &= \frac{1}{(\sqrt{2\pi}\hat{\sigma})^n} \exp\left\{-\frac{\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2}{2\hat{\sigma}^2}\right\} \\ &= \frac{1}{(\sqrt{2\pi}\hat{\sigma})^n} \exp\left\{-\frac{n \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2}{2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2}\right\} \\ &= \frac{1}{(\sqrt{2\pi}\hat{\sigma})^n} \exp\left\{-\frac{n}{2}\right\} \end{aligned}$$

Wilksov Λ je enak

$$\begin{aligned}\Lambda &= \frac{L_A}{L_0} = \frac{\frac{1}{(\sqrt{2\pi}\hat{\sigma}_A)^n} \exp\left\{-\frac{n}{2}\right\}}{\frac{1}{(\sqrt{2\pi}\hat{\sigma}_0)^n} \exp\left\{-\frac{n}{2}\right\}} \\ &= \frac{\hat{\sigma}_0^n}{\hat{\sigma}_A^n} \\ &= \left(\frac{\sum_{i=1}^n (y_i - \hat{\beta}_{00})^2}{\sum_{i=1}^n (y_i - \hat{\beta}_{0A} - x_i \hat{\beta}_{1A})^2} \right)^n\end{aligned}$$

Vrednost maksimuma pod alternativno domnevo izračunamo tako, da vstavimo ocenjene $\hat{\beta}_0$ in $\hat{\beta}_1$, za izračun vrednosti pod ničelno domnevo moramo oceniti še β_0 v ničelnem modelu. Dobljeni Wilksov Λ se porazdeljuje kot χ_1^2 .

Predlogi za vaje v R-u:

- Naj bo X enakomerno porazdeljena spremenljivka (med 0 in 20, zaokrožena navzdol), $\beta_0 = 15$, $\beta_1 = 4$, $\sigma = 10$. Generirajte vzorec velikosti 10, narišite podatke in vrišite populacijsko ter ocenjeno vrednost premice.

```
> set.seed(1)
> n <- 10                                #velikost vzorca
> beta0 <- 15
> beta1 <- 4
> sigma <- 10
> x <- floor(runif(n)*20)                 #navzdol zaokrožene vrednosti x
> x <- sort(x)                            #uredimo podatke po velikosti x
> y <- rnorm(n,mean=beta0+beta1*x,sd=sigma) #generiramo y-one
> plot(x,y)                               #narisemo tocke
> popul <- beta0 + beta1*x                 #populacijska vrednost premice
> lines(x,popul,col="grey",lwd=2)         #dodamo popul. vrednost premice
> fit <- lm(y~x)                           #ocenimo premico na podatkih
> summary(fit)                             #ogledamo si ocene koeficientov
> beta0h <- fit$coef[1]                    #ocenjena beta0
> beta1h <- fit$coef[2]                    #ocenjena beta1
> napoved <- beta0h + beta1h*x
> lines(x,napoved,lwd=2)                  #vrisemo ocenjeno premico na sliko
```

- Izračunajte posplošeni test razmerja verjetij v R-u

```

> fit0 <- lm(y~1) #pod nic. domnevo - le konstanta
> res0 <- y - fit0$coef #ostanki pod nicelno domnevo
> resA <- y - beta0h - beta1h*x #ostanki pod alternativno domnevo
#zanima nas razlika log verjetij - konstanto lahko izpustimo:
> logl0 <- -.5*n*log(sum(res0^2)) #loglik pod nicelno
> loglA <- -.5*n*log(sum(resA^2)) #loglik pod alternativno
> Lambda <- 2*(loglA-logl0) #Wilksov lambda
> 1-pchisq(Lambda,1) #likelihood ratio test
[1] 4.048e-05

```

5.2 Matrično računanje

Vrednosti neodvisnih spremenljivk združimo v matriko X (design matrix), vrednosti odvisne spremenljivke ter koeficientov predstavljajo vektorja Y in β :

$$X = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & \dots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix}; Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}; \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}$$

Matrika X je dimenzije $n \times (p + 1)$, kjer je p število spremenljivk. Če naš model ne bi vseboval konstante, bi prvi stolpec X izpustili.

