

11. I tabellen nedan finns för några år antalet skolor i Finland som gav undervisning på gymnasienivå:

År	2003	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013
Antal	483	479	471	461	449	449	441	439	433	428	421

Vi antar att (åtminstone approximativt) $Y_j = \beta_0 + \beta_1 x_j + \varepsilon_j$, där Y_j är antalet skolor som gav undervisning på gymnasienivå år x_j och slumpvariablerna ε_j är $N(0, \sigma^2)$ -fördelade och oberoende.

- Bestäm estimat b_0 och b_1 för koefficienterna i regressionsmodellen $Y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ med minsta kvadratmetoden.
- Beräkna estimatet s^2 för restvariansen och testa nollhypotesen $H_0 : \beta_1 \geq -5$ på signifikansnivån 1%.
- Vad är modellens förklaringsgrad?

I tabellen nedan finns data som underlättar räkningarna

\bar{y}	s_x^2	s_y^2	s_{xy}
450.36	11	430.85	-68

Lösning: (a) Medelvärdet av x -värdena är $\bar{x} = \frac{1}{11} \sum_{j=2003}^{2013} j = 2008$ så att estimaten för parametrarna β_1 och β_0 blir

$$b_1 = \frac{s_{xy}}{s_x^2} = \frac{-68}{11} = -6.1818$$

$$b_0 = \bar{y} - b_1 \bar{x} = 12863.$$

(b) Korrelationskoefficienten är

$$r_{xy} = \frac{s_{xy}}{\sqrt{s_x^2 \cdot s_y^2}} = \frac{-68}{\sqrt{11 \cdot 430.85}} = -0.98776$$

och då blir estimatet för restvariansen

$$s^2 = \frac{n-1}{n-2} s_y^2 (1 - r_{xy}^2) = \frac{10}{9} \cdot 430.85 \cdot (1 - (-0.98776)^2) = 11.647.$$

När vi testar nollhypotesen $\beta_1 \leq 5$ så använder vi testvariabeln

$$W_1 = \frac{B_1 - \beta_1}{\sqrt{\frac{s^2}{(n-1)s_x^2}}},$$

som är $t(n-2)$ -fördelad. Den här testvariabeln får värdet

$$w_1 = \frac{-6.1818 + 5}{\sqrt{\frac{11.647}{10 \cdot 11}}} = -3.6319$$

Eftersom alternativet till nollhypotesen är ensidigt blir p -värdet

$$p = F_{t(9)}(-3.6319) = 0.0027,$$

så vi kan förkasta nollhypotesen på signifikansnivån 0.01.

(c) Förklaringsgraden är enligt definitionen $r_{xy}^2 = 0.97567$.

I2. I Danmark var under åren 2004-2013 andelen (i procent) av befolkning som var 65 år eller äldre följande:

År	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013
Andel	14.9	15.0	15.2	15.3	15.6	15.9	16.3	16.8	17.3	17.8

Vi antar att (åtminstone approximativt) $Y_j = \beta_0 + \beta_1(x_j - 2015) + \varepsilon_j$, där Y_j är andelen år x_j av befolkningen i Danmark som äldre än 65 år och slumpvariablerna ε_j är $N(0, \sigma^2)$ -fördelade och oberoende.

(a) Bestäm estimat b_0 och b_1 för koefficienterna i regressionsmodellen $Y_i = \beta_0 + \beta_1(x_i - 2015) + \varepsilon_i$ med minsta kvadratmetoden. Vad är β_0 ?

(b) Beräkna estimatet s^2 för restvariansen och testa nollhypotesen $H_0 : \beta_0 = 18.8$ på signifikansnivån 1%

I tabellen nedan finns data som underlättar räkningarna

\bar{y}	s_x^2	s_y^2	r_{xy}
16.01	9.1667	1.0188	0.97260

När du byter ut talen x_j mot talen $x_j - 2015$ är det endast medelvärdet \bar{x} som måste räknas om.

