

## 2. Regression Analysis

Data set :  $(y_i, x_i) \quad i=1, \dots, 12$

- In an experiment to investigate the variation of the specific heat ( $Y$ ) of a certain chemical with the temperature ( $x$ ), two specific heat measurements were made for each of the 6 temperatures under study.

Temperature ( $^{\circ}C$ )	50	60	70	80	90	100
Specific heat ( $^{\circ}C/g$ )	1.60	1.63	1.67	1.70	1.71	1.71
	1.64	1.65	1.67	1.72	1.72	1.74

$n=12$

- Estimate the regression equation using the least squares method.
- Obtain a 95% confidence interval on the intercept  $\beta_0$ .
- Test the hypothesis  $H_0 : \beta_1 = 0$  against the hypothesis  $H_1 : \beta_1 \neq 0$ .
- Get an estimate of the specific heat corresponding to temperature  $75^{\circ}C$ . Determine the 95% confidence interval and prediction interval for this temperature value.

a) Model:  $\underline{y} = \underline{X} \underline{\beta} + \underline{\varepsilon}$

$$\hat{E}[Y | \underline{x}] = \hat{\mu}_Y | \underline{x} = \hat{\beta}_0 + \hat{\beta}_1 x = \underline{x}^T \hat{\underline{\beta}}$$

$$\underline{y} = \begin{bmatrix} 1.60 \\ 1.64 \\ 1.63 \\ 1.65 \\ 1.67 \\ 1.67 \\ \vdots \\ 1.71 \\ 1.74 \end{bmatrix}_{12 \times 1}$$

$$\underline{X} = \begin{bmatrix} 1 & 50 \\ 1 & 50 \\ 1 & 60 \\ 1 & 60 \\ \vdots & \vdots \\ 1 & 100 \\ 1 & 100 \end{bmatrix}_{12 \times 2}$$

$$\underline{x}^T \underline{y} = \begin{bmatrix} \sum_{i=1}^{12} y_i \\ \sum_{i=1}^{12} x_i y_i \end{bmatrix} = \begin{bmatrix} 20.16 \\ 1519.9 \end{bmatrix}$$

$$\hat{\underline{\beta}} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$$

$$\underline{C}^{-1} = \underline{X}^T \underline{X} = \begin{bmatrix} n & \sum_{i=1}^{12} x_i \\ \sum_{i=1}^{12} x_i & \sum_{i=1}^{12} x_i^2 \end{bmatrix} = \begin{bmatrix} 12 & 900 \\ 900 & 71000 \end{bmatrix}$$

$$\tilde{C} = (\tilde{X}^T \tilde{X})^{-1} = \frac{1}{12 \times 71000 - 900^2} \begin{bmatrix} 71000 & -900 \\ -900 & 12 \end{bmatrix}$$

$$= \begin{bmatrix} \frac{71}{42} & \frac{-3}{140} \\ \frac{-3}{140} & \frac{1}{3500} \end{bmatrix}$$

$$\hat{\beta} = \begin{matrix} \underbrace{\tilde{C}}_{2 \times 2} \underbrace{\tilde{X}^T \tilde{y}}_{2 \times 1} \end{matrix} = \begin{bmatrix} \frac{71}{42} & \frac{-3}{140} \\ \frac{-3}{140} & \frac{1}{3500} \end{bmatrix} \begin{bmatrix} 20.16 \\ 1519.9 \end{bmatrix} = \begin{bmatrix} 1.5107 \\ 0.002257 \end{bmatrix}$$

$$\hat{y} = \underbrace{1.5107}_{\hat{\beta}_0} + \underbrace{0.002257}_{\hat{\beta}_1} x, \quad \forall x \in (50; 100)$$

(b) Obtain a 95% confidence interval on the intercept  $\beta_0$ .

