

## Peatükk 4

### $F$ -test

#### 4.1 Mudelite võrdlemine $F$ -testi abil

Olgu meil vaatluse all kaks mudelit — lihtsam mudel  $\mathcal{M}_0$  mudeli maatriksiga  $\mathbf{X}_0$  ja keerukam mudel  $\mathcal{M}_1$  mudeli maatriksiga  $\mathbf{X}_1$ . Eeldame, et lihtsam mudel on erijuht keerukamast — keerukamast mudelist on võimalik lihtsamat mudelit saada fikseerides osade parameetrite väärtused (näiteks võrdsustades mõne parameetri nulliga). Seega on mudel  $\mathcal{M}_0 : y = c_0 + \varepsilon$  erijuht mudelist  $\mathcal{M}_1 : y = c_0 + c_1x + \varepsilon$ , sest keerukamast mudelist saaksime lihtsama, kui võtaksime  $c_1 = 0$ . Matemaatiliselt korrektsemalt kirja pandult: üks mudel on teise erijuht, kui  $\mathcal{C}(\mathbf{X}_0) \subset \mathcal{C}(\mathbf{X}_1)$ . Meie sooviks on testida, kas lihtsama mudeli kasutamine on õigustatud (nullhüpotees: lihtsam mudel on õige) või sunnivad andmed meid kasutama keerukamat mudelit:

$$\begin{aligned}H_0 : \mathbf{y} &\sim N(\mathbf{X}_0\boldsymbol{\beta}_0, \sigma^2\mathbf{I}) \\H_1 : \mathbf{y} &\sim N(\mathbf{X}_1\boldsymbol{\beta}_1, \sigma^2\mathbf{I})\end{aligned}$$

Paneme tähele: kui lihtsam mudel  $\mathcal{M}_0$  on õige, siis on õige ka keerukam mudel  $\mathcal{M}_1$ . Selgitus: Kui  $\mathcal{C}(\mathbf{X}_0) \subset \mathcal{C}(\mathbf{X}_1)$  siis on võimalik leida maatriks  $\mathbf{M}$  selliselt, et  $\mathbf{X}_0 = \mathbf{X}_1\mathbf{M}$ . Kui lihtsam mudel on korrektne, siis valiku  $\boldsymbol{\beta}_1 = \mathbf{M}\boldsymbol{\beta}_0$  korral on ka keerukam mudel õige:  $\mathbf{X}_1\boldsymbol{\beta}_1 = \mathbf{X}_1\mathbf{M}\boldsymbol{\beta}_0 = \mathbf{X}_0\boldsymbol{\beta}_0$ .

Antud peatükis eeldame, et keerukam mudel on igal juhul õige. Meie soovime vaid teada, kas ka lihtsam mudel on õige.

Kuna on üks mudel erijuht teisest? Vaatame paari näidet. Kahes esimeses näites on lihtsam mudel erijuht keerukamast. Viimases näites on tegemist aga mudelitega, milles üks pole vaadeldav teise mudeli lihtsustusena (seega ei saa neid mudeleid võrrelda siin peatükis toodud lähenemise abil).

**Näide 4.1** Kasutada prognoosimisel maakonda või valda?

Oletame, et oleme mõõtnud inimeste keskmist palka kolmes maakonnas — Harjumaal, Tartumaal ja Hiiumaal. Igas maakonnas oleme reaalselt küsitlenud inimesi kahes omavalitsuses (Harjumaal: Viimsi ja Saku vallas; Tartumaal Tartu linnas ja Konguta vallas; Hiiumaal Kärdla linnas ja Emmaste vallas). Kogutav andmestik võiks välja näha selline:

Palk	Maakond	Vald
$Y_1$	Harjumaa	Viimsi
$Y_2$	Harjumaa	Viimsi
$Y_3$	Harjumaa	Saku
$Y_4$	Harjumaa	Saku
$Y_5$	Tartumaa	Tartu
$Y_6$	Tartumaa	Tartu
$Y_7$	Tartumaa	Konguta
$Y_8$	Tartumaa	Konguta
$Y_9$	Hiiumaa	Kärdla
$Y_{10}$	Hiiumaa	Kärdla
$Y_{11}$	Hiiumaa	Emmaste
$Y_{12}$	Hiiumaa	Emmaste

