

# O LINEARNI REGRESIJI

Martin Raič

September 2014

# 1 Linearna regresija kot statistični model

Linearna regresija je statistični model, pri katerem opazimo slučajni vektor:

$$\mathbf{Y} = \mathbf{v} + \boldsymbol{\varepsilon}, \quad (1.1)$$

kjer je  $\mathbf{v}$  neopazljiv determinističen parameter,  $\boldsymbol{\varepsilon}$  pa neopazljiv slučajen šum. Pri tem vse skupaj živi v  $n$ -razsežnem evklidskem prostoru  $W$ . Za  $\mathbf{v}$  privzamemo, da pripada določenemu  $k$ -razsežnemu podprostoru  $V$ , za šum  $\boldsymbol{\varepsilon}$  pa privzamemo, da ima matematično upanje  $\mathbf{0}$  in kovariančno matriko  $\sigma^2 \mathbf{I}_W$ , kjer je  $\sigma^2$  še en neopazljiv determinističen parameter; le-temu bomo rekli *varianca modela*. Velja torej  $\mathbb{E} \boldsymbol{\varepsilon} = \mathbf{0}$  in  $\mathbb{E} \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^T = \mathbf{I}_W$ .

V praksi, ko podatke podajamo s številkami, je smiselno privzeti, da je  $W = \mathbb{R}^n$ ,  $V$  pa je zaloga vrednosti injektivne linearne preslikave  $\mathbf{X}: \mathbb{R}^k \rightarrow \mathbb{R}^n$ , ki jo bomo identificirali z matriko velikosti  $n \times k$ . Tako lahko model parametriziramo v obliki:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad (1.2)$$

kjer je  $\boldsymbol{\beta} \in \mathbb{R}^k$ . Ekvivalentno lahko zapišemo tudi:

$$\begin{aligned} Y_1 &= \beta_1 f_1(X_1) + \beta_2 f_2(X_1) + \cdots + \beta_k f_k(X_1) + \varepsilon_1, \\ Y_2 &= \beta_1 f_1(X_2) + \beta_2 f_2(X_2) + \cdots + \beta_k f_k(X_2) + \varepsilon_2, \\ &\vdots \\ Y_n &= \beta_1 f_1(X_n) + \beta_2 f_2(X_n) + \cdots + \beta_k f_k(X_n) + \varepsilon_n, \end{aligned} \quad (1.3)$$

kjer so  $X_1, \dots, X_n$  deterministične in opazljive spremenljivke. Tem pravimo *pojasnjevalne*, medtem ko so  $Y_1, \dots, Y_n$  *odvisne* spremenljivke. Tu je seveda:

$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} f_1(X_1) & f_2(X_1) & \cdots & f_k(X_1) \\ f_1(X_2) & f_2(X_2) & \cdots & f_k(X_2) \\ \vdots & \vdots & \ddots & \vdots \\ f_1(X_n) & f_2(X_n) & \cdots & f_k(X_n) \end{bmatrix}, \quad \boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}.$$

Obravnavali bomo naslednje naloge inferenčne statistike:

- točkasto in intervalsko ocenjevanje parametrov;
- napovedovanje spremenljivk, povezanih z modelom.
- testiranje hipotez v zvezi z modelom.

## 2 Točkasto ocenjevanje in napovedovanje

Oglejmo si najprej model, podan z zvezo (1.1). Iskali bomo čim boljšo *nepristransko* cenilko  $\hat{\mathbf{v}} = h(\mathbf{Y})$  parametra  $\mathbf{v}$ , t. j. tako, ki ima najmanjšo možno kovariančno matriko  $\text{var}(\hat{\mathbf{v}}) := \mathbb{E}[(\hat{\mathbf{v}} - \mathbf{v})(\hat{\mathbf{v}} - \mathbf{v})^T]$ . Seveda je to simetrična matrika, na teh pa bomo gledali delno urejenost, kjer je  $\mathbf{A} \leq \mathbf{B}$ , če je  $\mathbf{B} - \mathbf{A}$  pozitivno semidefinitna, t. j.  $\mathbf{a}^T \mathbf{A} \mathbf{a} \leq \mathbf{a}^T \mathbf{B} \mathbf{a}$  za vse  $\mathbf{a} \in V$ . Najprej se bomo omejili na *linearne* cenilke, t. j. take, kjer je  $h(\mathbf{y}) = \mathbf{A}\mathbf{y} + \mathbf{v}_0$ .

**Izrek 2.1** (Gauss, Markov). *Naj velja (1.1) skupaj s pripadajočimi predpostavkami. Tedaj je nepristranska linearna cenilka za  $\mathbf{v}$  z najmanjšo možno kovariančno matriko kar pravokotna projekcija opaženega vektorja  $\mathbf{Y}$  na podprostor  $V$ .*

Preden gremo dokazat izrek, uvedimo nekaj notacije za inkluzijo in pravokotno projekcijo. Naj bosta  $G \subseteq H$  evklidska prostora. Z  $\iota_{G \rightarrow H}$  bomo označevali inkluzijo iz  $G$  v  $H$ . Pri pravokotni projekciji pa bomo ločili dve niansi: označevali jo bomo tako s  $p_{H \rightarrow G}$  kot tudi s  $\mathbf{P}_{H \rightarrow G}$ , toda prvo bomo gledali kot preslikavo iz  $H$  na  $G$ , drugo pa kot preslikavo iz  $H$  v  $H$ , torej kot operator na  $H$ . Velja:

$$p_{H \rightarrow G}^T = \iota_{G \rightarrow H}, \quad \mathbf{P}_{H \rightarrow G}^T = \mathbf{P}_{H \rightarrow G}.$$

Pri zapisu cenilke lahko seveda uporabimo obe različici:

$$\hat{\mathbf{v}} = p_{W \rightarrow V} \mathbf{Y} = \mathbf{P}_{W \rightarrow V} \mathbf{Y}. \quad (2.1)$$

**DOKAZ IZREKA GAUSSA IN MARKOVA.** Najprej preverimo, kdaj je linearna cenilka  $\hat{\mathbf{v}} = \mathbf{A}\mathbf{Y} + \mathbf{v}_0$  nepristranska. Ker je  $\mathbb{E}\varepsilon = \mathbf{0}$ , mora veljati:

$$\mathbb{E}\hat{\mathbf{v}} = \mathbf{A}\mathbf{v} + \mathbf{v}_0 = \mathbf{v},$$

in sicer ne glede na dejansko vrednost neznanega parametra  $\mathbf{v}$ . Torej mora biti  $\mathbf{v}_0 = \mathbf{0}$ ,  $\mathbf{A}$  pa mora mirovati na  $V$ . Varianca cenilke  $\hat{\mathbf{v}}$  bo najmanjša natanko tedaj, ko bo izraz:

$$\mathbf{a}^T \text{var}(\hat{\mathbf{v}}) \mathbf{a} = \sigma^2 \mathbf{a}^T \mathbf{A} \mathbf{A}^T \mathbf{a} = \sigma^2 \|\mathbf{A}^T \mathbf{a}\|^2$$

minimalen za vse  $\mathbf{a} \in V$ . Ker  $\mathbf{A}$  miruje na  $V$ , za vse  $\mathbf{a}, \mathbf{b} \in V$  velja:

$$\mathbf{b}^T \mathbf{A}^T \mathbf{a} = (\mathbf{A}\mathbf{b})^T \mathbf{a} = \mathbf{b}^T \mathbf{a},$$

kar pomeni, da mora  $\mathbf{A}^T - \iota_{V \rightarrow W}$  slikati v ortogonalni komplement prostora  $V$ . Z drugimi besedami, lahko pišemo  $\mathbf{A}^T = \iota_{V \rightarrow W} + \mathbf{B}$ , kjer  $\mathbf{B}$  slika iz  $V$  v njegov ortogonalni komplement. Sedaj pa lahko izračunamo:

$$\|\mathbf{A}^T \mathbf{a}\|^2 = \|\mathbf{a} + \mathbf{B}\mathbf{a}\|^2 = \|\mathbf{a}\|^2 + \|\mathbf{B}\mathbf{a}\|^2$$

(upoštevali smo ortogonalnost). Slednji izraz pa je minimalen (za vse  $\mathbf{a}$ ), če je  $\mathbf{B} = \mathbf{0}$ , torej če je  $\mathbf{A}^T = \iota_{V \rightarrow W}$ , torej če je  $\mathbf{A} = p_{W \rightarrow W}$ . ■

**Opomba 2.2.** Iz dokaza izreka Gaussa in Markova lahko razberemo, da je, če je  $\hat{\mathbf{v}}$  tako kot v (2.1), kovariančna matrika te cenilke enaka:

$$\text{var}(\hat{\mathbf{v}}) = \sigma^2 \mathbf{I}_V. \quad (2.2)$$

Izrek Gaussa in Markova nam da tudi pripadajoči rezultat za linearne funkcionalne parametra  $\mathbf{v}$ : za poljuben  $\mathbf{a} \in V$  tudi  $\mathbf{a}^T \hat{\mathbf{v}}$  nepristranska cenilka karakteristike  $\mathbf{a}^T \mathbf{v}$  z najmanjšo možno varianco in slednja je enaka:

$$\text{var}(\mathbf{a}^T \hat{\mathbf{v}}) = \sigma^2 \|\mathbf{a}\|^2. \quad (2.3)$$

Privzemimo zdaj, da je model *Gaussov*, t. j. da je šum porazdeljen *večrazsežno normalno*, se pravi zvezno z gostoto:

$$f_{\varepsilon}(\mathbf{w} \mid \mathbf{v}) = e^{-\|\mathbf{w}\|^2/2}.$$

To pomeni, da je opažanje porazdeljeno zvezno z gostoto:

$$f_{\mathbf{Y}}(\mathbf{y} \mid \mathbf{v}) = e^{-\|\mathbf{y}-\mathbf{v}\|^2/2}.$$

Pri fiksnem  $\mathbf{y}$  je maksimum po  $\mathbf{v} \in V$  dosežen natanko tedaj, ko je  $\mathbf{v}$  pravokotna projekcija vektorja  $\mathbf{y}$  na  $V$ . To pomeni, da se cenilka  $\hat{\mathbf{v}} = p_{W \rightarrow V} \mathbf{Y}$  ujema s cenilko po metodi največjega verjetja.