b)  $CI_{95\%}(\beta_0) = ?$

Pivotal quantity:  $T = \frac{\hat{\beta}_0 - \beta_0}{\sqrt{\hat{\sigma}^2 c_{1,1}}} \sim t_{(12-2)} = t_{(10)}$   
 $n=12, p=2$

where,  $\hat{\sigma}^2 = \frac{SSE}{n-2} = \frac{SST-SSR}{n-2}$

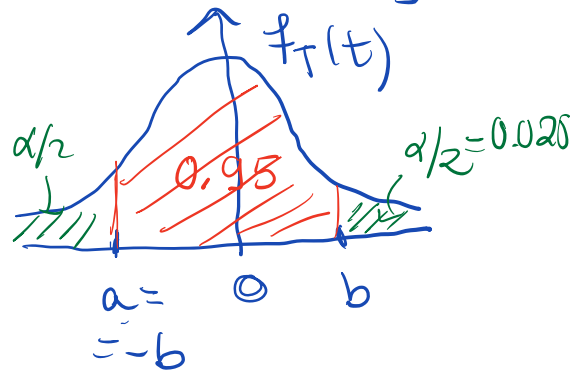
and  $\tilde{C} = \begin{bmatrix} \frac{71}{42} & \frac{-3}{140} \\ \frac{-3}{140} & \frac{1}{3500} \end{bmatrix} \rightarrow c_{1,1} = \frac{71}{42}$

Derivation of the C.I. (Random)

$$P(\underbrace{A \leq \beta_0 \leq B}_{\text{Random variables}}) = 0.95$$

to get this we take the Pivotal Quantity  
and we find the constants a and b :

$P(a \leq T \leq b) = 0.95$   
and taking the  $1-\alpha$   
symmetrical interval  
around 0  $\Rightarrow a = -b$



So,  $b = F_{t(10)}^{-1}(0.975) = t_{0.975}(10)$  (97.5% quantile of the  $t(10)$ )

$$P(-t_{0.975}(10) < \frac{\hat{\beta}_0 - \beta_0}{\sqrt{\hat{\sigma}^2 c_{11}}} < t_{0.975}(10)) = 0.95 \Leftrightarrow$$

$$P(-t_{0.975}(10) \times \sqrt{\hat{\sigma}^2 c_{11}} \leq \hat{\beta}_0 - \beta_0 \leq t_{0.975}(10) \times \sqrt{\hat{\sigma}^2 c_{11}}) = 0.95$$

$\Leftrightarrow$

( $\Rightarrow$ )

$$P(-\hat{\beta}_0 - t_{0.975(10)} \sqrt{\hat{\sigma}^2 c_{11}} \leq -\beta_0 \leq -\hat{\beta}_0 + t_{0.975(10)} \sqrt{\hat{\sigma}^2 c_{11}}) = 0.95$$

Multiplying by (-1), we obtain:

$$P(\underbrace{\hat{\beta}_0}_{\text{Random variable}} - \underbrace{t_{0.975(10)} \sqrt{\hat{\sigma}^2 c_{11}}}_A \leq \beta_0 \leq \underbrace{\hat{\beta}_0}_{\text{Random variable}} + \underbrace{t_{0.975(10)} \sqrt{\hat{\sigma}^2 c_{11}}}_B) = 0.95$$

random variable                      random variable                      random variable                      random variable

Concretization of the Random C.I.:

$$\hat{\beta}_0 = 1.5107 ; \quad t_{0.975(10)} = 2.2281 ; \quad c_{11} = \frac{71}{42}$$

$$\hat{\sigma}^2 = \frac{SSE}{n-2} = \frac{SST - SSR}{n-2} \quad SST = \sum y^2 - n\bar{y}^2$$

$$SST - SSR = \sum y^2 - \beta^T \sum x^T y = 33.8894 - [1.5107 \quad 0.002257] \begin{bmatrix} 20.16 \\ 1518.9 \end{bmatrix}$$

$$\sum y^2 = \sum_{i=1}^{12} y_i^2 = 33.8894 \quad = 0.002769$$

$$SSE = \left( \sum x^T y - n\bar{y}^2 \right)$$

$$\hat{\sigma}^2 = MSE = \frac{0.002769}{10} = 0.0002769$$

$$CI_{95\%}(\beta_0) = \left[ \hat{\beta}_0 \pm t_{0.975(10)} \sqrt{\hat{\sigma}^2 c_{11}} \right]$$

$$\therefore eI_{95\%}(\beta_0) = \left[ \underbrace{1.5107}_{\hat{\beta}_0} \pm 2.2281 \sqrt{\frac{0.0002769 \times 71}{42}} \right]$$



$$= [1.4625; 1.5589]$$

mean of C.I. =  $\hat{\beta}_0 = 1.5107$

(c) Test the hypothesis  $H_0: \beta_1 = 0$  against the hypothesis  $H_1: \beta_1 \neq 0$ .

