

Lösningar till tentamensskrivning för kursen
Linjära statistiska modeller

25 november 2020 9–16

Examinator: Ola Hössjer, tel. 070/672 12 18, ola@math.su.se

Uppgift 1

a) Låt \bar{Y}_j vara medelvärdet av Y_i för de personer som fått j mg av vaccinet. Vi noterar att

$$\begin{aligned}\bar{x} &= \sum_j \sum_{i;x_i=j} x_i / 20 = \sum_{j=0}^4 4j / 20 = \sum_{j=0}^4 j / 5 = 2, \\ \sum_{i=1}^{20} (x_i - \bar{x})^2 &= 4 \sum_{j=0}^4 (j - 2)^2 = 4 \cdot 10, \\ \sum_{i=1}^{20} (x_i - \bar{x}) Y_i &= 4 \sum_{j=0}^4 (j - 2) \bar{Y}_j \\ &= 4[(-2)2 + (-1)3.4 + 0 \cdot 4.0 + 1 \cdot 5.2 + 2 \cdot 6.1] \\ &= 4 \cdot 10.0.\end{aligned}$$

Det ger en minsta kvadrat-skattning

$$\hat{\beta} = \frac{\sum_i (x_i - \bar{x}) Y_i}{\sum_i (x_i - \bar{x})^2} = \frac{\sum_{j=0}^4 (j - 2) \bar{Y}_j}{\sum_{j=0}^4 (j - 2)^2} = \frac{10.0}{10} = 1.0.$$

b) Variansen för skattningen av effektparametern β ges av

$$\text{Var}(\hat{\beta}) = \frac{\sigma^2}{\sum_{i=1}^{20} (x_i - \bar{x})^2} = \frac{\sigma^2}{4 \cdot 10} = \frac{\sigma^2}{40}.$$

c) Låt s_j^2 vara stickprovsvariansen för den grupp av individer som får dosen j mg av vaccinet. Summan av kvadratavvikelsena $(Y_i - \bar{Y}_j)^2$ för individerna i i grupp j , ges av $(4 - 1)s_j^2 = 3s_j^2$. Genom att summera detta uttryck över alla fem grupper får vi den totala kvadratsumman

$$\text{Kvs(Inom dos)} = 3 \sum_{j=0}^4 s_j^2,$$

för variationskällan Inom dos, med $5(4 - 1) = 15$ frihetsgrader. Det ger en skattning

$$\begin{aligned}\hat{\sigma}^2 &= \text{Mkvs}(\text{Inom dos}) \\ &= \text{Kvs}(\text{Inom dos})/15 \\ &= \sum_{j=0}^4 3s_j^2/15 \\ &= \sum_{j=0}^4 s_j^2/5 \\ &= (1 + 1.8 + 2.4 + 2.0 + 2.8)/5 \\ &= 2.0\end{aligned}$$

av feltermensvariansen. Medelfelet för skattningen av β blir

$$d = \sqrt{\widehat{\text{Var}}(\hat{\beta})} = \sqrt{\frac{\hat{\sigma}^2}{40}} = \sqrt{\frac{2.0}{40}} = 0.2236.$$

d) Det följer av a) och c) ovan att ett tvåsidigt konfidensintervall för β med konfidensgrad 95.0% ges av

$$\begin{aligned}(\hat{\beta} - t_{0.025}(15)d, \hat{\beta} + t_{0.025}(15)d) &= (1.0 - 2.1314 \cdot 0.2236, 1.0 + 2.1314 \cdot 0.2236) \\ &= (0.523, 1.477),\end{aligned}$$

där $t_{0.025}(15) = \sqrt{F_{0.05}(1, 15)}$ fås ur tabell. Eftersom 0 inte ingår i detta intervall så kan vi förkasta nollhypotesen $\beta = 0$ på nivån 5.0%. Vaccinet har alltså signifikanta biverkningar.