- Zapišite vsoto vrednosti $\sum_{i=1}^n Y_i^2$ v matrični obliki.

$$Y^T Y = \begin{bmatrix} Y_1 & Y_2 & \dots & Y_n \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = Y_1^2 + Y_2^2 + \dots + Y_n^2$$

- Kaj dobimo, če matrično pomnožimo $X\beta$?

$$X\beta = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} = \begin{bmatrix} \beta_0 + \beta_1 x_1 \\ \beta_0 + \beta_1 x_2 \\ \vdots \\ \beta_0 + \beta_1 x_n \end{bmatrix} = \begin{bmatrix} E(Y_1) \\ E(Y_2) \\ \vdots \\ E(Y_n) \end{bmatrix} = E(Y)$$

- V matrični obliki oceno koeficientov po metodi najmanjših kvadratov (= po metodi največjega verjetja) zapišemo kot $\hat{\beta} = (X^T X)^{-1} X^T Y$. Pokažite, da za $p = 1$ dobite oceni:

$$\hat{\beta}_1 = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2}$$

$$\hat{\beta}_0 = \frac{\sum_{i=1}^n x_i^2 \sum_{i=1}^n y_i - \sum_{i=1}^n x_i \sum_{i=1}^n x_i y_i}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2}$$

Izračunajmo najprej $X^T X$:

$$X^T X = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} = \begin{bmatrix} n & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{bmatrix}$$

Inverz te 2×2 matrike je enak:

$$(X^T X)^{-1} = \frac{1}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2} \begin{bmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{bmatrix}$$

Izračunajmo še $X^T Y$:

$$X^T Y = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n Y_i \\ \sum_{i=1}^n x_i Y_i \end{bmatrix}$$

Velja torej

$$\begin{aligned}
(X^T X)^{-1} X^T Y &= \frac{1}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2} \begin{bmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{bmatrix} \begin{bmatrix} \sum_{i=1}^n Y_i \\ \sum_{i=1}^n x_i Y_i \end{bmatrix} \\
&= \frac{1}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2} \begin{bmatrix} \sum_{i=1}^n x_i^2 \sum_{i=1}^n Y_i - \sum_{i=1}^n x_i \sum_{i=1}^n x_i Y_i \\ -\sum_{i=1}^n x_i \sum_{i=1}^n Y_i + n \sum_{i=1}^n x_i Y_i \end{bmatrix}
\end{aligned}$$

Zgornja vrstica pri tem predstavlja oceno $\hat{\beta}_0$, spodnja pa $\hat{\beta}_1$.

- Izpeljite oceno po metodi najmanjših kvadratov še v matrični obliki. Pri tem boste potrebovali naslednje formule za matrično računanje:

$$\begin{aligned}
(A + B)^T &= A^T + B^T; \quad (A^T)^T = A; \quad (AB)^T = B^T A^T; \\
\frac{\partial \beta^T A}{\partial \beta} &= A; \quad \frac{\partial \beta^T A^T A \beta}{\partial \beta} = 2A^T A \beta
\end{aligned}$$

Namig: Kaj minimiziramo? Kako zapišemo vsoto kvadriranih ostankov v matrični obliki?

Če je v modelu ena spremenljivka, iščemo minimum funkcije

$$\sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_i)^2$$

V matrični obliki iščemo minimum funkcije

$$\begin{aligned}
(Y - X\beta)^T (Y - X\beta) &= (Y^T - \beta^T X^T)(Y - X\beta) \\
&= Y^T Y - \beta^T X^T Y - Y^T X \beta + \beta^T X^T X \beta \\
&= Y^T Y - 2\beta^T X^T Y + \beta^T X^T X \beta,
\end{aligned}$$

Pri čemer smo v zadnji vrstici uporabili, da je $(\beta^T X^T Y)^T = \beta^T X^T Y$, saj je matrika dimenzije 1×1 . Sedaj odvajamo po β

$$\frac{\partial}{\partial \beta} (Y^T Y - 2\beta^T X^T Y + \beta^T X^T X \beta) = -2X^T Y + 2X^T X \beta$$

in izenačimo z 0 (ter predpostavimo, da $X^T X$ ni singularna):

$$\begin{aligned} -2X^T Y + 2X^T X \hat{\beta} &= 0 \\ X^T X \hat{\beta} &= X^T Y \\ \hat{\beta} &= (X^T X)^{-1} X^T Y \end{aligned}$$

- Pokažite, da je ocena koeficientov nepristranska (vzemite, da so vrednosti x-ov dane in torej ne slučajne).

