



# Modelli lineari

## REGRESSIONE LINEARE



- Molti problemi dell'ingegneria e della scienza vogliono **studiare le relazioni tra due o più insiemi di variabili**
- Ad esempio, in un processo chimico è interessante studiare la dipendenza tra quantità di catalizzatore impiegato, temperatura e rendimento → obiettivo: predire rendimento per diversi valori di temperatura e quantità di catalizzatore



Variabile di risposta  $Y$  e variabili  $x_1, \dots, x_n$  di ingresso

risposta  
dipendente

esplicative  
(indipendenti)

Il modello suppone che la risposta sia funzione degli ingressi

$$Y = h(x_1, x_2, \dots, x_n)$$



$Y$  variabile dipendente,  $x_i$  variabili indipendenti



La relazione più semplice che è possibile immaginare  
è quella lineare

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \epsilon$$

$\beta_0, \beta_1, \dots, \beta_p$  costanti

$p = \text{n}^\circ$  di variabili  
esplicative

$\epsilon$  errore casuale

o residuo

part. zero da  $p=1$

$$Y = \beta_0 + \beta_1 x_1 + \epsilon$$

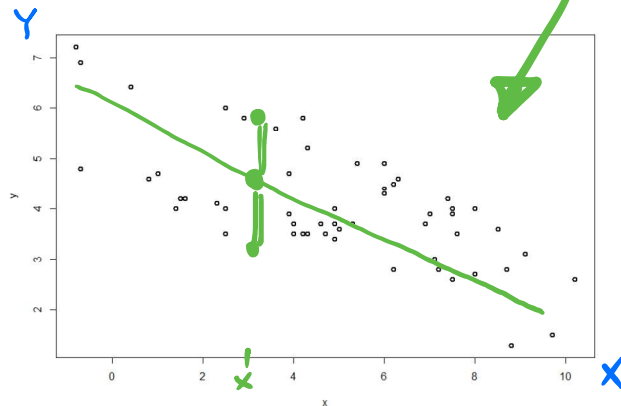


(Y, X)

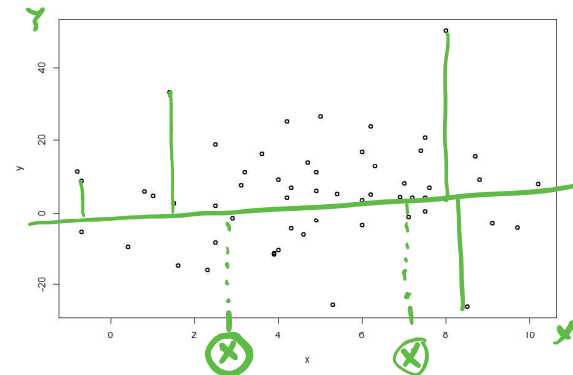
Partiamo dal caso più semplice (due variabili) e riprendiamo il concetto di correlazione lineare dalla statistica descrittiva

$$\frac{C_{XY}}{\sqrt{\sigma_x^2 \sigma_y^2}} = r = \frac{1}{n-1} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right) = \frac{\text{Cov}[X, Y]}{\sqrt{V[X] \cdot V[Y]}}$$

$r = -0.68$



$r = 0.1$



Per quale dei due esempi provereste ad impostare un modello lineare? Altre osservazioni?