Lösning: (a) Om vi byter ut talen x_j mot talen $\tilde{x}_j = x_j - 2015$ så blir $\bar{\tilde{x}} = -6.5$. För att räkna ut b_1 kan vi använda det faktum att r_{xy} är givet och att

$$b_1 = \frac{s_{xy}}{s_x^2} = \frac{s_{xy}}{\sqrt{s_x^2 \cdot s_y^2}} \cdot \sqrt{\frac{s_y^2}{s_x^2}} = r_{xy} \cdot \sqrt{\frac{s_y^2}{s_x^2}} = 0.9726 \cdot \sqrt{\frac{1.0188}{9.1667}} = 0.32424.$$

Estimatet b_0 för β_0 blir

$$b_0 = \bar{y} - b_1 \bar{\tilde{x}} = 16.01 - 0.32424 \cdot (-6.5) = 18.118$$

(b) Estimatet för restvariansen är

$$s^2 = \frac{n-1}{n-2} s_y^2 (1 - r_{xy}^2) = \frac{9}{8} \cdot 1.0188 \cdot (1 - 0.97260^2) = 0.061949.$$

När vi testar nollhypotesen $\beta_0 \leq 19$ så använder vi testvariabeln

$$W_0 = \frac{B_0 - \beta_0}{\sqrt{S^2 \left(\frac{1}{n} + \frac{\bar{\tilde{x}}^2}{(n-1)s_x^2} \right)}},$$

som är $t(n-2)$ -fördelad. Den här testvariabeln får värdet

$$w_1 = \frac{18.118 - 18.8}{\sqrt{0.061949 \cdot \left(\frac{1}{10} + \frac{(-6.5)^2}{9 \cdot 9.1667} \right)}} = -3.5023.$$

Eftersom alternativet till nollhypotesen är dubbelsidigt blir p -värdet

$$p = 2 \cdot F_{t(8)}(-3.5023) = 0.0081 < 0.01$$

så vi kan förkasta nollhypotesen på signifikansnivån 0.01.

I3. Utsläppen av växthusgaser i samband med avfallshantering i Finland var under åren 2001-2010 följande (i miljoner ton CO₂-ekv.):

2001	2002	2003	2004	2005	2006	2007	2008	2009	2010
3.14	2.92	2.75	2.61	2.40	2.46	2.38	2.28	2.19	2.19

Bestäm med stöd av dessa siffror ett 95% konfidensintervall för väntevärdet av utsläppen år 2014 om du antar att utsläppen kan beskrivas med modellen $Y_j = \beta_0 + \beta_1 x_j + \varepsilon_j$ där Y_j är utsläppen år x_j och slumpvariablerna ε_j är oberoende och $N(0, \sigma^2)$ -fördelade.

Kun kan använda följande data:

\bar{x}	\bar{y}	s_x^2	s_y^2	s_{xy}
2005.5	2.532	9.1667	0.10188	-0.92444

och kom ihåg att om du byter ut årtalen mot talen $x_j - 2014$ så är det endast \bar{x} som förändras och du skall räkna ett konfidensintervall för parametern β_0 .

Lösning: Om $\tilde{x}_j = x_j - 2014$ så är $\bar{\tilde{x}} = -8.5$ men $s_{\tilde{x}}^2 = s_x^2$ ja $s_{\tilde{x}y} = s_{xy}$.

Vi räknar först ut estimaten för parametrarna β_1 och β_0 och de blir

$$b_1 = \frac{s_{xy}}{s_x^2} = \frac{-0.92444}{9.1667} = -0.10085,$$

$$b_0 = \bar{y} - b_1 \bar{\tilde{x}} = 2.532 - (-0.10085) \cdot (-8.5) = 1.6748.$$

För att kunna räkna ett konfidensintervall behöver vi estimatet s^2 för restvariansen och först skall vi räkna ut stickprovskorrelationen som är

$$r_{xy} = \frac{s_{xy}}{\sqrt{s_x^2 \cdot s_y^2}} = \frac{-0.92444}{\sqrt{9.1667 \cdot 0.10188}} = -0.95659.$$

Estimatet för stickprovsvariansen är

$$s^2 = \frac{n-1}{n-2} s_y^2 (1 - r_{xy}^2) = \frac{9}{8} \cdot 0.10188 \cdot (1 - (-0.95659)^2) = 0.0097349.$$