Uppgift 2

a) Försöket kan beskrivas med en allmän linjär modell

$$\mathbf{Y} = \mathbf{A}\boldsymbol{\theta} + \boldsymbol{\varepsilon},$$

$\mathbf{Y} = (Y_1, \dots, Y_6)^T = (57.4, 58.4, 58.4, 30.2, 28.7, 27.7)^T$ är observationsvektorn, $\boldsymbol{\theta} = (\theta_1, \theta_2, \theta_3)^T$ parametervektorn,

$$\mathbf{A} = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \\ -1 & 1 & 1 \\ 1 & -1 & 1 \\ 1 & 1 & -1 \end{pmatrix}.$$

designmatrisen samt $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_6)^T$ en vektor av oberoende och $N(0, \sigma^2)$ -fördelade feltermer ε_i . Minsta kvadrat-skattningen av parametervektorn ges av

$$\hat{\boldsymbol{\theta}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y} = \frac{1}{5} \mathbf{A}^T \mathbf{Y} = \begin{pmatrix} 28.4 \\ 29.0 \\ 29.6 \end{pmatrix} = \begin{pmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \\ \hat{\theta}_3 \end{pmatrix},$$

där vi i andra ledet utnyttjade att designmatrisen har ortogonala kolumner med kvadratsumma 5 (dvs $\mathbf{A}^T \mathbf{A} = 5\mathbf{I}_3$). I de två sista leden användes

$$\begin{aligned}\hat{\theta}_1 &= (Y_1 + Y_2 - Y_4 + Y_5 + Y_6)/5 \\ &= (57.4 + 58.4 - 30.7 + 28.7 + 27.7)/5 \\ &= 28.4\end{aligned}$$

för skattningen av längden av Ark 1, och motsvarande för längderna av de andra två arken.

b) Nollhypotesen kan skrivas som $H_0 : \theta_1 = \theta_2 = \theta_3 = \lambda$. Det svarar mot sambandet

$$\boldsymbol{\theta} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \lambda = \mathbf{B}\lambda.$$

c) För att skatta λ så inför vi matrisen $\mathbf{C} = \mathbf{A}\mathbf{B} = (2, 2, 2, 1, 1, 1)^T$ och noterar att under nollhypotesen så gäller $\mathbf{Y} = \mathbf{C}\lambda + \boldsymbol{\varepsilon}$, dvs \mathbf{C} är designmatris under H_0 . Det ger minsta kvadrat-skattningen

$$\hat{\lambda} = (\mathbf{C}^T \mathbf{C})^{-1} \mathbf{C}^T \mathbf{Y} = \frac{1}{15} (2Y_1 + 2Y_2 + 2Y_3 + Y_4 + Y_5 + Y_6) = 29.0.$$

Vi kan nu skatta $\boldsymbol{\mu}$ under grund- och hypotesmodellen som

$$\begin{aligned}\hat{\boldsymbol{\mu}} &= \mathbf{A}\hat{\boldsymbol{\theta}} = (57.4, 58.0, 58.6, 30.2, 29.0, 27.8)^T, \\ \hat{\boldsymbol{\mu}} &= \mathbf{C}\hat{\lambda} = (58.0, 58.0, 58.0, 29.0, 29.0, 29.0)^T.\end{aligned}$$

Kvadratsumman för differensen mellan dessa två vektorer är

$$\|\hat{\boldsymbol{\mu}} - \hat{\boldsymbol{\mu}}\|^2 = (-0.6)^2 + 0^2 + 0.6^2 + 1.2^2 + 0^2 + (-1.2)^2 = 3.60.$$

Eftersom vi har $N = 6$ observationer, grundmodellen har $k = 3$ parametrar och hypotesmodellen $l = 1$ parameter får vi slutligen en

$$\text{F-kvot} = \frac{\|\hat{\boldsymbol{\mu}} - \hat{\boldsymbol{\mu}}\|^2 / (3 - 1)}{\text{Kvs(Residual)} / (6 - 3)} = \frac{3.60/2}{0.3/3} = 18 > F_{0.05}(2, 3) = 9.55.$$

Vi kan alltså på nivån 5% förkasta nollhypotesen att arken har samman längd.

Uppgift 3

a) Vi börjar med att fylla i antal frihetsgrader och medelkvadratsummor i variansanalystabellen, plus en ny kolumn med uttryck för $E(\text{Mkvs})$.

