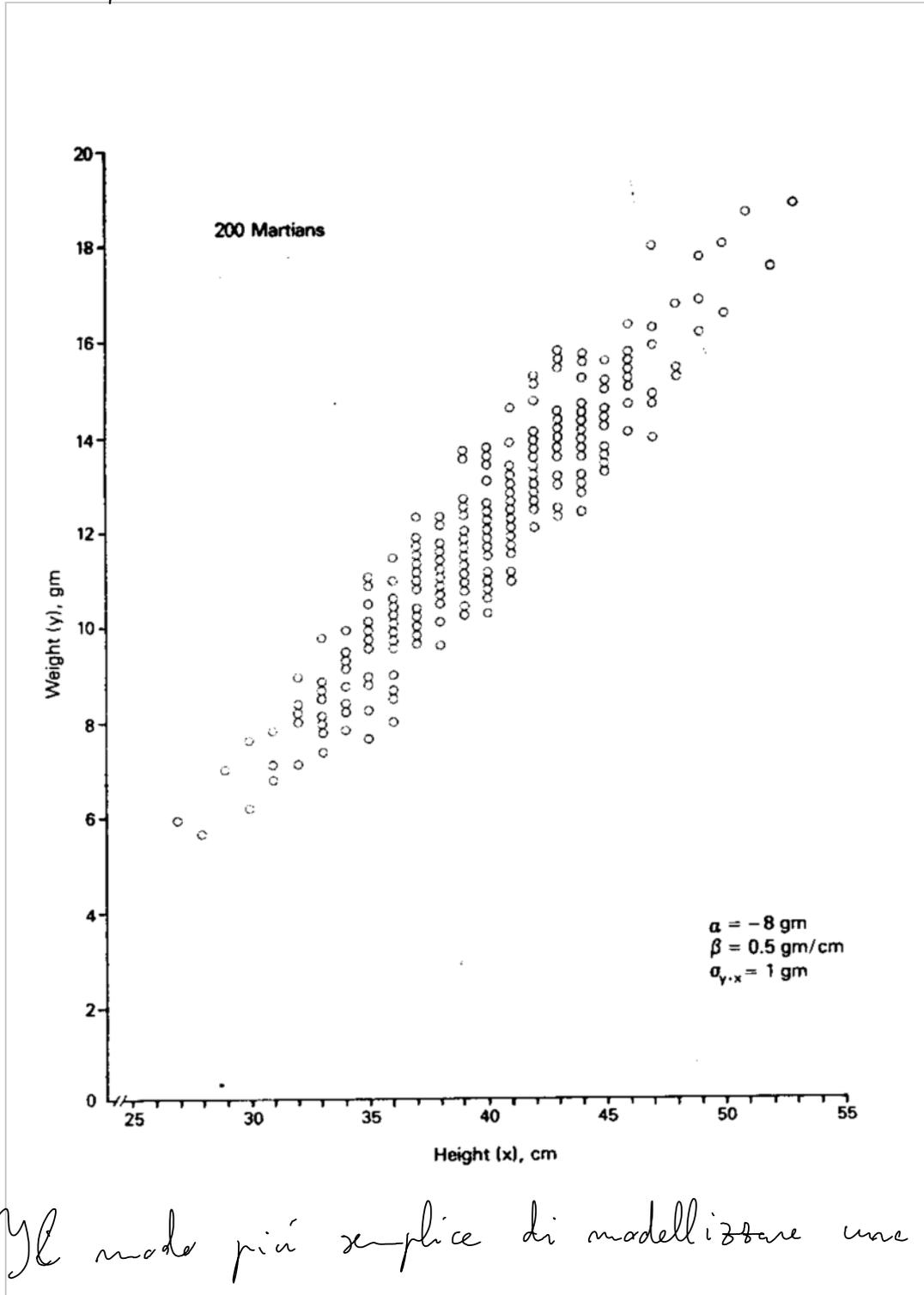


REGRESSIONE LINEARE

giovedì 5 giugno 2014 08.54

Yuppona di registrare altezze e pesi dei 200 marziani, con questi risultati:



Il modo più semplice di modellizzare una dipendenza tra altezza X_i dell' i -esimo marziano e peso Y_i è

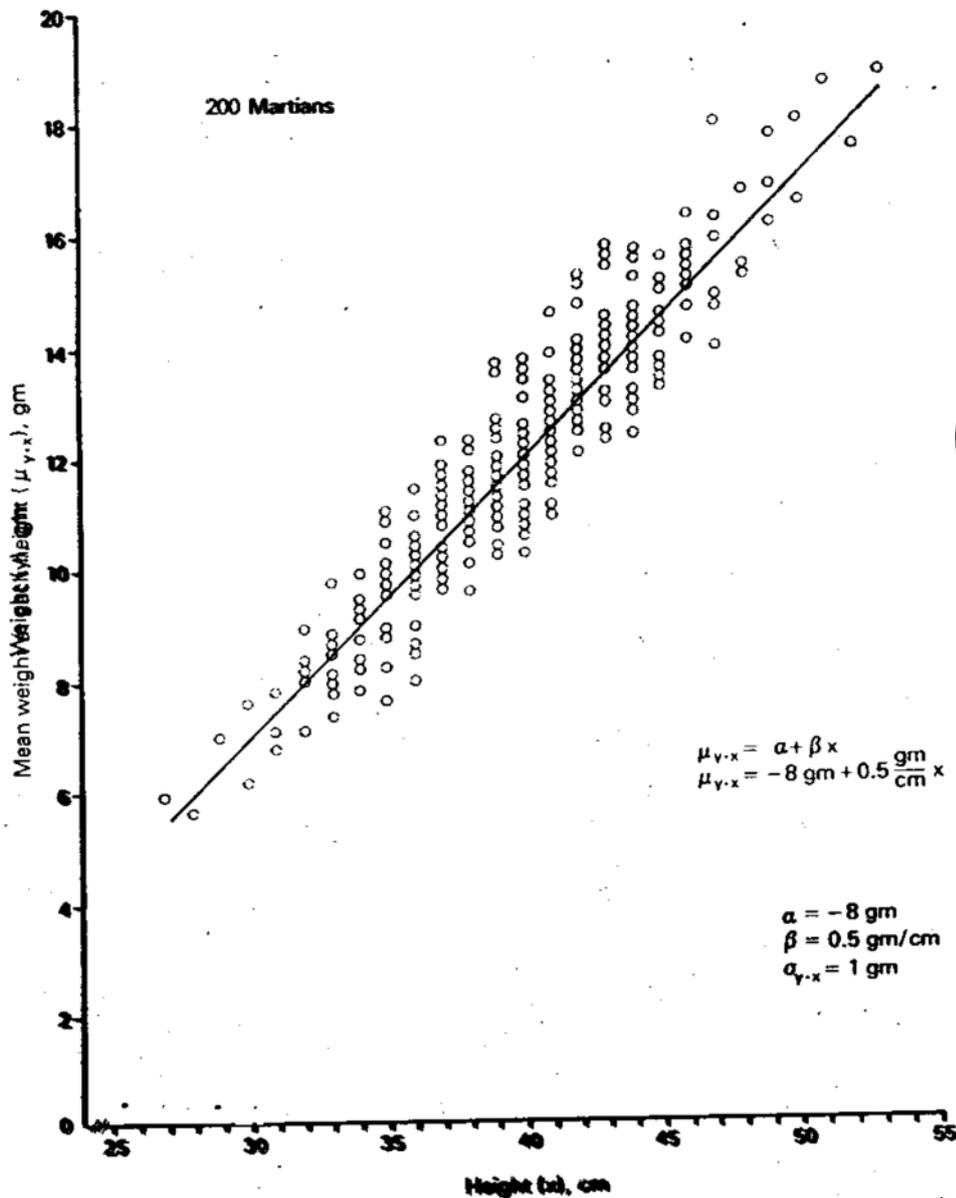
$$Y_i = \beta_1 X_i + \beta_0 + \varepsilon_i$$