Testvariabeln

$$W_0 = \frac{B_0 - \beta_0}{\sqrt{S^2 \left(\frac{1}{n} + \frac{\bar{\tilde{x}}^2}{(n-1)s_x^2} \right)}},$$

är $t(n-2)$ -fördelad så att

$$\Pr(-2.306 \leq W_0 \leq 2.306) = 0.95$$

eftersom $2.306 = F_{t(8)}^{-1}(0.975)$ av vilket följer att

$$\Pr \left(B_0 - 2.306 \sqrt{S^2 \left(\frac{1}{n} + \frac{\bar{\tilde{x}}^2}{(n-1)s_x^2} \right)} \leq \beta_0 \leq B_0 + 2.306 \sqrt{S^2 \left(\frac{1}{n} + \frac{\bar{\tilde{x}}^2}{(n-1)s_x^2} \right)} \right) = 0.95.$$

Konfidensintervallet är alltså

$$\left[B_0 - 2.306 \sqrt{S^2 \left(\frac{1}{n} + \frac{\bar{\tilde{x}}^2}{(n-1)s_x^2} \right)}, B_0 + 2.306 \sqrt{S^2 \left(\frac{1}{n} + \frac{\bar{\tilde{x}}^2}{(n-1)s_x^2} \right)} \right],$$

och om vi nu sätter in de värden vi har för de olika variablerna så får vi intervallet
 $[1.4501, 1.8995]$.

I4. Med hjälp av observationerna (x_j, y_j) , $j = 1, \dots, n$ har vi räknat ut koefficienterna i regressionslinjen $y = b_0 + b_1x$ i det fall att x är den ”förklarande” variabeln (dvs. då vi antar att $Y_j = \beta_0 + \beta_1x_j + \varepsilon_j$) och också för den ”inversa” regressionslinjen $x = a_0 + a_1y$ där y är den ”förklarande” variabeln (dvs. då vi antar att $X_j = \alpha_0 + \alpha_1y_j + \varepsilon_j$). Vad kan man säga om följande påståenden (då vi antar att $s_x > 0$ och $s_y > 0$):

- (a) Linjerna $y = b_0 + b_1x$ och $x = a_0 + a_1y$ sammanfaller om och endast om $|r_{xy}| = 1$.
- (b) Ifall nollhypotesen $H_0 : \beta_1 = 0$ förkastas på sigifikansnivån 0.01 så förkastas också nollhypotesen $H_0 : \alpha_1 = 0$ på signifikansnivån 0.01.

Motivera dina svar!

Ledning: (a) Regressionslinjerna kan skrivas i formen $y - \bar{y} = b_1(x - \bar{x})$ och $x - \bar{x} = a_1(y - \bar{y})$. (b) Kom ihåg hur man kan skriva testvariabeln då man testar $\beta_1 = 0$ med hjälp av korrelationskoefficienten och vad blir på motsvarande sätt testvariabeln då man testar nollhypotesen $\alpha_1 = 0$?

Lösning: (a) Båda linjerna går genom punkten (\bar{x}, \bar{y}) vilket betyder att linjerna sammanfaller om och endast om de har samma vinkelkoefficient. Nu är $b_1 = r_{xy} \frac{s_y}{s_x}$ och $a_1 = r_{xy} \frac{s_x}{s_y}$ och vinkelkoefficienterna är b_1 och $\frac{1}{a_1}$ så att dessa är lika då $b_1 a_1 = 1$, dvs. då $r_{xy}^2 = 1$ eller $r_{xy} = \pm 1$. Påståendet gäller.

(b) Testvariabeln kan i båda fallen eftersom vi testar hypoteserna $\beta_1 = 0$ och $\alpha_1 = 0$ skrivas i formen $\frac{r_{xy} \sqrt{n-2}}{\sqrt{1-r_{xy}^2}}$ så det är frågan om samma test och därför gäller även detta påstående.