**Izrek 2.3.** *Če je model Gaussov, ima cenilka  $\hat{\mathbf{v}} = p_{W \rightarrow V} \mathbf{Y}$  najmanjšo možno kovariacijsko matriko izmed vseh možnih nepristranskih cenilk.*

DOKAZ. Sklicali se bomo na *Lehmann–Schefféjev izrek* (glej npr. [1], izrek 5.5, str. 298). Najprej iz zapisa:

$$\begin{aligned} f_{\mathbf{Y}}(\mathbf{y} \mid \mathbf{v}) &= \exp\left(-\frac{\|\mathbf{v}\|^2}{2}\right) \exp\left(-\frac{\|\mathbf{y}\|^2}{2}\right) \exp(\mathbf{v}^T \mathbf{y}) = \\ &= \exp\left(-\frac{\|\mathbf{v}\|^2}{2}\right) \exp\left(-\frac{\|\mathbf{y}\|^2}{2}\right) \exp(\mathbf{v}^T p_{W \rightarrow V} \mathbf{y}) \end{aligned}$$

razberemo, da gre za neizrojeno eksponentno družino s pripadajočo zadostno statistiko  $p_{W \rightarrow V} \mathbf{Y}$ , ki je posledično kompletna (glej npr. [1], izrek 2.74, str. 108). Naj bo  $\mathbf{a} \in V$ . Ker je cenilka  $\mathbf{a}^T \hat{\mathbf{v}}$  karakteristike  $\mathbf{a}^T \mathbf{v}$  funkcija pripadajoče zadostne statistike in nepristranska, ima po Lehmann–Schefféjevem izreku najmanjšo možno varianco: za vsako drugo nepristransko cenilko  $h_{\mathbf{a}}(\mathbf{Y})$  velja  $\text{var}(h_{\mathbf{a}}(\mathbf{Y})) \geq \sigma^2 \|\mathbf{a}\|^2$ . Brž ko je torej  $h(\mathbf{Y})$  kakšna druga nepristranska cenilka za  $\mathbf{v}$ , je tudi  $\mathbf{a}^T h(\mathbf{Y})$  nepristranska cenilka za  $\mathbf{a}^T \mathbf{v}$  in po prejšnjem velja:

$$\mathbf{a}^T \text{var}(h(\mathbf{Y})) \mathbf{a} = \text{var}(\mathbf{a}^T h(\mathbf{Y})) \geq \sigma^2 \|\mathbf{a}\|^2 = \mathbf{a}^T \text{var}(\hat{\mathbf{v}}) \mathbf{a},$$

torej je  $\text{var}(h(\mathbf{Y})) \geq \text{var}(\hat{\mathbf{v}})$ . ■

Razliki  $\hat{\varepsilon} := \mathbf{Y} - \hat{\mathbf{v}}$  pravimo *rezidual* in je ocena za šum. Rezidual je pravokotna projekcija vektorja  $\mathbf{Y}$  na ortogonalni komplement prostora  $V$  v  $W$ . Če torej za evklidska prostora  $G \subseteq H$  s  $q_{H \rightarrow G}$  oziroma s  $\mathbf{Q}_{H \rightarrow G}$  označimo ortogonalni projektor s prostora  $G$  na ortogonalni komplement prostora  $G$  v  $H$  (pri čemer je  $q_{H \rightarrow G}$  definirana kot preslikava na ustrezen ortogonalni komplement,  $\mathbf{Q}_{H \rightarrow G}$  pa kot operator na  $H$ , tako da velja  $\mathbf{Q}_{H \rightarrow G} = \mathbf{I}_H - \mathbf{P}_{H \rightarrow G}$  in tudi  $\mathbf{P}_{H \rightarrow G} \mathbf{Q}_{H \rightarrow G} = \mathbf{Q}_{H \rightarrow G} \mathbf{P}_{H \rightarrow G} = \mathbf{0}$ ), velja:

$$\hat{\varepsilon} = q_{W \rightarrow V} \mathbf{Y} = \mathbf{Q}_{W \rightarrow V} \mathbf{Y}$$

in ker je  $q_{W \rightarrow V} \mathbf{v} = \mathbf{0}$  in  $\mathbf{Q}_{W \rightarrow V} \mathbf{v} = \mathbf{0}$ , je tudi:

$$\hat{\varepsilon} = q_{W \rightarrow V} \varepsilon = \mathbf{Q}_{W \rightarrow V} \varepsilon.$$

Pravokotna projekcija vektorja na določen podprostor je vektor v tem podprostoru, ki je izvornemu vektorju najbližje glede na evklidsko normo. V kontekstu našega modela (1.1) to pomeni, da je norma reziduala najmanjša možna, torej da je vsota kvadratov komponent reziduala minimalna. Zato pravimo, da je cenilka  $\hat{\mathbf{v}}$  za  $\mathbf{v}$  dobljena po *metodi najmanjših kvadratov*.

Rezidual  $\hat{\boldsymbol{\varepsilon}}$  je nekoreliran z  $\hat{\mathbf{v}}$ , saj je:

$$\begin{aligned} \text{cov}(\hat{\mathbf{v}}, \hat{\boldsymbol{\varepsilon}}) &= \text{cov}(\mathbf{P}_{W \rightarrow V} \boldsymbol{\varepsilon}, \mathbf{Q}_{W \rightarrow V} \boldsymbol{\varepsilon}) = \\ &= \mathbf{P}_{W \rightarrow V} \text{var}(\boldsymbol{\varepsilon}) \mathbf{Q}_{W \rightarrow V}^T = \\ &= \sigma^2 \mathbf{P}_{W \rightarrow V} \mathbf{I}_W \mathbf{Q}_{W \rightarrow V} = \\ &= \sigma^2 \mathbf{P}_{W \rightarrow V} \mathbf{Q}_{W \rightarrow V} = \\ &= \mathbf{0}. \end{aligned}$$

Izračunamo lahko še:

$$\begin{aligned} \mathbb{E} \|\hat{\boldsymbol{\varepsilon}}\|^2 &= \text{sl} \mathbb{E}[\hat{\boldsymbol{\varepsilon}} \hat{\boldsymbol{\varepsilon}}^T] = \\ &= \text{sl} \mathbb{E}[\mathbf{Q}_{W \rightarrow V} \hat{\boldsymbol{\varepsilon}} \hat{\boldsymbol{\varepsilon}}^T \mathbf{Q}_{W \rightarrow V}^T] = \\ &= \text{sl}(\mathbf{Q}_{W \rightarrow V} \mathbb{E}[\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^T] \mathbf{Q}_{W \rightarrow V}) = \\ &= \sigma^2 \text{sl}(\mathbf{Q}_{W \rightarrow V} \mathbf{I}_W \mathbf{Q}_{W \rightarrow V}) = \\ &= \sigma^2 \text{sl} \mathbf{Q}_{W \rightarrow V} = \\ &= (n - k) \sigma^2. \end{aligned}$$

Torej je:

$$\hat{\sigma}^2 = \frac{\|\hat{\boldsymbol{\varepsilon}}\|^2}{n - k} \quad (2.4)$$

nepristranska cenilka za  $\sigma^2$ . Količini  $\hat{\sigma}^2$  bomo zato rekli kar *ocenjena varianca modela*.

Zdaj si lahko ogledamo, kolikšno napako v grobem naredimo pri ocenjevanju linearnih funkcionalov parametrov. Če ocenjujemo določen linearni funkcional parametra  $\mathbf{v}$  iz zveze (1.1), t. j.  $\mathbf{a}^T \mathbf{v}$ , iz opombe 2.2 vemo, da je  $\mathbf{a}^T \hat{\mathbf{v}}$  nepristranska cenilka z varianco  $\sigma^2 \|\mathbf{a}\|^2$ . Grobo merilo za napako, ki smo jo naredili pri ocenjevanju, je *ocenjena standardna napaka*:

$$\text{SE}(\mathbf{a}^T \mathbf{v}) = \hat{\sigma} \|\mathbf{a}\|. \quad (2.5)$$

Oznaka SE tu pomeni operator, ki je definiran na linearnih funkcionalih na parametričnem prostoru  $V$  ob predpostavki, da je definirano tudi opažanje  $\mathbf{Y}$ . Pri tem dani linearni funkcional preslika v kvadratni koren nepristranske cenilke variance taistega funkcionala, uporabljenega na projekciji opažanja  $\mathbf{Y}$  na  $V$ . Pri tem cenilka variance izhaja iz (2.4), kar je vedno možno, saj je varianca funkcionala znan večkratnik variance modela. Pri pisavi  $\text{SE}(\mathbf{a}^T \mathbf{v})$  je nekaj zlorabe notacije: vektor  $\mathbf{v}$  smo identificirali z identiteto na  $V$ . Tovrstne zlorabe notacije so v matematiki zelo pogoste in se jih navadno niti ne zavedamo: tudi pri diferencialih pišemo spremenljivko, v resnici pa jo uporabimo kot preslikavo.