con $\varepsilon_i \sim N(0, \sigma_{Y|X}^2)$

La retta $y = \beta_1 x + \beta_0$ si chiama retta delle medie

In questo modello ci sono 3 parametri: $\beta_1, \beta_0, \sigma_{Y|X}^2$

$\beta_1 =$ coeff. angolare, $\beta_0 =$ ordinata all'origine



$$\beta_1 = 0,5 \text{ g/cm}$$

$$\beta_0 = -8 \text{ g}$$

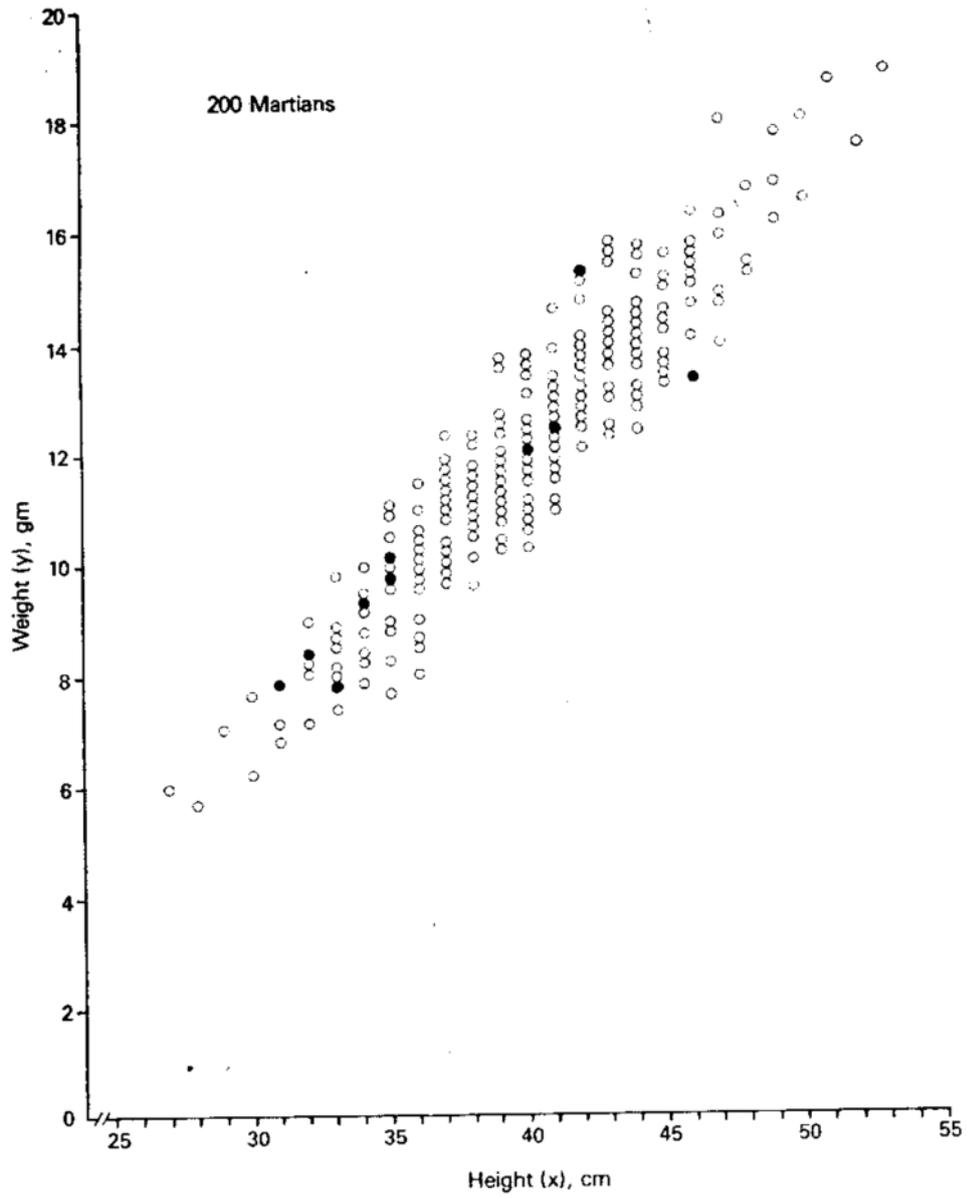
Supponiamo ora di avere un campione (X_i, Y_i) , $i = 1, \dots, n$