## Caso $p = 1$ : regressione lineare semplice

Modello:  $Y = \beta_0 + \beta_1 X + \varepsilon$

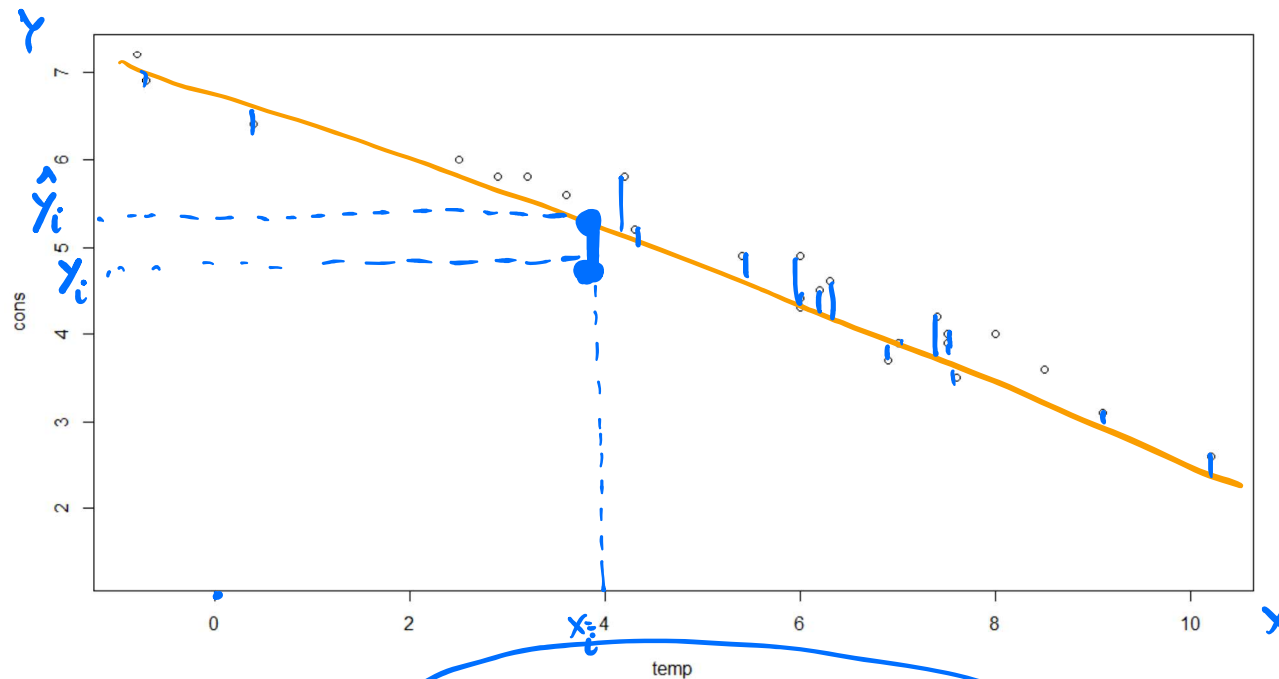
errore casuale

$\varepsilon$  è una v.c.  
 $E[\varepsilon] = 0$

Cosa prevede:  $E[Y] = \beta_0 + \beta_1 X$

Dati  $n$  punti sperimentali  $(x_i, y_i)$ , si vuole determinare una stima  $\hat{y} = b_0 + b_1 x$  della retta di regressione  $Y = \beta_0 + \beta_1 X + \varepsilon$ , minimizzando la somma dei quadrati degli errori casuali  $\varepsilon_i = Y_i - \hat{Y}_i$  tra i valori della variabile casuale  $Y$  e i valori previsti con la retta in corrispondenza del valore della variabile indipendente  $x_{i-}$

Esempio: relazione tra temperatura esterna (in gradi Celsius) X e consumo di gas Y (in ft<sup>3</sup>) per il riscaldamento di un appartamento.