very important test (see if the regression is significant)

if  $\beta_1 = 0$  then  $E(Y|x) = \beta_0 + \beta_1 x$   $x$  is important to explain  $E(Y)$ ?

$\beta_1 < 0 \vee \beta_1 > 0$

$H_0: \beta_1 = 0$  vs  $H_1: \beta_1 \neq 0$  (two sided test)

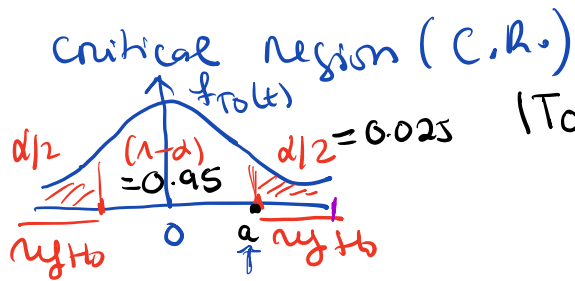
Pivotal quantity:  $T = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\hat{\sigma}^2 c_{2,2}}} \sim t_{(10)}$   $K=1$

Under  $H_0$ , we have the test statistic:

$(H_0 \text{ true})$   
 $T_{|H_0} = T_0 = \frac{\hat{\beta}_1 - 0}{\sqrt{\hat{\sigma}^2 c_{2,2}}} \sim t_{(10)}$

Assuming that  $\alpha = 0.05$

Note that  $\alpha = P(\text{'type I error'}) = P(\text{'reject } H_0 | H_0 \text{ is true}) = 0.05$ , we can obtain the critical region (where we reject  $H_0$ )



$$|T_0| > a \text{ and } P(|T_0| > a) = 0.05$$

$$\Rightarrow a = F_{t(10)}^{-1}(0.975) = t_{0.975}(10) = 2.2281$$

$$CR = ]-\infty, -2.2281[ \cup ]2.2281, +\infty[$$

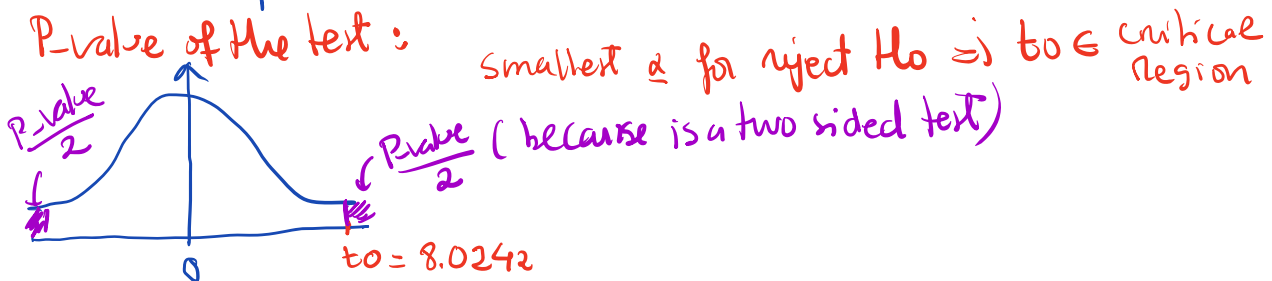
0.05

(critical region for  $\alpha = 0.05$ )

Observed value of the test statistic:

$$t_0 = \frac{0.002257}{\sqrt{\frac{0.0002769}{3500}}} \approx 8.0242 \in CR. \Rightarrow \text{ry } H_0 \text{ for } \alpha = 0.05,$$

so, for this dataset it seems that the temperature ( $x$ ) it is important to explain the Expected value of heat ( $Y$ )



$$P\text{-value} = P(|T_0| \geq t_0) = 2P(T_0 \geq 8.0242) = 2(1 - t_{t(10)}(8.0242))$$

$$= 2(1 - 0.9999943) = 1.1467 \times 10^{-5} \Rightarrow \text{ry } H_0 \forall \alpha \left. \begin{array}{l} 1\% \\ 5\% \\ \text{usual } 10\% \end{array} \right\}$$

