

Lösningar till tentamensskrivning för kursen Linjära statistiska modeller

20 oktober 2022 14–19

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

Uppgift 1

a) Eftersom $\bar{x} = 2016$ följer att $\sum_i (x_i - \bar{x})^2 = 2(1^2 + 2^2 + 3^2 + 4^2) = 60$.
Det ger minsta kvadrat-skattningar

$$\begin{aligned}\hat{\alpha} &= \bar{y} = 26.9, \\ \hat{\beta} &= \sum_i y_i (x_i - \bar{x}) / \sum_i (x_i - \bar{x})^2 = 27.3/60 = 0.455\end{aligned}$$

av intercept och lutning. Det förväntade kilopriset 2022 är

$$\mu(2022) = \tilde{\alpha} + \beta(2022 - \bar{x}) = \tilde{\alpha} + 6\beta,$$

och det skattas med

$$\hat{\mu}(2022) = \hat{\alpha} + 6\hat{\beta} = 26.9 + 6 \cdot 0.455 = 29.63.$$

b) Eftersom vi har $N = 9$ observationer följer att

$$\begin{aligned}\text{Var}(\hat{\mu}(2022)) &= \text{Var}(\hat{\alpha}) + 6^2 \cdot \text{Var}(\hat{\beta}) = \sigma^2 \left(\frac{1}{9} + \frac{6^2}{\sum_{i=1}^9 (x_i - \bar{x})^2} \right) \\ &= \sigma^2 \left(\frac{1}{9} + \frac{36}{60} \right) = 0.711\sigma^2,\end{aligned}$$

där felvariansen skattas med

$$\hat{\sigma}^2 = \frac{\text{Kvs(Residual)}}{N - 2} = \frac{0.9274}{7} = 0.1325.$$

Det ger ett medelfel

$$d = \sqrt{\widehat{\text{Var}}(\hat{\mu}(2022))} = \sqrt{0.711\hat{\sigma}^2} = \sqrt{0.711 \cdot 0.1325} = 0.3069.$$

c) Konfidensintervallet ges av

$$I = \hat{\mu}(2022) \pm t_{0.025}(N - 2)d = 29.63 \pm 2.3646 \cdot 0.3069 = (28.90, 30.36).$$

Uppgift 2

a) Matrisformuleringen av den linjära modellen $\mathbf{Y} = \mathbf{A}\boldsymbol{\theta} + \boldsymbol{\varepsilon}$ ges av

$$\begin{pmatrix} Y_{11} \\ Y_{12} \\ Y_{13} \\ Y_{14} \\ Y_{15} \\ Y_{21} \\ Y_{22} \\ Y_{23} \\ Y_{24} \\ Y_{25} \end{pmatrix} = \begin{pmatrix} 1 & x_1 & 0 & 0 \\ 1 & x_2 & 0 & 0 \\ 1 & x_3 & 0 & 0 \\ 1 & x_4 & 0 & 0 \\ 1 & x_5 & 0 & 0 \\ 0 & 0 & 1 & x_1 \\ 0 & 0 & 1 & x_2 \\ 0 & 0 & 1 & x_3 \\ 0 & 0 & 1 & x_4 \\ 0 & 0 & 1 & x_5 \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \beta_1 \\ \alpha_2 \\ \beta_2 \end{pmatrix} + \begin{pmatrix} \varepsilon_{11} \\ \varepsilon_{12} \\ \varepsilon_{13} \\ \varepsilon_{14} \\ \varepsilon_{15} \\ \varepsilon_{21} \\ \varepsilon_{22} \\ \varepsilon_{23} \\ \varepsilon_{24} \\ \varepsilon_{25} \end{pmatrix}.$$

b) Låt $\boldsymbol{\lambda} = (\lambda_1, \lambda_2)^T$ innehålla det gemensamma interceptet $\lambda_1 = \alpha_1 = \alpha_2$ och den gemensamma effektparametern $\lambda_2 = \beta_1 = \beta_2$ under nollhypotesen H_0 . Denna nollhypotes svarar mot $\boldsymbol{\theta} = \mathbf{B}\boldsymbol{\lambda}$, där

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

c) Eftersom antalet observationer $N = 10$, grundmodellen har $k = 4$ parametrar och hypotesmodellen $l = 2$ parametrar följer att

$$\begin{aligned} \text{Mkvs(Avv. från } H_0) &= \text{Kvs(Avv. från } H_0)/(k - l) = 6.8910/2 = 3.4455, \\ \text{Mkvs(Residual)} &= \text{Kvs(Residual)/(N - k)} = 1.1060/6 = 0.1843. \end{aligned}$$

Det ger en

$$\text{F-kvot} = \frac{\text{Mkvs(Avv. från } H_0)}{\text{Mkvs(Residual)}} = \frac{3.4455}{0.1843} = 18.69$$

som överstiger $F_{0.05}(k - l, N - k) = F_{0.05}(2, 6) = 5.1$. Vi kan alltså förkasta nollhypotesen att medicinerna har samma effekt.

