

Kazalo

1	UVOD	1
2	ENOSTAVNA LINEARNA REGRESIJA	7
2.1	Ocenjevanje parametrov modela po metodi najmanjših kvadratov	7
2.2	Statistične lastnosti cenilk parametrov	9
2.3	Cenilka za varianco napak	12
2.4	Koeficient determinacije	13
2.5	Diagnostika modela	14
2.6	Inferenca za parametre modela	16
2.6.1	Intervali zaupanja za β_0 in β_1	17
2.6.2	Testiranje domnev za β_0 in β_1	18
2.7	Analiza variance	19
2.8	Napovedovanje	21
2.9	Interpretacija parametrov enostavnega regresijskega modela	22
2.10	Hkratna inferenca in Bonferronijev popravek	23
2.11	Primer: SKT	25
2.12	Simulacija po predpostavkah modela enostavne linearne regresije	33
3	VAJE	39
3.1	Čas teka Collina Jacksona	39
3.2	Simulacije za enostavno linearno regresijo	39
3.3	Varianca odzivne spremenljivke	40

1 UVOD

V življenju se srečujemo z raznolikimi procesi, ki jih opisujemo z različnimi spremenljivkami, za katere pogosto verjamemo, da so medsebojno povezane oziroma, da je ena spremenljivka odvisna od drugih. V poglavjih, ki sledijo se bomo ukvarjali z **linearno odvisnostjo spremenljivk**. To odvisnost bomo opisali z linearnim modelom oziroma linearnim regresijskim modelom.

Primeri linearnih modelov, ki jih bomo obravnavali so na primer odvisnost zgornjega krvnega tlaka od starosti in spola osebe, odvisnost količine padavin od nadmorske višine in geografskega položaja, odvisnost koncentracije ozona v prizemni plasti zraka od temperature zraka, relativne važnost, sončnega obseva, hitrost vetra in podobno. V naštetih primerih zgornji krvni tlak, količina padavin in koncentracija ozona predstavljajo t. i. **odzivno spremenljivko**; starost, spol, nadmorska višina, geografska dolžina in širina, temperatura zraka, in ostale meteorološke spremenljivke pa t. i. **napovedne spremenljivke**.

Zanima nas linearna odvisnost povprečne/pričakovane vrednosti odzivne spremenljivke $\mathbb{E}(y)$ od napovedne spremenljivke x , kar ponazorimo z zvezo

$$\mathbb{E}(y) = \beta_0 + \beta_1 x.$$

β_0 in β_1 sta parametra modela.

Splošno o **statističnih modelih**:

- Kvalitativno lahko statistični model opišemo

$$\text{Odziv} = \text{Signal} + \text{Šum}.$$

Statistični model zapišemo

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i.$$

Signal ($\beta_0 + \beta_1 x_i$) je del odzivne spremenljivke, ki je pojasnjen z napovedno spremenljivko in ga imenujemo tudi sistematični oziroma deterministični del statističnega modela. Šum (ε_i) predstavlja nepojasneni del odzivne spremenljivke.

- V idealnem svetu bi bila linearna zveza eksaktna, kar pomeni, da bi pri isti vrednosti napovedne spremenljivke, dobili isto vrednost odzivne spremenljivke. Ker pa svet ni idealen, so vsi modeli samo aproksimacije izbranih procesov.

... all models are approximations. Essentially, all models are wrong, but some are useful. However, the approximate nature of the model must always be borne in mind. (Box and Draper (1987), str. 424)

- S statističnim modelom ocenimo signal in šum odzivne spremenljivke.
- Z **regresijsko analizo** poskušamo poiskati najboljšo povezavo med odzivno spremenljivko in napovednimi spremenljivkami, kar pomeni, čim bolje opisati signal.
- V končni fazi statistični model predstavlja redukcijo pogosto obsežnega nabora podatkov na majhno število parametrov.

Kaj je dober model?

- Dober statistični model podatke reducira tako, da lahko na podlagi interpretacije parametrov naredimo smiselne odločitve.
- Model se dobro prilega podatkom, če sistematični del modela dobro opiše variabilnost odzivne spremenljivke, kar pomeni, da je negotovost majhna.
- Model je dober, če je parsimoničen, kar pomeni, da vsebuje smiselno majhno število parametrov.
- Pri modeliranju je vedno treba narediti kompromis med kompleksnostjo in interpretabilnostjo modela.

Modeliranje je zaporedje treh korakov, ki jih ciklično ponavljamo, dokler ne pridemo do končnega modela:

- začasna formulacija modela
- ocenjevanje parametrov
- diagnostika modela

Za napovedne spremenljivke v splošnem velja:

- so vnaprej izbrane s strani načrtovalca raziskave;
- v načrtovanem poskusu igra izbira vrednosti napovednih spremenljivk zelo pomembno vlogo pri inferenci o vplivih napovednih spremenljivk na odzivno spremenljivko;
- če vrednosti napovednih spremenljivk niso izbrane vnaprej, predpostavimo, da so vsaj točne (brez merskih napak).

Predpostavke linearnega regresijskega modela

1. Imamo **odzivno spremenljivko** (*response variable*) y in m **napovednih/pojasnjevalnih spremenljivk** (*predictors/explanatory variables*).
2. Odzivna spremenljivka y je številska, njene vrednosti so medsebojno neodvisne.
3. Napovedne spremenljivke so lahko številske (npr. starost) in/ali opisne (npr. spol, ocena vrednotena na petmestni lestvici). Vse te spremenljivke določajo t. i. **regresorje**; to so številske spremenljivke, ki so vključene v model. Iz napovednih spremenljivk dobimo regresorje na različne načine:
 - številsko spremenljivko v model vključimo direktno kot en regresor; včasih je ta spremenljivka predhodno transformirana (npr. *log*). V določenih primerih je številska spremenljivka vključena v model z več regresorji (npr. polinomska regresija, zleпки);
 - za opisno spremenljivko z d vrednostmi se v model vključi $d - 1$ regresorjev z vrednostmi 0 in 1 (neme spremenljivke, *dummy variables*);
 - dodatne regresorje lahko dobimo z upoštevanjem interakcij med obstoječimi napovednimi spremenljivkami v modelu.

Torej, m napovednih spremenljivk generira k regresorjev, ki jih označimo x_1, \dots, x_k , $k \geq m$.

4. Pričakovana vrednost odzivne spremenljivke pogojno na regresorje je linearna funkcija regresorjev x_j :

$$\mathbb{E}(y|x_1, \dots, x_k) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k. \quad (1)$$

Pri tem so β_j , $j = 0, \dots, k$, parametri modela, ki so v linearnem odnosu z odzivno spremenljivko. Te parametre ocenjujemo v postopku linearnega modeliranja.

5. Varianca odzivne spremenljivke y pogojno na regresorje x_1, \dots, x_k je konstantna, ta lastnost se imenuje **homoskedastičnost**:

$$Var(y|x_1, \dots, x_k) = \sigma^2 > 0. \quad (2)$$

Opomba: v tem gradivu je oznaka za slučajno spremenljivko, ki predstavlja odzivno spremenljivko v modelih, označena z majhno črko y in ne z veliko Y , kot je v literaturi bolj običajno.

Vrednosti odzivne in napovednih spremenljivk so dobljene na vzorcu, ki ima n enot. Linearni regresijski model za i -to enoto, $i = 1, \dots, n$, zapišemo takole:

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \varepsilon_i, \quad (3)$$

ε_i se imenuje **napaka** (*error*). Privzete predpostavke se zrcalijo v njenih lastnostih:

$$\mathbb{E}(\varepsilon_i|x_1, \dots, x_k) = 0, \quad (4)$$

$$Var(\varepsilon_i|x_1, \dots, x_k) = \sigma^2 \quad \text{oziroma} \quad Var(\varepsilon_i|x_1, \dots, x_k) = \frac{\sigma^2}{w_i}, \quad i = 1, \dots, n, \quad (5)$$

ε_i so medsebojno neodvisni. Varianca napak je v splošnem lahko obratno sorazmerna z znanimi pozitivnimi utežmi w_i , $i = 1, \dots, n$ (apriorne uteži).

Najbolj pogosta dodatna predpostavka je, da je porazdelitev y pri x_1, \dots, x_k normalna; ta model imenujemo **normalni linearni model**:

$$y|x_1, \dots, x_k \sim N(\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k, \sigma^2), \quad (6)$$

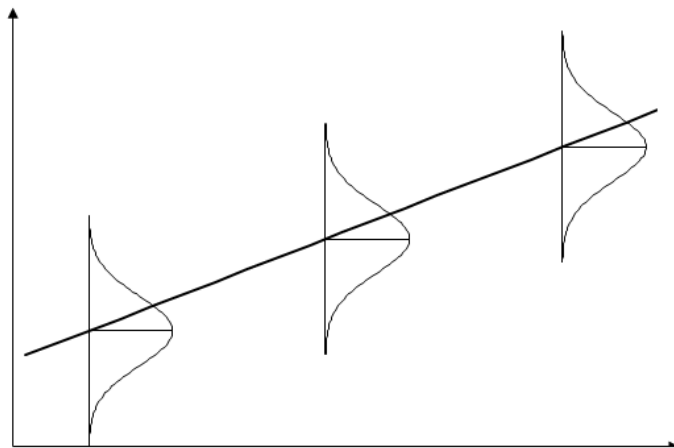
y je slučajna spremenljivka, njena porazdelitev je pri vsakem $\mathbf{x} = (x_1, \dots, x_k)^T$ normalna s povprečjem na regresijski ravnini in varianco σ^2 .

Posledično je porazdelitev napak $\varepsilon_i \sim iid N(0, \sigma^2)$. Oznaka *iid* pomeni **neodvisno enako porazdelitev** (*independent identically distributed*), kar pomeni:

- ε_i so porazdeljeni po normalni porazdelitvi;
- homogenost variance: $Var(\varepsilon_i|x_1, \dots, x_k) = \frac{\sigma^2}{w_i}$, varianca σ^2 je konstanta, prav tako v naprej postavljene uteži w_i ;
- ε_i so neodvisni.

Ilustracija predpostavk za normalni regresijski model z enim številskim regresorjem je na Sliki 1.

Normalni linearni model sodi v širšo skupino posplošenih linearnih modelov (*Generalized Linear Models*, GLM), kjer za odzivno spremenljivko privzamemo druge verjetnostne porazdelitve (logistična regresija, Poissonova regresija,...).



Slika 1: Predpostavke za normalni linearni regresijski model z eno napovedno številsko spremenljivko

V praksi parametre linearnega modela ocenjujemo na podlagi vrednosti odzivne spremenljivke in napovednih spremenljivk dobljenih na vzorcu n enot. Za ocenjevanje parametrov lahko uporabimo različne metode: **metodo najmanjših kvadratov** (*Ordinary Least Squares*, OLS), **tehtano metodo najmanjših kvadratov** (*Weighted Least Squares*, WLS), **posplošeno metodo najmanjših kvadratov** (*Generalized Least Squares*, GLS) ali **metodo največjega verjetja** (*Maximum Likelihood*, ML).

