


Correlazione e regressione

Correlazione

Correlazione

Come posso determinare il legame tra
due o più variabili?




COEFFICIENTE DI CORRELAZIONE
(r di Pearson)

- massimo consumo di ossigeno e prestazione nelle gare di resistenza
- indice di forza relativa e capacità di salto

Correlazione

COEFFICIENTE DI CORRELAZIONE



- $-1 < r < +1$
- $r=0$ indica assenza di *correlazione*
- se $r > 0$ le due variabili *covariano*
- se $r < 0$ le due variabili *controvariano*
- $r=1$ o $r=-1$ esiste una relazione matematica tra le due variabili

Correlazione

Calcolo di r

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{(\sum (x_i - \bar{x})^2)(\sum (y_i - \bar{y})^2)}}$$

Misura la forza della associazione tra due variabili

Correlazione

ESEMPIO: velocità di corsa e consumo di ossigeno

t (ore)	VO ₂ (ml/kg min)
2.2	78
2.3	78.5
2.5	76
2.8	75.8
2.9	60.1
3.0	59.2
3.1	59.0
3.1	59.2
3.3	58.7
3.4	58.0

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{(\sum (x_i - \bar{x})^2)(\sum (y_i - \bar{y})^2)}}$$

Correlazione

t	(t-t)	VO ₂	(VO ₂ -VO ₂)	(t-t) (VO ₂ -VO ₂)
2.2	-0.66	78	11.75	-7.755
2.3	-0.56	78.5	12.25	-6.86
2.5	-0.36	76	9.75	-3.51
2.8	-0.06	75.8	9.55	-0.573
2.9	0.04	60.1	-6.15	-0.246
3.0	0.14	59.2	-7.05	-0.987
3.1	0.24	59.0	-7.25	-1.74
3.1	0.24	59.2	-7.05	-1.692
3.3	0.44	58.7	-7.55	-3.322
3.4	0.54	58.0	-8.25	-4.45
				-31.14

NUMERATORE

Correlazione

t	(t-t)	(t-t) ²	VO ₂	(VO ₂ -VO ₂)	(VO ₂ -VO ₂) ²
2.2	-0.66	0.4356	78	11.75	138.06
2.3	-0.56	0.3136	78.5	12.25	150.06
2.5	-0.36	0.1296	76	9.75	95.06
2.8	-0.06	0.0036	75.8	9.55	91.20
2.9	0.04	0.0016	60.1	-6.15	37.82
3.0	0.14	0.0196	59.2	-7.05	49.70
3.1	0.24	0.0576	59.0	-7.25	52.56
3.1	0.24	0.0576	59.2	-7.05	49.70
3.3	0.44	0.1936	58.7	-7.55	57.00
3.4	0.54	0.2916	58.0	-8.25	68.06
		1.504			789.25