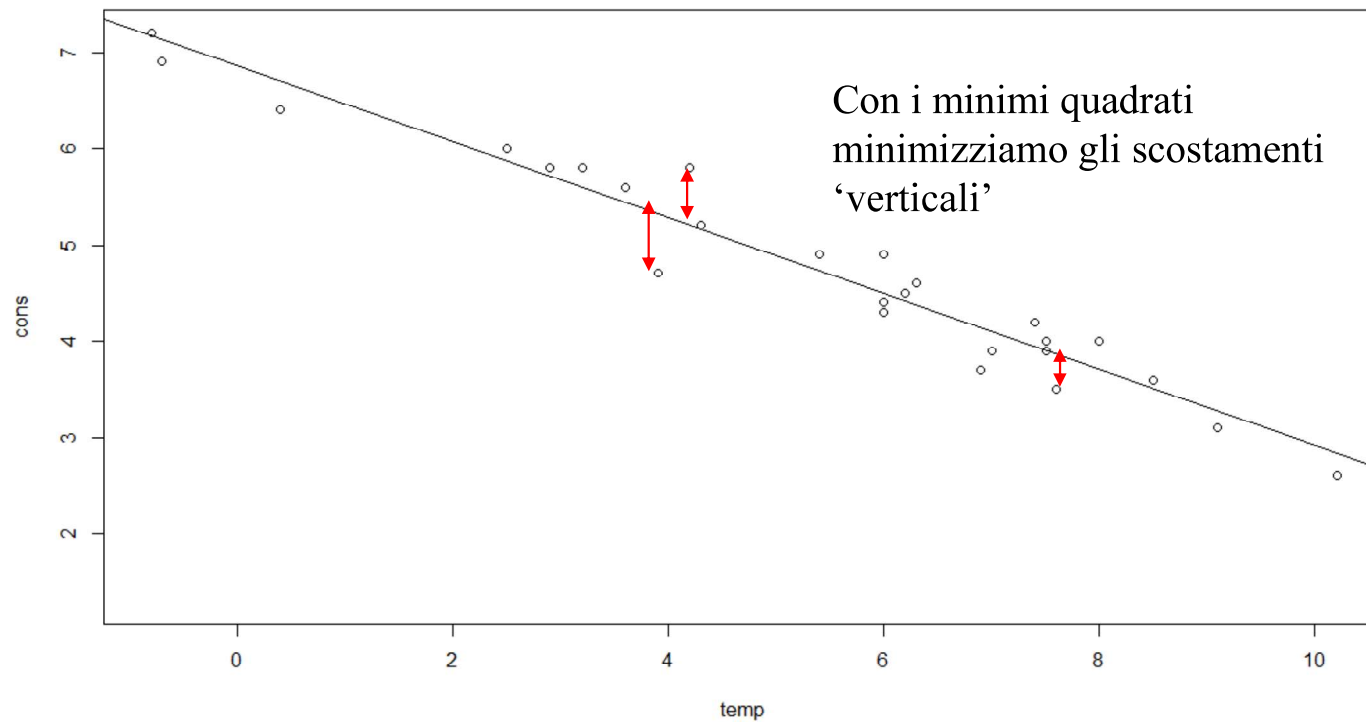
$$\min \sum_{i=1}^n (y_i - (b_0 + b_1 x_i))^2$$

$$r = -0.97$$

$$\min \sum_{i=1}^n (y_i - \hat{y}_i)^2$$



**Retta di regressione lineare**  **Metodo dei minimi quadrati**





## Stime dei parametri di regressione e proprietà

$$SS(\beta_0, \beta_1) = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

*SS Somma dei quadrati degli scarti tra risposte stimate e reali*

Il metodo dei minimi quadrati consiste nello scegliere come stimatori di  $\beta_0$  e  $\beta_1$  i due valori che minimizzano SS

$$\begin{cases} \frac{\partial SS}{\partial \beta_0} = -2 \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_i) = 0 \\ \frac{\partial SS}{\partial \beta_1} = -2 \sum_{i=1}^n x_i (Y_i - \beta_0 - \beta_1 x_i) = 0 \end{cases}$$





$$\begin{cases} \sum_{i=1}^n Y_i = nB_0 + B_1 \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i Y_i = B_0 \sum_{i=1}^n x_i + B_1 \sum_{i=1}^n x_i^2 \end{cases}$$

**Stimatori ai  
minimi quadrati**

$$\begin{cases} B_0 = \bar{Y} - B_1 \bar{x} \\ B_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{cases}$$

con:

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$



$$\begin{cases} b_0 = \bar{y} - b_1 \bar{x} \\ b_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{cases}$$

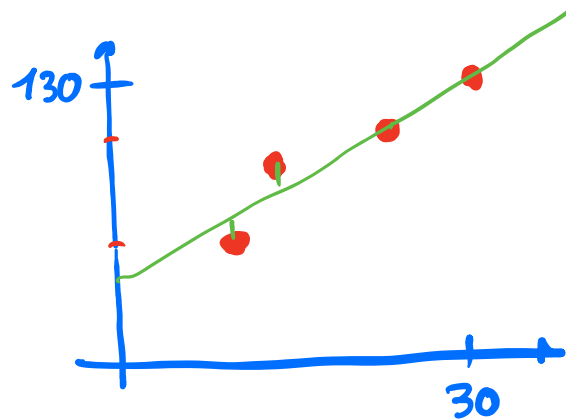
$$\text{con } \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

o in forma equivalente ma più adatta al calcolo:

$$\begin{cases} b_0 = \bar{y} - b_1 \bar{x} \\ b_1 = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2} \end{cases}$$