Variationskälla	Kvs	f	Mkvs	$E(\text{Mkvs})$
Tjocklek	18.0	4-1=3	18.0/3=6.0	$\sigma_\epsilon^2 + 2\sigma_\gamma^2 + 2 \cdot 3\sigma_\alpha^2$
Färg	18.4	3-1=2	18.4/2=9.2	$\sigma_\epsilon^2 + 2\sigma_\gamma^2 + 2 \cdot 4\sigma_\beta^2$
Samspel	7.2	(4-1)(3-1)=6	7.2/6 = 1.2	$\sigma_\epsilon^2 + 2\sigma_\gamma^2$
Inom celler	6.0	4 · 3(2 - 1) = 12	6.0/12=0.5	σ_ϵ^2
Totalt	52.8	4 · 3 · 2 - 1 = 23		

Genom att kombinera informationen från de två högra kolumnerna får vi väntevärdesriktiga skattningar

$$\begin{aligned}\hat{\sigma}_\epsilon^2 &= \text{Mkvs}(\text{Inom celler}) = 0.5, \\ \hat{\sigma}_\gamma^2 &= [\text{Mkvs}(\text{Samspel}) - \text{Mkvs}(\text{Inom celler})]/2 = (1.2 - 0.5)/2 = 0.35, \\ \hat{\sigma}_\alpha^2 &= [\text{Mkvs}(\text{Tjocklek}) - \text{Mkvs}(\text{Samspel})]/6 = (6.0 - 1.2)/6 = 0.80, \\ \hat{\sigma}_\beta^2 &= [\text{Mkvs}(\text{Färg}) - \text{Mkvs}(\text{Samspel})]/8 = (9.2 - 1.2)/8 = 1.00\end{aligned}$$

av de fyra varianskomponenterna.

b) Genom att utnyttja den givna definitionen av Y_{ijk} för en tvåvägs variansanalysmodell typ II, och medelvärdesbilda över alla 24 observationer ijk fås uttrycket

$$\hat{\mu} = \bar{Y}_{...} = \mu + \bar{\alpha}_{.} + \bar{\beta}_{.} + \bar{\gamma}_{..} + \bar{\epsilon}_{...} \quad (1)$$

för totalmedelvärdet av alla uppmätta färgnyanser. Eftersom alla stokastiska termer i (1) är oberoende så följer att

$$\begin{aligned}\text{Var}(\hat{\mu}) &= \text{Var}(\bar{\alpha}_{.}) + \text{Var}(\bar{\beta}_{.}) + \text{Var}(\bar{\gamma}_{..}) + \text{Var}(\bar{\epsilon}_{...}) \\ &= \sigma_\alpha^2/4 + \sigma_\beta^2/3 + \sigma_\gamma^2/(4 \cdot 3) + \sigma_\epsilon^2/(4 \cdot 3 \cdot 2).\end{aligned}$$

Av detta följer i sin tur att medelfelet ges av

$$\begin{aligned}d &= \sqrt{\hat{\sigma}_\alpha^2/4 + \hat{\sigma}_\beta^2/3 + \hat{\sigma}_\gamma^2/12 + \hat{\sigma}_\epsilon^2/24} \\ &= \sqrt{0.8/4 + 1.0/3 + 0.35/12 + 0.5/24} \\ &= 0.764.\end{aligned}$$

Uppgift 4

a) Vi börjar med utöka den givna tabellen med en till kolumn, där de fyra cellmedelvärdena $\bar{Y}_{ij.} = (Y_{ij1} + Y_{ij2})/2$ och totalmedelvärdet $\bar{Y}_{...}$ anges.

S	M	Y_{ij1}	Y_{ij2}	$\bar{Y}_{ij.}$
-	-	6.4	6.0	6.2
+	-	8.0	7.6	7.8
-	+	5.6	5.8	5.7
+	+	7.0	7.2	7.1
Medel				6.7

Minsta kvadrat-skattningar av de två huvudeffekterna, och av samspels-effekten, ges av

$$\begin{aligned}\hat{S} &= (-\bar{Y}_{--} + \bar{Y}_{+-} - \bar{Y}_{-+} + \bar{Y}_{++})/4 = 0.75, \\ \hat{M} &= (-\bar{Y}_{--} - \bar{Y}_{+-} + \bar{Y}_{-+} + \bar{Y}_{++})/4 = -0.30, \\ \widehat{SM} &= (+\bar{Y}_{--} - \bar{Y}_{+-} - \bar{Y}_{-+} + \bar{Y}_{++})/4 = -0.05.\end{aligned} \quad (2)$$