Uppgift 3

a) Antalet frihetsgrader för residualerna $Y_{ij} - \bar{Y}_i$ är $5(3 - 1) = 10$. Det ger en skattning

$$\begin{aligned} \hat{\sigma}_\varepsilon^2 &= \text{Mvs(Inom block)} \\ &= \frac{1}{10} \text{Kvs(Inom block)} \\ &= \frac{1}{10} \sum_{i=1}^5 \sum_{j=1}^3 (Y_{ij} - \bar{Y}_i)^2 \\ &= \frac{1}{10} \sum_{i=1}^5 2s_i^2 \\ &= \frac{1}{5} \sum_{i=1}^5 s_i^2 \\ &= \frac{1}{5} (0.035 + 0.031 + 0.020 + 0.046 + 0.023) \\ &= 0.031 \end{aligned}$$

av σ_ε^2 . Eftersom radmedelvärdena

$$\bar{Y}_i = \mu + \delta_i + \bar{\varepsilon}_i \sim N(\mu, \sigma_\delta^2 + \frac{\sigma_\varepsilon^2}{3})$$

är sinsemellan oberoende, kan vi använda deras stickprovsvarians för att skatta $\sigma_\delta^2 + \sigma_\varepsilon^2/3$:

$$\begin{aligned} \hat{\sigma}_\delta^2 + \frac{1}{3}\hat{\sigma}_\varepsilon^2 &= \frac{1}{5-1} \sum_{i=1}^5 (\bar{Y}_i - \bar{Y}_{..})^2 \\ &= \frac{1}{4} [(3.1 - 2.9)^2 + (2.6 - 2.9)^2 + (2.8 - 2.9)^2 + (3.3 - 2.9)^2 \\ &\quad + (2.7 - 2.9)^2]. \\ &= 0.085. \end{aligned}$$

Det ger i sin tur en skattning

$$\hat{\sigma}_\delta^2 = 0.085 - \frac{0.031}{3} = 0.075$$

av σ_δ^2 .

b) Minsta kvadrat-skattningen av μ är $\hat{\mu} = \bar{Y}_{..} = 2.9$. Eftersom

$$\text{Var}(\hat{\mu}) = \frac{\sigma_\delta^2}{5} + \frac{\sigma_\varepsilon^2}{15}$$

kan vi utnyttja räkningarna i deluppgift b) för att erhålla ett medelfel

$$\begin{aligned} d &= \sqrt{\widehat{\text{Var}}(\hat{\mu})} \\ &= \sqrt{\frac{1}{5}(\hat{\sigma}_\delta^2 + \frac{1}{3}\hat{\sigma}_\varepsilon^2)} \\ &= \sqrt{\frac{0.085}{5}} \\ &= 0.1304, \end{aligned}$$

och konfidensintervall

$$\begin{aligned} I &= \hat{\mu} \pm t_{0.025}(4)d \\ &= 2.9 \pm \sqrt{F_{0.05}(1, 4)} \cdot 0.1304 \\ &= 2.9 \pm \sqrt{7.7} \cdot 0.1304 \\ &= (2.54, 3.26) \end{aligned}$$

för μ med konfidensgrad 95%.

Uppgift 4

a) AR(1)-processen är stationär då $|\phi| < 1$.

b) Sätt $X_t = Y_t - \mu$, så att

$$X_t = \phi X_{t-1} + \varepsilon_t. \tag{1}$$

Eftersom X_{t-1} beror av feltermen $\varepsilon_{t-1}, \varepsilon_{t-2}, \dots$ så kommer X_{t-1} att vara oberoende av ε_t . Av detta och stationariteten hos $\{X_t\}$ följer att

$$\begin{aligned}\sigma^2 &= \text{Var}(Y_t) = \text{Var}(X_t) = \text{Var}(\phi X_{t-1} + \varepsilon_t) \\ &= \phi^2 \text{Var}(X_{t-1}) + \text{Var}(\varepsilon_t) = \phi^2 \sigma^2 + \sigma_\varepsilon^2.\end{aligned}$$