Za osnovne linearne modele sta najpogostejše uporabljeni metodi OLS in WLS. V tem primeru minimiramo vsoto kvadriranih odklonov oziroma tehtano vsoto kvadriranih odklonov vrednosti odzivne spremenljivke od njene pričakovane vrednosti, označimo jo $S(\beta_0, \beta_1, \dots, \beta_k)$. Funkcija $S(\beta_0, \beta_1, \dots, \beta_k)$ ima $k + 1$ parametrov:

$$S(\beta_0, \beta_1, \dots, \beta_k) = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik}))^2. \quad (7)$$

ali v primeru tehtane vsote kvadratov odklonov

$$S(\beta_0, \beta_1, \dots, \beta_k) = \sum_{i=1}^n w_i (y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik}))^2. \quad (8)$$

Funkcijo (7) ali (8) parcialno odvajamo po parametrih β_j , $j = 0, \dots, k$, in odvode izenačimo z 0. Dobimo t. i. normalni sistem $k + 1$ linearnih enačb. Rešitev tega sistema so cenilke parametrov, ki jih označimo b_j , $j = 0, \dots, k$.

Z modelom **napovedane vrednosti** (*fitted values*) označimo \hat{y}_i :

$$\hat{y}_i = b_0 + b_1 x_{i1} + \dots + b_k x_{ik}, \quad i = 1, \dots, n. \quad (9)$$

Razliko med dejansko vrednostjo y_i in napovedano vrednostjo \hat{y}_i imenujemo **ostanek** (*residual*) in ga označimo z e_i :

$$e_i = y_i - \hat{y}_i, \quad i = 1, \dots, n. \quad (10)$$

Matematična teorija pokaže, da so ostanki e_i nekorelirani z napovedanimi vrednostmi \hat{y}_i , kar uporabljamo pri analizi modela z grafičnimi prikazi (dokaz v poglavju 3.4).

Varianca ostanka je (dokaz v poglavju 3.4):

$$\text{Var}(e_i) = \sigma^2(1 - h_{ii}). \quad (11)$$

Pri izračunu moramo varianco σ^2 oceniti, cenilko variance označimo s^2 . Količina h_{ii} se imenuje vzvod (*leverage* ali *hat-value*); h_{ii} je vrednost med $1/n$ in 1, o tem kasneje.

Če v izrazu (11) uporabimo ocenjeno varianco s^2 , lahko izračunamo **standardizirane ostanke** e_{si} :

$$e_{si} = \frac{y_i - \hat{y}_i}{s\sqrt{1 - h_{ii}}}, \quad i = 1, \dots, n, \quad (12)$$

za katere pri pogoju, da je $n \gg k$ velja, da je njihova porazdelitev približno $N(0, 1)$.

Ali so predpostavke modela izpolnjene, ugotavljamo z analizo ostankov in standardiziranih ostankov.

Če lahko privzamemo normalni linearni model, matematična statistika pove, da so porazdelitve parametrov normalne z znanimi parametri, kar omogoča dodatne izračune:

- za vsako cenilko parametra lahko izračunamo njeno standardno napako;
- izračunamo interval zaupanja za vsak parameter modela;
- testiramo lahko domneve o parametrih modela;
- izračunamo napovedi in intervale zaupanja za povprečno napoved in za posamično napoved.

2 ENOSTAVNA LINEARNA REGRESIJA

Če je v modelu (3) samo ena napovedna spremenljivka x , ki je številska, ta model imenujemo enostavna linearna regresija:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, n. \quad (13)$$

ε_i so neodvisne slučajne spremenljivke s pričakovano vrednostjo $\mathbb{E}(\varepsilon_i) = 0$ in konstantno varianco $Var(\varepsilon_i) = \sigma^2$ za vsak $i = 1, \dots, n$. Ocenjujemo dva parametra modela β_0 in β_1 . V nadaljevanju bomo predstavili, kako ju ocenimo po metodi najmanjših kvadratov (OLS).

2.1 Ocenjevanje parametrov modela po metodi najmanjših kvadratov

Lema 1.1: Po metodi najmanjših kvadratov sta cenilki parametrov β_0 in β_1

$$b_0 = \bar{y} - b_1 \bar{x},$$

$$b_1 = \frac{SS_{xy}}{SS_{xx}}.$$

V zapisu smo uporabili izraza SS_{xx} in SS_{xy} . SS_{xx} je vsota kvadriranih odklonov za napovedno spremenljivko

$$SS_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n (x_i - \bar{x})x_i,$$

in SS_{xy} vsota produktov odklonov napovedne in odzivne spremenljivke

$$SS_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i(y_i - \bar{y}) = \sum_{i=1}^n y_i(x_i - \bar{x}).$$

Dokaz: Vsoto kvadriranih odklonov y_i od njene pričakovane vrednosti $(\beta_0 + \beta_1 x_i)$:

$$S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 \quad (14)$$

parcialno odvajamo po β_0 in β_1 :

$$\frac{\partial S}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i),$$

$$\frac{\partial S}{\partial \beta_1} = -2 \sum_{i=1}^n x_i (y_i - \beta_0 - \beta_1 x_i).$$

Ko odvoda izenačimo z 0, dobimo sistem dveh linearnih enačb:

$$\sum_{i=1}^n y_i = nb_0 + b_1 \sum_{i=1}^n x_i, \quad (15)$$

$$\sum_{i=1}^n x_i y_i = b_0 \sum_{i=1}^n x_i + b_1 \sum_{i=1}^n x_i^2. \quad (16)$$

Z rešitvijo sistema enačb dobimo ceniki za β_0 in β_1 , označimo ju b_0 in b_1 . Iz (15) sledi:

$$b_0 = \bar{y} - b_1 \bar{x}, \quad (17)$$

ko ta rezultat uporabimo v (16) dobimo

$$\begin{aligned} \sum_{i=1}^n (x_i y_i - x_i(\bar{y} - b_1 \bar{x}) - x_i b_1 x_i) &= \sum_{i=1}^n (x_i y_i - x_i \bar{y} + b_1 x_i \bar{x} - b_1 x_i^2) = 0 \\ \sum_{i=1}^n (x_i y_i - x_i \bar{y}) &= b_1 \sum_{i=1}^n (x_i^2 - x_i \bar{x}). \end{aligned}$$

Iz tega sledi:

$$\begin{aligned} b_1 &= \frac{\sum_{i=1}^n (x_i y_i - x_i \bar{y})}{\sum_{i=1}^n (x_i^2 - x_i \bar{x})} \\ &= \frac{\sum_{i=1}^n x_i (y_i - \bar{y})}{\sum_{i=1}^n x_i (x_i - \bar{x})} \\ &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ &= \frac{SS_{xy}}{SS_{xx}}, \end{aligned} \quad (18)$$

Model enostavne linearne regresije predstavlja premico, ki se po metodi najmanjših kvadratov najbolj prilaga podatkom. Cenilka b_0 predstavlja presečišče premice z ordinatno osjo, b_1 pa njen naklon.

Enačba regresijske premice je:

$$\hat{y} = b_0 + b_1 x. \quad (19)$$

Model velja na intervalu vrednosti napovedne spremenljivke $[x_{min}, x_{max}]$.

2.2 Statistične lastnosti cenilk parametrov

Ob danih predpostavkah za linearni regresijski model sta cenilki b_0 in b_1 funkciji y_i in posledično tudi funkciji ε_i . Statistične lastnosti cenilk opišemo na podlagi njihove pristraskosti (*bias*) in variance. Pri tem bomo potrebovali nekaj pravil matematične statistike.

Lema 1.2: Za slučajno spremenljivko X in konstanti $a, b \in \mathbb{R}$ velja

$$\mathbb{E}(a + bX) = a + b\mathbb{E}(X)$$

Lema 1.3: Za slučajni spremenljivki X, Y in konstanti $a, b \in \mathbb{R}$ velja

$$\mathbb{E}(aX + bY) = a\mathbb{E}(X) + b\mathbb{E}(Y)$$

$$\text{Var}(aX) = a^2\text{Var}(X)$$

$$\text{Var}(aX + bY) = a^2\text{Var}(X) + b^2\text{Var}(Y) - 2ab\text{Cov}(X, Y)$$

Definicija 1.1: Cenilka $\hat{\theta}$ za parameter θ je nepristranska (*unbiased*) če je pristraskost (*bias*) enaka nič:

$$\mathbb{E}(\hat{\theta}) - \theta = 0.$$

Lema 1.4: Cenilki b_0 in b_1 sta nepristranski cenilki parametrov β_0 in β_1 , velja $\mathbb{E}(b_i) = \beta_i$, $i = 0, 1$.

Dokaz: Najprej izračunajmo pristraskost za b_1 . pri tem bomo velikokrat uporabili Lemo 1.2.

$$\mathbb{E}(b_1) = \mathbb{E}\left(\frac{SS_{xy}}{SS_{xx}}\right) = \frac{1}{SS_{xx}}\mathbb{E}(SS_{xy})$$

Iz predpostavk modela vemo, da je $\mathbb{E}(y_i) = \beta_0 + \beta_1 x_i$ in

$$\begin{aligned}
 \mathbb{E}(SS_{xy}) &= \mathbb{E} \left(\sum_{i=1}^n (x_i - \bar{x}) y_i \right) \\
 &= \sum_{i=1}^n (x_i - \bar{x}) \mathbb{E}(y_i) \\
 &= \sum_{i=1}^n (x_i - \bar{x}) (\beta_0 + \beta_1 x_i) \\
 &= \sum_{i=1}^n (x_i - \bar{x}) \beta_0 + \sum_{i=1}^n (x_i - \bar{x}) \beta_1 x_i \\
 &= \sum_{i=1}^n (x_i - \bar{x}) \beta_1 x_i, \quad \text{ker je} \quad \sum_{i=1}^n (x_i - \bar{x}) = 0 \\
 &= \beta_1 \sum_{i=1}^n (x_i - \bar{x}) x_i \\
 &= \beta_1 S_{xx}.
 \end{aligned}$$

Za b_0 velja

$$\begin{aligned}
 \mathbb{E}(b_0) &= \mathbb{E}(\bar{y}) - \bar{x} \mathbb{E}(b_1) \\
 &= \mathbb{E}(\bar{y}) - \bar{x} \beta_1,
 \end{aligned}$$

velja $\mathbb{E}(\bar{y}) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(y_i)$

$$\begin{aligned}
 \mathbb{E}(\bar{y}) &= \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_i) \\
 &= \beta_0 + \beta_1 \bar{x},
 \end{aligned}$$

sledi, da je $\mathbb{E}(b_0) = \beta_0$.