Kas inimese palga prognoosimiseks piisab faktortunnuse „maakond“ kasutamisest või saame parema mudeli tunnust „vald“ kasutades? Prognoosimiseks maakonda kasutava mudeli mudelimaatriks  $\mathbf{X}_0$  ja valda kasutava mudeli (keerukama mudeli) mudelimaatriks  $\mathbf{X}_1$  näevad välja nii:

$$X_0 = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \end{bmatrix} \quad X_1 = \begin{bmatrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

Kuna  $\mathcal{C}(\mathbf{X}_0) \subset \mathcal{C}(\mathbf{X}_1)$  siis on üks vaadeldav mudel erijuht teisest.

**Näide 4.2** *Kaks mõõtmist — kas kasutada mõlemat või mõõtmiste keskmist?*

*Soovime prognoosida laste kaalu. Meil on  $i$ . lapse kaalu prognoosimiseks kasutada kaks tunnust:*

- lastearsti poolt mõõdetud lapse pikkus ( $x_{1i}$ )
- lapsevanema poolt väidetav lapse pikkus ( $x_{2i}$ )

*Valime kahe mudeli vahel. Keerukam mudel kasutab prognoosimiseks nii lapsevanema kui ka lastearsti mõõdetud pikkust:*

$$y_i = c_0 + c_1x_{1i} + c_2x_{2i} + \varepsilon_i.$$

*Lihtsam mudel kasutab aga kaalu prognoosimiseks kahe pikkusemõõtmise keskmist:*

$$y_i = c_0^* + c_1^*(x_{1i} + x_{2i})/2 + \varepsilon_i.$$

*Toodud mudelite puhul on lihtsam mudel keerukam erijuht, sest valides  $c_0 = c_0^*$ ,  $c_1 = c_1^*/2$  ja  $c_2 = c_1^*/2$  saame keerukamast mudelist lihtsama mudeli.*

**Näide 4.3** *Kaks mõõtmist — kas kasutada lastearsti või lapsevanema poolt mõõdetud lapse pikkust? Vaatame eelmises näites kirjeldatud olukorda — prognoosime lapse kaalu pikkuse järgi. Kas peaksime kasutama lastearsti või lapsevanema poolt mõõdetud lapse pikkust? Ehk kas peaksime eelistama mudelit*

$$y_i = c_0 + c_1x_{1i} + \varepsilon_i$$

*või mudelit*

$$y_i = c_0 + c_2x_{2i} + \varepsilon_i?$$

*Toodud mudelist ei saa kumbagi vaadelda kui teise mudeli erijuhtu (välja arvatud juhul, kui lapsevanema ja lastearsti mõõtmistulemused langevad täpselt kokku). Seega ei saa nende kahe mudeli seast sobivaimat leida käesoleva peatükis käsitlemist leidavate meetodite abil.*

Kuidas otsustada, kas kasutada lihtsamat või keerukamat mudelit? Üldine printsiip (Ockhami habemenoa printsiip) on järgmine — kui pole (piisavalt) tõendusmaterjali, mis räägiks keerukama mudeli kasuks, siis eelistame lihtsamat mudelit. Kuidas siis mudelite korral otsustada, kas andmed „sunivad“ meid keerukamat mudelit kasutama või mitte?

Vaatame mõlema mudeli prognoose vaatlusvektorile,  $\hat{y}_0 = \mathbf{X}_0\hat{\beta}_0$  ja  $\hat{y}_1 = \mathbf{X}_1\hat{\beta}_1$ . Kui lihtsam mudel on ka õige, siis peaks prognooside vahe  $\hat{y}_1 - \hat{y}_0$

tulema väike (nullilähedane), ehk prognooside erinevuste ruutude summa  $(\hat{\mathbf{y}}_1 - \hat{\mathbf{y}}_0)^T(\hat{\mathbf{y}}_1 - \hat{\mathbf{y}}_0)$  peaks tulema pisike. Seda ruutude summat saab kirja panna ka veidi teisel moel:

$$\begin{aligned} (\hat{\mathbf{y}}_1 - \hat{\mathbf{y}}_0)^T(\hat{\mathbf{y}}_1 - \hat{\mathbf{y}}_0) &= ((\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\mathbf{y})^T(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\mathbf{y} \\ &= \mathbf{y}^T(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\mathbf{y} \\ &= \mathbf{y}^T(\mathbf{P}_{\mathbf{X}_1}\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_1}\mathbf{P}_{\mathbf{X}_0} - \mathbf{P}_{\mathbf{X}_0}\mathbf{P}_{\mathbf{X}_1} + \mathbf{P}_{\mathbf{X}_0}\mathbf{P}_{\mathbf{X}_0})\mathbf{y}. \end{aligned}$$

