

Examen

Soluții

27 Iunie 2018

Exercițiul 1

12p

1. Reamintim că estimatorii lui α și β obținuți prin metoda celor mai mici pătrate sunt dați de

$$\hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \text{și} \quad \hat{\alpha} = \bar{y} - \hat{\beta}\bar{x}$$

iar varianțele lor sunt

$$Var(\hat{\beta}) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \text{și} \quad Var(\hat{\alpha}) = \frac{\sigma^2 \sum_{i=1}^n x_i^2}{n \sum_{i=1}^n (x_i - \bar{x})^2}.$$

4p

- a) Dacă β este cunoscut iar α este necunoscut atunci estimatorul $\tilde{\alpha}$ corespunde la valoarea care minimizează funcția

$$S(\alpha) = \sum_{i=1}^n [y_i - (\alpha + \beta x_i)]^2.$$

Rezolvând ecuația $S'(\alpha) = 0$ obținem

$$S'(\alpha) = -2 \sum_{i=1}^n [y_i - (\alpha + \beta x_i)] = 0 \iff \tilde{\alpha} = \hat{y} - \beta \bar{x}.$$

4p

- b) Pentru a calcula varianța lui $\tilde{\alpha}$ să observăm că, folosind relația $y_i = \alpha + \beta x_i$,

$$\tilde{\alpha} = \hat{y} - \beta \bar{x} = \underbrace{\left(\alpha + \beta \bar{x} + \frac{1}{n} \sum_{i=1}^n \varepsilon_i \right)}_{\hat{y}} - \beta \bar{x} = \alpha + \frac{1}{n} \sum_{i=1}^n \varepsilon_i$$

și cum variabilele aleatoare ε_i sunt centrate, necorelate și de varianță σ^2 găsim că

$$Var(\tilde{\alpha}) = \frac{\sigma^2}{n}.$$

Cum $\sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2 \leq \sum_{i=1}^n x_i^2$ deducem că $Var(\tilde{\alpha}) \leq Var(\hat{\alpha})$.

4p

- c) Presupunând acum că α este cunoscut și β este necunoscut avem că estimatorul $\tilde{\beta}$ corespunde la valoarea care minimizează funcția

$$S(\beta) = \sum_{i=1}^n [y_i - (\alpha + \beta x_i)]^2.$$

Rezolvând ecuația $S'(\beta) = 0$ găsim

$$S'(\beta) = -2 \sum_{i=1}^n x_i [y_i - (\alpha + \beta x_i)] = 0 \iff \tilde{\beta} = \frac{\sum_{i=1}^n x_i (y_i - \alpha)}{\sum_{i=1}^n x_i^2}.$$

În mod similar cu punctul anterior putem rescrie $\tilde{\beta}$ prin

$$\tilde{\beta} = \frac{\sum_{i=1}^n x_i (y_i - \alpha)}{\sum_{i=1}^n x_i^2} = \frac{\sum_{i=1}^n x_i \overbrace{(\alpha + \beta x_i + \varepsilon_i - \alpha)}^{y_i}}{\sum_{i=1}^n x_i^2} = \beta + \frac{\sum_{i=1}^n x_i \varepsilon_i}{\sum_{i=1}^n x_i^2},$$

ceea ce conduce la $Var(\tilde{\beta}) = \frac{\sigma^2}{\sum_{i=1}^n x_i^2}$.

Remarcăm că $\sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2 \leq \sum_{i=1}^n x_i^2$ ceea ce arată că $Var(\tilde{\beta}) \leq Var(\hat{\beta})$.