$$E(\hat{\beta}) = (X^T X)^{-1} X^T E(Y) = (X^T X)^{-1} X^T X \beta = \beta$$

- Izpeljite formulo za standardno napako ocenjenih koeficientov v matrični obliki. Intuitivno razložite od česa je odvisna standardna napaka koeficienta β_1 (za $p = 1$). Uporabite formulo: $\text{var}(cY) = c \text{var} Y c^T$.

$$\begin{aligned} \text{var}(\hat{\beta}) &= \text{var}[(X^T X)^{-1} X^T Y] = (X^T X)^{-1} X^T \text{var} Y [(X^T X)^{-1} X^T]^T \\ &= \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1} \end{aligned}$$

Izpeljali smo variančno-kovariančno matriko, variance so na diagonali. Standardna napaka SE_{β_1} je torej enaka

$$\begin{aligned} SE_{\beta_1} &= \sqrt{\frac{\sigma^2}{(X^T X)_{22}^{-1}}} \\ &= \sqrt{\frac{n\sigma^2}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}} \\ &= \sqrt{\frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \\ &= \sqrt{\frac{\sigma^2}{n\sigma_x^2}} = \frac{\sigma}{\sigma_x \sqrt{n}} \end{aligned}$$

Standardna napaka koeficienta je tako kot vedno odvisna od velikosti vzorca n ter razpršenosti podatkov. Vrednost σ je standardni odklon

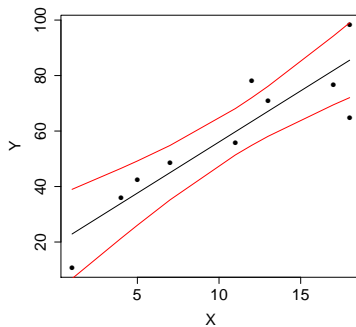
ostankov okrog premice - večja kot je, bolj se lahko zmotimo pri oceni premice. Vendar tu ni pomembna le absolutna velikost variance ostankov, zanima nas variabilnost ostankov glede na variabilnost neodvisne spremenljivke. Če je razpon x -ov majhen, je naša ocena pri isti variabilnosti ostankov manj natančna. Razložimo to na našem primeru - če bi v vzorec zajeli le posameznike, ki so se učili 3-5 ur, bi bila povezanost med spremenljivkama (npr. merjena s korelacijskim koeficientom) pri istih regresijskih koeficientih dosti manjša in zato možna večja odstopanja pri ocenjevanju.

- Kako bi izračunali interval zaupanja za napovedano premico v našem primeru ($p = 1$)? Kako bo tak interval izgledal na sliki? Izračunamo standardno napako za vsako točko posebej.

$$\text{var}(\hat{y}_i) = \text{var}(\hat{\beta}_0 + \hat{\beta}_1 x_i) = \text{var}(\hat{\beta}_0) + x_i^2 \text{var}(\hat{\beta}_1) + 2x_i \text{cov}(\hat{\beta}_0, \hat{\beta}_1)$$

Matrični izračun (za vse točke naenkrat):

$$\text{var}(\hat{Y}) = \text{var}(X\hat{\beta}) = X \text{var}(\hat{\beta}) X^T = \sigma^2 X(X^T X)^{-1} X^T$$



Slika 1: Točke na vzorcu in ocenjena premica z intervalom zaupanja.

- Recimo, da nas zanima, kako sta število ur učenja in spol (0=ženski, 1=moški) povezana z rezultatom na izpitu iz statistike. Predpostavimo model, ki vključuje interakcijo. Kako bi preverili, ali je število ur učenja pri moških povezano z rezultatom na izpitu? Model, ki ga prilagodimo podatkom, zapišemo kot:

$$Y = \beta_0 + \beta_1 \text{spol} + \beta_2 \text{ure} + \beta_3 (\text{ure} * \text{spol})$$

Ničelna domneva, ki jo želimo preveriti, je torej $H_0 : \beta_2 + \beta_3 = 0$. Zapišemo v matrični obliki. Naj bo vektor $a^T = [0,0,1,1]$, zanima nas $H_0 : a^T \beta = 0$. Varianca $a^T \beta$ je enaka $\text{var}(a^T \beta) = a^T \text{var} \beta a$, za preverjanje ničelne domneve uporabimo Waldov test.