Kuna  $\mathcal{C}(\mathbf{X}_0) \subset \mathcal{C}(\mathbf{X}_1)$  siis leidub maatriks  $\mathbf{M}$  nii, et  $\mathbf{X}_0 = \mathbf{X}_1\mathbf{M}$ . Seega

$$\begin{aligned} \mathbf{P}_{\mathbf{X}_1}\mathbf{P}_{\mathbf{X}_0} &= \mathbf{P}_{\mathbf{X}_1}\mathbf{X}_1\mathbf{M}(\mathbf{X}_0^T\mathbf{X}_0)^{-1}\mathbf{X}_0^T \\ &= \mathbf{X}_1\mathbf{M}(\mathbf{X}_0^T\mathbf{X}_0)^{-1}\mathbf{X}_0^T \\ &= \mathbf{X}_0(\mathbf{X}_0^T\mathbf{X}_0)^{-1}\mathbf{X}_0^T \\ &= \mathbf{P}_{\mathbf{X}_0} \end{aligned} \tag{4.1}$$

ja järelikult

$$(\hat{\mathbf{y}}_1 - \hat{\mathbf{y}}_0)^T(\hat{\mathbf{y}}_1 - \hat{\mathbf{y}}_0) = \mathbf{y}^T(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\mathbf{y}.$$

Kui suurt prognooside erinevust võiksime näha, kui mõlemad mudelid on õiged? Sellele küsimusele aitab vastata meile juba tuttav valem (vaata lemmat 3.1):  $E\mathbf{y}^T\mathbf{A}\mathbf{y} = \text{tr}(\mathbf{A}\mathbf{V}) + \boldsymbol{\mu}^T\mathbf{A}\boldsymbol{\mu}$ . Seda valemit ja tähistusi  $p_0 := \text{rank}(\mathbf{X}_0)$  ja  $p_1 := \text{rank}(\mathbf{X}_1)$  kasutades võime kirjutada:

$$\begin{aligned} E[\mathbf{y}^T(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\mathbf{y}] &= \text{tr}[(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\sigma^2] + \boldsymbol{\beta}_0^T\mathbf{X}_0^T(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\mathbf{X}_0\boldsymbol{\beta} \\ &= \sigma^2(p_1 - p_0) + 0, \end{aligned}$$

Sest  $\mathbf{X}_0^T\mathbf{P}_{\mathbf{X}_1} = \mathbf{X}_0^T$  (meenuta:  $\mathbf{X}_0^T = \mathbf{M}^T\mathbf{X}_1^T$ ) ja  $\mathbf{X}_0^T\mathbf{P}_{\mathbf{X}_0} = \mathbf{X}_0^T$ . Teisisõnu, nullhüpoteesi kehtides:

$$E\left[\frac{\mathbf{y}^T(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\mathbf{y}}{p_1 - p_0}\right] = \sigma^2.$$

Meenutame, et oli teinegi nihketa hinnang jääkide dispersioonile, MSE:

$$E\left[\frac{\mathbf{y}^T(\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})\mathbf{y}}{n - p_1}\right] = \sigma^2.$$

Viimane hinnang on õige ka siis, kui lihtsam mudel ei kehti (keerukam mudel on aga siiski õige). Seega võiks nende kahe hinnangu suhe nullhüpoteesi kehtides (kui mõlemad mudelid on õiged) olla ligikaudu 1:

$$\frac{\mathbf{y}^T(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\mathbf{y}}{p_1 - p_0} / \frac{\mathbf{y}^T(\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})\mathbf{y}}{n - p_1} \stackrel{H_0}{\approx} 1.$$

Aga kui kaugelt ühest võib see suhe puhtalt juhuse tõttu kalduda?