Če pa je model podan z zvezo (1.2), lahko to prevedemo na (1.1), tako da postavimo  $\mathbf{v} = \mathbf{X}\boldsymbol{\beta}$ . Tedaj se pravokotni projektor izraža v obliki  $\mathbf{P}_{W \rightarrow V} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ . Pišemo lahko  $\boldsymbol{\beta} = \mathbf{X}^{-1} \mathbf{v}$ , kjer je inverz  $\mathbf{X}^{-1}$  definiran kot preslikava iz  $V$  na  $\mathbb{R}^k$  (medtem ko je  $\mathbf{X}$  preslikava iz  $\mathbb{R}^k$  v  $\mathbb{R}^n = W$ ). Velja:

$$\mathbf{X}^{-1} p_{W \rightarrow V} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T. \quad (2.6)$$

Od tod dobimo še cenilko za  $\boldsymbol{\beta}$ :

$$\hat{\boldsymbol{\beta}} = \mathbf{X}^{-1} \hat{\mathbf{v}} = \mathbf{X}^{-1} p_{W \rightarrow V} \hat{\mathbf{v}} = \mathbf{X}^{-1} p_{W \rightarrow V} \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}.$$

Tudi  $\hat{\boldsymbol{\beta}}$  je nepristranska cenilka za  $\boldsymbol{\beta}$  z najmanjšo možno kovariančno matriko: to sledi iz ustrezne lastnosti cenilke  $\hat{\mathbf{v}}$  za  $\mathbf{v}$ . Ta kovariančna matrika je enaka:

$$\text{var}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{I}_W \left( (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \right)^T = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}. \quad (2.7)$$

Spet lahko izpeljemo tudi pripadajoče rezultate za linearne funkcionalne oblike  $\mathbf{c}^T \boldsymbol{\beta}$ , kjer je  $\mathbf{c} \in \mathbb{R}^k$ . S pomočjo zveze (2.7) dobimo varianco:

$$\text{var}(\mathbf{c}^T \hat{\boldsymbol{\beta}}) = \mathbf{c}^T \text{var}(\hat{\boldsymbol{\beta}}) \mathbf{c} = \sigma^2 \mathbf{c}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c} \quad (2.8)$$

To pa lahko izpeljemo tudi tako, da zvezo (1.2) prevedemo na (1.1), tako da postavimo  $\mathbf{v} := \mathbf{X}\boldsymbol{\beta}$ . V skladu z zvezo (2.6) potem dobimo:

$$\mathbf{c}^T \boldsymbol{\beta} = \mathbf{c}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{v},$$

torej lahko postavimo  $\mathbf{a} := \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c}$  in po (2.3) velja:

$$\text{var}(\mathbf{c}^T \hat{\boldsymbol{\beta}}) = \sigma^2 \|\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c}\|^2$$

in ni se težko prepričati, da je to isto kot v (2.8). Iz variance dobimo ocenjeno standardno napako:

$$\text{SE}(\mathbf{c}^T \boldsymbol{\beta}) = \hat{\sigma} \sqrt{\mathbf{c}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c}} = \hat{\sigma} \|\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c}\|.$$

Ocenjena varianca modela  $\hat{\sigma}^2$ , definirana v (2.4), se izraža z rezidualom  $\hat{\boldsymbol{\varepsilon}}$ . Iz prevedbe zapisa (1.2) na zapis (1.1) dobimo, da ima rezidual zdaj obliko:

$$\hat{\boldsymbol{\varepsilon}} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}.$$

Z drugimi besedami, rezidual je razlika med *opaženimi vrednostmi*  $\mathbf{Y}$  in *teoretičnimi vrednostmi*  $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ .

Videli smo, da je, če je model podan z zvezo (1.1), rezidual nekoreliran z  $\hat{\mathbf{v}}$ . No, če je podan z zvezo (1.2), pa je nekoreliran z  $\hat{\boldsymbol{\beta}}$ , saj je:

$$\text{cov}(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\varepsilon}}) = \text{cov}(\mathbf{X}^{-1} p_{W \rightarrow V} \hat{\mathbf{v}}, \hat{\boldsymbol{\varepsilon}}) = \mathbf{X}^{-1} p_{W \rightarrow V} \text{cov}(\hat{\mathbf{v}}, \hat{\boldsymbol{\varepsilon}}) = \mathbf{0}.$$

Včasih pa nas zanimajo tudi količine, ki nastanejo kot vsota določenega funkcionala na parametričnem prostoru in dodatnega šuma. Aproximaciji teh količin z opaženimi količinami pravimo *napovedovanje*. Če privzamemo, da je pričakovana vrednost dodatnega šuma enaka nič, lahko dobimo nepristransko napoved v smislu, da je  $\mathbb{E}(\hat{S}) = \mathbb{E}(S)$ , kjer je  $S$  količina, ki jo napovedujemo,  $\hat{S}$  pa njena napoved. Če želimo zdaj oceniti napako, ki jo naredimo, potrebujemo *srednjo kvadratično napako*  $\mathbb{E}((\hat{S} - S)^2)$ . Kvadratni koren ocene srednje kvadratične napake bo groba ocena za napako, ki jo bomo naredili. Temu bomo prav tako rekli ocenjena standardna napaka.

Ocenjeno standardno napako bo možno določiti pod pogojem, da je dodatni šum nekoreliran s šumom, ki se nanaša na opažanja, t. j.  $\boldsymbol{\varepsilon}$ , ter da je varianca dodatnega šuma znan večkratnik variance modela. Tako bomo zapisali:

$$S = \mathbf{a}^T \mathbf{v} + \lambda \eta \quad \text{ozioroma} \quad S = \mathbf{c}^T \boldsymbol{\beta} + \lambda \eta \quad (2.9)$$

kjer je slučajna spremenljivka  $\eta$  nekorelirana z  $\boldsymbol{\varepsilon}$  ter velja  $\mathbb{E} \eta = 0$  in  $\text{var}(\eta) = \sigma^2$ ,  $\lambda$  pa je znan koeficient. Nepristranska napoved za  $S$  bo seveda:

$$\hat{S} = \mathbf{a}^T \hat{\mathbf{v}} \quad \text{ozioroma} \quad \hat{S} = \mathbf{c}^T \hat{\boldsymbol{\beta}}, \quad (2.10)$$

srednja kvadratična napaka pa bo:

$$\mathbb{E}((\hat{S} - S)^2) = \text{var}(\hat{S} - S) = \text{var}(\mathbf{a}^T \hat{\mathbf{v}}) + \lambda^2 \sigma^2 = \sigma^2 (\|\mathbf{a}\|^2 + \lambda^2)$$

ozioroma:

$$\mathbb{E}((\hat{S} - S)^2) = \text{var}(\hat{S} - S) = \text{var}(\mathbf{c}^T \hat{\boldsymbol{\beta}}) + \lambda^2 \sigma^2 = \sigma^2 (\mathbf{c}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c} + \lambda^2).$$

Tako bo ocenjena standardna napaka:

$$\text{SE}(S) := \frac{\hat{\sigma}}{\sigma} \sqrt{\mathbb{E}((\hat{S} - S)^2)} \quad (2.11)$$

vedno opazljiva, in sicer enaka:

$$\text{SE}(S) = \hat{\sigma} \sqrt{\|\mathbf{a}\|^2 + \lambda^2} \quad (2.12)$$

ozioroma:

$$\text{SE}(S) = \hat{\sigma} \sqrt{\mathbf{c}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c} + \lambda^2} = \hat{\sigma} \sqrt{\|\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c}\|^2 + \lambda^2}. \quad (2.13)$$

### 3 Intervalsko in območno ocenjevanje in napovedovanje

Pri intervalskem ocenjevanju se bomo omejili na *linearne* funkcionalov vektorjev  $\mathbf{v}$  ozioroma  $\boldsymbol{\beta}$ , ki jim bo lahko prištet še neodvisen *šum*. V slednjem primeru bomo temu rekli *napovedovanje*. Poleg tega se bomo ukvarjali tudi z linearnimi preslikavami vektorjev  $\mathbf{v}$  in  $\boldsymbol{\beta}$

(v večrazsežne prostore), spet lahko s prištetim šumom. Tovrstne količine bomo ocenjevali oz. napovedovali *območno*.

Pri intervalskem oz. območnem ocenjevanju oz. napovedovanju moramo seveda vedeti nekaj več o porazdelitvah. Privzeli bomo *Gaussov model*, kar pomeni, da je šum porazdeljen *večrazsežno normalno*.

Potrebovali bomo naslednja dejstva o večrazsežni normalni porazdelitvi in njenih izpeljankah:

- Če sta  $V \subseteq W$  evklidska prostora in je slučajni vektor  $\mathbf{Z}$  porazdeljen standardno normalno na  $W$ , je  $p_{W \rightarrow V} \mathbf{Z}$  porazdeljen standardno normalno na  $V$ .
- Če sta slučajna vektorja  $\mathbf{Z}_1$  in  $\mathbf{Z}_2$  nekorelirana in je njuna skupna porazdelitev večrazsežna normalna, sta ta dva slučajna vektorja tudi neodvisna.
- Če je  $\mathbf{Z}$  porazdeljen standardno normalno na evklidskem prostoru  $V$ , ima slučajna spremenljivka  $\|\mathbf{Z}\|^2$  porazdelitev hi kvadrat s toliko prostostnimi stopnjami, kolikor je dimenzija prostora  $V$ ; pišemo  $\|\mathbf{Z}\|^2 \sim \chi^2(\dim V)$ .
- Če je slučajna spremenljivka  $Z$  porazdeljena standardno enorazsežno normalno in je  $U \sim \chi^2(n)$  od nje neodvisna slučajna spremenljivka, ima slučajna spremenljivka  $\frac{Z}{\sqrt{U}} \sqrt{n}$  Studentovo porazdelitev z  $n$  prostostnimi stopnjami; pišemo  $\frac{Z}{\sqrt{U}} \sqrt{n} \sim \text{Student}(n)$ .
- Če sta  $U \sim \chi^2(m)$  in  $V \sim \chi^2(n)$  neodvisni slučajni spremenljivki, ima  $\frac{U/m}{V/n}$  Fisherjevo porazdelitev z  $(m, n)$  prostostnimi stopnjami; pišemo  $\frac{U/m}{V/n} \sim \text{Fisher}(m, n)$ .