DENOMINATORE

Correlazione

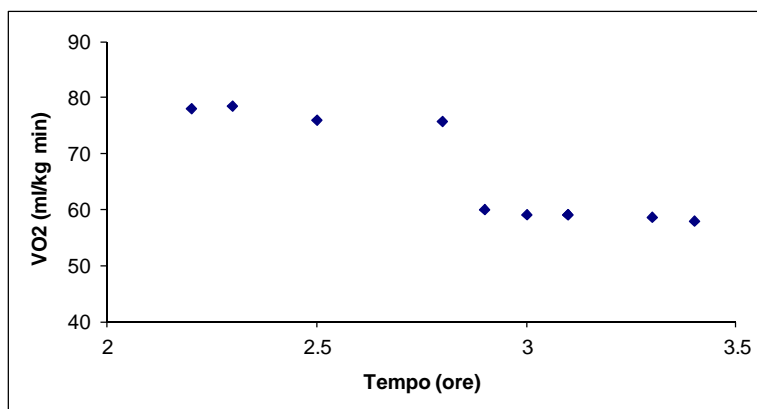
Sostituisco i valori

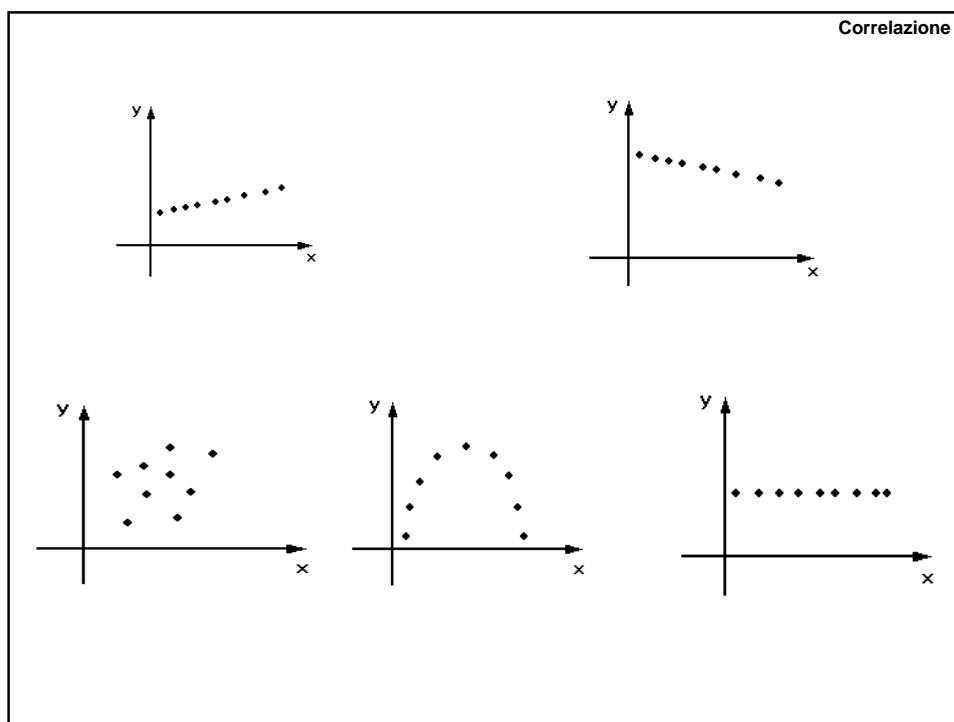
$$r = \frac{-31.14}{\sqrt{1.504 \times 789.25}} = -0.904$$

le due variabili sono inversamente correlate

(all'aumentare del tempo corrisponde una diminuzione del consumo di ossigeno)

Correlazione





La regressione lineare semplice

Regressione

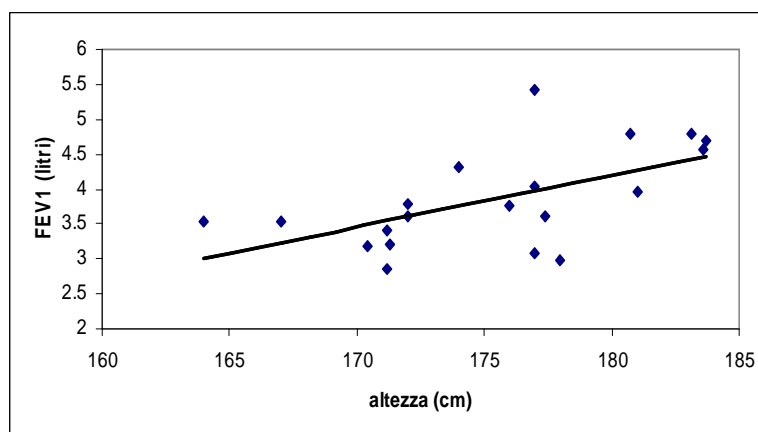
REGRESSIONE

Esempio: relazione tra FEV₁ (Volume espiratorio forzato) e altezza

altezza (cm)	FEV ₁ (litri)	altezza (cm)	FEV ₁ (litri)	altezza (cm)	FEV ₁ (litri)
164.0	3.54	172.0	3.78	178.0	2.98
167.0	3.54	174.0	4.32	180.7	4.80
170.4	3.19	176.0	3.75	181.0	3.96
171.2	2.85	177.0	3.09	183.1	4.78
171.2	3.42	177.0	4.05	183.6	4.56
171.3	3.20	177.0	5.43	183.7	4.68
172.0	3.60	177.4	3.60		