Alternativt minsta kvadrat-skattas parametervektorn $\boldsymbol{\theta} = (\mu, \bar{S}, \bar{M}, \widehat{SM})^T$ med hjälp av formeln $\hat{\boldsymbol{\theta}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y}$, där $\mathbf{Y} = (\bar{Y}_{--}, \bar{Y}_{+-}, \bar{Y}_{-+}, \bar{Y}_{++})^T$

är observationsvektorn för den reducerade modellen med cellmedelvärden, och där

$$\mathbf{A} = \begin{pmatrix} 1 & -1 & -1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & 1 & 1 \end{pmatrix}$$

är designmatrisen.

b) Vi har att

$$\begin{aligned} \text{Kvs(Inom celler)} &= \sum_{i,j} [(Y_{ij1} - \bar{Y}_{ij\cdot})^2 + (Y_{ij2} - \bar{Y}_{ij\cdot})^2] \\ &= \sum_{i,j} (Y_{ij1} - Y_{ij2})^2 / 2 \\ &= [(6.4 - 6.0)^2 + (8.0 - 7.6)^2 + (5.6 - 5.8)^2 + (7.0 - 7.2)^2] / 2 \\ &= 0.2. \end{aligned}$$

Eftersom antal frihetsgrader för variationskällan Inom celler är $2 \cdot 2(2 - 1) = 4$, så följer att

$$\hat{\sigma}^2 = \text{Mkvs(Inom celler)} = \frac{\text{Kvs(Inom celler)}}{4} = 0.05$$

är en väntevärdesriktig skattning av feltermsvariansen σ^2 .

c) Låt $\Delta = 2\bar{M}$ vara den parameter vi söker konfidensintervall för. Det följer av (2) att $\hat{\Delta} = 2(-0.30) = -0.60$ och

$$\begin{aligned} \text{Var}(\hat{\Delta}) &= 2^2 [\text{Var}(\bar{Y}_{--}) + \text{Var}(\bar{Y}_{+-}) + \text{Var}(\bar{Y}_{-+}) + \text{Var}(\bar{Y}_{++})] / 4^2 \\ &= 4[4(\sigma^2/2)] / 16 \\ &= \sigma^2/2. \end{aligned}$$

Det ger ett medelfel $d = \sqrt{\hat{\sigma}^2/2} = \sqrt{0.05/2} = 0.1581$. Eftersom antalet frihetsgrader för variationskällan Inom celler är $4(2 - 1) = 4$ får vi ett konfidensintervall

$$\begin{aligned} I_{\Delta} &= (\hat{\Delta} - t_{0.025}(4)d, \hat{\Delta} + t_{0.025}(4)d) \\ &= (-0.6 - 2.7764 \cdot 0.158, -0.6 + 2.7764 \cdot 0.158) \\ &= (-1.039, -0.161) \end{aligned}$$

med täckningsgrad 0.95. Eftersom 0 inte tillhör detta intervall så sänks rehabiliteringstiden signifikant då fler skruvar används för att fästa protesen.

Uppgift 5

a) Hattmatrisen definieras enligt $\mathbf{H} = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$.

b) Eftersom minsta kvadrat-skattningen av parametervektorn ges av $\hat{\boldsymbol{\theta}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y}$ kan vi skriva om residualvektorn som

$$\begin{aligned} \mathbf{e} &= \mathbf{Y} - \mathbf{A}\hat{\boldsymbol{\theta}} \\ &= \mathbf{Y} - \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y} \\ &= \mathbf{Y} - \mathbf{H}\mathbf{Y} \\ &= (\mathbf{I}_N - \mathbf{H})\mathbf{Y}, \\ &= (\mathbf{I}_N - \mathbf{H})(\mathbf{A}\boldsymbol{\theta} + \boldsymbol{\varepsilon}) \\ &= (\mathbf{I}_N - \mathbf{H})\boldsymbol{\varepsilon}, \end{aligned}$$

där vi i fjärde ledet införde enhetsmatrisen \mathbf{I}_N av ordning N och i sjätte ledet utnyttjade att $\mathbf{H}\mathbf{A}\boldsymbol{\theta} = \mathbf{A}\boldsymbol{\theta}$, vilket följer av definitionen av hattmatrisen. Kovariansmatrisen för residualvektorn ges av