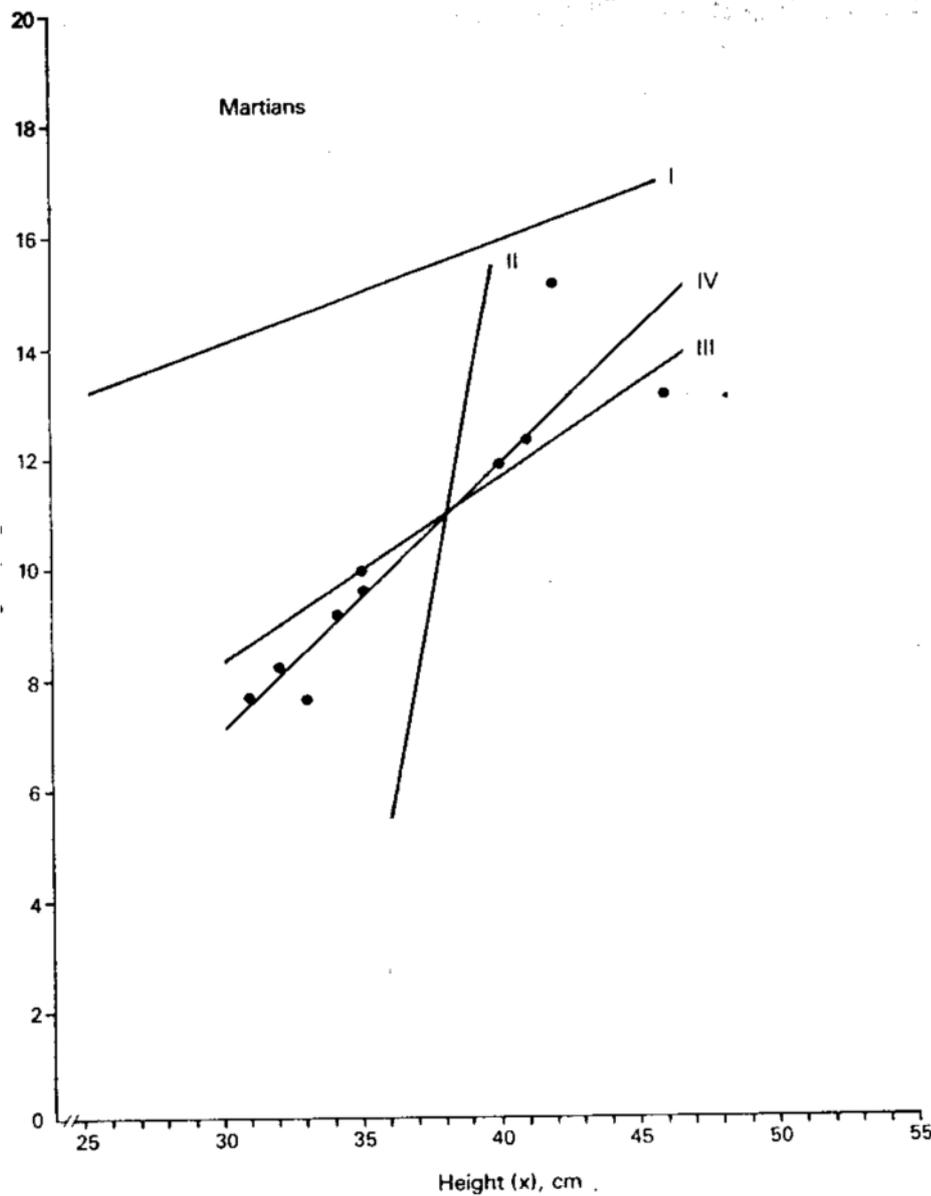
con $Y_i = \beta_1 X_i + \beta_0 + \epsilon_i$, con ϵ_i i.i.d. $\sim N(0, \sigma_{Y|X}^2)$

