

Probability and Statistics

7 – Regression

Stefan Heiss

Technische Hochschule Ostwestfalen-Lippe
Dep. of Electrical Engineering and Computer Science

January 26, 2024

Theorem (7.5)

If there exists some $\sigma > 0$, such that $Y_x \sim \mathcal{N}(\alpha + \beta x, \sigma)$ for all $x \in \mathbb{R}$, then:

(i)

$$B \sim \mathcal{N}(\beta, \sigma/\sqrt{s_x})$$

B is an unbiased estimator for β

(ii)

$$A \sim \mathcal{N}\left(\alpha, \sigma \cdot \sqrt{\frac{\sum_{i=1}^n x_i^2}{n s_x}}\right)$$

A is an unbiased estimator for α

(iii) For any x_0 :

$$A + B x_0 \sim \mathcal{N}\left(\alpha + \beta x_0, \sigma \cdot \sqrt{\frac{s_x + n(x_0 - \bar{x})^2}{n s_x}}\right)$$

(iii) \Rightarrow (ii) : (iii) with $x_0 = 0$

$$\sum x_i^2 = s_x + n \bar{x}^2$$

see Q.2)

Proof of Theorem (7.5) (i)

$$\underline{B} = \frac{S_{XY}}{s_X} = \frac{1}{s_X} \cdot \left(\sum_{i=1}^n (x_i - \bar{x}) Y_i - \bar{Y} \sum_{i=1}^n (x_i - \bar{x}) \right) = \frac{1}{s_X} \cdot \underline{\sum_{i=1}^n (x_i - \bar{x}) Y_i}$$

and B is therefore a linear combination of the independent normally distributed random variables Y_1, \dots, Y_n . Hence B itself has a normal distribution with:

$$\begin{aligned} E(B) &= \frac{1}{s_X} \cdot \sum_{i=1}^n (x_i - \bar{x}) E(Y_i) = \frac{1}{s_X} \cdot \sum_{i=1}^n (x_i - \bar{x}) (\alpha + \beta x_i) \\ &= \frac{1}{s_X} \left(\alpha \cdot \sum_{i=1}^n (x_i - \bar{x}) + \beta \cdot \sum_{i=1}^n x_i (x_i - \bar{x}) \right) = \frac{1}{s_X} (\beta \cdot s_X) = \beta \end{aligned}$$

$$Var(B) = \frac{1}{s_X^2} \cdot \sum_{i=1}^n (x_i - \bar{x})^2 Var(Y_i) = \frac{\sigma^2 s_X}{s_X^2} = \frac{\sigma^2}{s_X}$$

Proof of Theorem (7.5) (iii)

$$B = \frac{S_{xY}}{s_x} = \frac{1}{s_x} \cdot \left(\sum_{i=1}^n (x_i - \bar{x}) Y_i - \bar{Y} \sum_{i=1}^n (x_i - \bar{x}) \right) = \frac{1}{s_x} \cdot \underbrace{\sum_{i=1}^n (x_i - \bar{x}) Y_i}_{\cdot (x_0 - \bar{x})}$$

From the proof of (i) given above, we have:

$$A = \bar{Y} - B \cdot \bar{x} \quad \checkmark$$

$$A + B x_0 = \underbrace{\bar{Y}} + B (x_0 - \bar{x}) = \sum_{i=1}^n \left(\underbrace{\frac{1}{n}} + \underbrace{\frac{(x_0 - \bar{x})(x_i - \bar{x})}{s_x}} \right) \underbrace{Y_i}$$

Proof of Theorem (7.5) (iii)

From the proof of (i) given above, we have:

$$A + Bx_0 = \bar{Y} + B(x_0 - \bar{x}) = \sum_{i=1}^n \left(\frac{1}{n} + \frac{(x_0 - \bar{x})(x_i - \bar{x})}{s_x} \right) Y_i$$

Therefore, $A + Bx_0$ is a sum of the independent normally distributed random variables Y_i .