Meenutame üht-teist jaotuste kohta.

Kui  $\mathbf{y} \sim N(\boldsymbol{\mu}; \mathbf{I})$  ja  $\mathbf{A}$  on sümmeetriline ja idempotentne maatriks ( $\mathbf{A}\mathbf{A} = \mathbf{A}$ ), siis  $\mathbf{y}^T \mathbf{A} \mathbf{y} \sim \chi^2_{df=\text{rank}(\mathbf{A}); \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu}}$ . Seega, nullhüpoteesi kehtides ( $\mathbf{X}\boldsymbol{\beta} = \mathbf{X}_0\boldsymbol{\beta}_0$ ):

$$\begin{aligned} (\mathbf{y}^T/\sigma)(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})(\mathbf{y}/\sigma) &= \mathbf{y}^T(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\mathbf{y}/\sigma^2 \\ &\sim \chi^2_{p_1-p_0; \boldsymbol{\beta}_0^T \mathbf{X}_0^T (\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0}) \mathbf{X}_0 \boldsymbol{\beta}_0 / \sigma^2} \\ &\sim \chi^2_{p_1-p_0; 0} \end{aligned}$$

ja

$$\mathbf{y}^T(\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})\mathbf{y}/\sigma^2 \sim \chi^2_{n-p_1}.$$

Nüüd veel meenutust: kui  $X \sim \chi_k^2$  ja  $Y \sim \chi_l^2$  on sõltumatud juhuslikud suurused, siis

$$\frac{X/k}{Y/l} \sim F_{k,l}.$$

Järelikult, kui lugeja ja nimetaja oleksid sõltumatud (hetke pärast kontrollime), oleks:

$$\frac{\mathbf{y}^T(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\mathbf{y} / [\sigma^2(p_1 - p_0)]}{\mathbf{y}^T(\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})\mathbf{y} / [\sigma^2(n - p_1)]} \stackrel{H_0}{\sim} F_{p_1-p_0; n-p_1}.$$

Kas lugeja ja nimetaja on sõltumatud? Selle näitamiseks piisab, kui saaksime näidata  $(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\mathbf{y}$  ja  $(\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})\mathbf{y}$  sõltumatust (sest  $f(X) \perp g(Y)$  kui  $X \perp Y$ ). Kuna  $((\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0}) : (\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})^T) \mathbf{y}$  on mitmemõõtmelise normaaljaotusega (kui  $\mathbf{y}$ -on normaaljaotusega) siis piisab sõltumatuse näitamiseks vaid sellest, kui suudame näidata võrdust  $\text{cov}[(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\mathbf{y}, (\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})\mathbf{y}] = 0$ :

$$\begin{aligned} \text{cov}[(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\mathbf{y}, (\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})\mathbf{y}] &= (\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\text{cov}(\mathbf{y}, \mathbf{y})(\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})^T \\ &= \sigma^2(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})(\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})^T \\ &= \sigma^2(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0} - \mathbf{P}_{\mathbf{X}_1} + \mathbf{P}_{\mathbf{X}_0}\mathbf{P}_{\mathbf{X}_1}) \\ &= 0. \end{aligned}$$

Seega võime järeldada, et

$$\frac{\mathbf{y}^T(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\mathbf{y}/(p_1 - p_0)}{\mathbf{y}^T(\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})\mathbf{y}/(n - p_1)} \stackrel{H_0}{\sim} F_{p_1-p_0; n-p_1}.$$

Muuseas, mõnes kontekstis võib olla kasu teadmisest, et

$$\begin{aligned} \mathbf{y}^T(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\mathbf{y} &= \mathbf{y}^T(\mathbf{I} - \mathbf{P}_{\mathbf{X}_0} - (\mathbf{I} - \mathbf{P}_{\mathbf{X}_1}))\mathbf{y} \\ &= \mathbf{y}^T(\mathbf{I} - \mathbf{P}_{\mathbf{X}_0})\mathbf{y} - \mathbf{y}^T(\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})\mathbf{y} \\ &= SSE_0 - SSE_1 \end{aligned}$$

ehk prognooside erinevuste ruutude summa on sama, mis lihtsama mudeli jääkide ruutude summa ( $SSE_0$ ) miinus keerukama mudeli jääkide ruutude summa ( $SSE_1$ ).