Lema 1.5: varianci in kovarianca cenilk parametrov so

$$\begin{aligned}
 Var(b_0) &= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{SS_{xx}} \right), \\
 Var(b_1) &= \frac{\sigma^2}{SS_{xx}}, \\
 Cov(b_0, b_1) &= -\frac{\sigma^2 \bar{x}}{SS_{xx}}.
 \end{aligned} \tag{20}$$

Dokaz: Varianca cenilke b_1 je

$$Var(b_1) = Var\left(\frac{SS_{xy}}{SS_{xx}}\right) = \frac{Var(SS_{xy})}{SS_{xx}^2},$$

ugotoviti moramo, kakšna je $Var(SS_{xy})$

$$\begin{aligned} Var(SS_{xy}) &= Var\left(\sum_{i=1}^n (x_i - \bar{x})y_i\right) \\ &= \sum_{i=1}^n (x_i - \bar{x})^2 Var(y_i) \\ &= \sum_{i=1}^n (x_i - \bar{x})^2 \sigma^2 \\ &= \sigma^2 SS_{xx} \end{aligned}$$

in

$$Var(b_1) = \frac{\sigma^2 SS_{xx}}{SS_{xx}^2} = \frac{\sigma^2}{SS_{xx}}.$$

Varianca cenilke b_0 je

$$\begin{aligned} Var(b_0) &= Var(\bar{y} - b_1 \bar{x}) \\ &= Var(\bar{y}) + Var(b_1 \bar{x}) - 2Cov(\bar{y} b_1 \bar{x}), \end{aligned}$$

poglejmo, kako izrazimo $Cov(\bar{y} b_1)$.

$$\begin{aligned} Cov(\bar{y} b_1) &= \frac{1}{SS_{xx}} Cov(\bar{y}, SS_{xy}) \\ &= \frac{1}{n SS_{xx}} \sum_{i=1}^n Cov(y_i, SS_{xy}) \\ &= \frac{1}{n SS_{xx}} \sum_{i=1}^n Cov(y_i, \sum_{j=1}^n (x_j - \bar{x}) y_j) \\ &= \frac{1}{n SS_{xx}} \sum_{i=1}^n \sum_{j=1}^n (x_j - \bar{x}) Cov(y_i, y_j) \\ &= \frac{n}{n SS_{xx}} \sum_{j=1}^n (x_j - \bar{x}) \sigma^2 \\ &= \frac{\sigma^2}{SS_{xx}} \sum_{j=1}^n (x_j - \bar{x}) \\ &= 0. \end{aligned}$$

Sledi

$$\begin{aligned} Var(b_0) &= Var(\bar{y}) + \bar{x}^2 Var(b_1) \\ &= \frac{\sigma^2}{n} + \bar{x}^2 \frac{\sigma^2}{SS_{xx}} \\ &= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{SS_{xx}} \right). \end{aligned}$$

Za izpeljavo kovariance $Cov(b_0, b_1)$ uporabimo rešitev $b_0 = \bar{y} - b_1 \bar{x}$:

$$\begin{aligned} Cov(b_0, b_1) &= Cov(\bar{y} - b_1 \bar{x}, b_1) \\ &= Cov(\bar{y}, b_1) - \bar{x} Var(b_1) \\ &= -\bar{x} \frac{\sigma^2}{SS_{xx}}. \end{aligned}$$

Iz enačb (20) vidimo, da sta varianci in kovarianca cenilk parametrov enostavnega linearnega modela odvisni od vrednosti napovedne spremenljivke x_i , $i = 1, \dots, n$, in od variance napake σ^2 . Vrednosti x_i so znane, varianco napake pa moramo v praksi oceniti. Glede na predpostavko linearnega regresijskega modela (5) je σ^2 varianca napak ε_i gledano pogojno na regresorje.

2.3 Cenilka za varianco napak

Napake $\varepsilon_i = y_i - \beta_0 - \beta_1 x_i$ ocenimo z t. i. **ostankom** $e_i = y_i - \hat{y}_i = y_i - b_0 - b_1 x_i$, $i = 1, \dots, n$ (*residual*) in definiramo **vsoto kvadriranih ostankov** ($SS_{residual}$):

$$SS_{residual} = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2. \quad (21)$$

Matematična statistika pokaže, da je nepristranska cenilka za σ^2 :

$$s^2 = \frac{SS_{residual}}{n - 2}. \quad (22)$$

V imenovalcu (22) delimo z $n - 2$ namesto z n , saj smo dve stopinji prostosti porabili za oceno parametrov modela. Uporabna formula za cenilko variance napak je

$$s^2 = \frac{SS_{yy} - b_1^2 SS_{xx}}{n - 2}. \quad (23)$$

Glede odzivne spremenljivke smo do zdaj povedali:

- odzivno spremenljivko označimo z y_i in na podlagi predpostavke linearnega modela velja

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

- povprečje oziroma pričakovana vrednost odzivne spremenljivke je

$$\mathbb{E}(y_i) = \beta_0 + \beta_1 x_i$$

- vzorčno povprečje odzivne spremenljivke je

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

- z modelom prilagojene vrednosti odzivne spremenljivke so

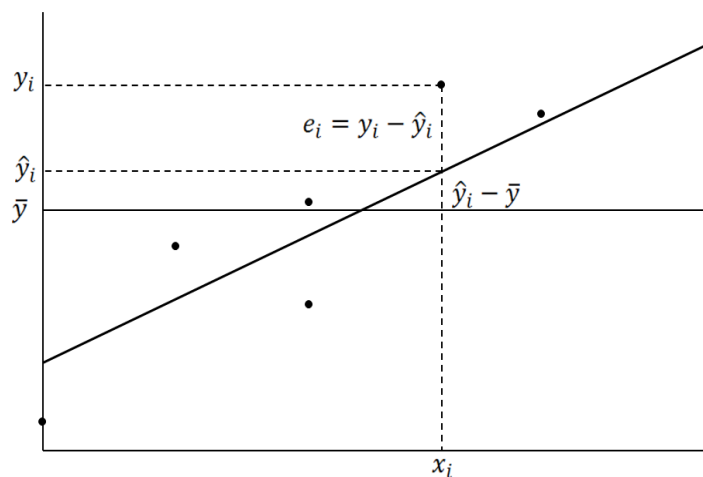
$$\hat{y}_i = b_0 + b_1 x_i$$

- v nadaljevanju bomo videli, da se variabilnost odzivne spremenljivke razdeli na dva dela

$$SS_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2 = SS_{model} + SS_{residual}$$

2.4 Koeficient determinacije

Koeficient determinacije je enostavna mera, ki opredeljuje kakovost linearnega regresijskega modela; za njen izračun ne potrebujemo nobenih predpostavk. Izhodišče za njegov izračun je razvidno iz Slike 2.



Slika 2: Izhodišče za izračun koeficienta determinacije: $y_i - \bar{y} = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})$

Odklon od povprečja zapišemo $y_i - \bar{y} = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})$. To enačbo kvadriramo in seštejemo po vseh enotah, $i = 1, \dots, n$. Z upoštevanjem formul za b_0 in za b_1 (17,18) se pokaže, da je $2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) = 0$ in ostane:

$$\begin{aligned} \sum_{i=1}^n (y_i - \bar{y})^2 &= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2, \\ SS_{yy} &= SS_{model} + SS_{residual}. \end{aligned} \quad (24)$$

Vsota kvadriranih odklonov (*Sum of Squares*) za odzivno spremenljivko y , SS_{yy} , se razdeli na dva dela: na del, ki ga pojasni regresijski model, SS_{model} , in na del, ki ostane z regresijskim modelom nepojasnen, $SS_{residual}$.

Koeficient determinacije R^2 je delež variabilnosti za y , ki je pojasnjen z regresijskim modelom:

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{SS_{model}}{SS_{yy}}. \quad (25)$$

Lastnosti koeficienta determinacije:

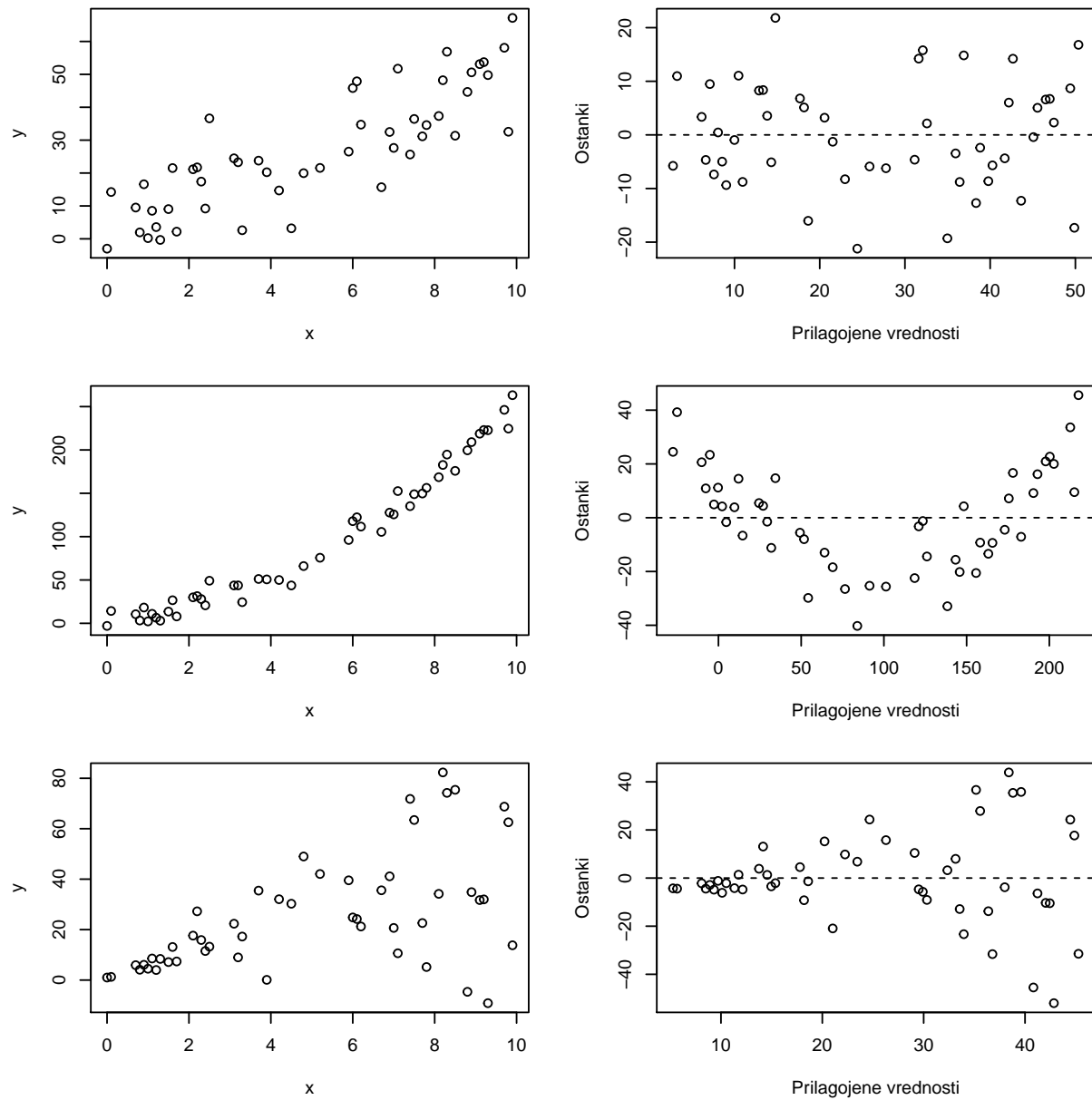
- je nenegativna vrednost;
- je manjši ali enak 1; ima vrednost 1, če je $SS_{model} = SS_{yy}$, ko so vse točke na premici;
- R^2 je odvisen od zaloge vrednosti napovedne spremenljivke;
- previdni moramo biti pri uporabi R^2 v kontekstu multiple regresije, saj vsak dodani regresor poveča vrednost R^2 , tudi če je vpliv tega regresorja na odzivno spremenljivko statistično nepomemben.

