

Regression

Modell: Vi antar ett linjärt samband

$$Y = \beta_0 + \beta_1 X \text{ som modelleras av}$$

$$\sum_i = \beta_0 + \beta_1 x_i + E_i \text{ där}$$

i) x_1, x_2, \dots, x_n är givna kända tal. De kallas för förklarande variabler,

ii) E_1, E_2, \dots, E_n är obero. s.v. (störningar)

Sådana att $E_i \sim N(0, \sigma^2)$

Vi antar här att $E_i \sim N(0, \sigma^2)$,

iii) Vi vill skatta β_0 = intercept och β_1 = lutningskoeff. / slope.

iv) \sum_1, \dots, \sum_n kallas för beroende variabler.

v) Linjen $y = \beta_0 + \beta_1 x$ kallas för den teoretiska regressionslinjen.

vi) Mätningar kommer resultera i data $(x_1, y_1), \dots, (x_n, y_n)$ där

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

Och ε_i är en realisation av E_i .

Vi skattar β_0, β_1 genom att välja $\hat{\beta}_0$ och $\hat{\beta}_1$, så att

$$Q(b_0, b_1) = \sum_{k=1}^n (y_k - (b_0 + b_1 x_k))^2 = \sum_{k=1}^n e_k^2 \text{ minimeras.}$$

<u>Ex:</u>	F-halt	0.0	0.3	
Elasticitet	12.86	13.10		

"visa bild 1"

Om $b_0 = 14$ och $b_1 = 0.5$ så blir

$$Q(14, 0.5) \approx 68.1$$

Vi minimerar $Q(b_0, b_1)$ genom att lösa

$$\frac{\partial Q}{\partial b_0} = \frac{\partial Q}{\partial b_1} = 0.$$

Om vi läter $s_{xy} = \sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y})$,

$s_{xx} = \sum_{k=1}^n (x_k - \bar{x})^2$ och $s_{yy} = \sum_{k=1}^n (y_k - \bar{y})^2$ får vi

att minimum uppnås av

$$\hat{\beta}_1 = \frac{s_{xy}}{s_{xx}} \quad \text{och} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}.$$

Vår anpassade regressionslinje blir

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

där $\hat{\beta}_1 \approx 1.74$ och $\hat{\beta}_0 \approx 12.76$

Som ett mått på anpassning beräknas

förhörlingsgraden (coefficient of determination)

$$R^2 = 1 - \frac{\sum_{k=1}^n (y_k - \hat{y}_k)^2}{\sum_{k=1}^n (y_k - \bar{y})^2} = \frac{s_{xy}^2}{s_{xx} s_{yy}}$$

där $y_k - \hat{y}_k = y_k - (\hat{\beta}_0 + \hat{\beta}_1 x_k)$

- Vi har att $R^2 \in [0, 1]$ och att $R^2 \approx 1$ är bra, men $R^2 \approx 0$ är dåligt.

Under antagandet att $E_i \sim N(0, \sigma^2)$ är i.i.d. kan man visa att

$$\hat{\beta}_1 = \hat{\beta}_1(E_1, \dots, E_n) \sim N(\beta_1, \frac{\sigma^2}{s_{xx}}) \text{ och}$$

$$\hat{\beta}_0 = \hat{\beta}_0(E_1, \dots, E_n) \sim N(\beta_0, \sigma^2 \frac{\sum_{k=1}^n x_k^2}{n s_{xx}})$$

- OBS! $\mathbb{E}[\hat{\beta}_0] = \beta_0$ och $\mathbb{E}[\hat{\beta}_1] = \beta_1$ så båda är VVR,

K.I. för β_1

Oftast är σ^2 okänd. I så fall skattar vi

$$\sigma^2 \text{ med } S_r^2 = \frac{1}{n-2} \sum_{k=1}^n (y_k - \hat{y}_k)^2 = \frac{1}{n-2} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k)^2$$

$$= \frac{1}{n-2} SSE$$

Man kan visa att $\mathbb{E}[S_r^2] = \sigma^2$ så
att S_r^2 är en VVR sättbar av σ^2