Hence, $A + Bx_0$ has a normal distribution with:

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i, \quad Y_i \sim \mathcal{N}(\alpha + \beta x_i, \sigma^2)$$

$$\begin{aligned} E(A + Bx_0) &= E(\bar{Y} + B(x_0 - \bar{x})) = E(\bar{Y}) + E(B(x_0 - \bar{x})) \\ &= \frac{1}{n} \sum_{i=1}^n (\alpha + \beta x_i) + \beta(x_0 - \bar{x}) \quad \xleftarrow{(i)} \\ &= \alpha + \beta x_0 \end{aligned}$$

Proof of Theorem (7.5) (iii)

$$\begin{aligned}
 \text{Var}(A + Bx_0) &= \sigma^2 \sum_{i=1}^n \left(\frac{1}{n} + \frac{(x_0 - \bar{x})(x_i - \bar{x})}{s_x} \right)^2 \\
 &= \frac{\sigma^2}{n^2 s_x^2} \sum_{i=1}^n (s_x + n(x_0 - \bar{x})(x_i - \bar{x}))^2 \\
 &= \frac{\sigma^2}{n^2 s_x^2} \left(ns_x^2 + 2s_x n(x_0 - \bar{x}) \underbrace{\sum_{i=1}^n (x_i - \bar{x})}_0 + n^2 (\bar{x} - x_0)^2 \underbrace{\sum_{i=1}^n (x_i - \bar{x})^2}_{s_x} \right) \\
 &= \frac{\sigma^2}{n^2 s_x^2} (ns_x^2 + n^2 (x_0 - \bar{x})^2 s_x) \\
 &= \frac{\sigma^2}{n s_x} (s_x + n(x_0 - \bar{x})^2) = \frac{\sigma^2}{n} + \frac{\sigma^2}{s_x} (x_0 - \bar{x})^2
 \end{aligned}$$

Theorem (7.5)

 σ not known

If there exists some $\sigma > 0$, such that $Y \sim \mathcal{N}(\alpha + \beta x, \sigma)$ for all $x \in \mathbb{R}$, then:

(iv) S_R is independent of A and B and:

without proof

$$S_R / \sigma^2 \sim \chi_{n-2}^2$$

(omitting the proof)

(v)

$$\sqrt{\frac{(n-2)s_x}{S_R}} \cdot (B - \beta) \sim t_{n-2}$$

(vi)

$$\sqrt{\frac{n(n-2)s_x}{S_R \cdot \sum_{i=1}^n x_i^2}} \cdot (A - \alpha) \sim t_{n-2}$$

Proof of Theorem (7.5) (v)

From (i) and (iv) it follows that

$$\frac{B - \beta}{\sigma / \sqrt{s_x}} \sim \mathcal{N}(0, 1)$$

and:

$$\frac{\frac{B - \beta}{\sigma / \sqrt{s_x}}}{\sqrt{\frac{S_R}{\sigma^2 (n-2)}}} = \sqrt{\frac{(n-2) s_x}{S_R}} \cdot (B - \beta) \sim t_{n-2}$$

(Handwritten red annotations: a circle around the denominator's fraction, a green checkmark on the sigma, and the text 'chi^2_{n-2}' below the denominator)

Proof of Theorem (7.5) (vi)

From (ii) and (iv) it follows that

$$\frac{A - \alpha}{\frac{\sigma \sqrt{\sum_{i=1}^n x_i^2}}{\sqrt{n s_x}}} \sim \mathcal{N}(0, 1)$$

and:

$$\frac{\sqrt{n s_x} \cdot \frac{A - \alpha}{\sigma \sqrt{\sum_{i=1}^n x_i^2}}}{\sqrt{\frac{S_R}{\sigma^2 (n-2)}}} = \sqrt{\frac{n(n-2) s_x}{S_R \sum_{i=1}^n x_i^2}} \cdot (A - \alpha) \sim t_{n-2}$$

Theorem (7.5)

If there exists some $\sigma > 0$, such that $Y \sim \mathcal{N}(\alpha + \beta x, \sigma)$ for all $x \in \mathbb{R}$, then:

(vii) For any x_0 :

$$\sqrt{\frac{n(n-2)s_x}{S_R(s_x + n(x_0 - \bar{x})^2)}} \cdot ((A + Bx_0) - (\alpha + \beta x_0)) \sim t_{n-2}$$

(viii) If Y_0 denotes the response for the input value x_0 , then $Y_0 \sim \mathcal{N}(\alpha + \beta x_0, \sigma)$ and:

$$\sqrt{\frac{n(n-2)s_x}{S_R((n+1)s_x + n(x_0 - \bar{x})^2)}} \cdot (Y_0 - (A + Bx_0)) \sim t_{n-2}$$

Proof of Theorem (7.5) (vii)

From (iii) and (iv) it follows that

$$\frac{A + B x_0 - (\alpha + \beta x_0)}{\sigma \sqrt{\frac{s_x + n(x_0 - \bar{x})^2}{n s_x}}} \sim \mathcal{N}(0, 1)$$

and:

$$\frac{\frac{A + B x_0 - (\alpha + \beta x_0)}{\sigma \sqrt{\frac{s_x + n(x_0 - \bar{x})^2}{n s_x}}}}{\sqrt{\frac{S_R}{\sigma^2 (n-2)}}} = \sqrt{\frac{n(n-2) s_x}{S_R (s_x + n(x_0 - \bar{x})^2)}} \cdot (A + B x_0 - (\alpha + \beta x_0)) \sim t_{n-2}$$

Proof of Theorem (7.5) (viii)

$$V_{\sigma}(A + Bx_0) = \frac{\sigma^2}{n s_x} (s_x + n(x_0 - \bar{x})^2) + \sigma^2$$

$$Y_0 \sim \mathcal{N}(\alpha + \beta x_0, \underline{\sigma})$$

is independent of Y_1, \dots, Y_n (and consequently independent of $A + Bx_0$). It follows from (iii):

$$\begin{aligned} Y_0 - (A + Bx_0) &\sim \mathcal{N}\left(0, \sqrt{\frac{\sigma^2}{n s_x} (ns_x + s_x + n(x_0 - \bar{x})^2)}\right) \\ &= \mathcal{N}\left(0, \sigma \cdot \sqrt{\frac{(n+1)s_x + n(x_0 - \bar{x})^2}{n s_x}}\right) \end{aligned}$$

and:

$$\frac{\frac{Y_0 - (A + Bx_0)}{\sigma \sqrt{\frac{(n+1)s_x + n(x_0 - \bar{x})^2}{n s_x}}}}{\sqrt{\frac{S_R}{\sigma^2(n-2)}}} = \sqrt{\frac{n(n-2)s_x}{S_R((n+1)s_x + n(x_0 - \bar{x})^2)}} \cdot (Y_0 - (A + Bx_0)) \sim t_{n-2}$$

Notation (7.6)

Since the Greek character α is already used to denote one of the constants in our simple linear regression equation, we will use $1 - \gamma$ and γ in the following to denote confidence and significance levels, respectively.

(7.7) Confidence interval for β with confidence level $1 - \gamma$

Let

$$\delta_t := F_{t_{n-2}}^{-1} \left(1 - \frac{\gamma}{2} \right)$$

Then:

$$\begin{aligned} \Pr \left(\left| \sqrt{\frac{(n-2)s_x}{S_R}} \cdot (B - \beta) \right| \leq \delta_t \right) &= F_{t_{n-2}}(\delta_t) - F_{t_{n-2}}(-\delta_t) \\ &= 1 - \frac{\gamma}{2} - \left(1 - \left(1 - \frac{\gamma}{2} \right) \right) = 1 - \gamma \end{aligned}$$

