

STAT-1302; Lecture 17; π -day, 2024

Ch. 13

The Simple Linear Regression Model

$$Y_i = A + BX_i + \epsilon_i, \quad i=1, \dots, n$$

where $\epsilon_i \stackrel{\text{independent}}{\sim} N(0, \sigma^2)$

ϵ_i → Epsilon

Parameters:

A = intercept

B = Slope

σ^2 = noise or error variance

Notes:

1. X_i (independent variable) is a fixed variable (i.e. no randomness). It is part of the data collected.
2. Y_i (dependent or response variable) is a

random variable. It is also part of the data collected.

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

3. ϵ_i is a random variable that is meant to capture any variation in Y_i that has not been captured by x_i .

Interpretation of B:

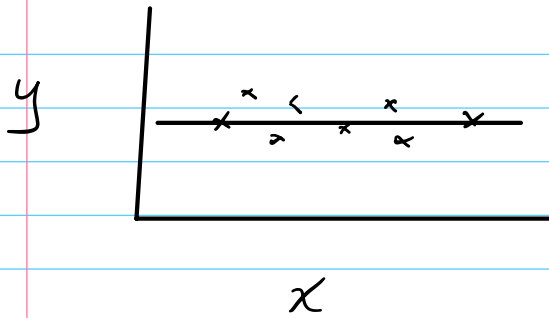
$B > 0 \Rightarrow$ a positive linear relationship between X and Y .

(ie. as X increases, then so does y .)

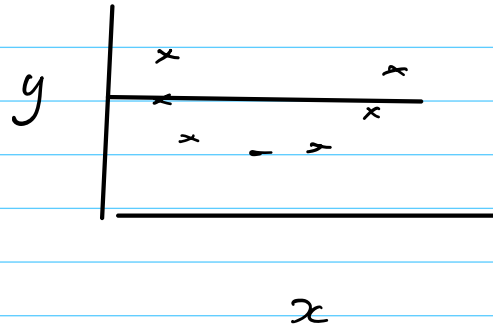
$B < 0 \Rightarrow$ a negative linear relationship between X and Y .

(ie. as X increases, then y decreases.)

$B = 0 \Rightarrow$ no linear relationship between X and y .



$$B=0$$



$$B=0$$

no linear relationship between X and y

More on Interpretation of B :

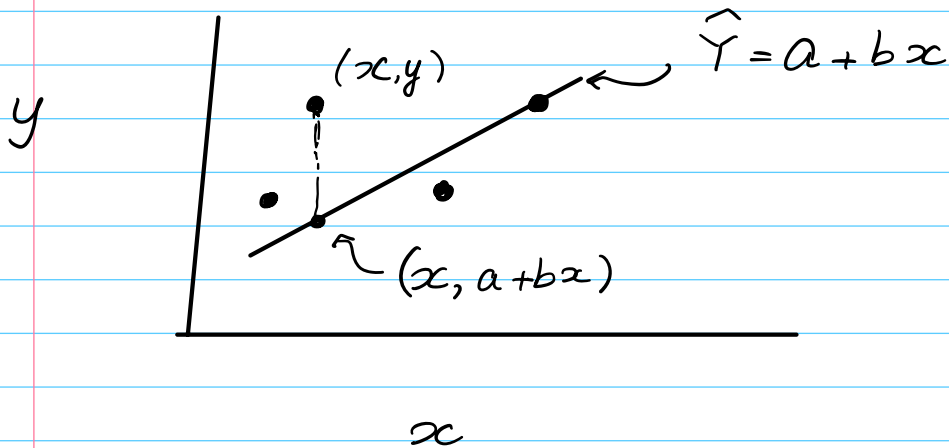
MCQ

- B is the change in the mean of Y when X is increased by one unit.

Fitted/Estimate Simple linear regression model:

$$\hat{Y} = a + bX$$

where "a" and "b" are the least squares estimators of A and B , respectively.



The least squares regression line is found by

minimizing $\sum_{i=1}^n (Y_i - (a + bX_i))^2$ with respect to a and b . (Calculus Problem).

$$b = \frac{\sum_{i=1}^n (Y_i - \bar{Y})(X_i - \bar{X})}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

$$a = \bar{Y} - b\bar{X}. \quad \leftarrow \text{Given on formula sheet.}$$

To compute "b" use the short-cut formula given on formula sheet:

$$b = \frac{SS_{xy}}{SS_{xx}} \quad \text{where} \quad \text{Given}$$

$$SS_{xx} = \sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i\right)^2}{n}$$

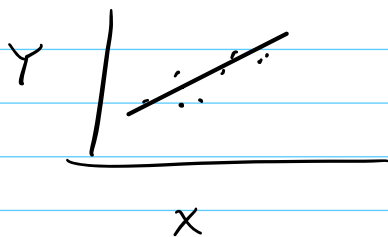
$$SS_{xy} = \sum_{i=1}^n x_i y_i - \frac{\left(\sum_{i=1}^n x_i\right) \left(\sum_{i=1}^n y_i\right)}{n}$$

n = total no. of observations :