OBS: $S^2 = \frac{1}{n-1} \sum_{k=1}^n (y_k - \bar{y})^2$ är ej bra här

"tar ej hänsyn till lutningen"

○ Liknande tidigare i kursen får vi att

$$\frac{\hat{\beta}_1 - \beta_1}{S_r / \sqrt{S_{xx}}} \sim t(n-2)$$

Vi kan då hitta ett I.L.I. för β_1 .

$$1-\alpha = P\left(-t_{\alpha/2}(n-2) \leq \frac{\hat{\beta}_1 - \beta_1}{S_r / \sqrt{S_{xx}}} \leq t_{\alpha/2}(n-2)\right)$$

$$= = = P\left(\hat{\beta}_1 - t_{\alpha/2}(n-2) \frac{S_r}{\sqrt{S_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{\alpha/2}(n-2) \frac{S_r}{\sqrt{S_{xx}}}\right)$$

○ så att

$$I_{\beta_1} = \hat{\beta}_1 \pm t_{\alpha/2}(n-2) \frac{S_r}{\sqrt{S_{xx}}}$$

är ett I.L.I. med kont. grad $(1-\alpha)100\%$ för β_1 .

Tal: För exemplet ovan, bestäm ett 95% I.L.I. för β_1 .

L: Data ger oss att

$$S_r^2(x) = \frac{1}{n-2} \sum_{k=1}^n (y_k - \hat{y}_k)^2 \approx 0.239$$

Tabell ger oss att $t_{0.025}(6) \approx 2.45$
 så ett 95% numeriskt I.L.T. för β_1
 blir

$$I_{\beta_1} = \hat{\beta}_1 \pm 2.45 \cdot \frac{\sqrt{0.239}}{\sqrt{S_{xx}}} \approx [1.31, 2.16].$$

- Utifrån vårt I.L.T. ovan verkar det
 troligt att $\beta_1 \neq 0$. tolkningen av detta är
 att halten F påverkar elasticiteten!
- Vi vill göra hypotestest på detta.

Vi betraktar följande situationer:

I	$H_0: \beta_1 = 0$	II	$H_0: \beta_1 = 0$	III	$H_0: \beta_1 = 0$
	$H_1: \beta_1 > 0$		$H_1: \beta_1 < 0$		$H_1: \beta_1 \neq 0$

- Tal: Testa om halten F påverkar
 elasticiteten på signifikansnivåen
 $\alpha = 0.05$.

L: Vi shall alltså testa

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

Vi gör detta genom att beräkna
 p-värde. Testet är tvåsidigt så att
 $p\text{-värde} = 2 P(T(E) \geq |T(\epsilon)|)$

Här väljer vi teststatistiken

$$T(E_1, \dots, E_n) = \frac{\hat{\beta}_1(E_1, \dots, E_n)}{S_r(E_1, \dots, E_n) / \sqrt{S_{xx}}}$$

så att

$$T(\varepsilon_1, \dots, \varepsilon_n) = \frac{\hat{\beta}_1(\varepsilon_1, \dots, \varepsilon_n)}{S_r(\varepsilon_1, \dots, \varepsilon_n) / \sqrt{S_{xx}}} \approx \frac{1.74}{\sqrt{0.239} / \sqrt{S_{xx}}} \approx 10.$$

Vi ser från tabell (sid 699) att

$$2P(T(E) \geq |T(\varepsilon)|) \approx 2P(T(E) \geq 10) \leq 2P(T(E) \geq 5.96)$$

$$\leq 2 \cdot 0.0005 = 0.001 \text{ där } T(E) \sim t(6)$$

Då p-värdet < α så förhas tar vi H_0 på

signifikansnivån $\alpha = 0.05$,

Kort rimlighetsanalys

Om förklaringsgraden är låg ($R^2 \leq 0.8$)

så bör vi vara skeptiska till vårt modell-
antagande. Vi kan då betrakta residualerna

$$e_k = y_k - \hat{y}_k$$

och leta efter mönster.

"Visa figurer"

I detta exempel blev $R^2 \approx 0.62$,

kanse är

$$\Sigma_i = \beta_0 + \beta_1(x_i)^2 + e_i$$

en bättre modell?