2.5 Diagnostika modela

Diagnostika je namenjena preverjanju predpostavk linearnega modela. V praksi na podlagi podatkov ocenimo parametre modela, za tem pa je potrebno preveriti, ali je bilo tako modeliranje upravičeno. Preveriti moramo sledeče:

- linearnost odvisnosti odzivne spremenljivke od napovedne spremenljivke (razsevni grafikon y glede na x);
- varianca napak oziroma varianca odzivne spremenljivke pogojno na napovedne spremenljivke je konstantna (slika ostankov glede na napovedane vrednosti);
- pričakovana vrednost napak je 0 (slika ostankov glede na napovedane vrednosti);
- napake so medsebojno neodvisne (težko preveriti, verjamemo, da so bili podatki pridobljeni z ustreznim načinom vzorčenja).

Osnova za diagnostiko modela so grafični prikazi ostankov $e_i = y_i - \hat{y}_i$.



Slika 3: Razsevni grafikoni za y in x (levo) in slike ostankov (desno), slučajna porazdelitev ostankov s pričakovano vrednostjo 0 (zgoraj), porazdelitev ostankov v obliki črke U izraža nelinearnost odnosa (sredina), porazdelitev ostankov v obliki klina izraža nekonstantno varianco

2.6 Inferenca za parametre modela

Ob predpostavki $\varepsilon_i \sim iid N(0, \sigma^2)$ in posledičnem dejstvu, da sta cenilki za presečišče in naklon premice linearni funkciji normalno porazdeljenih spremenljivk, velja, da je tudi njuna porazdelitev normalna:

$$\begin{aligned} b_0 &\sim N\left(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{SS_{xx}}\right)\right), \\ b_1 &\sim N\left(\beta_1, \frac{\sigma^2}{SS_{xx}}\right). \end{aligned} \tag{26}$$

Za porazdelitev cenilke variance napak s^2 velja (dokaz pri OTS)

$$\frac{(n-2)s^2}{\sigma^2} \sim \chi_{n-2}^2 \tag{27}$$

Cenilki varianc parametrov modela $s_{b_0}^2$ in $s_{b_1}^2$ izračunamo tako, da v (20) σ^2 zamenjamo z s^2 .

Lema 1.6: Če sta $X \sim N(0, 1)$ in $Y \sim \chi_n^2$ neodvisno slučajni spremenljivki, potem velja

$$\frac{X}{\sqrt{Y/n}} \sim t_n$$

(dokaz pri OTS)

Lema 1.7: Če uporabimo predstavljene lastnosti cenilk in Lemo 1.6, lahko pod predpostavko normalne porazdelitve napak ε_i , za $i = 1, \dots, n$ izpeljemo

$$\begin{aligned} \frac{b_0 - \beta_0}{s\sqrt{\left(\frac{1}{n} + \frac{\bar{x}^2}{SS_{xx}}\right)}} &\sim t_{n-2} \quad \text{in} \quad \frac{b_1 - \beta_1}{\frac{s}{\sqrt{SS_{xx}}}} \sim t_{n-2}. \\ \frac{b_0 - \beta_0}{sb_0} &\sim t_{n-2} \quad \text{in} \quad \frac{b_1 - \beta_1}{sb_1} \sim t_{n-2}. \end{aligned}$$

Dokaz: Izraz

$$\frac{b_0 - \beta_0}{s\sqrt{\left(\frac{1}{n} + \frac{\bar{x}^2}{SS_{xx}}\right)}},$$

v števcu in imenovalcu delimo z $Var(b_0)$ in izraz v imenovalcu postavimo pod koren ter števec in imenovalec pomnožomo z $(n-2)$:

$$\frac{\frac{b_0 - \beta_0}{\sigma\sqrt{\left(\frac{1}{n} + \frac{\bar{x}^2}{SS_{xx}}\right)}}}{\sqrt{\frac{(n-2)s^2}{\sigma^2}}}$$

V števcu je slučajna spremenljivka porazdeljena $N(0, 1)$, v imenovalcu pa je koren slučajne spremenljivke porazdeljene χ_{n-2}^2 .

Na enak način lahko to pokažemo za

$$\frac{b_1 - \beta_1}{\frac{s}{\sqrt{SS_{xx}}}}.$$

2.6.1 Intervali zaupanja za β_0 in β_1

Za slučajni vektor $Y = (Y_1, \dots, Y_n)$ s porazdelitvijo, ki je odvisna od nepoznanega parametra θ iščemo interval zaupanja za θ .

- Iščemo interval $(L(Y), U(Y))$ za katerega velja

$$\mathbb{P}(L(Y) < \theta < U(Y)) \geq 1 - \alpha.$$

- Spodnja in zgornja meja intervala zaupanja sta slučajni spremenljivki, θ pa je konstanta.
- Interval $(L(Y), U(Y))$ pokrije vrednost θ z verjetnostjo $1 - \alpha$.
- Če poskus ponovimo velikokrat, in vsakič izračunamo $(L(Y), U(Y))$, bomo dobili vsaj $100(1 - \alpha)\%$ slučajnih intervalov, ki vsebujejo pravo vrednost θ .

Pri izračunu intervalov zaupanja za β_0 in β_1 izhajamo iz porazdelitve statistik

$$\frac{b_0 - \beta_0}{s_{b_0}} \sim t_{n-2},$$

$$\frac{b_1 - \beta_1}{s_{b_1}} \sim t_{n-2}.$$

Naj bo $t_{\alpha/2; n-2}$ $(1 - \alpha/2)$ -ti kvantil t -porazdelitve s stopinjami prostosti $n - 2$.

$$P(-t_{\alpha/2; n-2} \leq \frac{b_1 - \beta_1}{s_{b_1}} \leq t_{\alpha/2; n-2}) = 1 - \alpha,$$

ko preuredimo zgornji zapis, dobimo

$$P(b_1 - t_{\alpha/2; n-2}s_{b_1} \leq \beta_1 \leq b_1 + t_{\alpha/2; n-2}s_{b_1}) = 1 - \alpha,$$

$100(1 - \alpha)\%$ intervala zaupanja (IZ) za β_0 in β_1 sta:

$$\begin{aligned} & (b_0 - t_{\alpha/2; n-2}s_{b_0}, \quad b_0 + t_{\alpha/2; n-2}s_{b_0}), \\ & (b_1 - t_{\alpha/2; n-2}s_{b_1}, \quad b_1 + t_{\alpha/2; n-2}s_{b_1}), \end{aligned}$$

$t_{\frac{\alpha}{2}; n-2}$ je $(1 - \alpha/2)100$ -ti centil t -porazdelitve s stopinjami prostosti $n - 2$.

2.6.2 Testiranje domnev za β_0 in β_1

Testiramo ničelno domnevo za β_1

$$H_0 : \beta_1 = \beta$$

nasproti alternativni domnevi

$$H_1 : \beta_1 \neq \beta.$$

Uporabimo testno statistiko

$$T = \frac{b_1 - \beta}{\frac{s}{\sqrt{SS_{xx}}}},$$

ki je pod ničelno domnevo porazdeljena t_{n-2} . Ničelno domnevo zavrnemo pri stopnji značilnosti α , če je

$$T < -t_{\alpha/2;n-2} \quad \text{ali} \quad T > t_{\alpha/2;n-2}$$

Ob tem je verjetnost, da T leži v območju zavrnitve ničelne domneve, ko je ta pravilna, α .

Za dani vzorec podatkov izračunamo vrednost testne statistike $T = t$. Verjetnost

$$p = P(|T| \geq |t| | \beta_1 = \beta)$$

predstavlja t. i. **p -vrednost**. Če je $p < \alpha$ ničelno domnevo zavrnemo in če je $p \geq \alpha$ ničelne domneve ne moremo zavrniti.

Enako kot ničelno domnevo z dvostransko alternativno domnevo, lahko testiramo tudi ničelno domnevo z **enostransko alternativno domnevo** $H_1 : \beta_1 < \beta$ ali $H_1 : \beta_1 > \beta$. V tem primeru je p -vrednost polovica p -vrednosti pri testiranju dvostranske alternativne domneve.

Podobno lahko testiramo ničelno domnevo za β_0

$$H_0 : \beta_0 = \beta$$

nasproti alternativni domnevi

$$H_1 : \beta_0 \neq \beta.$$

Uporabimo testno statistiko

$$T = \frac{b_0 - \beta}{\sigma \sqrt{\left(\frac{1}{n} + \frac{\bar{x}^2}{SS_{xx}}\right)}},$$

ki je pod ničelno domnevo porazdeljena t_{n-2} . Ničelno domnevo zavrnamo pri stopnji značilnosti α , če je

$$T < -t_{\alpha/2;n-2} \quad \text{ali} \quad T > t_{\alpha/2;n-2}$$

Ob tem je verjetnost, da T leži v območju zavrnitve ničelne domneve, ko je ta pravilna, α .

Namen linearnega modela je, da ugotovimo ali obstaja pomemben odnos med y in x . Če je odvisnost obstaja, potem jo želimo opisati (naraščajoča/padajoča, tesna/šibeka). Če y ni odvisna od x , potem nam vrednosti spremenljivke x ne pomagajo napovedati y in najboljša napoved za y je vodoravna premica $y = \bar{y}$. Za vodoravno premico velja $\beta_1 = 0$ in ničelna domneva, ki pravi, da y ni odvisen od x je

$$H_0 : \beta_1 = 0.$$

Alternativna domneva je

$$H_1 : \beta_1 \neq 0.$$

Osnovni povzetek regresijskega modela v R vsebuje rezultat testiranja te domneve na podlagi T -statistike. Prav tako v povzetku najdemo tudi rezultat testiranja ničelne domneve

$$H_0 : \beta_0 = 0.$$

nasproti alternativni domnevi

$$H_1 : \beta_0 \neq 0.$$

Ta ničelna domneva je redko vsebinsko zanimiva in jo v rezultatih ignoriramo.

2.7 Analiza variance

Videli smo, da se vsota kvadriranih odklonov za odzivno spremenljivko y razdeli na dva dela: na del, ki ga pojasni regresijski model (SS_{model}), in na del, ki ostane z regresijskim modelom nepojasnen ($SS_{residual}$). Vsak od teh členov ima pripadajoče stopinje prostosti Df (*Degrees of freedom*). Stopinje prostosti za skupno variabilnost Df_{yy} so $n - 1$, kjer je n število enot, ki so vključene v model. Stopinje prostosti za regresijski model Df_{model} z eno napovedno spremenljivko so $k = 1$, kjer je $k + 1 = 2$ število ocenjenih parametrov v modelu. Stopinje prostosti za ostanek, $Df_{residual}$ predstavljajo razliko, torej $n - 2$.

Opredelimo še srednji kvadrirani odklon (*Mean Square*) za posamezno komponento, $MS = SS/Df$. Vse omenjene količine po virih variabilnosti uredimo v t. i. tabelo analize variance (ANOVA).