$(\sigma_{Y|X}^2 = \underline{\underline{\text{varianza del rumore}}})$

Vogliamo stimare i parametri β_1 , β_0 e $\sigma_{Y|X}^2$

A





Per decidere quali sono gli stimatori "migliori", prendiamo quelli che minimizzano il quadrato delle distanze delle y :

$$f(b_0, b_1) = \sum_{i=1}^n \left(\underset{\substack{\uparrow \\ \text{dei dati}}}{Y_i} - \underbrace{b_1 X_i + b_0}_{\substack{\text{della retta delle medie}}} \right)^2$$

Gli stimatori migliori sono i b_0, b_1 che minimizzano f :

$$\sum X_i Y_i = n \bar{X} \bar{Y}$$

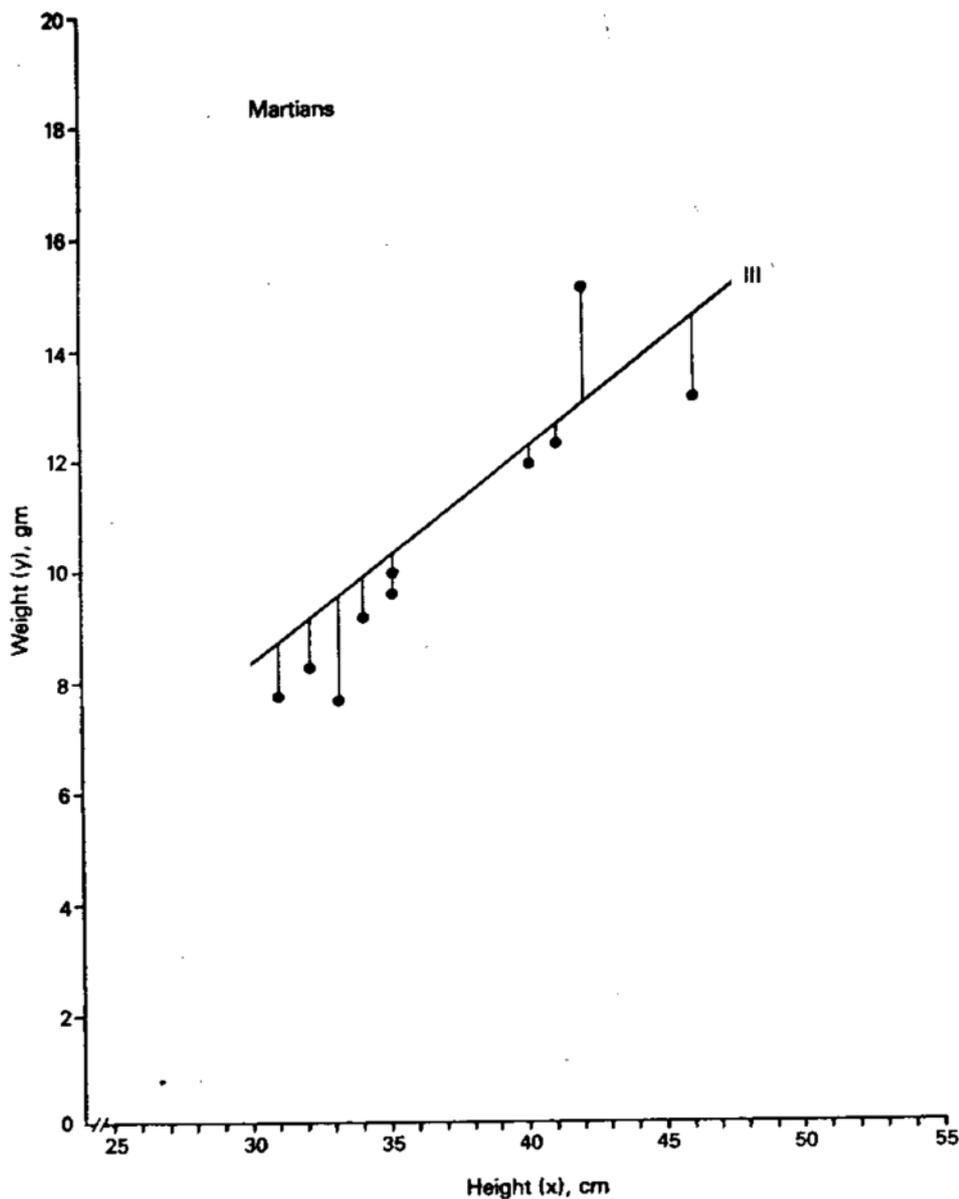
$$b_1 = \frac{\sum X_i Y_i - n \bar{X} \bar{Y}}{\sum X_i^2 - n \bar{X}^2}, \quad b_0 = \bar{Y} - b_1 \bar{X}$$