Regressione

La retta è la migliore rappresentazione della relazione tra le due variabili



In questo esempio, vogliamo sapere quale è la media (valore atteso) del FEV_1 per gli studenti di una certa altezza e quale è l'incremento del FEV_1 all'aumento unitario dell'altezza

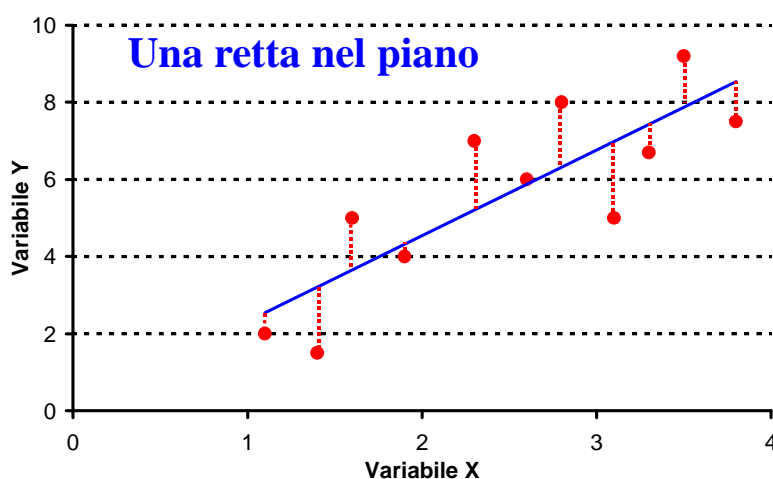
FEV_1 è la *variabile di risposta* o *dipendente*
 altezza è la *variabile esplicativa* o *indipendente*

$$FEV_1 = a + b \times \text{altezza}$$

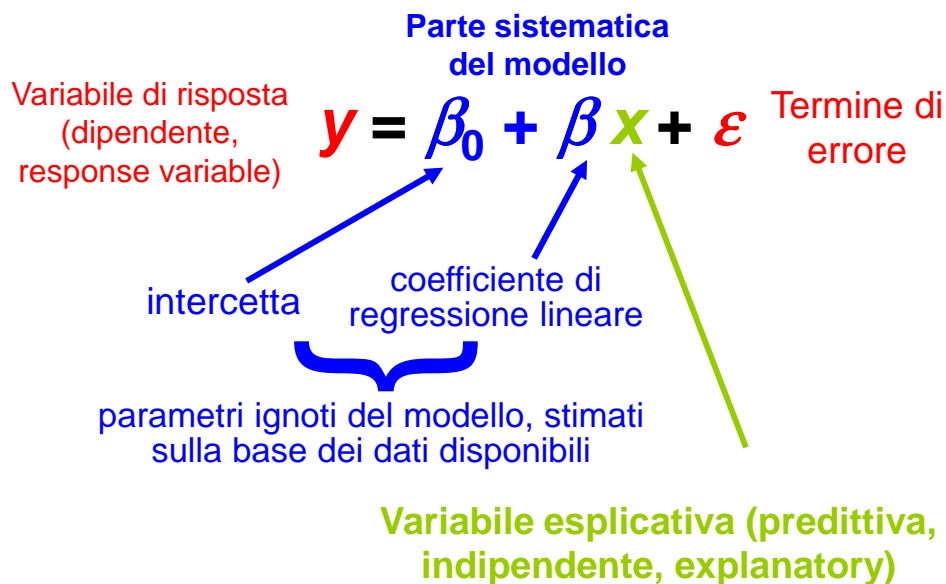
REGRESSIONE LINEARE SEMPLICE

Il modello di regressione lineare semplice - 1

$$y = \beta_0 + \beta x + \varepsilon$$



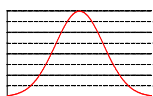
Il modello di regressione lineare semplice - 2



Il modello di regressione lineare semplice - 3

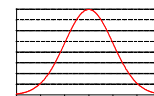
$$y = \beta_0 + \beta x + \varepsilon$$

Variabile di risposta
(dipendente)