- (d) Get an estimate of the specific heat corresponding to temperature  $75^\circ\text{C}$ . Determine the 95% confidence interval and prediction interval for this temperature value.

$$x = 75^\circ\text{C} \quad x \in [50, 100]$$

$$\underline{x}_0 = \begin{bmatrix} 1 \\ 75 \end{bmatrix}$$

$$\hat{E}[Y | \underline{x}_0] = \hat{\mu}_{Y | \underline{x}_0} = 1.5107 + 0.002257 \times 75 = 1.68$$

$$CI_{95\%}(\mu | \underline{x}_0) = ?$$

$$\text{Pivotal Quantity: } T = \frac{\hat{\mu}_{Y | \underline{x}_0} - \mu_{Y | \underline{x}_0}}{\sqrt{\hat{\sigma}^2 \underline{x}_0^T (\underline{X}^T \underline{X})^{-1} \underline{x}_0}} \sim t(10)$$

$$CI_{95\%}(\mu_{Y | \underline{x}_0}) = \left[ \hat{\mu}_{Y | \underline{x}_0} \pm t_{0.975}(10) \sqrt{\hat{\sigma}^2 \underline{x}_0^T (\underline{X}^T \underline{X})^{-1} \underline{x}_0} \right]$$

$$t_{0.975}(10) = 2.2281$$

$$\underline{x}_0^T (\underline{X}^T \underline{X})^{-1} \underline{x}_0 = (1 \ 75) \begin{bmatrix} \frac{71}{42} & -\frac{3}{140} \\ -\frac{3}{140} & \frac{1}{3500} \end{bmatrix} \begin{bmatrix} 1 \\ 75 \end{bmatrix} =$$

$$= \left[ \frac{71}{42} - \frac{3 \times 75}{140}, \frac{-3}{140} + \frac{75}{3500} \right] \begin{bmatrix} 1 \\ 75 \end{bmatrix} \approx 0.083333$$

$$\hat{\sigma}^2 \underline{x}_0^T (\underline{X}^T \underline{X})^{-1} \underline{x}_0 = 0.0002769 \times 0.083333 \approx 2.3075 \times 10^{-5}$$

$$\sqrt{\hat{\sigma}^2 \underline{x}_0^T (\underline{X}^T \underline{X})^{-1} \underline{x}_0} \approx 0.0048$$

$$CI_{95\%}(\mu_{Y | \underline{x}_0}) = [1.68 \pm 2.2281 \times 0.0048] \\ = [1.669 ; 1.691]$$

$$PI_{95\%}(Y_0) = ?$$

$$\text{Pivotal quantity: } T = \frac{\hat{Y}_0 - Y_0}{\sqrt{\hat{\sigma}^2 (1 + \underline{x}_0^T (\underline{X}^T \underline{X})^{-1} \underline{x}_0)}} \sim t_{(10)}$$

$$PI_{95\%}(Y_0) = \left[ \hat{Y}_0 \pm t_{0.975(10)} \sqrt{\hat{\sigma}^2 (1 + \underline{x}_0^T (\underline{X}^T \underline{X})^{-1} \underline{x}_0)} \right]$$

$$\hat{Y}_0 = \hat{E}[Y | \underline{x}_0] = 1.68$$

$$t_{0.975(10)} = 2.2281$$

$$\sqrt{\hat{\sigma}^2 (1 + \underline{x}_0^T (\underline{X}^T \underline{X})^{-1} \underline{x}_0)} = \sqrt{0.0002769 + 2.3075 \times 10^{-5}} \\ \approx 0.01732$$

$$PI_{95\%}(Y_0) = [1.68 \pm 2.2281 \times 0.01732]$$

$$= [1.641; 1.719] \quad \text{more wide than the C.I.}$$