$$y = b_1 x + b_0 \quad \underline{\text{retta di regressione}}$$

nota: $\sum_1^n X_i^2 - n \bar{X}^2 = (n-1) S_x^2$. Allo stesso modo,

$$\text{definiamo } S_{xy} = \frac{1}{n-1} \left(\sum_1^n X_i Y_i - n \bar{X} \bar{Y} \right)$$

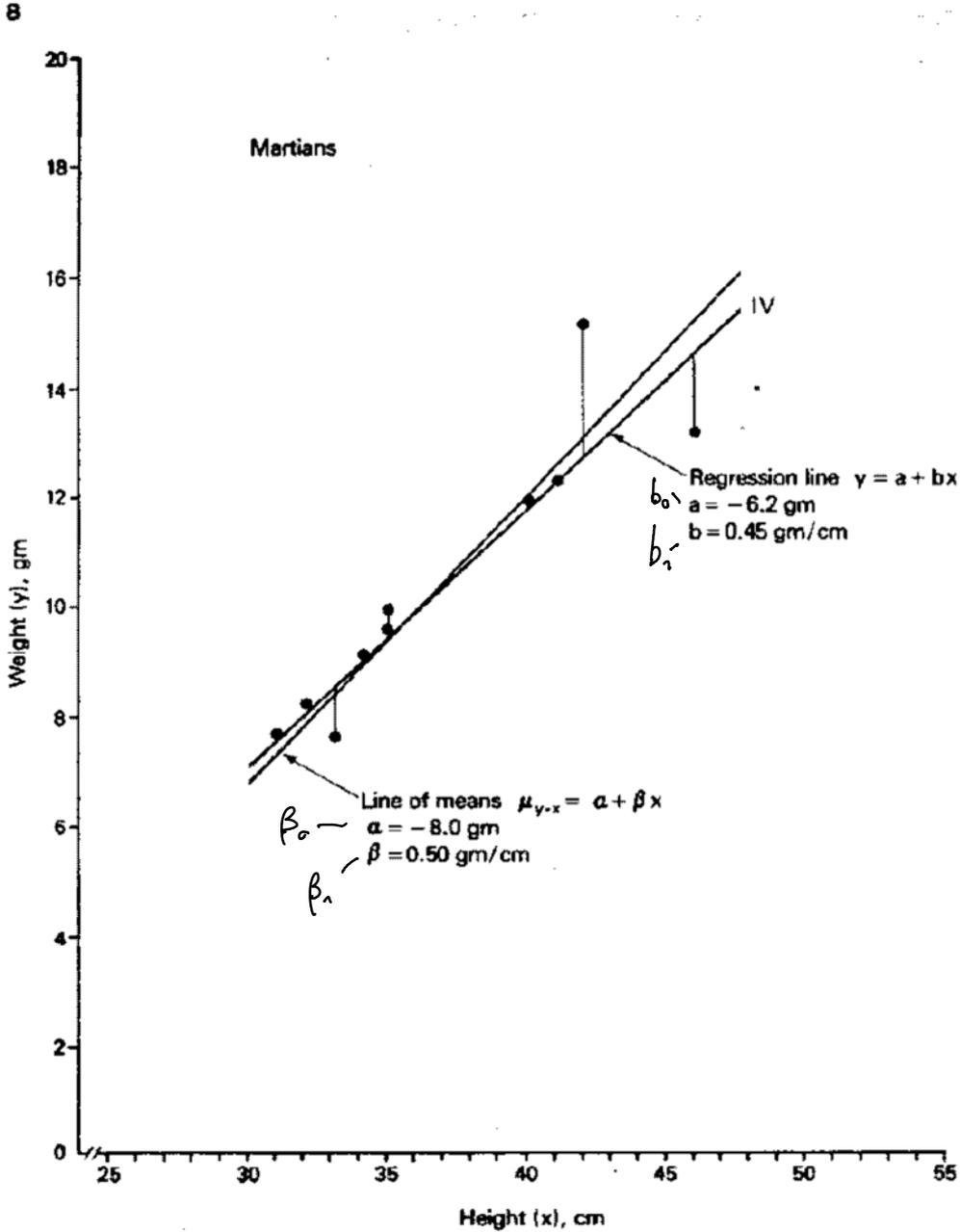
covarianza campionaria: è uno stimatore di $\text{Cov}(X, Y)$