8p

2. Avem modelele

$$\begin{aligned} y_i &= \alpha + \beta x_i + \varepsilon_i \\ y_i &= \alpha' + \beta' (x_i - \bar{x}) + \varepsilon_i \\ &= \alpha' + \beta' z_i + \varepsilon_i \end{aligned}$$

4p

a) Deoarece $\bar{z} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}) = 0$ deducem că

$$\hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n z_i (y_i - \bar{y})}{\sum_{i=1}^n z_i^2} = \hat{\beta}'.$$

De asemenea avem $\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x}$ și respectiv $\hat{\alpha}' = \bar{y} - \hat{\beta}'\bar{z} = \bar{y}$ ceea ce arată că $\hat{\alpha}' \neq \hat{\alpha}$. Mai mult observăm că $\hat{\alpha}' \sim \mathcal{N}\left(\alpha + \beta\bar{z}, \frac{\sigma^2}{n}\right) = \mathcal{N}\left(\alpha, \frac{\sigma^2}{n}\right)$.

4p

b) Pentru a verifica necorelarea dintre $\hat{\alpha}'$ și $\hat{\beta}'$ să observăm că din

$$\hat{\alpha}' = \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i, \quad \hat{\beta}' = \frac{\sum_{i=1}^n z_i (y_i - \bar{y})}{\sum_{i=1}^n z_i^2} = \sum_{j=1}^n \left(\frac{z_j}{\sum_{i=1}^n z_i^2} \right) y_j$$

avem

$$Cov(\hat{\alpha}', \hat{\beta}') = -\sigma^2 \sum_{j=1}^n \frac{1}{n} \left(\frac{z_j}{\sum_{i=1}^n z_i^2} \right) = 0.$$

Altfel se poate consulta exercițiul 1.4 din Seminarul 2.

Exercițiul 2

5p

1. Pentru a determina valorile lipsă vom folosi proprietatea de simetrie a matricei $\mathbf{X}^\top \mathbf{X}$, i.e. $\mathbf{X}^\top \mathbf{X} = (\mathbf{X}^\top \mathbf{X})^\top$:

$$\mathbf{X}^\top \mathbf{X} = \begin{pmatrix} ? & ? & 9791.6 \\ ? & 3306476 & ? \\ ? & 471237.9 & 67660 \end{pmatrix} = \underbrace{\begin{pmatrix} ? & ? & ? \\ ? & 3306476 & 471237.9 \\ 9791.6 & ? & 67660 \end{pmatrix}}_{(\mathbf{X}^\top \mathbf{X})^\top}$$

ceea ce conduce la $\mathbf{X}^\top \mathbf{X} = \begin{pmatrix} ? & ? & 9791.6 \\ ? & 3306476 & 471237.9 \\ 9791.6 & 471237.9 & 67660 \end{pmatrix}$. De asemenea putem observa că matricea de design este $\mathbf{X} = (\mathbf{1}, \mathbf{x}, \mathbf{z})$ unde $\mathbf{z} = \sqrt{\mathbf{x}} = (\sqrt{x_1}, \dots, \sqrt{x_n})$ ceea ce conduce la

$$\mathbf{X}^\top \mathbf{X} = \begin{pmatrix} \mathbf{1}^\top \\ \mathbf{x}^\top \\ \mathbf{z}^\top \end{pmatrix} (\mathbf{1} \quad \mathbf{x} \quad \mathbf{z}) = \begin{pmatrix} \mathbf{1}^\top \mathbf{1} & \mathbf{1}^\top \mathbf{x} & \mathbf{1}^\top \mathbf{z} \\ \mathbf{x}^\top \mathbf{1} & \mathbf{x}^\top \mathbf{x} & \mathbf{x}^\top \mathbf{z} \\ \mathbf{z}^\top \mathbf{1} & \mathbf{z}^\top \mathbf{x} & \mathbf{z}^\top \mathbf{z} \end{pmatrix} = \begin{pmatrix} n & n\bar{x} & n\bar{z} \\ n\bar{x} & \mathbf{x}^\top \mathbf{x} & \mathbf{x}^\top \mathbf{z} \\ n\bar{z} & \mathbf{z}^\top \mathbf{x} & \mathbf{z}^\top \mathbf{z} \end{pmatrix}$$