Tabela 1: Shema tabele ANOVA za enostavni linearni regresijski model

Vir variabilnosti	Df	SS	$MS = SS/df$	F
Model	1	SS_{model}	MS_{model}	$MS_{model}/MS_{residual}$
Ostanek (<i>Residual</i>)	$n - 2$	$SS_{residual}$	$MS_{residual}$	
Skupaj	$n - 1$	SS_{yy}		

Ob predpostavki $\varepsilon_i \sim iid N(0, \sigma^2)$ za F -statistiko velja, da je njena ničelna porazdelitev F -porazdelitev s stopinjami prostosti $SP_{model} = 1$ in $SP_{residual} = n - 2$ (Dokaz: Rice, 1995, str. 448).

Iz tabele ANOVA dobimo:

- cenilko za varianco σ^2 , ki jo označimo s^2 . Teorija pove, da je po metodi najmanjših kvadratov ta cenilka enaka $MS_{residual}$. Količino s imenujemo **standardna napaka regresije** (*Residual standard error*).
- F -statistika testira domnevo o ničelnem vplivu napovedne spremenljivke:
 $H_0 : \beta_1 = 0$,
 $H_1 : \beta_1 \neq 0$.

F -test dobi večji pomen v primeru, ko imamo k napovednih spremenljivk v modelu, ker testira ničelno domnevo o hkratni ničnosti vseh parametrov v modelu

$$H_0 : \beta_1 = \dots = \beta_k = 0$$

nasproti alternativni domnevi

$$H_1 : vsaj en \beta_i \neq 0, \quad i = 1, \dots, k.$$

Povezava med T in F statistiko:

Lema 1.8: Če je slučajna spremenljivka X porazdeljena po t -porazdelitvi s stopinjami prostosti v , $X \sim t_v$, potem je slučajna spremenljivka X^2 porazdeljena po F -porazdelitvi s stopinjami prostosti 1 in v , $X^2 \sim F_{1,v}$.

- za testiranje domneve $\beta_1 = 0$ je T -statistika pod ničelno domnevo porazdeljena po t -porazdelitvi s stopinjami prostosti $n - 2$

$$T = \frac{b_1}{\frac{s}{\sqrt{SS_{xx}}}} \sim t_{n-2}$$

- če zgornji izraz kvadriramo, dobimo F -statistiko

$$F = \frac{b_1^2 SS_{xx}}{s^2} \sim F_{1, n-2}$$

- če upoštevamo, da je $SS_{model} = b_1^2 S_{xx}$ in $s^2 = SS_{residual}/(n-2)$, F -statistiko zapišemo še drugače:

$$F = \frac{SS_{model}/1}{SS_{residual}/(n-2)} \sim F_{1,n-2}.$$

Vidimo, da je F -statistika skalirano razmerje vsote kvadriranih odklonov modela in ostanka. Če je SS_{model} veliko večja od $SS_{residual}$, bo F -statistika velika in ničelno domnevo, ki pravi, da regresorji niso uporabni pri napovedovanju odzivne spremenljivke, bomo zavrnili.

2.8 Napovedovanje

Pogosto naredimo linearni model z namenom napovedovanja vrednosti odzivne spremenljivke pri izbrani vrednosti napovedne spremenljivke x_0 . Poleg napovedi nas zanima tudi njena natančnost. Vemo, da na podlagi ocenjenih parametrov modela lahko izračunamo **povprečno napoved** $\hat{y}(x_0)$

$$\hat{y}(x_0) = b_0 + b_1 x_0,$$

prava napoved je

$$\mathbb{E}(y(x_0)) = \beta_0 + \beta_1 x_0,$$

z y_0 pa označimo eno izmed možnih vrednosti odzivne spremenljivke pri x_0

$$y_0 = \beta_0 + \beta_1 x_0 + \varepsilon,$$

kjer je $\varepsilon(0, \sigma^2)$.

Pokažemo lahko, da je porazdelitev povprečne napovedi pri x_0 normalna:

$$\hat{y}(x_0) \sim N\left(\beta_0 + \beta_1 x_0, \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_{xx}}\right)\right).$$

in velja, podobno kot pri intervalih zaupanja za parametre modela, da je statistika

$$\frac{\hat{y}(x_0) - \beta_0 - \beta_1 x_0}{s^2 \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_{xx}}}} \sim t_{n-2}.$$

$100(1 - \alpha)\%$ interval zaupanja za povprečno napoved $\beta_0 + \beta_1 x_0$ pri x_0 je

$$\left(\hat{y}(x_0) - t_{\alpha/2; n-2} s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_{xx}}}, \hat{y}(x_0) + t_{\alpha/2; n-2} s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_{xx}}} \right).$$

Napoved nove vrednosti y_0 imenujemo **posamična napoved** in je enaka povprečni napovedi $\hat{y}(x_0)$:

$$\mathbb{E}(\hat{y}(x_0) - y_0) = \beta_0 + \beta_1 x_0 - \beta_0 + \beta_1 x_0 = 0$$

Ker sta y_0 in $\hat{y}(x_0)$ neodvisni, velja tudi:

$$\begin{aligned} \text{Var}(\hat{y}(x_0) - y_0) &= \text{Var}(\hat{y}(x_0)) + \text{Var}(y_0) \\ &= \sigma^2 + \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_{xx}} \right) \\ &= \sigma^2 \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_{xx}} \right). \end{aligned}$$

Varianca posamične napovedi je za člen σ^2 večja od variance povprečne napovedi.

Podobno kot pri intervalu zaupanja za posamično napoved lahko definiramo T -statistiko:

$$\frac{\hat{y}(x_0) - y_0}{s^2 \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_{xx}}}} \sim t_{n-2},$$

in $100(1 - \alpha) \%$ interval zaupanja za posamično napoved je

$$\left(\hat{y}(x_0) - t_{\alpha/2; n-2} s \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_{xx}}}, \hat{y}(x_0) + t_{\alpha/2; n-2} s \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_{xx}}} \right).$$

Širina intervalov zaupanja za povprečno kot za posamično napoved je odvisna od velikosti vzorca n , variance napak, ocenjene z s^2 in od položaja x_0 najmanjša je pri povprečju \bar{x} in narašča s kvadratom razdalje od povprečja. Meje pripadajočih intervalov zaupanja so na hiperbolah (Slika 7). Omenjena natančnost napovedi velja na intervalu $[x_{min}, x_{max}]$, ekstrapolacija je nevarna.

2.9 Interpretacija parametrov enostavnega regresijskega modela

Parameter β_1 v modelu enostavne linearne regresije predstavlja naklon premice, kar pomeni, če se vrednost napovedne spremenljivke x poveča za 1, se pričakovana vrednost odzivne spremenljivke poveča za β_1 :

Vzemimo točko na premici pri x_0 , $\hat{y}(x_0) = \beta_0 + \beta_1 x_0$, in točko na premici pri vrednosti $x_0 + 1$:

$$\hat{y}(x_0 + 1) = \beta_0 + \beta_1(x_0 + 1) = \beta_0 + \beta_1 x_0 + \beta_1 = \hat{y}(x_0) + \beta_1,$$

Iz tega sledi, da je

$$\beta_1 = \hat{y}(x_0 + 1) - \hat{y}(x_0).$$

2.10 Hkratna inferenca in Bonferronijev popravek

S simulacijami lahko pokažemo, da posamična intervala zaupanja za parametra β_0 in β_1 ne določata tudi skupnega območja zaupanja za oba parametra hkrati.

Podobno velja tudi za intervala zaupanja za povprečno ali posamično napoved modela v točki x_0 . Označimo z $\mathbb{E}(x_0; \alpha)$ $100(1 - \alpha)$ % interval zaupanja za povprečno napoved:

$$\mathbb{E}(x_0; \alpha) = \left(\hat{y}(x_0) - t_{\alpha/2; n-2} s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_{xx}}}, \hat{y}(x_0) + t_{\alpha/2; n-2} s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_{xx}}} \right).$$

Vemo, da velja

$$\mathbb{P}(\beta_0 + \beta_1 x_0 \in \mathbb{E}(x_0; \alpha)) = 1 - \alpha.$$

za napovedi v dveh različnih točkah $x_{0,1}$ in $x_{0,2}$ pa ugotovimo, da v splošnem

$$\mathbb{P}(\beta_0 + \beta_1 x_{0,1} \in \mathbb{E}(x_0; \alpha) \text{ in } \beta_0 + \beta_1 x_{0,2} \in \mathbb{E}(x_0; \alpha)) \neq 1 - \alpha.$$

Zaradi tega moramo uporabiti t. i. popravke pri izračunu hkratnih intervalov zaupanja za napovedi v več točkah hkrati. Najbolj tradicionalen je Bonferronijev popravek ozirom Bonferronijeva prilagoditev.

Bonferronijev popravek

Recimo, da imamo k parametrov $\theta_1, \dots, \theta_k$.

Pokažemo lahko, da če za parametre $\theta_1, \dots, \theta_k$ izračunamo posamične $100(1 - \frac{\alpha}{k})$ % intervale zaupanja, imajo ti intervali skupno pokritje vsaj $100(1 - \alpha)$ %.

Ta popravek deluje tudi za hkratno testiranje ničelnih domnev. Pri testiranju vrednost α zamenjamo z $\frac{\alpha}{k}$.

Intuicija:

- Če izvedemo en statistični test pri $\alpha = 0.05$ in je ničelna domneva pravilna, je verjetnost, da zavrnemo pravilno ničelno domnevo $\alpha = 0.05$.
- Če izvedemo 15 neodvisnih testov in je ničelna domneva v vseh 15 primerjavah pravilna, verjetnost, da zavrnemo pravilno ničelno domnevo ni več $\alpha = 0.05$, ampak ena minus verjetnost, da nobena od njih ni statistično značilna $1 - (1 - \alpha)^{15} = 0.5367$. Če uporabimo Bonferronijev popravek je ta verjetnost $1 - (1 - \alpha/15)^{15} = 0.0489$.
- Iz tega sledi, da brez Bonferronijevega popravka lahko v več kot polovico primerov zavrnemo ničelno domnevo, ki je pravilna.

- Pokažemo lahko, da Bonferronijev popravek velja tudi, če testi niso neodvisni.
- Bonferronijev popravek zagotavlja veljavno inferenco, ker zagotovi, da je skupna verjetnost za napako I. vrste $\leq \alpha$. Pri tem pa lahko plačamo visoko ceno, ker je zelo konzervativen. Zelo se zmanjša moč testa, saj je verjetnost, da zavrnemo ničelno domnevo, ko je ta napačna manjša oziroma napaka II. vrste večja.
- Obstajajo tudi druge prilagoditve za hkratno testiranje domnev, o teh bomo govorili kasneje.