Oglejmo si zdaj, kaj pomenijo ta dejstva pri predpostavkah iz prejšnjega razdelka. Najprej se spomnimo, da sta slučajna vektorja  $\hat{\mathbf{v}}$  in  $\hat{\boldsymbol{\epsilon}}$  pri predpostavkah, navedenih pod zvezo (1.1), nekorelirana. Če je model Gaussov, sta torej neodvisna. Podobno sta, če privzamemo zvezo (1.2) in vse potrebne predpostavke, neodvisna slučajna vektorja  $\hat{\boldsymbol{\beta}}$  in  $\hat{\boldsymbol{\epsilon}}$ .

V prejšnjem razdelku smo izračunali matematično upanje slučajne spremenljivke  $\|\hat{\boldsymbol{\epsilon}}\|^2$ . Če je model Gaussov, poznamo tudi njeno porazdelitev. Slučajni vektor  $\boldsymbol{\epsilon}/\sigma$  je namreč porazdeljen standardno normalno. Slučajni vektor  $\hat{\boldsymbol{\epsilon}}/\sigma = q_{W \rightarrow V} \boldsymbol{\epsilon}/\sigma$  pa je pravokotna projekcija tega vektorja na ortogonalni komplement prostora  $V$  v  $W$ , torej je prav tako porazdeljen standardno normalno na tem ortogonalnem komplementu, ki je  $(n-k)$ -razsežen prostor. To pa pomeni, da je:

$$\frac{\|\hat{\boldsymbol{\epsilon}}\|^2}{\sigma^2} \sim \chi^2(n-k). \quad (3.1)$$

Lotimo se sedaj iskanja napovednega intervala za količino  $S$ , definirano tako kot v (2.9), t. j.  $S = \mathbf{a}^T \mathbf{v} + \lambda \eta$ , če je model podan z zvezo (1.1), oziroma  $S = \mathbf{c}^T \boldsymbol{\beta} + \lambda \eta$ , če je model podan z zvezo (1.2). Tu je slučajna spremenljivka  $\eta$  porazdeljena normalno  $N(0, \sigma^2)$  in



neodvisna od  $\varepsilon$ ,  $\lambda$  pa je znan koeficient. Naj bo  $\hat{S}$  nepristranska napoved za  $S$ , definirana v (2.10). Slučajna spremenljivka:

$$\frac{\hat{S} - S}{\sqrt{\mathbb{E}((\hat{S} - S)^2)}}$$

je porazdeljena standardno normalno in je neodvisna od  $\varepsilon$ . Zato je:

$$\frac{\frac{\hat{S} - S}{\sqrt{\mathbb{E}((\hat{S} - S)^2)}}}{\frac{\|\varepsilon\|}{\sigma}} \sqrt{n - k} = \frac{\sigma}{\hat{\sigma}} \frac{\hat{S} - S}{\sqrt{\mathbb{E}((\hat{S} - S)^2)}} = \frac{\hat{S} - S}{\text{SE}(S)} \sim \text{Student}(n - k) \quad (3.2)$$

(uporabili smo še (2.4) in (2.11)). Če torej s  $t_p(m)$  označimo kvantil Studentove porazdelitve z  $m$  prostostnimi stopnjami za verjetnost  $p$ , je napovedni interval za  $S$  pri stopnji zaupanja  $\beta$  možno zapisati v obliki:

$$\hat{S} - \text{SE}(S) t_{(1+\beta)/2}(n - k) < S < \hat{S} + \text{SE}(S) t_{(1+\beta)/2}(n - k).$$

Oglejmo si še primer, ko napovedujemo slučajni vektor  $\mathbf{S} := \mathbf{A}\mathbf{v} + \mathbf{\Lambda}\boldsymbol{\eta}$  z vrednostmi v vektorskem prostoru  $U$ , kjer je dodatni šum  $\boldsymbol{\eta} \sim \text{N}(\mathbf{0}, \sigma^2 \mathbf{I}_U)$  neodvisen od  $\varepsilon$ . V tem primeru lahko določimo *napovedno območje* – množico v  $U$ . Velja:

$$\mathbf{A}\hat{\mathbf{v}} - \mathbf{S} = \mathbf{A}(\hat{\mathbf{v}} - \mathbf{v}) - t\boldsymbol{\eta} \sim \text{N}(\mathbf{0}, \sigma^2(\mathbf{A}\mathbf{A}^T + \mathbf{\Lambda}\mathbf{\Lambda}^T)),$$

torej je:

$$\frac{1}{\sigma^2} \left\| (\mathbf{A}\mathbf{A}^T + \mathbf{\Lambda}\mathbf{\Lambda}^T)^{-1/2} (\mathbf{A}\hat{\mathbf{v}} - \mathbf{S}) \right\|^2 = \frac{1}{\sigma^2} (\mathbf{A}\hat{\mathbf{v}} - \mathbf{S})^T (\mathbf{A}\mathbf{A}^T + \mathbf{\Lambda}\mathbf{\Lambda}^T)^{-1} (\mathbf{A}\hat{\mathbf{v}} - \mathbf{S}) \sim \chi^2(m).$$

kjer je  $m = \dim U$ . Če se spomnimo še na (3.1) in dejstvo, da sta slučajna vektorja  $\hat{\mathbf{v}}$  in  $\hat{\boldsymbol{\varepsilon}}$  neodvisna, dobimo:

$$\frac{n - k}{m} \frac{(\mathbf{A}\hat{\mathbf{v}} - \mathbf{S})^T (\mathbf{A}\mathbf{A}^T + \mathbf{\Lambda}\mathbf{\Lambda}^T)^{-1} (\mathbf{A}\hat{\mathbf{v}} - \mathbf{S})}{\|\hat{\boldsymbol{\varepsilon}}\|^2} \sim \text{Fisher}(m, n - k). \quad (3.3)$$

Če torej z  $F_p(m, n)$  označimo kvantil Fisherjeve porazdelitve z  $(m, n)$  prostostnimi stopnjami za verjetnost  $p$ , je iskano napovedno območje pri stopnji zaupanja  $\beta$  določeno s pogojem:

$$(\mathbf{S} - \mathbf{A}\hat{\mathbf{v}})^T (\mathbf{A}\mathbf{A}^T + \mathbf{\Lambda}\mathbf{\Lambda}^T)^{-1} (\mathbf{S} - \mathbf{A}\hat{\mathbf{v}}) < \frac{m}{n - k} F_\beta(m, n - k) \|\hat{\boldsymbol{\varepsilon}}\|^2.$$

Če pa je model podan z zvezo (1.2), lahko poiščemo napovedno območje za slučajni vektor  $\mathbf{S} := \mathbf{C}\boldsymbol{\beta} + \lambda\boldsymbol{\eta}$ . V tem primeru spet postavimo  $\mathbf{v} = \mathbf{X}\boldsymbol{\beta}$  in  $\mathbf{A} = \mathbf{C}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$ . Dobimo napovedno območje:

$$(\mathbf{S} - \mathbf{C}\hat{\boldsymbol{\beta}})^T \left( \mathbf{C}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{C}^T + \mathbf{\Lambda}\mathbf{\Lambda}^T \right)^{-1} (\mathbf{S} - \mathbf{C}\hat{\boldsymbol{\beta}}) < \frac{m}{n - k} F_\beta(m, n - k) \|\hat{\boldsymbol{\varepsilon}}\|^2.$$

## 4 Testiranje hipotez

V celotnem razdelku se bomo omejili na primer, ko je model Gaussov. Testiranje enodimenzionalnih ničelnih hipotez oblike  $s = s_0$  kjer je  $s = \mathbf{a}^T \mathbf{v}$  (če je model podan z zvezo (1.1)) oziroma  $s = \mathbf{c}^T \boldsymbol{\beta}$  (če je model podan z zvezo (1.2)) je enostavno in temelji na relaciji (3.2). Karakteristiko  $s$  najprej točkasto ocenimo tako kot v (2.10), torej:

$$\hat{s} = \mathbf{a}^T \hat{\mathbf{v}} \quad \text{oziroma} \quad \hat{s} = \mathbf{c}^T \hat{\boldsymbol{\beta}}.$$

Za testiranje ničelne hipoteze  $s = s_0$  izračunamo testno statistiko:

$$T = \frac{\hat{s} - s_0}{\text{SE}(s)},$$

kjer je ocenjena standardna napaka  $\text{SE}(s)$  definirana tako kot v (2.5) oziroma (2). Testna statistika je torej *razmerje med opaženo razliko in ocenjeno standardno napako*. V skladu z (3.2) je  $T \sim \text{Student}(n - k)$ , zato bomo količini  $T$  rekli *Studentova statistika*; na njej torej izvedemo  $T$ -test z  $n - k$  prostostnimi stopnjami (eno- ali dvostranskega).

Podobno bi lahko na podlagi relacije (3.3) (pri  $\boldsymbol{\Lambda} = 0$ ) testirali hipoteze o vektorjih  $\mathbf{A}\mathbf{v}$  ali  $\mathbf{C}\boldsymbol{\beta}$ : le-te bi trdile, da se ustrezni vektor nahaja na nekem afinem prostoru. Vendar pa se izkaže priročneje hipoteze formulirati nekoliko drugače.