(7.7) Confidence interval for β with confidence level $1 - \gamma$

Given sampled data values $y = (y_1, \dots, y_n)$ from (Y_1, \dots, Y_n) , the value $B(y)$ and $S_R(y)$ can be calculate (substituting the Y_i 's in the definition of B and S_R by the sampled values in y) and with a confidence of level $1 - \gamma$:

$$\left| \sqrt{\frac{(n-2)s_x}{S_R(y)}} \cdot (B(y) - \beta) \right| \leq \delta_t = F_{t_{n-2}}^{-1} \left(1 - \frac{\gamma}{2} \right)$$

$$\Leftrightarrow |B(y) - \beta| \leq \sqrt{\frac{S_R(y)}{(n-2)s_x}} \cdot \delta_t = \sqrt{\frac{S_R(y)}{(n-2)s_x}} \cdot F_{t_{n-2}}^{-1} \left(1 - \frac{\gamma}{2} \right)$$

$$\Leftrightarrow \beta = B(y) \pm \sqrt{\frac{S_R(y)}{(n-2)s_x}} \cdot F_{t_{n-2}}^{-1} \left(1 - \frac{\gamma}{2} \right)$$

(7.8) Hypothesis test of $H_0 : \beta = \beta_0$

Given a significance level γ and sampled data values $y = (y_1, \dots, y_n)$ from (Y_1, \dots, Y_n) :

- H_0 is rejected if $\sqrt{\frac{(n-2)s_x}{S_R(y)}} \cdot |B(y) - \beta_0| > F_{t_{n-2}}^{-1}\left(1 - \frac{\gamma}{2}\right)$
- H_0 is accepted if $\sqrt{\frac{(n-2)s_x}{S_R(y)}} \cdot |B(y) - \beta_0| \leq F_{t_{n-2}}^{-1}\left(1 - \frac{\gamma}{2}\right)$

The p -value is:

$$\gamma_y := 2 \left(1 - F_{t_{n-2}} \left(\sqrt{\frac{(n-2)s_x}{S_R(y)}} \cdot |B(y) - \beta_0| \right) \right)$$

(7.9) Confidence interval for $\alpha + \beta x_0$

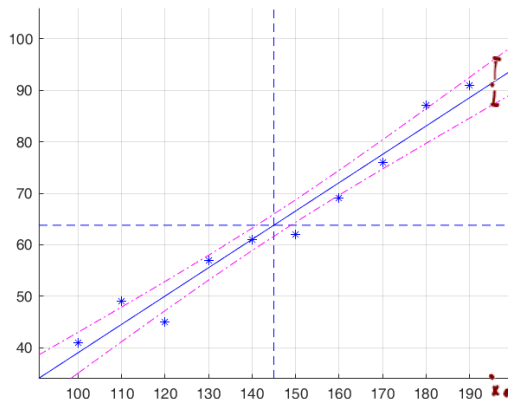
Given an input value x_0 and sampled data values $y = (y_1, \dots, y_n)$ from (Y_1, \dots, Y_n) , we have with a confidence of $100(1 - \gamma)\%$:

$$\alpha + \beta x_0 = A(y) + B(y) x_0 \pm \sqrt{\frac{S_R(y) \cdot (s_x + n \underline{x_0 - \bar{x}}^2)}{n(n-2) s_x}} \cdot F_{t_{n-2}}^{-1} \left(1 - \frac{\gamma}{2} \right)$$

Example

The following table contains 10 data pairs relating the yield of a laboratory experiment y_i to the temperature x_i at which the experiment was run.

i	x_i	y_i
1	100	41
2	110	49
3	120	45
4	130	57
5	140	61
6	150	62
7	160	69
8	170	76
9	180	87
10	190	91



conf. level
of 80%

(7.10) Prediction interval for a response at input value x_0

For a given input value x_0 and sampled data values $y = (y_1, \dots, y_n)$ from (Y_1, \dots, Y_n) a $100(1 - \gamma)$ percent confidence interval for the response value is given by:

$$A(y) + B(y)x_0 \pm \sqrt{\frac{S_R(y) \cdot ((n+1)s_x + n(\underline{x_0 - \bar{x}})^2)}{n(n-2)s_x}} \cdot F_{t_{n-2}}^{-1} \left(1 - \frac{\gamma}{2} \right)$$

Example

The following table contains 10 data pairs relating the yield of a laboratory experiment y_i to the temperature x_i at which the experiment was run.

i	x_i	y_i
1	100	41
2	110	49
3	120	45
4	130	57
5	140	61
6	150	62
7	160	69
8	170	76
9	180	87
10	190	91

