

2 Oct - Lec 2

- Descriptive properties of the estimated linear model

1) the estimated regression line passes for the point (\bar{x}, \bar{y})

i.e. $\bar{y} = \hat{\beta}_1 + \hat{\beta}_2 \bar{x}$

compute \hat{y} at \bar{x} : $\hat{y} = \hat{\beta}_1 + \hat{\beta}_2 \bar{x} = \bar{y} - \cancel{\hat{\beta}_2 \bar{x}} + \cancel{\hat{\beta}_2 \bar{x}} = \bar{y}$

2) the mean of the response at the observed locations (x_1, \dots, x_n) is equal to the mean of the predicted values at those locations

$$\frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n \hat{y}_i \quad (\bar{y} = \bar{\hat{y}})$$

$$\begin{aligned} \bar{y} &= \frac{1}{n} \sum_{i=1}^n \hat{y}_i = \frac{1}{n} \sum_{i=1}^n (\hat{\beta}_1 + \hat{\beta}_2 x_i) = \frac{1}{n} \cancel{n \hat{\beta}_1} + \hat{\beta}_2 \bar{x} = \hat{\beta}_1 + \hat{\beta}_2 \bar{x} \\ &= \underbrace{\bar{y} - \hat{\beta}_2 \bar{x}}_{\hat{\beta}_1} + \hat{\beta}_2 \bar{x} = \bar{y} \end{aligned}$$

3) the sample mean of the residuals is equal to zero

i.e. $\bar{e} = \frac{1}{n} \sum_{i=1}^n e_i = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) = 0$

$$\bar{e} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_i) = \bar{y} - \hat{\beta}_1 - \hat{\beta}_2 \bar{x} = \bar{y} - \bar{y} + \hat{\beta}_2 \bar{x} - \hat{\beta}_2 \bar{x} = 0$$

- Inferential properties of the estimated linear model

(now we need those assumptions on the errors)

$\hat{\beta}_1$ and $\hat{\beta}_2$ are the estimates (they are not random variables, they are numbers).

we study the properties of the corresponding estimators $\hat{\beta}_1 = \hat{\beta}_1(Y)$, $\hat{\beta}_2 = \hat{\beta}_2(Y)$

(the random variable is $Y = (Y_1, \dots, Y_n)$, and the estimators $\hat{\beta}_1$ and $\hat{\beta}_2$ a transformation of Y .)

$$\begin{aligned} \hat{\beta}_1 &= \bar{Y} - \hat{\beta}_2 \bar{x} \\ \hat{\beta}_2 &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{aligned}$$

EXPECTED VALUE AND VARIANCE

- Let's start with $\hat{\beta}_2$

$$\begin{aligned} \hat{\beta}_2 &= \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \left\{ \sum_{i=1}^n (x_i - \bar{x}) Y_i - \bar{Y} \sum_{i=1}^n (x_i - \bar{x}) \right\} \\ &\quad = 0 \text{ since } \sum_{i=1}^n (x_i - \bar{x}) = \sum_{i=1}^n x_i - n\bar{x} \\ &\quad = n\bar{x} - n\bar{x} = 0 \\ &= \sum_{i=1}^n \underbrace{\frac{(x_i - \bar{x})}{\sum_{h=1}^n (x_h - \bar{x})^2}}_{w_i} \cdot Y_i \end{aligned}$$

w_i is a linear combination of Y_1, \dots, Y_n

$$\begin{aligned} E[\hat{\beta}_2] &= E\left[\sum_{i=1}^n w_i Y_i\right] = \sum_{i=1}^n w_i E[Y_i] = \sum_{i=1}^n w_i (\beta_1 + \beta_2 x_i) = \\ &= \underbrace{\beta_1 \sum_{i=1}^n w_i}_{A} + \underbrace{\beta_2 \sum_{i=1}^n w_i x_i}_{B} = \beta_2 \end{aligned}$$

$$A = \sum_{i=1}^n w_i = \sum_{i=1}^n \frac{x_i - \bar{x}}{\sum_{h=0}^n (x_h - \bar{x})^2} = 0$$

$$B = \sum_{i=1}^n w_i x_i = \sum_{i=1}^n \frac{x_i (x_i - \bar{x})}{\sum_{h=0}^n (x_h - \bar{x})^2} = \frac{\sum_{i=1}^n x_i^2 - \bar{x} \sum_{i=1}^n x_i}{\sum_{h=0}^n (x_h - \bar{x})^2} = \frac{\sum_{i=1}^n x_i^2 - n\bar{x}^2}{\sum_{h=0}^n (x_h - \bar{x})^2} = 1$$

$$\text{since } \text{cov}(Y_i, Y_k) = 0$$

$$\text{var}(\hat{\beta}_2) = \text{var}\left(\sum_{i=1}^n w_i Y_i\right) = \sum_{i=1}^n w_i^2 \text{var}(Y_i) = \sum_{i=1}^n w_i^2 \sigma^2 =$$

$$= \sigma^2 \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\left[\sum_{h=0}^n (x_h - \bar{x})^2\right]^2} = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

REMARKS

- $\hat{\beta}_1$ and $\hat{\beta}_2$ are UNBIASED estimators (i.e. $E[\hat{\beta}_1] = \beta_1$; $E[\hat{\beta}_2] = \beta_2$)

- if σ^2 increases $\Rightarrow \text{var}(\hat{\beta}_1)$ and $\text{var}(\hat{\beta}_2)$ increase (fixing the rest)

- if $\sum(x_i - \bar{x})^2$ increases $\Rightarrow \text{var}(\hat{\beta}_1)$ and $\text{var}(\hat{\beta}_2)$ decrease

↓
I can estimate the line better

if the x_i 's are well spread out

- $\text{var}(\hat{\beta}_1)$ and $\text{var}(\hat{\beta}_2)$ depend on σ^2 (unknown) \rightarrow can we estimate it?

ESTIMATE of σ^2

Recall that σ^2 is the (common) variance of the errors e_1, \dots, e_n . However, these are not observable quantities.

The corresponding sample quantities (observable) are the RESIDUALS $e_i = y_i - \hat{y}_i$, $i=1, \dots, n$.

(note: they are not an estimate of the errors)

Idea to estimate σ^2 : we estimate it using the sample variance of the residuals, i.e.

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (e_i - \bar{e})^2$$

we have seen that $\bar{e} = 0 \Rightarrow \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n e_i^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_i)^2$

We can consider the corresponding estimator $\hat{\sigma}^2$ to study its properties

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_i)^2$$

It can be shown that $E[\hat{\sigma}^2] = \frac{n-2}{n} \sigma^2$ it is a BIASED estimator of σ^2

- if n is large, the bias is small

indeed, $\lim_{n \rightarrow \infty} E[\hat{\sigma}^2] = \sigma^2$ asymptotically unbiased

- we can define an unbiased version $S^2 = \frac{n}{n-2} \hat{\sigma}^2 = \frac{n}{n-2} \cdot \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_i)^2$

$$= \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_i)^2$$

- ④ once we compute the estimate of σ^2 , S^2 ,

we can plug it into $\text{var}(\hat{\beta}_1)$ and $\text{var}(\hat{\beta}_2)$ to obtain an estimate of these quantities

$$\hat{\text{var}}(\hat{\beta}_1) = S^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)$$

$$\hat{\text{var}}(\hat{\beta}_2) = \frac{S^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$