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

Ex. Y = food expenditure

X = monthly income (in hundreds of dollars)



Summary Statistics:

$$\sum x = 212, \quad \sum y = 64, \quad \sum xy = 2150, \quad \sum x^2 = 7222.$$

$$n = 7$$

a) Find the least squares regression line where food expenditure (y) is regressed against monthly income (x).

$$b = \frac{SS_{xy}}{SS_{xx}}$$

$$SS_{xy} = \sum_{i=1}^7 x_i y_i - \frac{\left(\sum_{i=1}^7 x_i\right) \left(\sum_{i=1}^7 y_i\right)}{7}$$

$$= 2150 - \frac{212(64)}{7} = 211.7143$$

$$SS_{xx} = \sum x^2 - \frac{\left(\sum x\right)^2}{7} \quad n=7$$

$$= 7222 - \frac{212^2}{7} = 801.4286$$

$$b = \frac{SS_{xy}}{SS_{xx}} = \frac{211.7143}{801.4286} = 0.2642$$

$$a = \bar{Y} - b\bar{X} = \frac{\sum y}{7} - 0.2642 \times \frac{\sum x}{7}$$

$$= \frac{64}{7} - 0.2642 \times \frac{212}{7} = 1.1414$$

The estimated least squares regression line is

$$\hat{Y} = 1.1414 + 0.2642 \times X$$

must write " $\hat{\quad}$ " in least squares regression line

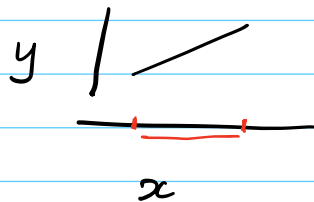
Or

$$\hat{\text{Expenditure}} = 1.1414 + 0.2642 \times \text{Income}$$

b) What is the predicted monthly food expenditure for a household whose monthly income is 10 (i.e. \$1,000) ?

$$\begin{aligned}\hat{\text{Expenditure}} &= 1.1414 + 0.2642 \times 10 \\ &= 3.7834\end{aligned}$$

i.e. \$378.34. (x is in hundreds of dollars).



Residual:

$$e_i = \underset{\substack{\downarrow \\ \text{Observed response}}}{Y_i} - \underset{\substack{\downarrow \\ \text{predicted response} \\ \text{for regression line.}}}{\hat{Y}_i}$$

Ex. c) What is the residual when $X=35$ and $y=9$?

$$\hat{Y} = 1.1414 + 0.2642(35) = 10.3884$$

$$e = Y - \hat{Y} = 9 - 10.3884 = -1.3884.$$

That is, the expenditure is overestimated by \$138.84 when monthly income is \$3500.

[Extra: STAT-3103.

We use e_i 's to assess the adequacy to the fitted regression line.

$$\sum e_i \approx 0 \Rightarrow E(e_i) = 0 \text{ assumption holds.}]$$

Interpretation of a :

" a " is the mean of the response variable when $X=0$.

Assumptions of the Simple linear regression model:

$$Y = A + BX + E$$

1. E is a random variable that has mean zero at each x .
2. E_1, \dots, E_n , the errors associated with each of x_1, \dots, x_n , respectively are independent.
3. For any given x , the errors E_1, \dots, E_n , are normally distributed.
4. E_1, \dots, E_n have constant variance σ^2 .

$$\left[\begin{array}{l} \text{i.e. } Y_i = A + BX_i + E_i \text{ where} \\ E_i \stackrel{\text{ind}}{\sim} N(0, \sigma^2) \text{ for } i=1, \dots, n. \end{array} \right]$$

Ex. Data for understanding how the Price of a particular Car model depreciates with age are as follows.

Age	8	3	6	9	2	5	6	3
Price	45	210	100	33	267	134	109	235

Age (x): in years

Price (y): in hundreds of dollars.

a) Determine the estimated least squares regression line when price is regressed against age.

Given: $\sum xy = 4549$, $\sum x^2 = 264$, $\sum y = 1144$,
 $\sum x = 42$; $n = 8$.

$$b = \frac{SS_{xy}}{SS_{xx}}$$

$$\begin{aligned} SS_{xy} &= \sum xy - \frac{(\sum x)(\sum y)}{n} \\ &= 4549 - \frac{42(1144)}{8} = -1457 \end{aligned}$$

$$SS_{xx} = \sum x^2 - \frac{(\sum x)^2}{n} = 264 - \frac{42^2}{8} = 43.5$$

$$b = \frac{SS_{xy}}{SS_{xx}} = \frac{-1457}{43.5} = -33.49$$

$$a = \bar{Y} - b\bar{X} = \frac{1144}{8} - (-33.49)\frac{42}{8}$$
$$= 318.83$$

$$\widehat{\text{Price}} = 318.83 - 33.49 \times \text{Age}$$

b) Give a brief interpretation of "a" and "b" in part (a).