I5. Vi har ett observerat stickprov (x_i, y_i) $i = 1, \dots, 25$ och vi har räknat ut stickprovskorrelationen som blev $r_{xy} = -0.47$. Bilda ett (approximativt) symmetriskt 99% konfidensintervall för korrelationen ρ_{XY} med hjälp av Fischers transformation som säger att $\frac{1}{2} \ln \left(\frac{1+R_{xy}}{1-R_{xy}} \right) \sim_a N \left(\frac{1}{2} \ln \left(\frac{1+\rho_{xy}}{1-\rho_{xy}} \right), \frac{1}{n-3} \right)$ tex. på följande sätt:

- (a) Bestäm talet z så att $\Pr(-z \leq Z \leq z) \approx 0.99$ om $Z \sim_a N(0, 1)$.

- (b) Välj $Z = \frac{\frac{1}{2} \ln \left(\frac{1+r_{xy}}{1-r_{xy}} \right) - \frac{1}{2} \ln \left(\frac{1+\rho_{xy}}{1-\rho_{xy}} \right)}{\sqrt{\frac{1}{n-3}}}$ (i enlighet med Fishers transformation) och bestäm tal a och b så att

$$-z \leq Z \leq z \quad \Leftrightarrow \quad a \leq \ln \left(\frac{1+\rho_{xy}}{1-\rho_{xy}} \right) \leq b.$$

- (c) Bestäm talen r_L och r_U så att

$$a \leq \ln \left(\frac{1+\rho_{xy}}{1-\rho_{xy}} \right) \leq b \quad \Leftrightarrow \quad r_L \leq \rho_{XY} \leq r_U.$$

Konfidensintervallet är nu $[r_L, r_U]$.

Lösning: (a) Eftersom $0.005 + 0.005 = 0.01 = 1 - 0.99$ så är $z = z_{0.005} = F_{N(0,1)}^{-1}(0.995) = 2.5758$.

(b) Eftersom $r_{xy} = -0.47$ så är $\frac{1}{2} \ln \left(\frac{1+r_{xy}}{1-r_{xy}} \right) = -0.51007$ och $\sqrt{\frac{1}{25-3}} = 0.2132$ så att

$$-z \leq Z \leq z \Leftrightarrow -2.5758 \leq \frac{1}{0.2132} \cdot \left(-0.51007 - \frac{1}{2} \ln \left(\frac{1+\rho_{xy}}{1-\rho_{xy}} \right) \right) \leq 2.5758$$

$$\Leftrightarrow 0.2132 \cdot (-2.5758) + 0.51007 \leq -\frac{1}{2} \ln \left(\frac{1+\rho_{xy}}{1-\rho_{xy}} \right) \leq 0.2132 \cdot 2.5758 + 0.51007$$

$$\Leftrightarrow -2 \cdot (0.2132 \cdot 2.5758 + 0.51007) \leq \ln \left(\frac{1+\rho_{xy}}{1-\rho_{xy}} \right) \leq 2 \cdot (0.2132 \cdot (-2.5758) + 0.51007),$$

så att $a = -2 \cdot (0.2132 \cdot 2.5758 + 0.51007) = -2.1185$ och $b = -2 \cdot (0.2132 \cdot (-2.5758) + 0.51007) = 0.078181$.

(c) Vi har

$$-2.1185 \leq \ln \left(\frac{1+\rho_{xy}}{1-\rho_{xy}} \right) \leq 0.078181 \Leftrightarrow e^{-2.1185} \leq \frac{1+\rho_{xy}}{1-\rho_{xy}} \leq e^{0.078181}$$

$$\Leftrightarrow e^{-2.1185} \cdot (1-\rho_{xy}) \leq 1+\rho_{xy} \leq e^{0.078181} \cdot (1-\rho_{xy})$$

$$\Leftrightarrow \frac{e^{-2.1185} - 1}{e^{-2.1185} + 1} \leq \rho_{XY} \leq \frac{e^{0.078181} - 1}{e^{0.078181} + 1}$$

så att

$$r_L = \frac{e^{-2.1185} - 1}{e^{-2.1185} + 1} = -0.78538,$$

$$r_U = \frac{e^{0.078181} - 1}{e^{0.078181} + 1} = 0.039071.$$