## Ülesanne

Soovitakse võrrelda kahte mudelit:

Mudel 1:  $Y_{ij} = \mu + \varepsilon_{ij}$ ,  $i = 1 \dots 4$

Mudel 2:  $Y_{ij} = \mu + \alpha_i + \varepsilon_{ij}$ ,  $i = 1 \dots 4$

Kaks erinevat teadlast on omaenda juhuslike valimite põhjal (mõlemad valimid on sõltumatud, kuid võetud samast populatsioonist) kirjeldanud ühte neist mudelitest. Esimene teadlane on raporteerinud esimese mudeli jääkide ruutude summa  $SSE_1$ , oma vaatluste arvu  $n_1$  ja andnud ka parameetri  $\mu$  hinangu. Lisaks teame ka tema hinnatud parameetrite arvu,  $p_1 = 1$ . Teine teadlane on aga oma valimi põhjal hinnanud keerukama mudeli ja raporteerinud keerukama mudeli (mudel 2) jääkide ruutude summa  $SSE_2$ , keerukama mudeli parameetrite hinnangud (st teame ka parameetrite arvu,  $p_2$ ) ja vaatluste arvud igas grupis.

Kas me võime kasutada mudelite võrdlemiseks F-testi kujul:

$$F = \frac{SSE_1 - SSE_2}{p_2 - p_1} / \frac{SSE_2}{n_2 - p_2}?$$

Põhjenda oma otsust!

Kui arvad, et antud teststatistikut võib kasutada, siis millise jaotusega on teststatistik nullhüpoteesi kehtides?

## 4.2 F-testi geomeetrilisest interpretatsioonist

Vaatlusvektor  $\mathbf{y}$  „elab“ vektorruumis  $\mathbb{R}^n$ . Jagame selle vektorruumi kolmeks (ortogonaalseks) osaks:

$$\mathbb{R}^n = \mathcal{C}(\mathbf{X}_1)^\perp \oplus \mathcal{C}(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0}) \oplus \mathcal{C}(\mathbf{X}_0),$$

kus  $\oplus$  tähistab vektorruumide otsesummat.

Pane tähele:  $\mathcal{C}(\mathbf{X}_1) = \mathcal{C}(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0}) \oplus \mathcal{C}(\mathbf{X}_0)$  ehk vektorruum  $\mathcal{C}(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})$  kirjeldab keerukama mudeli poolt lisatavat osa (või täiendavat vabadust) kus prognoosivektoril  $\hat{\mathbf{y}}$  lubatakse viibida.

F-test kontrollib, kas vaatlusvektori vari (matemaatilises keeles projektsioon) ruumis  $\mathcal{C}(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})$  on sama hüplev kui ruumis  $\mathcal{C}(\mathbf{X}_1)^\perp$ , kus kogu varieeruvus on meile teadaolevalt põhjustatud vaid juhusest. Kui kahtlusaluses osas vaatlusvektori varieeruvus on samasugune nagu taustavarieeruvus, siis pole meil otseselt vajadust oma lihtsama mudeli mudelimaatriksi poolt genereeritavat vektorruumi  $\mathcal{C}(\mathbf{X}_0)$  täiendada.

Märkus: vektori  $\mathbf{y}$  ortogonaalse projektsiooni vektorruumidesse  $\mathcal{C}(\mathbf{X}_1)^\perp$ ,  $\mathcal{C}(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})$  ja  $\mathcal{C}(\mathbf{X}_0)$  saame leida vastavalt projektsioonimaatriksite  $\mathbf{I} - \mathbf{P}_{\mathbf{X}_1}$ ,  $\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0}$  ja  $\mathbf{P}_{\mathbf{X}_0}$  abil, kusjuures on lihtne näha, et mistahes pikusega  $n$  vektori korral  $\mathbf{y} = (\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})\mathbf{y} + (\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\mathbf{y} + \mathbf{P}_{\mathbf{X}_0}\mathbf{y}$ . Lisaks on valemite 4.1 kasutades lihtne järeldada, et kõik saadud projektsioonid ( $(\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})\mathbf{y}$ ;  $(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\mathbf{y}$  ja  $\mathbf{P}_{\mathbf{X}_0}\mathbf{y}$ ) on teineteisega risti ehk ortogonaalsed.