Če delamo z zvezo (1.1), problem formuliramo tako, da v prototipu  $\mathbf{Y} = \mathbf{v} + \boldsymbol{\varepsilon}$  testiramo ožji model  $\mathbf{v} \in U$  proti širšemu modelu  $\mathbf{v} \in V$ , kjer je  $U$  podprostor prostora  $V$ . Lahko bi vzeli, da je  $U$  afin podprostor, a problem lahko brez težav premaknemo, tako da bomo privzeli, da je  $U$  kar vektorski podprostor. Spomnimo se oznak  $n = \dim W$  in  $k = \dim V$  in označimo še  $l := \dim U$ . Test temelji na razmerju projekcij  $q_{V \rightarrow U} \mathbf{Y}$  in  $q_{W \rightarrow V} \mathbf{Y} = \hat{\boldsymbol{\varepsilon}}$ . Vemo, da sta ti dve projekciji nekorelirani, če je model Gaussov, pa tudi neodvisni. Ker je:

$$\frac{1}{\sigma^2} \|q_{V \rightarrow U} \mathbf{Y}\|^2 \sim \chi^2(k - l) \quad \text{in} \quad \frac{1}{\sigma^2} \|q_{W \rightarrow V} \mathbf{Y}\|^2 \sim \chi^2(n - k),$$

je:

$$F := \frac{n - k}{k - l} \frac{\|q_{V \rightarrow U} \mathbf{Y}\|^2}{\|q_{W \rightarrow V} \mathbf{Y}\|^2} = \frac{n - k}{k - l} \frac{\|\mathbf{Q}_{V \rightarrow U} \mathbf{Y}\|^2}{\|\mathbf{Q}_{W \rightarrow V} \mathbf{Y}\|^2} \sim \text{Fisher}(k - l, n - k).$$

Glede na porazdelitev bomo količini  $F$  rekli *Fisherjeva statistika*; na njej izvedemo enostranski  $F$ -test. Na to Fisherjevo statistiko pa lahko pogledamo še malo drugače. Glede na to, da je:

$$\mathbf{Q}_{W \rightarrow U} = \mathbf{Q}_{W \rightarrow V} + \mathbf{Q}_{V \rightarrow U},$$

rezidual v ožjem modelu  $\mathbf{Q}_{W \rightarrow U} \mathbf{Y}$  razpade na *pojasnjeni rezidual*  $\mathbf{Q}_{V \rightarrow U} \mathbf{Y}$  in *nepojasnjeni rezidual*  $\mathbf{Q}_{W \rightarrow V} \mathbf{Y}$ . Izraz ‘pojasnjen’ tu pomeni ‘pojasnjen z razširitvijo modela’: ta rezidual izgine, ko ožji model nadomestimo s širšim, nepojasnjeni rezidual pa je tisti, ki ostane tudi potem, ko model že razširimo.

Pojasnjeni in nepojasnjeni rezidual sta pravokotna, zato se kvadrata njunih norm seštejeta v kvadrat norme celotnega reziduala. Tako v skladu s formulo (2.4) razpade tudi

varianca, ocenjena na podlagi ožjega modela:

$$\hat{\sigma}^2 = \frac{\|\mathbf{Q}_{W \rightarrow U} \mathbf{Y}\|^2}{n-l} = \frac{\|\mathbf{Q}_{V \rightarrow U} \mathbf{Y}\|^2}{n-l} + \frac{\|\mathbf{Q}_{W \rightarrow V} \mathbf{Y}\|^2}{n-l}.$$

Prvemu členu bomo rekli *pojasnjena ocenjena varianca modela*, drugemu pa *nepojasnjena ocenjena varianca modela*.

Formula (2.4) pa pove še nekaj drugega: če v splošnem za vektorska prostora  $G \subseteq H$  s:

$$\hat{\sigma}_{H \rightarrow G}^2 = \frac{\|\mathbf{Q}_{H \rightarrow G} \mathbf{Y}\|^2}{\dim H - \dim G} \quad (4.1)$$

označimo varianco, ocenjeno na podlagi razkoraka med  $G$  in  $H$ , dobimo, da je Fisherjeva statistika enaka:

$$F = \frac{\hat{\sigma}_{V \rightarrow U}^2}{\hat{\sigma}_{W \rightarrow V}^2}. \quad (4.2)$$

To je torej razmerje med varianco, ocenjeno s pomočjo pojasnjenega reziduala, in varianco, ocenjeno s pomočjo nepojasnjene reziduala.

Če pa delamo z zvezo (1.2), zadevo zastavimo kot test ožjega modela:

$$\mathbf{Y} = \mathbf{Z}\boldsymbol{\gamma} + \boldsymbol{\eta}; \quad \boldsymbol{\gamma} \in \mathbb{R}^l \quad (4.3)$$

proti širšemu modelu:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}; \quad \boldsymbol{\beta} \in \mathbb{R}^k; \quad (4.4)$$

veljati mora  $\text{im } \mathbf{Z} \subset \text{im } \mathbf{X}$ , torej  $l < k$ . Če postavimo  $U := \text{im } \mathbf{Z}$  in  $V := \text{im } \mathbf{X}$  (in se spomnimo, da je  $W = \mathbb{R}^n$ ), velja:

$$\mathbf{Q}_{W \rightarrow V} = \mathbf{I}_W - \mathbf{P}_{W \rightarrow V} = \mathbf{I}_W - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T, \quad (4.5)$$

$$\mathbf{Q}_{V \rightarrow U} = \mathbf{P}_{W \rightarrow V} - \mathbf{P}_{W \rightarrow U} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T - \mathbf{Z}(\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T. \quad (4.6)$$

Glede na to, da za vsak ortogonalni projektor  $\mathbf{Q}$  velja  $\|\mathbf{Q}\mathbf{Y}\|^2 = \mathbf{Y}^T \mathbf{Q}^T \mathbf{Q} \mathbf{Y} = \mathbf{Y}^T \mathbf{Q} \mathbf{Y}$ , se ocenjeni varianci, potrebni za izračun Fisherjeve statistike, izražata na naslednji način:

$$\begin{aligned} \hat{\sigma}_{W \rightarrow V}^2 &= \frac{\|\mathbf{Q}_{W \rightarrow V} \mathbf{Y}\|^2}{n-k} = \frac{\mathbf{Y} \left( \mathbf{I}_W - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \right) \mathbf{Y}^T}{n-k}, \\ \hat{\sigma}_{V \rightarrow U}^2 &= \frac{\|\mathbf{Q}_{V \rightarrow U} \mathbf{Y}\|^2}{k-l} = \frac{\mathbf{Y} \left( \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T - \mathbf{Z}(\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \right) \mathbf{Y}^T}{k-l}. \end{aligned}$$

## 5 Posebni primeri

### 5.1 Pričakovana vrednost

Denimo, da opazimo vrednosti nekoreliranih slučajnih spremenljivk  $Y_1, Y_2, \dots, Y_n$ , ki imajo vse matematično upanje  $\mu$  in varianco  $\sigma^2$ . Želeli bi oceniti  $\mu$  in  $\sigma$ . To lahko predstavimo

kot  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , kjer je:

$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}, \quad \boldsymbol{\beta} = \mu, \quad \boldsymbol{\varepsilon} = \begin{bmatrix} Y_1 - \mu \\ Y_2 - \mu \\ \vdots \\ Y_n - \mu \end{bmatrix}.$$

Z drugimi besedami, to je zveza (1.3) za  $k = 1$  in  $f_1(x) = 1$ . Ni težko videti, da so izpolnjene vse predpostavke, navedene k (1.2).

Za ocenjevanje parametra  $\mu$  izračunamo  $\mathbf{X}^T \mathbf{X} = n$ , torej je točkasta ocena:

$$\hat{\mu} = \frac{1}{n} [1 \quad 1 \quad \cdots \quad 1] \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \frac{Y_1 + Y_2 + \cdots + Y_n}{n} =: \bar{Y}.$$

Rezidual je enak:

$$\hat{\boldsymbol{\varepsilon}} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \bar{Y} = \begin{bmatrix} Y_1 - \bar{Y} \\ Y_2 - \bar{Y} \\ \vdots \\ Y_n - \bar{Y} \end{bmatrix}$$

in nepristranska cenilka za varianco je:

$$\hat{\sigma}^2 = \frac{\|\hat{\boldsymbol{\varepsilon}}\|^2}{n-1} = \frac{1}{n-1} \left[ (Y_1 - \bar{Y})^2 + (Y_2 - \bar{Y})^2 + \cdots + (Y_n - \bar{Y})^2 \right].$$

To je *popravljen vzorčna varianca*. Izračunajmo še ocenjeno standardno napako za  $\mu$ . Pišemo lahko  $\mu = \mathbf{c}^T \boldsymbol{\beta}$ , kjer je, kot že rečeno,  $\boldsymbol{\beta} = \mu$  in  $\mathbf{c} = 1$ . Torej je  $\mathbf{c}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c} = 1/n$  in iz (2) dobimo ocenjeno standardno napako:

$$\text{SE}(\mu) = \frac{\hat{\sigma}}{\sqrt{n}}.$$

## 5.2 Enostavna linearna regresija

Enostavna linearna regresija je model *linearnega trenda*, ki pomeni, da je odvisna spremenljivka linearna funkcija pojasnjevalne spremenljivke z določenimi napakami. Naredimo več meritev: pri  $i$ -ti meritvi naj ima pojasnjevalna spremenljivka vrednost  $X_i$ , odvisna pa  $Y_i$ . Iščemo premico  $y = a + bx$ , ki se tem meritvam najboljše prilega, merilo za prileganje pa naj bo vsota kvadratov odstopanj *v smeri y*. Iskani premici bomo rekli *regresijska premica*. Glede na kriterije, ki smo jih postavili, je vrednosti pojasnjevalne spremenljivke  $X_i$  smiselno proglasiti za deterministične, vrednosti odvisne spremenljivke  $Y_i$  pa za slučajne,

tako da so odstopanja od premice v smeri  $y$  šumi. Tako lahko model predstavimo z zvezo (1.3), kjer je:

$$k = 2, \quad f_1(x) = 1, \quad f_2(x) = x, \quad \beta_1 = a, \quad \beta_2 = b.$$

Ekvivalentno lahko model predstavimo z zvezo (1.2), kjer je:

$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} 1 & X_1 \\ 1 & X_2 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} a \\ b \end{bmatrix}.$$

Vektor  $\boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$  je vektor odstopanj od premice v smeri  $y$ . Seveda privzamemo, da so

$\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$  nekorelirane z  $\mathbb{E} \varepsilon_i = 0$  in  $\text{var}(\varepsilon_i) = \sigma^2$ .