Height (x), cm

Passiamo ricreare

$$b_1 = \frac{S_{xy}}{S_x^2}$$



stima delle varianze del rumore:

$$\sigma_{y|x}^2 = \text{Var}[\varepsilon_i] = \text{Var}[y_i - \beta_1 x_i - \beta_0]$$

$$\sigma_{Y|X}^2 = \text{Var}[\varepsilon_i] = \text{Var}[Y_i - \beta_1 X_i - \beta_0]$$

può essere stimato da

$$S_{Y|X}^2 := \frac{1}{n-2} \sum_{i=1}^n (Y_i - b_1 X_i - b_0)^2$$

che è una stima bene corretta, cioè $E[S_{Y|X}^2] = \sigma_{Y|X}^2$

Si può dimostrare che $S_{Y|X}^2 = \frac{n-1}{n-2} (s_y^2 - b_1^2 s_x^2)$

esempio (es. 9.1) rette di regressione per

| X_i | Y_i | X_i^2 | Y_i^2 | $X_i Y_i$ |
|-------|-------|---------|---------|-----------|
| 3 | 2,2 | 9 | 4,84 | 6,6 |
| 4 | 3,6 | 16 | 12,96 | 14,4 |
| 1 | 2,1 | 1 | 4,41 | 2,1 |
| 2 | 2,6 | 4 | 6,76 | 5,2 |

$n=4$

TOT. 10 | 11,0 | 30 | 31,42 | 29,8

$$\bar{X} = \frac{\sum X_i}{n} = \frac{10}{4} = 2,5$$

$$\bar{Y} = \frac{\sum Y_i}{n} = \frac{11}{4} = 2,75$$

$$S_x^2 = \frac{1}{n-1} \left(\sum X_i^2 - n \bar{X}^2 \right) = \frac{1}{3} \left(30 - 4 \cdot 2,5^2 \right) = 1,67$$

$$S_y^2 = \frac{1}{n-1} \left(\sum_{i=1}^n Y_i^2 - n \bar{Y}^2 \right) = \frac{1}{3} \left(31,42 - 4 \cdot 2,75^2 \right) = 0,39$$

$$S_{xy} = \frac{1}{n-1} \left(\sum X_i Y_i - n \bar{X} \bar{Y} \right) = \frac{1}{3} \left(29,8 - 4 \cdot 2,5 \cdot 2,75 \right) = 0,27$$

$$b_1 = \frac{S_{xy}}{S_x^2} = \frac{0,27}{1,67} = 0,46$$

$$b_0 = \bar{Y} - b_1 \bar{X} = 2,75 - 0,46 \cdot 2,5 = 1,60$$

$$s_{y|x}^2 = \frac{n-1}{n-2} (s_y^2 - b_1^2 s_x^2) = \frac{3}{2} (0,39 - 0,46^2 \cdot 1,67) = 0,056$$

$$s_{y|x} = \sqrt{s_{y|x}^2} = \sqrt{0,056} = 0,236$$