2.11 Primer: SKT

Analizirajmo odvisnost sistoličnega krvnega tlaka (SKT) od starosti oseb (datoteka SKT.txt).

```
> tlak<-read.table(file="SKT.txt", header = TRUE)
> head(tlak)
```

```
    spol SKT starost
1     m 158     41
2     m 185     60
3     m 152     41
4     m 159     47
5     m 176     66
6     m 156     47
```

```
> str(tlak)
```

```
'data.frame':      69 obs. of  3 variables:
 $ spol   : chr  "m" "m" "m" "m" ...
 $ SKT    : int  158 185 152 159 176 156 184 138 172 168 ...
 $ starost: int   41 60 41 47 66 47 68 43 68 57 ...
```

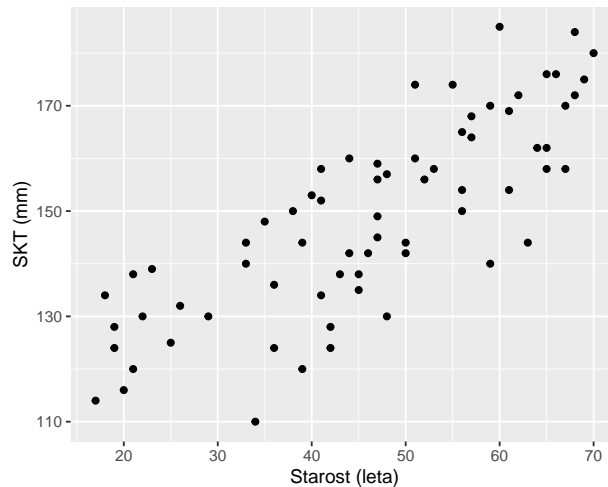
Analiziramo opisne statistike za SKT in starost.

```
> summary(tlak)
```

spol	SKT	starost
Length:69	Min. :110.0	Min. :17.00
Class :character	1st Qu.:135.0	1st Qu.:36.00
Mode :character	Median :149.0	Median :47.00
	Mean :148.7	Mean :46.14
	3rd Qu.:162.0	3rd Qu.:59.00
	Max. :185.0	Max. :70.00

Ustrezni grafični prikaz za analizo odvisnosti SKT od starost je razsevni diagram. Slika 4 kaže, da je linearna zveza primerna za te podatke.

```
> library(ggplot2)
> ggplot(data=tlak) +
+   geom_point(mapping=aes(x=starost, y=SKT)) +
+   xlab("Starost (leta)") +
+   ylab("SKT (mm)")
```



Slika 4: Odvisnost SKT od starosti za vzorec 69 oseb

Naredimo linearni regresijski model za odvisnost SKT od starosti, poimenovali ga bomo `model.SKT`.

```
> model.SKT <- lm(SKt~starost, data=tlak)
> # Objekt model.SKT ima naslednje komponente
> names(model.SKT)

[1] "coefficients" "residuals"      "effects"         "rank"
[5] "fitted.values" "assign"          "qr"              "df.residual"
[9] "xlevels"       "call"           "terms"          "model"

> model.SKT$coef # ni potrebno napisati celega imena "coefficients"

(Intercept)      starost
103.3490547      0.9833276
```

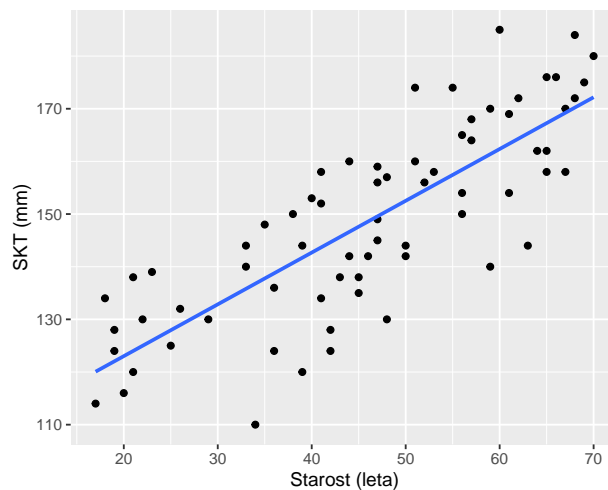
Najbolj osnovni rezultati v objektu `model.SKT` sta cenilka za presečišče (`Intercept`) in cenilka naklona SKT glede na `starost`. Enačba regresijske premice je (Slika 5):

$$\hat{y} = 103.35 + 0.98 x.$$

- Pričakovan SKT pri novorojenčku je 103.35; ta interpretacija ni vsebinsko smiselna, ker nimamo podatkov za starost oseb mlajših od 17 let.

- Z vsakim letom starosti se SKT v povprečju poveča za 0.98 mm. Bolje povedano: na deset let se SKT v povprečju poveča za 9.8 mm.

```
> ggplot(data=tlak, mapping=aes(x=starost, y=SKT)) +  
+   geom_point() +  
+   geom_smooth(method="lm", se=FALSE) +  
+   xlab("Starost (leta)") +  
+   ylab("SKT (mm)")
```

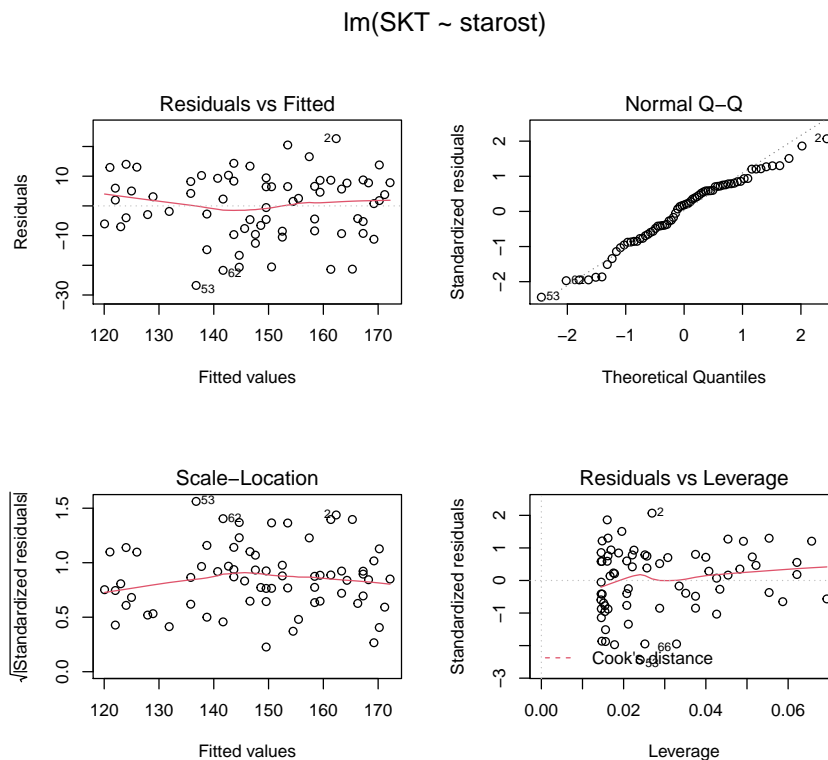


Slika 5: Odvisnost SKT od starosti za vzorec 69 oseb in regresijska premica

Kakovost prileganja regresijske premice podatkom in izpolnjevanje predpostavk linearnega modela ocenimo na podlagi slik ostankov modela. Z ukazom `plot(model.SKT)` dobimo 4 slike ostankov:

- Graf 1: *Residuals vs Fitted*;
- Graf 2: *Normal Q-Q plot*;
- Graf 3: *Scale-Location plot of $\sqrt{|residuals|}$ against fitted values*;
- Graf 4: *Plot of residuals against leverages, and a plot of Cook's distances against leverage/(1-leverage)*.

```
> par(mfrow = c(2, 2), oma=c(0,0,3,0))
> plot(model.SKT)
```



Slika 6: Grafični prikaz ostankov glede na napovedano vrednost (zgoraj levo) ter QQ graf standardiziranih ostankov (zgoraj desno), kvadratni koren absolutne vrednosti standardiziranih ostankov glede na napovedano vrednost (spodaj levo) in standardizirani ostanki glede na vzvod (spodaj desno)

Zaenkrat pogledjmo prva dva grafa. Če model ustreza podatkom, morajo biti točke na Graf 1 razporejene slučajno okoli vrednosti 0, kar pomeni, da je gladilnik približno na osi x . Če naj bi veljal normalni linearni model, morajo biti standardizirani ostanki porazdeljeni **približno** po $N(0, 1)$, to pomeni, da bi morale biti točke na Graf 2 **približno** na črtkani črti. Slika 6 kaže, da model dovolj dobro ustreza podatkom, ni pa idealen.

Povzetek modela dobimo z ukazom `summary(model.SKT)`. Izpis je sestavljen iz treh delov:

- informacija o ostankih;
- informacija o koeficientih regresijskega modela: cenilka, standardna napaka, t -vrednost, p -vrednost;
- koeficient determinacije in informacija iz tabele ANOVA.

```
> summary(model.SKT)
```

Call:

```
lm(formula = SKT ~ starost, data = tlak)
```

Residuals:

Min	1Q	Median	3Q	Max
-26.782	-7.632	1.968	8.201	22.651

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	103.34905	4.33190	23.86	<2e-16 ***
starost	0.98333	0.08929	11.01	<2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 11.1 on 67 degrees of freedom

Multiple R-squared: 0.6441, Adjusted R-squared: 0.6388

F-statistic: 121.3 on 1 and 67 DF, p-value: < 2.2e-16

```
> names(summary(model.SKT)) # za dodatne ali posamezne izpise
```

[1] "call"	"terms"	"residuals"	"coefficients"
[5] "aliased"	"sigma"	"df"	"r.squared"
[9] "adj.r.squared"	"fstatistic"	"cov.unscaled"	

V povzetku modela so standardne napake cenilk parametrov (*Std. Error*) izračunane na podlagi variančno-kovariančne matrike cenilk parametrov, ki jo vrne funkcija `vcov`. Na njeni diagonali so pripadajoče variance, izven diagonale je kovarianca:

$$\begin{aligned} \text{Var}(b_0) &= s^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{SS_{xx}} \right), \\ \text{Var}(b_1) &= \frac{s^2}{SS_{xx}}, \\ \text{Cov}(b_0, b_1) &= -\frac{s^2 \bar{x}}{SS_{xx}}, \end{aligned}$$

$s^2 = SS_{\text{residual}} / (n - 2)$ je cenilka za σ^2 .

```
> vcov(model.SKT)
```

	(Intercept)	starost
(Intercept)	18.7653188	-0.367938365
starost	-0.3679384	0.007973539

```
> # standardne napake cenilk parametrov
> sqrt(diag(vcov(model.SKT)))
```

```
(Intercept)      starost
 4.33189553    0.08929467
```

Za posamezen parameter modela testiramo ničelno domnevo, da je njegova vrednost enaka 0. Ilustrirajmo postopek za β_0 in β_1 .

$H_0 : \beta_0 = 0.$ Regresijska premica gre skozi izhodišče.
 $H_1 : \beta_0 \neq 0.$

```
> t.b0 <- model.SKT$coef[1]/sqrt(vcov(model.SKT)[1,1]); t.b0
```

```
(Intercept)
 23.8577
```

ali pa

```
> (t.b0 <- coef(summary(model.SKT))[1,1]/coef(summary(model.SKT))[1,2])
```

```
[1] 23.8577
```

$$t = \frac{b_0 - 0}{s_{b_0}} = \frac{103.34905}{4.3319} = 23.86.$$

```
> # p-vrednost
```

```
> p.b0 <- 2*pt(abs(t.b0), df=model.SKT$df.residual, lower.tail=FALSE); p.b0
```

```
[1] 1.836799e-34
```

Ničelno domnevo, da gre premica skozi izhodišče, zavrnamo ($p < 0.0001$). Opomba: rezultat ni vsebinsko smiseln. Testiranje te domneve je vedno v izpisu, pogosto pa ta domneva ni raziskovalno zanimiva.