și cum $\mathbf{z}^\top \mathbf{z} = n\bar{x}$, iar $n = 1429$ rezultă că

$$\mathbf{X}^\top \mathbf{X} = \begin{pmatrix} 1429 & 67660 & 9791.6 \\ 67660 & 3306476 & 471237.9 \\ 9791.6 & 471237.9 & 67660 \end{pmatrix}.$$

5p

2. Valoarea diametrului mediu este dată de

$$\bar{x} = \frac{(\mathbf{X}^\top \mathbf{X})_{1,2}}{n} = \frac{67660}{1429} \approx 47.3 \text{ cm}.$$

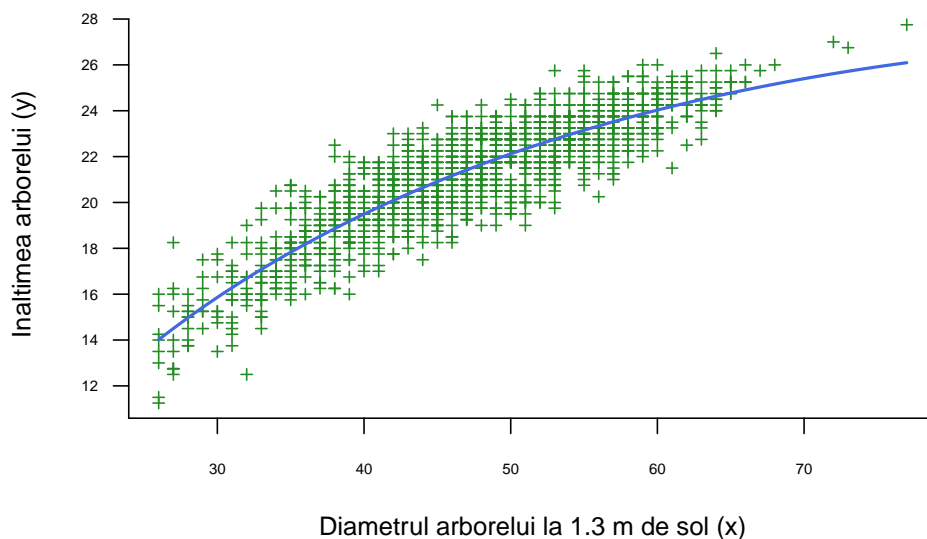
Media empirică a observațiilor y_i se poate deduce din prima componentă a vectorului $\mathbf{X}^\top \mathbf{Y}$

$$\bar{y} = \frac{(\mathbf{X}^\top \mathbf{Y})_1}{n} = \frac{30312.5}{1429} = 21.21 \text{ m}.$$

5p

3. Estimatorul lui $\boldsymbol{\beta} = (\beta_0, \beta_1, \beta_2)^\top$ obținut prin metoda celor mai mici pătrate este:

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} = \begin{pmatrix} -24.352 \\ -0.482 \\ 9.986 \end{pmatrix}$$



5p

4. Estimatorul nedeplasat $\hat{\sigma}^2$ a lui σ^2 este (a se vedea exercițiul 2.5 din Seminarul 2 și figura din secțiunea 2.3) dat de expresia

$$\hat{\sigma}^2 = \frac{\|\mathbf{Y} - \mathbf{X}\hat{\beta}\|^2}{n - (p + 1)} = \frac{\|\mathbf{Y}\|^2 - \|\mathbf{X}\hat{\beta}\|^2}{n - (p + 1)}$$

și cum

$$\|\mathbf{X}\hat{\beta}\|^2 = \hat{\beta}^\top \mathbf{X}^\top \mathbf{X} \hat{\beta} = \hat{\beta}^\top \mathbf{X}^\top \mathbf{X} \underbrace{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}}_{\hat{\beta}} = \hat{\beta}^\top \mathbf{X}^\top \mathbf{Y}$$

găsim că

$$\hat{\sigma}^2 = \frac{\mathbf{Y}^\top \mathbf{Y} - \hat{\beta}^\top \mathbf{X}^\top \mathbf{Y}}{n - (p + 1)} = \frac{651857.9 - 650017.2}{1429 - (2 + 1)} \approx 1.29.$$