Za točkasto oceno parametrov  $a$  in  $b$  izračunamo:

$$\begin{aligned} \mathbf{X}^T \mathbf{X} &= \begin{bmatrix} n & \sum_{i=1}^n X_i \\ \sum_{i=1}^n X_i & \sum_{i=1}^n X_i^2 \end{bmatrix}, \\ (\mathbf{X}^T \mathbf{X})^{-1} &= \frac{1}{n \sum_{i=1}^n X_i^2 - (\sum_{i=1}^n X_i)^2} \begin{bmatrix} \sum_{i=1}^n X_i^2 & -\sum_{i=1}^n X_i \\ -\sum_{i=1}^n X_i & n \end{bmatrix}, \\ \hat{\boldsymbol{\beta}} = \begin{bmatrix} \hat{a} \\ \hat{b} \end{bmatrix} &= \frac{1}{n \sum_{i=1}^n X_i^2 - (\sum_{i=1}^n X_i)^2} \begin{bmatrix} \sum_{i=1}^n X_i^2 \sum_{i=1}^n Y_i - \sum_{i=1}^n X_i \sum_{i=1}^n X_i Y_i \\ n \sum_{i=1}^n X_i Y_i - \sum_{i=1}^n X_i \sum_{i=1}^n Y_i \end{bmatrix}. \end{aligned}$$

Vektor rezidualov pa ima pomembno lastnost, da je vsota njegovih komponent enaka nič. Za vsak vektor  $\mathbf{a}$  je namreč vsoto komponent mogoče izraziti kot  $[1 \ 0] \mathbf{X}^T \mathbf{a}$ . Vsota komponent reziduala je tako enaka:

$$[1 \ 0] \mathbf{X}^T (\mathbf{Y} - \mathbf{X} \boldsymbol{\beta}) = [1 \ 0] \left( \mathbf{X}^T \mathbf{Y} - \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \right) = 0.$$

Sledi  $\sum_{i=1}^n (Y_i - \hat{a} - \hat{b} X_i) = 0$ . Če z  $\bar{X} := \frac{1}{n} \sum_{i=1}^n X_i$  in  $\bar{Y} := \frac{1}{n} \sum_{i=1}^n Y_i$  označimo povprečji, velja tudi  $\bar{Y} = \hat{a} + \hat{b} \bar{X}$ . Točka  $(\bar{X}, \bar{Y})$  torej vedno leži na regresijski premici.

S pomočjo povprečij pa se dasta tudi alternativno izraziti cenilki  $\hat{a}$  in  $\hat{b}$  regresijskih koeficientov. Velja namreč:

$$\begin{aligned} \hat{b} &= \frac{\sum_{i=1}^n (X_i - \bar{X}) Y_i}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\sum_{i=1}^n (X_i - \bar{X}) (Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}, \\ \hat{a} &= \bar{Y} - \hat{b} \bar{X} = \sum_{i=1}^n \left( \frac{1}{n} - \frac{\bar{X} (X_i - \bar{X})}{\sum_{i=1}^n (X_i - \bar{X})^2} \right) Y_i. \end{aligned}$$

Iz reziduala dobimo še nepristransko cenilko variance modela:

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{a} - \hat{b}X_i)^2.$$

V nadaljevanju bomo izračunali ocenjene standardne napake za  $a$ ,  $b$  in  $aX_0 + b + \lambda\varepsilon_0$ . Slednje je mišljeno kot nova izmerjena vrednost odvisne spremenljivke pri vrednosti pojasnjevalne spremenljivke  $X_0$ , pri čemer dopuščamo možnost, da je šum drugačne velikosti kot pri že znanih meritvah. Seveda privzamemo, da je novi šum  $\varepsilon_0$  nekoreliran z  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$  ter da velja  $\mathbb{E}\varepsilon_0 = 0$  in  $\text{var}(\varepsilon_0) = \sigma^2$ .

Ocenjeno standardno napako bomo najprej izračunali za strmino  $b$ . Za ta namen postavimo  $\mathbf{c} := \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ : tako bo  $\mathbf{c}^T \boldsymbol{\beta} = b$ . Izračunajmo:

$$\mathbf{c}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c} = \frac{n}{n \sum_{i=1}^n X_i^2 - (\sum_{i=1}^n X_i)^2} = \frac{1}{\sum_{i=1}^n (X_i - \bar{X})^2}.$$

Iz (2) zdaj dobimo:

$$\text{SE}(b) = \frac{\hat{\sigma}}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2}}.$$

Za statistično sklepanje o koeficientu  $a$  moramo postaviti  $\mathbf{c} := \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ . Dobimo:

$$\begin{aligned} \mathbf{c}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c} &= \frac{\sum_{i=1}^n X_i^2}{n \sum_{i=1}^n X_i^2 - (\sum_{i=1}^n X_i)^2} = \frac{\sum_{i=1}^n (X_i - \bar{X})^2 + n\bar{X}^2}{n \sum_{i=1}^n (X_i - \bar{X})^2} = \\ &= \frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2}. \end{aligned}$$

Sledi:

$$\text{SE}(a) = \hat{\sigma} \sqrt{\frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}.$$

Zdaj se pa spomnimo, da je  $a$  koordinata  $y$  na regresijski premici pri  $x = 0$ . Novo meritev odvisne spremenljivke pri  $X_0 = 0$  torej predstavimo kot  $a + \lambda\varepsilon_0$ . Podobno kot prej iz (2.11) dobimo:

$$\text{SE}(a + \lambda\varepsilon_0) = \hat{\sigma} \sqrt{\lambda^2 + \frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}.$$

Za izračun ocenjene standardne napake za novo izmerjeno vrednost pri splošni vrednosti pojasnjevalne spremenljivke  $X_0$  pa lahko model preprosto premaknemo. Naš izvirni model:

$$Y_i = a + bX_i + \varepsilon_i$$

prepišimo v obliki:

$$Y_i = \tilde{a} + b\tilde{X}_i + \varepsilon_i,$$

kjer je  $\tilde{X}_i = X_i - X_0$  in  $\tilde{a} = a + bX_0$ . Slednje je spet enostavna linearna regresija. Pomudimo se najprej pri oceni variance modela. To lahko ocenimo v obeh modelih. Toda zaradi linearnosti cenilk je  $\hat{\tilde{a}} = \hat{a} + \hat{b}X_0$  in zato:

$$Y_i - \hat{\tilde{a}} - \hat{b}\tilde{X}_i = Y_i - \hat{a} - \hat{b}X_0 - \hat{b}(X_i - X_0) = Y_i - \hat{a} - \hat{b}X_i.$$

Torej je ocena za  $\sigma^2$  v obeh modelih enaka:

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{\tilde{a}} - \hat{b}\tilde{X}_i)^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{a} - \hat{b}X_i)^2.$$

Zdaj lahko izračunamo želeno ocenjeno standardno napako:

$$\text{SE}(a + bX_0 + \lambda\varepsilon_0) = \text{SE}(\tilde{a} + \lambda\varepsilon_0) = \hat{\sigma} \sqrt{\frac{1}{n} + \frac{\tilde{X}^2}{\sum_{i=1}^n (\tilde{X}_i - \tilde{X})^2}} = \hat{\sigma} \sqrt{\frac{1}{n} + \frac{(\bar{X} - X_0)^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}.$$

### 5.3 Primerjava po parih

Denimo, da vsakič isto stvar izmerimo v dveh okoliščinah, recimo pred določeno akcijo in po njej. To napravimo na  $m$  enotah. Meritev na posamezni enoti v prvih okoliščinah (recimo pred akcijo) označimo z  $Y'_i$ , v drugih okoliščinah (recimo po akciji) pa z  $Y''_i$ . Zanima nas, ali spremenjene okoliščine (npr. izvedena akcija) kaj vplivajo na rezultate naših meritev. Privzamemo, da so meritve nekorelirane ter da je:

$$\mathbb{E}(Y'_i) = \mu_i, \quad \mathbb{E}(Y''_i) = \mu_i + \delta, \quad \text{var}(Y'_i) = (\sigma')^2, \quad \text{var}(Y''_i) = (\sigma'')^2.$$

Pričakovana vrednost se torej lahko spreminja od enote do enote, za pričakovani vpliv spremembe okoliščin (akcije) pa privzamemo, da je ves čas enak. Varianca je lahko odvisna od okoliščin, ne pa tudi od enote.

Če želimo prevesti v jezik linearne regresije, najprej zapišimo:

$$Y'_i = \mu_i + \varepsilon'_i, \quad Y''_i = \mu_i + \delta + \lambda\varepsilon''_i,$$

kjer so  $\varepsilon'_i$  in  $\varepsilon''_i$  nekorelirane z  $\mathbb{E}\varepsilon'_i = \mathbb{E}\varepsilon''_i = 0$  in  $\text{var}(\varepsilon'_i) = \text{var}(\varepsilon''_i) = (\sigma')^2$ . Torej je  $\lambda = \sigma''/\sigma$ . To je neznan determinističen parameter, ki pa ne bo konstitutivni del regresijskega modela. Model bo deloval, če  $\lambda$  v statističnem sklepanju ne bo zajet.