Löser vi ut denna ekvation med avseende på σ^2 så erhålls

$$\sigma^2 = \frac{\sigma_\varepsilon^2}{1 - \phi^2}.$$

Vidare har vi, givet $k \geq 1$, att

$$\begin{aligned}\gamma_k &= \text{Cov}(Y_t, Y_{t+k}) = \text{Cov}(X_t, X_{t+k}) = \text{Cov}(X_t, \phi X_{t+k-1} + \varepsilon_{t+k}) \\ &= \phi \text{Cov}(X_t, X_{t+k-1}) + \text{Cov}(X_t, \varepsilon_{t+k}) = \phi \gamma_{k-1}.\end{aligned}\quad (2)$$

Eftersom $\gamma_0 = \sigma^2$ så medför (2) (med induktionsbevis) att $\gamma_k = \gamma_0 \phi^k$. Det ger

$$\rho_k = \frac{\text{Cov}(Y_t, Y_{t+k})}{\sqrt{\text{Var}(Y_t)}\sqrt{\text{Var}(Y_{t+k})}} = \frac{\text{Cov}(X_t, X_{t+k})}{\sqrt{\text{Var}(X_t)}\sqrt{\text{Var}(X_{t+k})}} = \frac{\gamma_k}{\sqrt{\gamma_0}\sqrt{\gamma_0}} = \frac{\gamma_k}{\gamma_0} = \phi^k.$$

c) Med hjälp av ledningen kan vi skriva (den approximativa) log-likelihooden som

$$l(\boldsymbol{\psi}) = l(\alpha, \phi, \sigma_\varepsilon^2) = C_1 - C_2 \sum_{t=2}^T (y_t - \alpha - \phi y_{t-1})^2, \quad (3)$$

där $C_1 = -(T-1) \log(\sqrt{2\pi}\sigma_\varepsilon)$ och $C_2 = (2\sigma_\varepsilon^2)^{-1}$ är positiva och beror av σ_ε^2 . Oavsett värdet på σ_ε^2 så maximeras $l(\boldsymbol{\psi})$ genom att minimera kvadratsumman i (3). Detta är ekvivalent med att skatta (icke-centrerat) intercept α och lutning ϕ med minsta kvadratmetoden i en enkel linjär regression med data $\{(y_{t-1}, y_t)\}_{t=2}^T$, där $x_t = y_{t-1}$ och y_t svarar mot förklarande variabel och responsvariabel för observation $(x_t, y_t) = (y_{t-1}, y_t)$. Med $\bar{x} = \sum_{t=2}^T x_t / (T-1) = \sum_{t=1}^{T-1} y_t / (T-1)$ och $\bar{y} = \sum_{t=2}^T y_t / (T-1)$ får vi alltså

$$\begin{aligned}\hat{\phi} &= \sum_{t=2}^T (y_{t-1} - \bar{x})y_t / \sum_{t=2}^T (y_{t-1} - \bar{x})^2, \\ \hat{\alpha} &= \bar{y} - \hat{\phi}\bar{x}.\end{aligned}$$

Eftersom $\mu = \alpha / (1 - \phi)$ ger det ML-skattningen

$$\hat{\mu} = \hat{\alpha} / (1 - \hat{\phi}) = (\bar{y} - \hat{\phi}\bar{x}) / (1 - \hat{\phi}) \approx \bar{y}$$

av μ , där vi i sista ledet utnyttjade att $\bar{x} \approx \bar{y}$ om T är stor.

Uppgift 5

a) Låt $\hat{\boldsymbol{\theta}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y}$ vara minsta kvadrat-skattningen av parametervektorn $\boldsymbol{\theta}$. Härur följer att

$$\hat{\boldsymbol{\mu}} = \mathbf{A} \hat{\boldsymbol{\theta}} = \mathbf{A} (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y} = \mathbf{H} \mathbf{Y}.$$