X	Y	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$	$(y_i - \bar{y})^2$	$\hat{y}_i$
100	27	10	2	20	100	4	
60	18	-30	-7	210	900	21	
130	30	40	5	200	1600	25	
70	25	-20	0	0	400	0	
medie	90	25	0	107.5	8850	19.5	

$$y = b_0 + b_1 x$$



$$C_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n} = 107.5 \quad \frac{\sum (x_i - \bar{x})^2}{n} = \sigma_x^2 = 8850 \quad \sigma_y^2 = 19.5$$

$$r = \frac{C_{xy}}{\sqrt{\sigma_x^2 \cdot \sigma_y^2}} = \frac{107.5}{\sqrt{19.5 \cdot 8850}} \approx 0.89$$

$$b_1 = \frac{C_{xy}}{\sigma_x^2} = \frac{107.5}{8850} \approx 0.14$$

$$b_0 = \bar{y} - b_1 \bar{x} = 25 - 0.14 \cdot 90 \approx 12$$

$$Y = 0.14 \cdot X + 12$$

$$X_0 = 80 \rightarrow Y_0 = 0.14 \cdot 80 + 12 \approx 11.2 + 12 \approx 23$$

$$Y_0 = 23 + \epsilon \quad \epsilon = ?$$



## Distribuzione degli stimatori

$\beta_0, \beta_1, \varepsilon = ?$

$\beta_0, \beta_1$

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

per  $x = x_i$

$\varepsilon_i$  (errore indiv.  $i$ -esimo)

➤  $E[\varepsilon_i] = 0$  e  $var[\varepsilon_i] = \sigma^2$

➤  $cov[\varepsilon_i, \varepsilon_j] = 0$  per ogni  $i \neq j$



➤  $E[Y_i] = E[\beta_0 + \beta_1 x_i + \varepsilon_i] = \beta_0 + \beta_1 x_i$  per ogni  $i = 1, 2, \dots, n$

➤  $var[Y_i] = var[\beta_0 + \beta_1 x_i + \varepsilon_i] = \sigma^2$  per ogni  $i = 1, 2, \dots, n$

➤ le  $Y_i$  sono non correlate, cioè  $cov[Y_i, Y_j] = 0$  per ogni  $i \neq j$

➤ la variabile  $x$  è nota senza errore ed è osservata per almeno due valori distinti

+ Normalità

$\longleftrightarrow \varepsilon_i \sim N(0, \sigma^2)$



$$B_0, B_1 \sim N$$

$$\Rightarrow \mathbb{E}[B_0] = \beta_0 \quad \text{Var}(B_0) = \frac{\sigma^2 \sum_{i=1}^n x_i^2}{n(\sum_{i=1}^n x_i^2 - n\bar{x}^2)}$$

$$\Rightarrow \mathbb{E}[B_1] = \beta_1 \quad \text{var}(B_1) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\Rightarrow \text{cov}[B_0, B_1] = -\frac{\sigma^2 \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2}$$



Sia

$$SS_R = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (Y_i - B_0 - B_1 x_i)^2$$

*(Handwritten blue notes:  $= \sum (\epsilon_i - 0)^2$  with an arrow pointing to  $\epsilon_i^2$ ; blue circles around  $B_0$  and  $B_1$  with arrows pointing to them)*

Si può dimostrare che

$$\frac{SS_R}{\sigma^2} \sim \chi_{n-2}^2$$

$$E \left[ \frac{SS_R}{\sigma^2} \right] = n - 2$$

$$\frac{(1-r^2) \sum (y_i - \hat{y})^2}{n-2} = \hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

stimatore non distorto di  $\sigma^2$

Questa è la stima della varianza dell'errore. La sua radice è chiamata  
**ERRORE STANDARD (o deviazione standard dell'errore)**





Alcune definizioni utili per dopo

$$S_{xY} = \sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y}) = \sum_{i=1}^n x_i Y_i - n\bar{x}\bar{Y}$$

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2$$

$$S_{YY} = \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n Y_i^2 - n\bar{Y}^2$$

## Inferenza statistica sui parametri di regressione

$$\frac{B_0 - E[\beta_0]}{\sqrt{\text{var}[B_0]}} = \frac{B_0 - \beta_0}{\hat{\sigma} / \sqrt{\frac{nS_{xx}}{\sum_{i=1}^n x_i^2}}} \sim N(0,1) \quad \text{t-student}$$

Intervallo di fiducia

Test di ipotesi

$$\frac{B_1 - E[\beta_1]}{\sqrt{\text{var}[B_1]}} = \frac{B_1 - \beta_1}{\hat{\sigma} / \sqrt{S_{xx}}} \sim N(0,1)$$

Intervallo di fiducia

Test di ipotesi