Če označimo:

$$\mathbf{Y}' = \begin{bmatrix} Y'_1 \\ Y'_2 \\ \vdots \\ Y'_m \end{bmatrix}, \quad \mathbf{Y}'' = \begin{bmatrix} Y''_1 \\ Y''_2 \\ \vdots \\ Y''_m \end{bmatrix}, \quad \boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_m \end{bmatrix}, \quad \boldsymbol{\varepsilon}' = \begin{bmatrix} \varepsilon'_1 \\ \varepsilon'_2 \\ \vdots \\ \varepsilon'_m \end{bmatrix}, \quad \boldsymbol{\varepsilon}'' = \begin{bmatrix} \varepsilon''_1 \\ \varepsilon''_2 \\ \vdots \\ \varepsilon''_m \end{bmatrix}$$

ter če še z  $\mathbf{0}_m$  označimo stolpec iz  $m$  ničel, z  $\mathbf{1}_m$  stolpec iz  $m$  enic, z  $\mathbf{I}_m$  pa identiteto na  $\mathbb{R}^m$ , lahko model zapišemo v obliki (1.2) z:

$$\mathbf{Y} = \begin{bmatrix} \mathbf{Y}' \\ \frac{1}{\lambda} \mathbf{Y}'' \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} \mathbf{I}_m & \mathbf{0}_m \\ \frac{1}{\lambda} \mathbf{I}_m & \frac{1}{\lambda} \mathbf{1}_m \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \boldsymbol{\mu} \\ \delta \end{bmatrix}, \quad \boldsymbol{\varepsilon} = \begin{bmatrix} \boldsymbol{\varepsilon}' \\ \boldsymbol{\varepsilon}'' \end{bmatrix}.$$

Velja še  $n = 2m$  in  $k = m + 1$ . Izračunajmo:

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} (1 + \frac{1}{\lambda^2}) \mathbf{I}_m & \frac{1}{\lambda^2} \mathbf{1}_m \\ \mathbf{1}_m^T & n \lambda^{-2} \end{bmatrix},$$

$$(\mathbf{X}^T \mathbf{X})^{-1} = \begin{bmatrix} \frac{\lambda^2}{1 + \lambda^2} \mathbf{I}_m + \frac{1}{n(1 + \lambda^2)} \mathbf{1}_m \mathbf{1}_m^T & -\frac{1}{n} \mathbf{1}_m \\ -\frac{1}{n} \mathbf{1}_m^T & \frac{1}{n} (1 + \lambda^2) \end{bmatrix},$$

Od tod po krajšem računu dobimo cenilko za  $\boldsymbol{\beta}$ :

$$\hat{\boldsymbol{\beta}} = \begin{bmatrix} \frac{\lambda^2}{1 + \lambda^2} \mathbf{Y}' + \frac{1}{1 + \lambda^2} \mathbf{Y}'' - \frac{1}{n(1 + \lambda^2)} \mathbf{1}_m \mathbf{1}_m^T (\mathbf{Y}'' - \mathbf{Y}') \\ \frac{1}{n} \mathbf{1}_m^T (\mathbf{Y}'' - \mathbf{Y}') \end{bmatrix}.$$

Po komponentah lahko to zapišemo tudi takole:

$$\hat{\mu}_i = \frac{\lambda^2}{1 + \lambda^2} Y'_i + \frac{1}{1 + \lambda^2} Y''_i + \frac{1}{(1 + \lambda^2)} (\bar{Y}'' - \bar{Y}'),$$

$$\hat{\delta} = \bar{Y}'' - \bar{Y}'.$$

Posebej opazimo, da je cenilka za  $\delta$  neodvisna od dodatnega parametra  $\lambda$  (in zelo intuitivna). Izračunajmo še ocenjeno standardno napako. Najprej opazimo:

$$[0 \ \dots \ 0 \ 1] (\mathbf{X}^T \mathbf{X})^{-1} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} = \frac{1 + \lambda^2}{n}.$$

Zdaj moramo dobiti še cenilko variance modela, ta pa temelji na rezidualu. Najprej izračunajmo:

$$\mathbf{X} \hat{\boldsymbol{\beta}} = \begin{bmatrix} \frac{\lambda^2}{1 + \lambda^2} \mathbf{Y}' + \frac{1}{1 + \lambda^2} \mathbf{Y}'' - \frac{1}{n(1 + \lambda^2)} \mathbf{1}_m \mathbf{1}_m^T (\mathbf{Y}'' - \mathbf{Y}') \\ \frac{\lambda}{1 + \lambda^2} \mathbf{Y}' + \frac{1}{\lambda(1 + \lambda^2)} \mathbf{Y}'' + \frac{\lambda}{n(1 + \lambda^2)} \mathbf{1}_m \mathbf{1}_m^T (\mathbf{Y}'' - \mathbf{Y}') \end{bmatrix}.$$

Rezidual je tako enak:

$$\hat{\boldsymbol{\varepsilon}} = \frac{1}{1 + \lambda^2} \begin{bmatrix} \mathbf{I}_m \\ \lambda \mathbf{I}_m \end{bmatrix} (\mathbf{I}_m - \frac{1}{n} \mathbf{1}_m \mathbf{1}_m^T) (\mathbf{Y}'' - \mathbf{Y}'),$$

kvadrat njegove norme pa je enak:

$$\|\hat{\boldsymbol{\varepsilon}}\|^2 = \frac{1}{1 + \lambda^2} \|(\mathbf{I}_m - \frac{1}{n} \mathbf{1}_m \mathbf{1}_m^T) (\mathbf{Y}'' - \mathbf{Y}')\|^2 = \frac{1}{1 + \lambda^2} \sum_{i=1}^n ((Y''_i - Y'_i - (\bar{Y}'' - \bar{Y}'))^2).$$



Iz formule (2.4) dobimo:

$$\hat{\sigma}^2 = \frac{1}{(n-1)(1+\lambda^2)} \sum_{i=1}^n ((Y_i'' - Y_i' - (\bar{Y}'' - \bar{Y}'))^2$$

in končno:

$$\text{SE}(\delta) = \sqrt{\frac{1}{n(n-1)} \sum_{i=1}^n ((Y_i'' - Y_i' - (\bar{Y}'' - \bar{Y}'))^2}.$$

Tudi ocenjena standardna napaka ni odvisna od  $\lambda$ . Opazimo pa še nekaj: natanko isto cenilko in ocenjeno standardno napako bi dobili, če bi uporabili rezultate iz podrazdelka 5.1 na razlikah  $Y_i'' - Y_i'$ , ki seveda ustrezajo tamkajšnjemu regresijskemu modelu. To ne velja le pod predpostavkami iz tega razdelka. Lahko jih namreč oslabimo tako, da ne privzamemo več nekoreliranosti slučajnih spremenljivk  $Y_1', \dots, Y_m', Y_1'', \dots, Y_m''$ , marveč le nekoreliranost parov  $(Y_1', Y_1''), \dots, (Y_m', Y_m'')$ : znotraj para  $(Y_i', Y_i'')$  je dopustna tudi koreliranost.

## 5.4 Primerjava dveh skupin

Privzemimo, da isto količino izmerimo na dveh skupinah. Meritve na prvi skupini so  $Y_1', Y_2', \dots, Y_{m'}'$ , na drugi skupini pa  $Y_1'', Y_2'', \dots, Y_{m''}''$ . Za vse meritve privzamemo, da so neodvisne. Nadalje naj bo  $\mathbb{E}(Y_i') = \mu'$ ,  $\mathbb{E}(Y_i'') = \mu''$  in  $\text{var}(Y_i') = \text{var}(Y_i'') = \sigma^2$ . V nasprotju s prejšnjim podrazdelkom sta torej lahko skupini različno veliki, zato pa privzamemo popolno neodvisnost in homoskedastičnost, poleg tega pa se pričakovana vrednost ne sme spremenjati od enote do enote. Podobno kot prej pa nas bo zanimalo predvsem, ali so razlike med skupinama plod njunih različnih ustrojov. Zanimala nas bo torej razlika  $\mu'' - \mu'$ .

Regresijski model zapišimo v obliki:

$$Y_i' = \mu' + \varepsilon_i', \quad Y_i'' = \mu'' + \varepsilon_i'',$$

kjer so  $\varepsilon_i'$  in  $\varepsilon_i''$  nekorelirane z  $\mathbb{E} \varepsilon_i' = \mathbb{E} \varepsilon_i'' = 0$  in  $\text{var}(\varepsilon_i') = \text{var}(\varepsilon_i'') = \sigma^2$ . Če označimo:

$$\mathbf{Y}' = \begin{bmatrix} Y_1' \\ Y_2' \\ \vdots \\ Y_{m'}' \end{bmatrix}, \quad \mathbf{Y}'' = \begin{bmatrix} Y_1'' \\ Y_2'' \\ \vdots \\ Y_{m''}'' \end{bmatrix}, \quad \boldsymbol{\varepsilon}' = \begin{bmatrix} \varepsilon_1' \\ \varepsilon_2' \\ \vdots \\ \varepsilon_{m'}' \end{bmatrix}, \quad \boldsymbol{\varepsilon}'' = \begin{bmatrix} \varepsilon_1'' \\ \varepsilon_2'' \\ \vdots \\ \varepsilon_{m''}'' \end{bmatrix},$$

lahko model zapišemo v obliki (1.2) z:

$$\mathbf{Y} = \begin{bmatrix} \mathbf{Y}' \\ \mathbf{Y}'' \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} \mathbf{1}_{m'} & \mathbf{0}_{m'} \\ \mathbf{0}_{m''} & \mathbf{1}_{m''} \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \mu' \\ \mu'' \end{bmatrix}, \quad \boldsymbol{\varepsilon} = \begin{bmatrix} \boldsymbol{\varepsilon}' \\ \boldsymbol{\varepsilon}'' \end{bmatrix}.$$