Deoarece

$$T_2 = \frac{\hat{\beta}_2 - \beta_2}{\hat{\sigma} \sqrt{[(\mathbf{X}^\top \mathbf{X})^{-1}]_{3,3}}} = \frac{\hat{\beta}_2 - \beta_2}{\hat{\sigma}_{\hat{\beta}_2}} \sim t_{n-(p+1)} = t_{n-3}$$

găsim că un interval de încredere de nivel de încredere $\alpha = 95\%$ este

$$I(\beta_2) = \left[\hat{\beta}_2 - t_{n-3} \left(1 - \frac{\alpha}{2} \right) \hat{\sigma} \sqrt{[(\mathbf{X}^\top \mathbf{X})^{-1}]_{3,3}}, \hat{\beta}_2 + t_{n-3} \left(1 - \frac{\alpha}{2} \right) \hat{\sigma} \sqrt{[(\mathbf{X}^\top \mathbf{X})^{-1}]_{3,3}} \right]$$

și ținând cont că $t_{n-3}(0.975) = t_{1426}(0.975) \approx 1.96$ (din Teorema Limită Centrală) obținem

$$I(\beta_2) \approx [8.456, 11.517].$$

5p

5. Vrem să testăm ipotezele $H_0 : \beta_1 = 0$ versus $H_1 : \beta_1 \neq 0$ la un nivel de semnificație de 10%. Sub ipoteza nulă avem că

$$\frac{\hat{\beta}_1}{\hat{\sigma}_{\hat{\beta}_1}} = \frac{\hat{\beta}_1}{\hat{\sigma} \sqrt{[(\mathbf{X}^\top \mathbf{X})^{-1}]_{2,2}}} \sim t_{n-(p+1)} = t_{n-3} \approx \mathcal{N}(0, 1)$$

prin urmare este suficient să comparăm valoarea absolută a statisticii de test cu cea a cuantilei de ordin 0.95 pentru repartiția normală (adică 1.645):

$$\left| \frac{\hat{\beta}_1}{\hat{\sigma}_{\hat{\beta}_1}} \right| = \frac{|-0.482|}{\sqrt{1.29} \times \sqrt{0.002}} \approx 8.33 > 1.645$$

În consecință respingem ipoteza nulă $H_0 : \beta_1 = 0$.

5p

6. Notând cu $\mathbf{x}_{n+1}^\top = (1, x_{n+1}, \sqrt{x_{n+1}}) = (1, 49, 7)$, valoarea prezisă y_{n+1} a variabilei răspuns este

$$y_{n+1} = \mathbf{x}_{n+1}^\top \hat{\beta} \approx 21.89 m$$

iar un interval de predicție este

$$I(y_{n+1}) = \left[\mathbf{x}_{n+1}^\top \hat{\beta} - t_{n-(p+1)} \left(1 - \frac{\alpha}{2} \right) \hat{\sigma} \sqrt{1 + \mathbf{x}_{n+1}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{n+1}}, \mathbf{x}_{n+1}^\top \hat{\beta} + t_{n-(p+1)} \left(1 - \frac{\alpha}{2} \right) \hat{\sigma} \sqrt{1 + \mathbf{x}_{n+1}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{n+1}} \right]$$

care revine la $I(y_{n+1}) \approx [19.66, 24.12]$.

În mod similar găsim un interval de predicție și pentru y_{n+1} atunci când $\mathbf{x}_{n+1}^T = (1, x_{n+1}, \sqrt{x_{n+1}}) = (1, 25, 5)$:
 $I(y_{n+1}) \approx [11.25, 15.76]$.