$H_0 : \beta_1 = 0.$ Regresijska premica je vzporedna osi x ; ni odvisnosti SKT od starosti.
 $H_1 : \beta_1 \neq 0.$

```
> # t.b1<-model.SKT$coef[2]/sqrt(vcov(model.SKT)[2,2]); t.b1
```

```
> (t.b1 <- coef(summary(model.SKT))[2,1]/coef(summary(model.SKT))[2,2])
```

```
[1] 11.01216
```

$$t = \frac{b_1 - 0}{s_{b_1}} = \frac{0.98333}{0.08929} = 11.01.$$

```
> p.b1 <- 2*pt(t.b1, df=model.SKT$df.residual, lower.tail=FALSE); p.b1
[1] 1.116594e-16
```

Ničelno domnevo zavrnemo v korist alternativne domneve; zveza je statistično značilna. SKT je statistično značilno odvisen od starosti ($p < 0.0001$).

V tretjem delu povzetka modela je standardna napaka regresije s , koeficient determinacije R^2 in povzetek tabele analize variance modela (F -statistika s pripadajočo p -vrednostjo).

```
> summary(model.SKT)$r.squared
[1] 0.6441239

> (F <- summary(model.SKT)$fstatistic)

      value      numdf      dendif
121.2678    1.0000    67.0000

> F <- as.numeric(F)
> (p.F <- pf(F[1], df1=F[2], df2=F[3], lower.tail=FALSE))
[1] 1.116594e-16
```

Za vajo pogledajmo izračun koeficienta determinacije na naših podatkih:

```
> SS_model<-sum((model.SKT$fitted-mean(tlak$SKT))^2);SS_model
[1] 14951.25

> SS_res<-sum(model.SKT$residual^2);SS_res
[1] 8260.514

> R2<-SS_model/(SS_model+SS_res); R2
[1] 0.6441239
```

Približno dve tretjini variabilnosti SKT pojasni starost, ostala variabilnost ostane nepojasnjena. Le ugibamo lahko, kaj poleg starosti še vpliva na variabilnost SKT, morda spol, prehrana, genetika, fizična aktivnost, itd.

Vsote kvadriranih odklonov SS_{model} in $SS_{residual}$ dobimo v tabeli analize variance s funkcijo `anova()`:

```
> anova(model.SKT)
```

Analysis of Variance Table

Response: SKT

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
starost	1	14951.3	14951.3	121.27	< 2.2e-16 ***
Residuals	67	8260.5	123.3		

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Cenilka za skupno varianco σ^2 iz (6) je $s^2 = 123.29$, pripadajoča standardna napaka regresije je $s = 11.1$. V tem primeru F -statistika testira domnevo $H_0 : \beta_1 = 0$, ki je ekvivalentna domnevi, ki jo testira t -statistika; velja $F = t^2 = 121.27$.

Intervala zaupanja za parametra modela dobimo s funkcijo `confint`:

```
> confint(model.SKT)
```

	2.5 %	97.5 %
(Intercept)	94.7025550	111.995554
starost	0.8050947	1.161561

- Prvi interval zaupanja ni vsebinsko smiseln.
- Če se starost poveča za 10 let, pri 95 % zaupanju pričakujemo, da se bo SKT povečal na intervalu 8.1 mm do 11.6 mm.

Izračunajmo napovedi in njihove 95 % IZ za SKT za osebe stare 30 in 60 let.

```
> starost.napovedi<-data.frame(starost=c(30,60))
> ## povprečne napovedi
> povp.napovedi.SKT<-predict(model.SKT, starost.napovedi, interval="confidence")
> data.frame(cbind(starost.napovedi,povp.napovedi.SKT ))
```

	starost	fit	lwr	upr
1	30	132.8489	128.9247	136.7731
2	60	162.3487	158.7132	165.9842

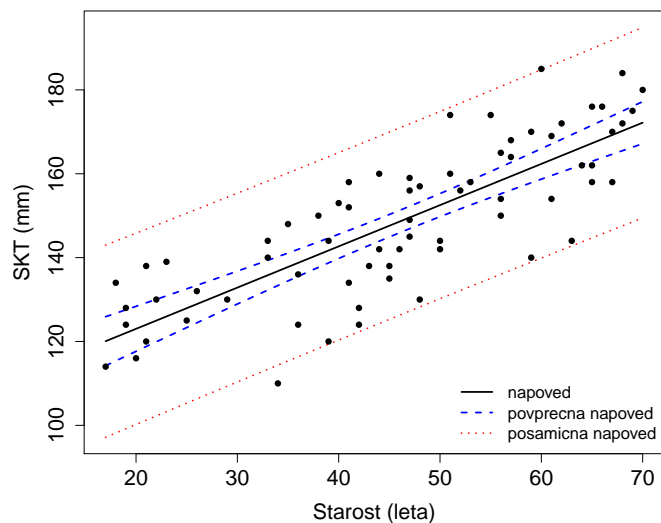
```
> ## posamične napovedi
> pos.napovedi.SKT<-predict(model.SKT, starost.napovedi, interval="prediction")
> data.frame(cbind(starost.napovedi,pos.napovedi.SKT ))
```

	starost	fit	lwr	upr
1	30	132.8489	110.3412	155.3566
2	60	162.3487	139.8895	184.8079

Interpretacija za osebe stare 30 let: napovedana vrednost za tlak 132.8 mm. Pripadajoč 95 % IZ za povprečje oseb s starostjo 30 let je (128.9 mm, 136.8 mm). Za posamezno osebo staro 30 let je pripadajoč 95 % IZ (110.3 mm, 155.4 mm).

Ali bi bila utemeljena napoved za osebo staro 95 let? Ali bi bila utemeljena napoved za osebo staro 10 let? Ne, ker gre za ekstrapolacijo.

Grafični prikaz 95 % intervalov zaupanja za povprečne in za posamične napovedi je na Sliki 7.



Slika 7: Napovedi za SKT, 95 % IZ za povprečne napovedi (notranji hiperboli) in za posamične napovedi (zunanji hiperboli)

2.12 Simulacija po predpostavkah modela enostavne linearne regresije

Simulacijo bomo izvedli, da bi ilustrirali statistične lastnosti cenilk parametrov linearnega modela ter velikost testa in moč testa pri testiranju domnev o parametrih. Za okvir simulacije vzemimo model odvisnosti SKT od `starost` (`model.SKT`). V tem modelu testiranje domneve $H_0 : \beta_0 = 0$ ni vsebinsko zanimivo, osredotočili se bomo na testiranje domneve $H_0 : \beta_1 = 0$.

V vseh simulacijah obdržimo iste vrednosti za napovedno spremenljivko `starost`. Za vrednosti odzivne spremenljivke SKT upoštevamo, da so pogojno na vrednosti napovedne spremenljivke porazdeljene normalno s pričakovano vrednostjo $\beta_0 + \beta_1 \text{starost}$ in varianco σ^2 , pri čemer za vrednosti β_0, β_1 in σ^2 vzamemo cenilke iz `model.SKT`: $SKT_i = 103 + 0.98 \text{starost}_i + \varepsilon_i$; napake $\varepsilon_i, i = 1, \dots, 69$, generiramo s funkcijo `rnorm()` za porazdelitev $N(0, \sigma^2 = 11^2)$.

Za vsak generirani vzorec vrednosti `SKT` izračunamo cenilke parametrov enostavnega linear-

nega modela b_0 in b_1 ter 95 % interval zaupanja za β_1 . Izpišemo tudi p -vrednost pri testiranju domneve $H_0 : \beta_1 = 0$.

```
> # vrednosti za starost vzamemo iz podatkovnega okvira tlak
> starost<-tlak$starost
> # velikost vzorca
> n<-length(tlak$starost)
> # standardni odklon napak
> sigma<-11
> # izbrana parametra modela
> beta0<-103
> beta1<-0.98
> # teoretični izračun variance cenilke b0
> var.b0<-sigma^2*(1/n+mean(starost)^2/sum((starost-mean(starost))^2))
> # teoretična standardna napaka za b0
> sqrt(var.b0)
```

```
[1] 4.291455
```

```
> # teoretični izračun variance cenilke b1
> var.b1<-sigma^2/sum((starost-mean(starost))^2)
> # teoretična standardna napaka za b1
> sqrt(var.b1)
```

```
[1] 0.08846106
```

```
> set.seed(77) # za ponovljivost rezultatov
> # vrednosti slučajnih napak
> epsilon<-rnorm(n, mean=0, sd=sigma)
> # simulirane vrednosti odzivne spremenljivke
> SKT<-beta0 + beta1*starost + epsilon

> # model enostavne linearne regresije na simuliranih podatkih
> mod<-lm(SKT~starost)
> (b0<-coef(mod)[1])
```

```
(Intercept)
105.9049
```

```
> (b1<-coef(mod)[2])
```

```
starost
0.9183129
```

```
> (p.b1<-coefficients(summary(mod))[2,4]) # za H0: beta1=0
```

```
[1] 1.054677e-14
```

```
> (sp.meja.b1<-confint(mod)[2,1])
```

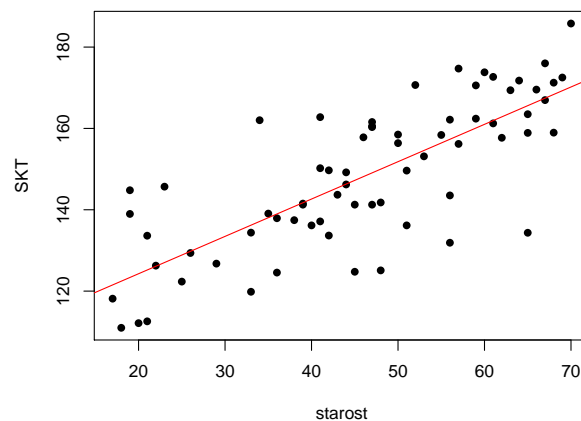
```
[1] 0.7327401
```

```
> (zg.meja.b1<-confint(mod)[2,2])
```

```
[1] 1.103886
```

```
> plot(starost, SKT, pch=16)
```

```
> abline(reg=mod, col="red")
```



Slika 8: Primer simuliranih podatkov in pripadajoča regresijska premica, `set.seed(77)`

Simulacijo bomo velikokrat ponovili. Na podlagi rezultatov simulacij želimo pogledati:

- kakšni sta porazdelitvi cenilk parametrov enostavnega linearnega modela;
- kolikšen delež intervalov zaupanja za β_1 ne vsebuje prave vrednosti parametra (velikost testa);
- kolikšna je moč testa pri testiranju ničelne domneve $H_0 : \beta_1 = 0$.