$$\begin{aligned} \text{Var}(\mathbf{e}) &= \text{Var}[(\mathbf{I}_N - \mathbf{H})\boldsymbol{\varepsilon}] = (\mathbf{I}_N - \mathbf{H})\text{Var}(\boldsymbol{\varepsilon})(\mathbf{I}_N - \mathbf{H})^T \\ &= \sigma^2(\mathbf{I}_N - \mathbf{H})(\mathbf{I}_N - \mathbf{H})^T = \sigma^2(\mathbf{I}_N - \mathbf{H}), \end{aligned} \quad (3)$$

där vi i tredje ledet utnyttjade att $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}_N$ och i fjärde ledet att hattmatrisen är idempotent ($\mathbf{H}\mathbf{H} = \mathbf{H}$) och symmetrisk ($\mathbf{H}^T = \mathbf{H}$). Notera att diagonalelement i av (3) svarar mot identiteten $\text{Var}(e_i) = \sigma^2(1 - h_{ii})$. Eftersom $E(e_i) = 0$ följer av detta att

$$E[\text{Kvs}(\text{Residual})] = \sum_{i=1}^N E(e_i^2) = \sum_{i=1}^N \text{Var}(e_i) = \sigma^2 \sum_{i=1}^N (1 - h_{ii}) = \sigma^2(N - k),$$

där vi i sista ledet utnyttjade $\sum_{i=1}^N h_{ii} = k$.

c) Prediktionsfelsevektorn kan skrivas som

$$\mathbf{Z} - \hat{\boldsymbol{\mu}} = \mathbf{A}\boldsymbol{\theta} + \boldsymbol{\varepsilon} - \mathbf{H}\mathbf{Y} = \mathbf{A}\boldsymbol{\theta} + \boldsymbol{\varepsilon} - \mathbf{H}(\mathbf{A}\boldsymbol{\theta} + \boldsymbol{\varepsilon}) = \boldsymbol{\varepsilon} - \mathbf{H}\boldsymbol{\varepsilon}.$$

Dess kovariansmatris är

$$\begin{aligned} \text{Var}(\mathbf{Z} - \hat{\boldsymbol{\mu}}) &= \text{Var}(\boldsymbol{\varepsilon} - \mathbf{H}\boldsymbol{\varepsilon}) = \text{Var}(\boldsymbol{\varepsilon}) + \text{Var}(\mathbf{H}\boldsymbol{\varepsilon}) \\ &= \sigma^2 \mathbf{I}_N + \sigma^2 \mathbf{H}\mathbf{H}^T = \sigma^2(\mathbf{I}_N + \mathbf{H}), \end{aligned} \quad (4)$$

där vi andra ledet utnyttjade att feltermvektorerna $\boldsymbol{\varepsilon}$ och $\boldsymbol{\varepsilon}$ är oberoende, i femte ledet att $\text{Var}(\boldsymbol{\varepsilon}) = \text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}_N$ och i sista ledet att $\mathbf{H}\mathbf{H}^T = \mathbf{H}$. Genom att identifiera diagonalelementen i höger- och vänsterleden av (4), ser vi att

$$E[\text{Kvs}(\text{Prediktion})] = \sum_{i=1}^N E[(Z_i - \hat{\mu}_i)^2] = \sigma^2 \sum_{i=1}^N (1 + h_{ii}) = \sigma^2(N + k).$$

d) Den predikterade kvadratsumman är väntevärdesriktig

$$\begin{aligned} E[\widehat{\text{Kvs}}(\text{Prediktion})] &= E[\text{Kvs}(\text{Residual})] + CE(\hat{\sigma}^2) \\ &= \sigma^2(N - k) + C\sigma^2 \\ &= \sigma^2(N + k) \end{aligned}$$

enligt resultatet i c), om $C = 2k$. Anta att vi använder $\widehat{K}_{vs}(\text{Prediktion})$ för att välja mellan olika modeller, genom att minimera detta uttryck. Vi tolkar $C\hat{\sigma}^2 = 2k\hat{\sigma}^2$ som en term som bestraffar modeller med många parametrar k för att undvika överanpassning, eftersom $K_{vs}(\text{Residual})$ är mindre ju större en modell är.