Näide.

Tehakse seitse mõõtmist. Kolme esimese mõõtmise ajal oleme kindlad, et signaali pole — kõik mida näeme, on müra:  $Y_i = \varepsilon_i, i = 1..3$ . Kaks viimast mõõtmist sisaldavad kindlalt signaali,  $Y_i = \mu_i + \varepsilon_i, i = 6; 7$ . Aga vahepealsed mõõtmised — nende puhul pole me kindlad. Võib olla, et vaatluste 4 ja 5 korral kehtib lihtsam mudel,  $Y_i = \varepsilon_i$ , aga võib olla algas signaali edastamine varem, ja õigem oleks keerukam variant —  $Y_i = \mu_i + \varepsilon_i$ . Ehk teisisõnu — me ei tea, kas sobivaks mudeliks on

$$\mathbf{X}_0 = \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{pmatrix} \quad \text{või} \quad \mathbf{X}_1 = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

Antud juhul on lihtsam mudel erijuht keerukamast.

Vaatame, millise varju heidab vaatlusvektor hüpoteeside testimise seisukohast olulistesse alamruumidesse.

$$\begin{array}{ccccccc}
 \begin{pmatrix} 2 \\ -10 \\ 5 \\ -4 \\ 8 \\ 89 \\ 60 \end{pmatrix} & = & \begin{pmatrix} 2 \\ -10 \\ 5 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} & + & \begin{pmatrix} 0 \\ 0 \\ 0 \\ -4 \\ 8 \\ 0 \\ 0 \end{pmatrix} & + & \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 89 \\ 60 \end{pmatrix} \\
 \mathbf{y} & & (\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})\mathbf{y} & & (\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\mathbf{y} & & \mathbf{P}_{\mathbf{X}_0}\mathbf{y} \\
 \text{rank} = 6 & & \text{rank} = 3 & & \text{rank} = 2 & & \text{rank} = 2 \\
 \text{vaatlusvektor} & & \text{kindlasti müra} & & \text{ebaselge} & & \text{kindlasti midagi}
 \end{array}$$

Hindame jääkide dispersiooni kahel viisil.

Dispersioonihinnangu  $\mathbf{y}^T(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})\mathbf{y}/\text{rank}(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0}) = ((-4)^2 + 8^2)/2$  puhul jagame jääkide ruutude summat kahega, sest (sisuliselt) kasutame esialgsetest vaatlustest vaid kahte.

Kui hindame jääkide dispersiooni valemiga  $\mathbf{y}^T(\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})\mathbf{y}/\text{rank}(\mathbf{I} - \mathbf{P}_{\mathbf{X}_1}) = (2^2 + (-10)^2 + 5^2)/3$  siis jagame ruutude summa läbi 3-ga, sest sisuliselt kasutame dispersiooni hindamiseks vaid kolme vaatlust — teisi ruumisruumi me ignoreerime.

Kahe mudeli võrdlemiseks kasutatava *F*-statistiku väärtuseks saame  $F = \frac{40}{43} \approx 0,93$  mis annab *F*-testi olulisustõenäosuseks 0,48... Seega jääb nullhüpotees kummutamata — on võimalik, et vaatluste 4 ja 5 ajal signaali edastamist ei toimunud.

## Ülesanne

Hindame regressioonanalüüsi mudelit,  $y_i = c_0 + c_1x_i + \varepsilon_i$ . Sõltumatu tunnuse  $x_i$  väärtused on 1; 2; 3; 4 ja sõltuva tunnuse  $y_i$  väärtusteks on 1; 3; 2; 4. Soovime testida sirge tõusu olulisust,  $H_0 : c_1 = 0$  (ehk nullhüpoteesi arvates kõlbab kasutada ka mudelit  $y_i = c_0 + \varepsilon_i$ ). Pane kirja kuidas näeb välja vaatlusvektori projektsioon ruumides  $\mathcal{C}(\mathbf{X}_1)^\perp$ ;  $\mathcal{C}(\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0})$  ja  $\mathcal{C}(\mathbf{X}_0)$ . Kontrolli, kas nende projektsioonide summa annab tagasi esialgse vaatlusvektori!