V ta namen bomo izvedli 1000 simulacij (`Nsim=1000`). Napišimo funkcijo, ki bo izvedla simulacije in za vsako simulacijo shranila v podatkovni okvir cenilke parametrov, spodnjo in zgornjo mejo intervala zaupanja za β_1 in p -vrednost pri testiranju ničelne domneve $H_0 : \beta_1 = 0$.

```
> Nsim <- 1000
> reg.sim <- function(x, beta0, beta1, sigma, n, Nsim) {
+ # pripravimo prazne vektorje za rezultate simulacij, cenilki parametrov b0 in b1,
+ # p-vrednost za testiranje domneve beta1=0,
+ # spodnjo in zgornjo mejo intervala zaupanja za beta1
+   b0 <- numeric(Nsim)
+   b1 <- numeric(Nsim)
+   p.b1 <- numeric(Nsim)
+   sp.meja.b1 <- numeric(Nsim)
+   zg.meja.b1 <- numeric(Nsim)
+
+   for (i in 1:Nsim) {
+     epsilon<-rnorm(n, mean=0, sd=sigma)
+     y<-beta0+beta1*x+epsilon
+     mod<-lm(y~x)
+     b0[i]<-coef(mod)[1]
+     b1[i]<-coef(mod)[2]
+     p.b1[i]<-coefficients(summary(mod))[2,4]
+     sp.meja.b1[i]<-confint(mod)[2,1]
+     zg.meja.b1[i]<-confint(mod)[2,2]
+   }
+   return(data.frame(b0,b1,p.b1,sp.meja.b1,zg.meja.b1))
+ }
> rez.1000<-reg.sim(x=starost, beta0=103, beta1=0.98, sigma=11, n=30, Nsim=1000)
> head(rez.1000)
```

	b0	b1	p.b1	sp.meja.b1	zg.meja.b1
1	106.90334	0.9055551	1.168895e-11	0.6844261	1.126684
2	104.78995	0.9275064	7.847449e-17	0.7607452	1.094268
3	99.62207	1.0249480	6.718667e-18	0.8505561	1.199340
4	100.33065	0.9441515	3.556979e-19	0.7933889	1.094914
5	105.65246	0.9421287	1.368567e-22	0.8138692	1.070388
6	93.07662	1.1370228	2.988673e-23	0.9867236	1.287322

```
> # standardna napaka b0 in b1 na podlagi porazdelitve Nsim cenilk parametrov
> (sd(rez.1000$b0))

[1] 4.336861

> (sd(rez.1000$b1))

[1] 0.08443493

> # 2.5 in 97.5 centil za b1 na podlagi simulacij
> (centili<-quantile(rez.1000$b1, probs=c(0.025, 0.975)))
```

```
      2.5%      97.5%
0.8207757 1.1485645
```

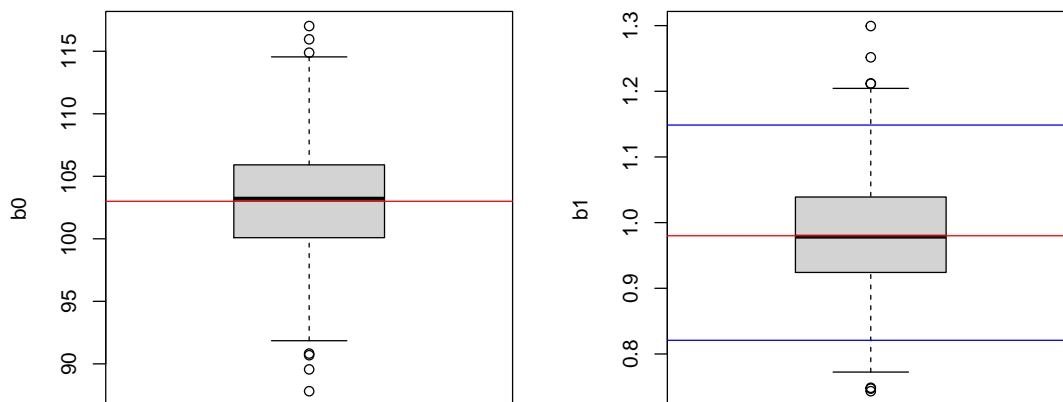
```
> # ocena verjetnosti za napako II. vrste za H0: beta1=0
> sum(rez.1000$p>0.05)/Nsim
```

```
[1] 0
```

```
> # ocena moči testa na podlagi Nsim simulacij
> (moc.testa<-1-sum(rez.1000$p>0.05)/Nsim)
```

```
[1] 1
```

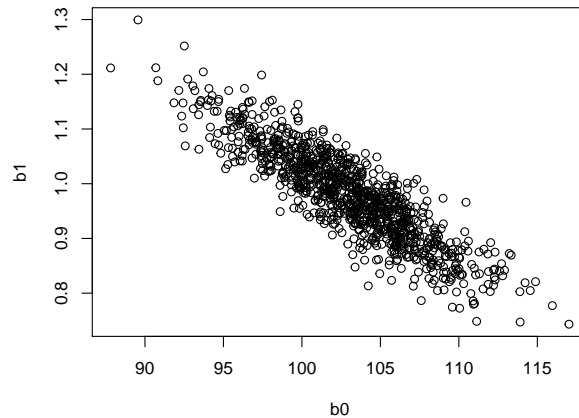
```
> par(mfrow=c(1,2))
> boxplot(rez.1000$b0, ylab="b0");
> abline(h=beta0, col="red")
> boxplot(rez.1000$b1, ylab="b1");
> abline(h=beta1, col="red");
> abline(h=centili, col="blue")
```



Slika 9: Porazdelitev cenilk parametrov b_0 (levo) in b_1 (desno) za $\sigma = 11$ in $n = 69$, `set.seed(77)`, rdeča črta kaže pravo vrednost za parameter, modri črti predstavljata 2.5 in 97.5 centil za b_1

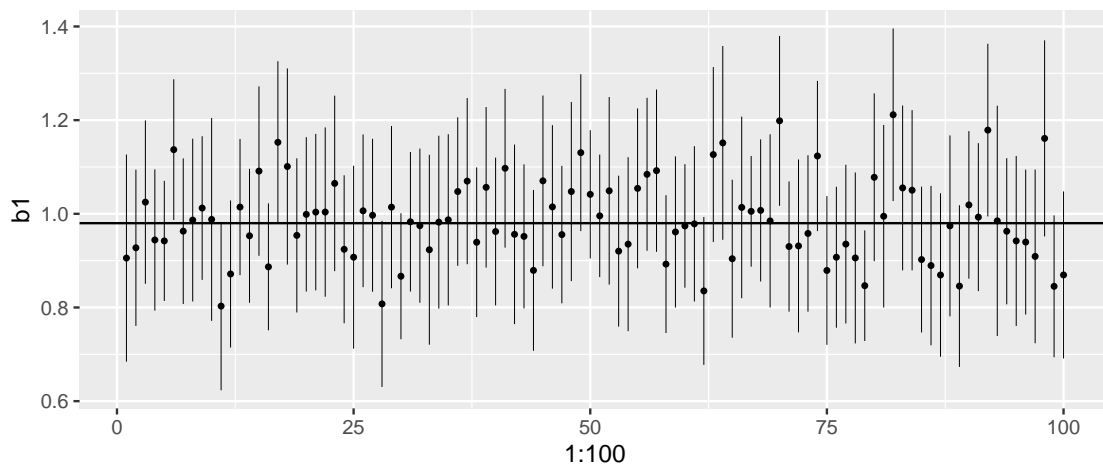
Slika 10 prikazuje negativno povezanost cenilk b_0 in b_1 , ki je teoretično izražena s kovarianco $Cov(b_0, b_1) = -\sigma^2 \bar{x} / S_{xx}$.

```
> plot(rez.1000$b0, rez.1000$b1, pch=1, xlab="b0", ylab="b1")
```



Slika 10: Povezanost cenilk b_0 in b_1 , `set.seed(77)`

```
> ggplot(rez.1000[1:100,], aes(x=1:100,y=b1,ymin=sp.meja.b1,ymax=zg.meja.b1)) +  
+ geom_pointrange(size=0.2, shape=16) +  
+ geom_hline(yintercept=beta1)
```



Slika 11: 100 intervalov zaupanja za β_1 za $\sigma = 11$ in $n = 69$, `set.seed(77)`

```
> # delež intervalov zaupanja, ki ne vsebujejo prave vrednosti parametra beta1,  
> # to je ocena velikosti testa  
> sum(rez.1000$sp.meja.b1>beta1 | rez.1000$zg.meja.b1<beta1)/Nsim
```

```
[1] 0.033
```

3 VAJE

3.1 Čas teka Collina Jacksona

V datoteki COLLIN.txt so podatki za 21 tekov na 110 m čez ovire tekača Collina Jacksona: hitrost vetra = `windspeed` (m/s) in čas teka = `time` (s) (Vir: Daly et al., str. 525). Podatki so bili dobljeni v poskusu v zaprtem prostoru, hitrost vetra je bila izbrana za vsak tek posebej vnaprej. Negativne vrednosti hitrosti vetra pomenijo, da je veter pihal v prsi tekača. Kako hitrost vetra vpliva na čas teka na 110 m čez ovire?

- Grafično prikažite podatke.
- Ocenite parametra linearnega regresijskega modela za odvisnost časa teka od hitrosti vetra.
- Analizirajte ostanke modela na podlagi grafičnih prikazov.
- Obrazložite cenilke parametra modela in njuna intervala zaupanja.
- Obrazložite koeficient determinacije.
- Izračunajte povprečno in posamično napoved časa teka ter pripadajoče 95 % intervale zaupanja za naslednje hitrosti vetra: -1 m/s, 0 m/s, 1 m/s in 4 m/s. Ali so vse napovedi upravičene? Zakaj?

3.2 Simulacije za enostavno linearno regresijo

Za izbrani vrednosti parametrov enostavne linearne regresije $\beta_0 = 100$ in $\beta_1 = 1$ izvedite simulacije, ki bodo ilustrirale vpliv velikosti vzorca n in vrednosti variance napak σ^2 na porazdelitev cenilk parametrov in na moč testa pri testiranju domneve $H_0 : \beta_1 = 0$.

Za vsako izbrano velikost vzorca n najprej generirajte vrednosti napovedne spremenljivke x na intervalu 15 do 70. Pri tem uporabite funkcijo `sample` z argumentom `replace=TRUE`: `x<-sample(c(17:70), size=n, replace=TRUE)`. Za tako določene vrednosti napovedne spremenljivke nadaljujte z simulacijami pri izbranih vrednostih σ^2 .

Navodilo:

- Izberite naslednje vrednosti za n : 7, 15, 30, 100 in naslednje vrednosti za σ : 5, 11, 22. Za vsako kombinacijo n in σ naredite 1000 simulacij.
- Grafično prikažite:
 - odvisnost širine intervala zaupanja za β_1 od n , za vsako vrednost σ ;
 - odvisnost širine intervala zaupanja za β_1 od σ , za vsak n ;
 - odvisnost moči testa od n , za vsako vrednost σ ;
 - odvisnost moči testa od σ , za vsak n .
- Napišite kratek povzetek vaših ugotovitev.

3.3 Varianca odzivne spremenljivke

Pri izrazu za razdelitev vsote kvadratov odklonov odzivne spremenljivke SS_{yy} na dva dela, tistega, ki ga pojasni linearni model SS_{model} in nepojasneni del $SS_{residual}$:

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2,$$

smo upoštevali, da je $2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) = 0$. Pokažite, da to velja.