Eftersom detta samband ska gälla för alla \mathbf{Y} så ges hattmatrisen av

$$\mathbf{H} = \mathbf{A} (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T. \quad (4)$$

b) Vi kan skriva residualvektorn som

$$\mathbf{e} = \mathbf{Y} - \hat{\boldsymbol{\mu}} = \mathbf{Y} - \mathbf{H} \mathbf{Y} = (\mathbf{I}_N - \mathbf{H}) \mathbf{Y},$$

där \mathbf{I}_N är identitetsmatrisen av ordning N . Av detta, och det faktum att $\text{Var}(\mathbf{Y}) = \text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}_N$, följer att residualvektorn har kovariansmatrisen

$$\begin{aligned} \text{Var}(\mathbf{e}) &= (\mathbf{I}_N - \mathbf{H}) \text{Var}(\mathbf{Y}) (\mathbf{I}_N - \mathbf{H})^T \\ &= \sigma^2 (\mathbf{I}_N - \mathbf{H}) (\mathbf{I}_N - \mathbf{H})^T \\ &= \sigma^2 (\mathbf{I}_N - \mathbf{H}^T - \mathbf{H} + \mathbf{H} \mathbf{H}^T) \\ &= \sigma^2 (\mathbf{I}_N - \mathbf{H} - \mathbf{H} + \mathbf{H}) \\ &= \sigma^2 (\mathbf{I}_N - \mathbf{H}). \end{aligned} \quad (5)$$

I det fjärde steget av (5) utnyttjades att hattmatrisen är symmetrisk ($\mathbf{H} = \mathbf{H}^T$) och idempotent ($\mathbf{H} = \mathbf{H} \mathbf{H}$), vilket följer ur (4) och räkneregler för matriser. Vidare ges väntevärdet för residualvektorn av

$$E(\mathbf{e}) = E[(\mathbf{I}_N - \mathbf{H}) \mathbf{Y}] = (\mathbf{I}_N - \mathbf{H}) \boldsymbol{\mu} = \boldsymbol{\mu} - \boldsymbol{\mu} = \mathbf{0}, \quad (6)$$

där vi i näst sista steget utnyttjade att $\mathbf{H} \boldsymbol{\mu} = \mathbf{H} \mathbf{A} \boldsymbol{\theta} = \mathbf{A} \boldsymbol{\theta} = \boldsymbol{\mu}$. Eftersom $\text{Var}(e_i)$ fås ur i :te diagonalelementet av kovariansmatrisen i (5), och $E(e_i)$ från komponent i av nollvektorn i (6), följer att

$$E(e_i^2) = E(e_i)^2 + \text{Var}(e_i) = 0^2 + \sigma^2 (\mathbf{I}_N - \mathbf{H})_{ii} = \sigma^2 (1 - h_{ii}). \quad (7)$$

c) Om den förklarande variabeln centreras i den enkla linjära regressionsmodellen fås

$$Y_i = \tilde{\alpha} + \beta(x_i - \bar{x}) + \varepsilon_i, \quad i = 1, \dots, N,$$

där $\bar{x} = \sum_{i=1}^N x_i / N$. Av detta följer att designmatrisen ges av

$$\mathbf{A} = \begin{pmatrix} 1 & x_1 - \bar{x} \\ \vdots & \vdots \\ 1 & x_N - \bar{x} \end{pmatrix} \implies \mathbf{A}^T \mathbf{A} = \begin{pmatrix} N & 0 \\ 0 & \sum_{j=1}^N (x_j - \bar{x})^2 \end{pmatrix}.$$

Sedan använder vi (4) för att räkna ut hattmatrisen enligt

$$\mathbf{H} = \mathbf{A} \begin{pmatrix} N & 0 \\ 0 & \sum_{j=1}^N (x_j - \bar{x})^2 \end{pmatrix}^{-1} \mathbf{A}^T = \mathbf{A} \begin{pmatrix} \frac{1}{N} & 0 \\ 0 & \frac{1}{\sum_{j=1}^N (x_j - \bar{x})^2} \end{pmatrix} \mathbf{A}^T.$$

Speciellt ges hattmatrisens i :te diagonalelement av

$$h_{ii} = \frac{1}{N} + \frac{(x_i - \bar{x})^2}{\sum_{j=1}^N (x_j - \bar{x})^2}. \quad (8)$$

En inflytelserik punkt (x_i, Y_i) drar i större utsträckning till sig den skattade regressionslinjen, så att residualen e_i blir liten. Ju större värdet på h_{ii} är, desto mindre blir variansen för denna residual enligt (7). Det följer alltså av (8) att en observationspunkt kan förväntas vara mer inflytelserik ju längre dess förklarande variabel x_i ligger ifrån centrum \bar{x} i förhållande till övriga observationer. Detta svarar mot den så kallade hävstångseffekten (leverage points).