Velja še  $n = m' + m''$  in  $k = 2$ . Izračunajmo:

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} m' & 0 \\ 0 & m'' \end{bmatrix}, \quad (\mathbf{X}^T \mathbf{X})^{-1} = \begin{bmatrix} \frac{1}{m'} & 0 \\ 0 & \frac{1}{m''} \end{bmatrix}.$$

Od tod po krajšem računu dobimo cenilko za  $\beta$ :

$$\hat{\beta} = \begin{bmatrix} \frac{1}{m'} \mathbf{1}_{m'}^T \mathbf{Y}' \\ \frac{1}{m''} \mathbf{1}_{m''}^T \mathbf{Y}'' \end{bmatrix}.$$

Po komponentah je to kar:

$$\hat{\mu}' = \bar{Y}', \quad \hat{\mu}'' = \bar{Y}'',$$

kar je spet zelo intuitivno. Izračunajmo zdaj standardno napako za  $\mu'' - \mu' = \mathbf{c}^T \beta$ , kjer je  $\mathbf{c} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$ . Velja:

$$\mathbf{c}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c} = \frac{1}{m'} + \frac{1}{m''}.$$

Lotimo se sedaj reziduala. Velja:

$$\mathbf{X} \hat{\beta} = \begin{bmatrix} \frac{1}{m'} \mathbf{1}_{m'} \mathbf{1}_{m'}^T \mathbf{Y}' \\ \frac{1}{m''} \mathbf{1}_{m''} \mathbf{1}_{m''}^T \mathbf{Y}'' \end{bmatrix},$$

torej je rezidual enak:

$$\hat{\boldsymbol{\varepsilon}} = \begin{bmatrix} (\mathbf{I}_{m'} - \frac{1}{m'} \mathbf{1}_{m'} \mathbf{1}_{m'}^T) \mathbf{Y}' \\ (\mathbf{I}_{m''} - \frac{1}{m''} \mathbf{1}_{m''} \mathbf{1}_{m''}^T) \mathbf{Y}'' \end{bmatrix}$$

kvadrat njegove norme pa je:

$$\|\hat{\boldsymbol{\varepsilon}}\|^2 = \sum_{i=1}^{m'} (Y_i' - \bar{Y}')^2 + \sum_{i=1}^{m''} (Y_i'' - \bar{Y}'')^2$$

in iz formule (2.4) dobimo:

$$\hat{\sigma}^2 = \frac{1}{m' + m'' - 2} \left( \sum_{i=1}^{m'} (Y_i' - \bar{Y}')^2 + \sum_{i=1}^{m''} (Y_i'' - \bar{Y}'')^2 \right),$$

tako da je končno:

$$\text{SE}(\mu'' - \mu') = \sqrt{\frac{1}{m' + m'' - 2} \left( \frac{1}{m'} + \frac{1}{m''} \right) \left( \sum_{i=1}^{m'} (Y_i' - \bar{Y}')^2 + \sum_{i=1}^{m''} (Y_i'' - \bar{Y}'')^2 \right)}.$$

## 5.5 Primerjava več skupin: analiza variance z enojno klasifikacijo

Recimo sedaj, da isto količino večkrat opazimo na več skupinah in želimo ovrednotiti razlike med skupinami. Naj  $Y_{ij}$ , kjer je  $i = 1, 2, \dots, r$  in  $j = 1, 2, \dots, m_i$ , označuje  $j$ -to opažanje na  $i$ -ti skupini. Testiramo ožji model:

$$Y_{ij} = \mu + \eta_{ij}$$

proti širšemu modelu:

$$Y_{ij} = \mu_i + \epsilon_{ij},$$

kjer so tako šumi  $\eta_{ij}$  kot tudi  $\epsilon_{ij}$  neodvisni in porazdeljeni normalno  $N(0, \sigma^2)$ . Označimo:

$$\mathbf{Y}_i = \begin{bmatrix} Y_{i1} \\ Y_{i2} \\ \vdots \\ Y_{im_i} \end{bmatrix}, \quad \mathbf{Y} = \begin{bmatrix} \mathbf{Y}_1 \\ \mathbf{Y}_2 \\ \vdots \\ \mathbf{Y}_r \end{bmatrix}$$

Označimo spet z  $\mathbf{1}_s$  stolpec iz  $s$  enic in še  $m := m_1 + m_2 + \cdots + m_r$ . To je tudi dimenzija opažanja, torej  $n = m$ . Ožji model je isti kot v podrazdelku 5.1, torej ustreza (4.3) z:

$$\mathbf{Z} = \mathbf{1}_m, \quad \boldsymbol{\gamma} = \mu, \quad l = 1.$$

Če označimo skupno povprečje:

$$\bar{Y} := \frac{1}{m} \mathbf{1}_m^T \mathbf{Y} = \frac{1}{m} \sum_{i=1}^r \sum_{j=1}^{m_i} Y_{ij},$$

neposredno iz podrazdelka 5.1 dobimo  $\hat{\mu} = \bar{Y}$  in:

$$\mathbf{P}_{W \rightarrow U} = \bar{Y} \mathbf{1}_m.$$

Širši model pa ustreza (4.4) z:

$$\mathbf{X} = \begin{bmatrix} \mathbf{1}_{m_1} & \mathbf{0}_{m_1} & \cdots & \mathbf{0}_{m_1} \\ \mathbf{0}_{m_2} & \mathbf{1}_{m_2} & \cdots & \mathbf{0}_{m_2} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0}_{m_r} & \mathbf{0}_{m_r} & \cdots & \mathbf{1}_{m_r} \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_r \end{bmatrix}, \quad k = r.$$

Izračunajmo:

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} m_1 & 0 & \cdots & 0 \\ 0 & m_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & m_r \end{bmatrix}, \quad (\mathbf{X}^T \mathbf{X})^{-1} = \begin{bmatrix} \frac{1}{m_1} & 0 & \cdots & 0 \\ 0 & \frac{1}{m_2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{1}{m_r} \end{bmatrix}.$$

Če označimo še povprečja posameznih skupin:

$$\bar{Y}_i = \frac{1}{m_i} \mathbf{1}_{m_i}^T \mathbf{Y}_i = \frac{1}{m_i} \sum_{j=1}^{m_i} Y_{ij},$$

velja:

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \begin{bmatrix} \bar{Y}_1 \\ \bar{Y}_2 \\ \vdots \\ \bar{Y}_r \end{bmatrix}, \quad \mathbf{P}_{W \rightarrow V} \mathbf{Y} = \begin{bmatrix} \bar{Y}_1 \mathbf{1}_{m_1} \\ \bar{Y}_2 \mathbf{1}_{m_2} \\ \vdots \\ \bar{Y}_r \mathbf{1}_{m_r} \end{bmatrix}.$$

Pojasneni rezidual je torej v skladu s formulo (4.6) enak:

$$\mathbf{Q}_{V \rightarrow U} \mathbf{Y} = \mathbf{P}_{W \rightarrow V} \mathbf{Y} - \mathbf{P}_{W \rightarrow U} \mathbf{Y} = \begin{bmatrix} (\bar{Y}_1 - \bar{Y}) \mathbf{1}_{m_1} \\ (\bar{Y}_2 - \bar{Y}) \mathbf{1}_{m_2} \\ \vdots \\ (\bar{Y}_r - \bar{Y}) \mathbf{1}_{m_r} \end{bmatrix}$$

ocena variance na njegovi podlagi pa je:

$$\hat{\sigma}_{W \rightarrow V}^2 = \frac{1}{r-1} \sum_{i=1}^r m_i (\bar{Y}_i - \bar{Y})^2.$$

Nepojasneni rezidual pa je v skladu s formulo (4.5) enak:

$$\mathbf{Q}_{W \rightarrow V} \mathbf{Y} = \mathbf{Y} - \mathbf{P}_{W \rightarrow V} \mathbf{Y} = \begin{bmatrix} \mathbf{Y}_1 - \bar{Y}_1 \mathbf{1}_{m_1} \\ \mathbf{Y}_2 - \bar{Y}_2 \mathbf{1}_{m_2} \\ \vdots \\ \mathbf{Y}_r - \bar{Y}_r \mathbf{1}_{m_r} \end{bmatrix}$$

in ocena variance na njegovi podlagi je:

$$\hat{\sigma}_{W \rightarrow V}^2 = \frac{1}{m-r} \sum_{i=1}^r \sum_{j=1}^{m_i} (Y_{ij} - \bar{Y}_i)^2.$$

Ožji model testiramo proti širšemu z  $F$ -testom na testni statistiki  $F = \hat{\sigma}_{V \rightarrow U}^2 / \hat{\sigma}_{W \rightarrow V}^2$  z  $(m-k, k-1)$  prostostnimi stopnjami.

Za konec si to oglejmo še v kontekstu varianc. Pojasnjena varianca:

$$\frac{1}{m-1} \|\mathbf{Q}_{V \rightarrow U} \mathbf{Y}\|^2 = \sum_{i=1}^r \frac{m_i}{m-1} (\bar{Y}_i - \bar{Y})^2$$

je *popravljen* varianca med skupinami, medtem ko je nepojasnjena varianca:

$$\frac{1}{m-1} \|\mathbf{Q}_{W \rightarrow V} \mathbf{Y}\|^2 = \frac{1}{m-1} \sum_{i=1}^r \sum_{j=1}^{m_i} (Y_{ij} - \bar{Y}_i)^2$$

*popravljen* varianca znotraj skupin.

## Literatura

- [1] M. J. Schervish: *Theory of Statistics*. Druga izdaja. Springer-Verlag, New York, 1997.