Predlogi za vaje v R-u:

- V R-u v matrični obliki zapišite podatke, izračunajte oceno po metodi najmanjših kvadratov ter jo primerjajte z izpisom R-ove funkcije `lm`.
- Ocenite tudi standardno napako ter jo primerjajte z izpisom
- Narišite sliko napovedane premice ter ji dodajte interval zaupanja.

```
> X <- cbind(1,x)
> sigma <- summary(fit)$sigma
> inv <- solve(t(X)%*%X)
> mat <- X%*%inv%*%t(X)
> se <- sigma*sqrt(diag(mat))
> betah <- c(beta0h,beta1h)
> plot(x,y)
> lines(x,X%*%betah)
> t8 <- qt(.975,8)
> lines(x,X%*%betah - t8*se,col=2)
> lines(x,X%*%betah + t8*se,col=2)
```

5.3 Predpostavke linearne regresije

Z osnovnim modelom linearne regresije naredimo štiri predpostavke:

- Ostanke so okrog premice porazdeljeni normalno
- Varianca ostankov ni odvisna od vrednosti neodvisne spremenljivke (homoskedastičnost)
- Ostanke so med seboj neodvisni.
- Povezanost med X in Y je linearna

Kaj se zgodi z ocenami koeficientov, njihovo pričakovano vrednostjo, standardno napako in z intervali zaupanja, če je katera izmed prvih treh predpostavk kršena?

- Kaj se spremeni v izpeljavah, če ostanki okrog premice niso porazdeljeni normalno?

V tem primeru ocena koeficientov po metodi največjega verjetja ni enaka oceni po metodi najmanjših kvadratov. Ocena po metodi najmanjših kvadratov bo identična kot do sedaj, enaka bo tudi ocena standardne napake. Prav tako bo ocena po metodi najmanjših kvadratov nepristranska ocena populacijskih vrednosti. Ne moremo pa o populacijskih vrednostih sklepati ničesar več, saj ne poznamo porazdelitve ocene okrog prave vrednosti.

- Recimo, da je varianca ostankov odvisna od x . Če varianca ostankov ni enaka za vsak x , moramo varianco pisati v matrični obliki, npr.

$$\Sigma = \sigma \begin{bmatrix} w_1 & 0 & \dots & 0 & 0 \\ 0 & w_2 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & w_n \end{bmatrix}$$

Varianca sicer ne vpliva na oceno koeficientov po metodi najmanjših kvadratov, vendar pa se spremeni ocena po metodi največjega verjetja, saj moramo maksimizirati funkcijo $-2(Y - X\beta)^T \Sigma^{-1} (Y - X\beta)$:

$$\begin{aligned} (Y - X\beta)^T \Sigma^{-1} (Y - X\beta) &= \\ &= (Y^T - \beta^T X^T) \Sigma^{-1} (Y - X\beta) \\ &= Y^T \Sigma^{-1} Y - \beta^T X^T \Sigma^{-1} Y - Y^T \Sigma^{-1} X \beta + \beta^T X^T \Sigma^{-1} X \beta \\ &= Y^T \Sigma^{-1} Y - 2\beta^T X^T \Sigma^{-1} Y + \beta^T X^T \Sigma^{-1} X \beta \end{aligned}$$

in zato

$$\begin{aligned} -2X^T \Sigma^{-1} Y + 2X^T \Sigma^{-1} X \hat{\beta} &= 0 \\ X^T \Sigma^{-1} X \hat{\beta} &= X^T \Sigma^{-1} Y \\ \hat{\beta} &= (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y \end{aligned}$$

Ustrezno se spremeni tudi varianca ocene:

$$\begin{aligned}\text{var}(\hat{\beta}) &= \text{var}[(X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y] \\ &= (X^T \Sigma^{-1} X)^{-1} X^T \text{var} Y [(X^T \Sigma^{-1} X)^{-1} X^T]^T \\ &= \sigma^2 (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} X (X^T \Sigma^{-1} X)^{-1} \\ &= \sigma^2 (X^T \Sigma^{-1} X)^{-1}\end{aligned}$$

Če so vrednosti w_i znane, se v ocenjevanje koeficientov in standardne napake torej le vrine diagonalna matrika. Statistično sklepanje je enako kot do sedaj.

- Kaj pa če ostanki med seboj niso neodvisni?

Potem variančna matrika Σ ni več diagonalna (je npr. bločno diagonalna). Rezultati bodo podobni tistim v prejšnji točki, bo pa seveda ocenjevanje odvisno od tega, kaj vemo o elementih Σ .