

## Linjär regression

Betrakta följande modell:

$$Y(x) \text{ är } N(m(x), \sigma)$$

det vill säga  $Y$  är en normalfördelad stokastisk variabel vars väntevärde beror på en underliggande parameter  $x$ . Vi kommer i huvudsak att betrakta funktioner

$$m(x) = \alpha' + \beta x,$$

det vill säga att väntevärdet för  $Y(x)$  är linjär i  $x$ . Problemet nu är hur vi med observationer

$$y_1(x_1), y_2(x_2), \dots, y_n(x_n)$$

skall skatta  $\alpha'$  och  $\beta$ .

För vår observationsserie, låt oss byta parametrering till

$$m(x) = \alpha' + \beta x = \alpha' + \beta x - \beta \bar{x} + \beta \bar{x} = \underbrace{\alpha' + \beta \bar{x}}_{=\alpha} + \beta(x - \bar{x}) = \alpha + \beta(x - \bar{x}),$$

det vill säga vi har infört  $\alpha$  istället för  $\alpha'$ . Med minsta kvadratmetoden får vi att skattningarna av  $\alpha$  och  $\beta$  är de värden som minimerar

$$Q(\alpha, \beta) = \sum_{i=1}^n (y_i - m(x_i))^2 = \sum_{i=1}^n (y_i - (\alpha + \beta(x_i - \bar{x})))^2.$$

Derivering ger ekvationen

$$\begin{aligned} \frac{\partial}{\partial \alpha} Q(\alpha, \beta) &= -2 \sum_{i=1}^n (y_i - (\alpha + \beta(x_i - \bar{x}))) = -2 \left( \sum_{i=1}^n y_i \right) + 2n\alpha + \beta \underbrace{\sum_{i=1}^n (x_i - \bar{x})}_{=0} \\ &= -2n\bar{y} + 2n\alpha = 0, \end{aligned}$$

eller  $\alpha^* = \bar{y}$ . Vidare,

$$\begin{aligned} \frac{\partial}{\partial \beta} Q(\alpha, \beta) &= -2 \sum_{i=1}^n (y_i - \bar{y} - \beta(x_i - \bar{x}))(x_i - \bar{x}) \\ &= -2 \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) + 2\beta \sum_{i=1}^n (x_i - \bar{x})^2 = 0, \end{aligned}$$

eller

$$\beta^* = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}}.$$

Skattningarna beskrivs av motsvarande stokastiska variabler:

$$\alpha^* = \bar{Y} \quad \beta^* = \frac{S_{XY}}{S_{XX}}$$

som båda är linjärkombinationer av normalfördelade variabler och alltså är normalfördelade. Slutligen får vi en skattning av  $\sigma^2$  som

$$s^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - (\alpha^* + \beta^*(x_i - \bar{x})))^2 = \frac{1}{n-2} Q(\alpha^*, \beta^*)$$

## Sammanfattning

$$\begin{aligned}\alpha^* &= \bar{Y} \text{ är } N\left(\alpha, \frac{\sigma}{\sqrt{n}}\right) \\ \beta^* &= \frac{S_{xY}}{S_{xx}} \text{ är } N\left(\beta, \frac{\sigma}{\sqrt{S_{xx}}}\right) \\ S^2 &= \frac{1}{n-2}(S_{YY} - \beta^* S_{xY}) = \frac{1}{n-2}Q(\alpha^*, \beta^*) \\ \frac{(n-2)S^2}{\sigma^2} &= \frac{Q(\alpha^*, \beta^*)}{\sigma^2} \text{ är } \chi^2(n-2)\end{aligned}$$

Fixera en punkt  $x_0$ . Ett konfidensintervall för  $m_0 = m(x_0) = \mathbf{E}[Y(x_0)] = \alpha + \beta(x_0 - \bar{x})$  erhålls som följer. Vi skattar  $\alpha + \beta(x_0 - \bar{x})$  med den normalfördelade variabeln

$$m_0^* = \alpha^* + \beta^*(x_0 - \bar{x})$$

som har väntevärde

$$\mathbf{E}[m_0^*] = \mathbf{E}[\alpha^* + \beta^*(x_0 - \bar{x})] = \alpha + \beta(x_0 - \bar{x}) = m_0$$

och varians

$$\mathbf{V}(m_0^*) = \mathbf{V}(\alpha^*) + \mathbf{V}(\beta^*)(x_0 - \bar{x})^2 = \frac{\sigma^2}{n} + \sigma^2 \frac{(x_0 - \bar{x})^2}{S_{xx}}.$$

(Här utnyttjade vi att  $\alpha^*$  och  $\beta^*$  är oberoende!) Alltså är

$$\frac{m_0^* - m_0}{S\sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}} \text{ är } t(n-2).$$

Ett  $1 - \gamma$  konfidensintervall ges av

$$m_0 \in m_0^* \pm t_{\gamma/2} S \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}.$$