Predittore lineare,
parte deterministica del modello,
senza variabilità casuale

Termine di errore,
parte probabilistica



**L'errore, e quindi la variabile di risposta,
si distribuisce NORMALMENTE**

Il modello di regressione lineare semplice - 4

Il FEV₁ (Y) dipende dalla statura (X₁)

$$E(y) = \beta_0 + \beta_1 x_1$$

E(y) = valore atteso (media) del FEV1 degli individui che hanno quella determinata statura

$$y = \beta_0 + \beta_1 x_1 + \varepsilon$$

y = FEV₁ di un determinato individuo, che dipende dalla statura, (parte sistematica del modello), ma anche da altre caratteristiche individuali (ε , parte probabilistica)

Il modello di regressione lineare semplice - 5

- **Modello teorico (ignoto)**

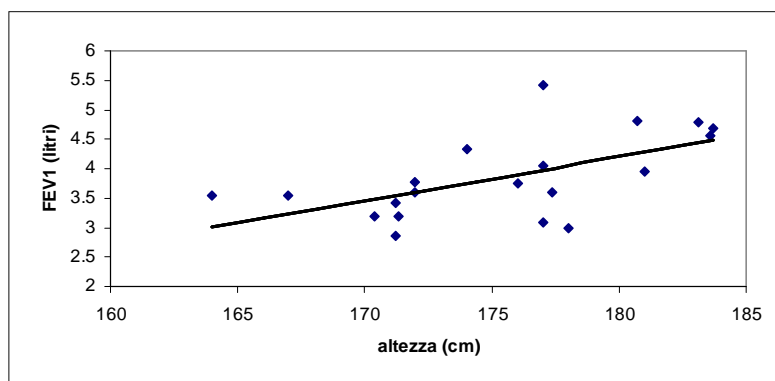
$$y = \beta_0 + \beta_1 x + \varepsilon$$

- **Regressione Lineare stimata**

$$\hat{y} = b_0 + b_1 x$$

Regressione

Come costruire la retta?
Metodo dei MINIMI QUADRATI



$$FEV_1 (\text{litri}) = -9.19 + 0.0744 \times \text{altezza (cm)}$$

Se uno studente è alto 170 cm, il suo FEV_1 è 3.458 litri

Regressione

$$FEV_1 (\text{litri}) = -9.19 + 0.0744 \times \text{altezza (cm)}$$

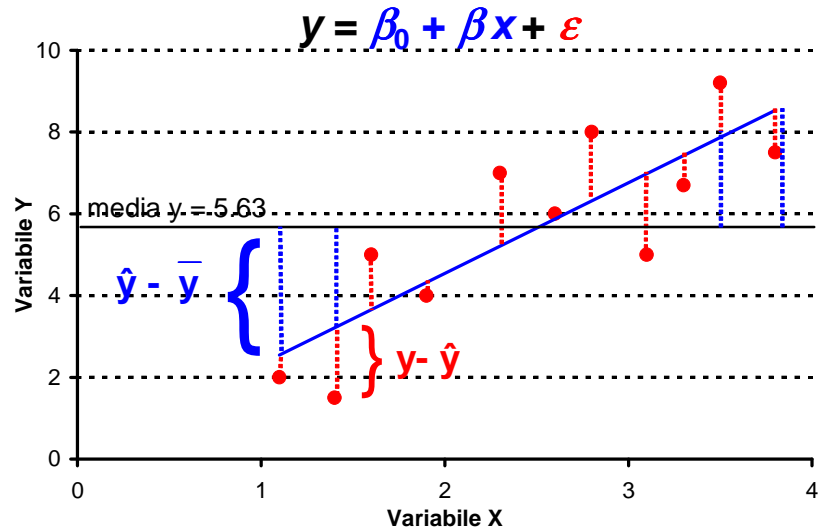
coefficiente di regressione

intercetta

$$b = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$a = \bar{y} - b\bar{x}$$

SCOMPOSIZIONE DELLA DEVIANZA nella Regressione lineare semplice - 1



$$(y - \bar{y}) = (\hat{y} - \bar{y}) + (y - \hat{y})$$

SCOMPOSIZIONE DELLA DEVIANZA nella Regressione lineare semplice - 2

Variabilità totale

$(y - \bar{y}) = (\hat{y} - \bar{y}) + (y - \hat{y})$

Variabilità spiegata dalla regressione

Variabilità residua

Si può dimostrare che:

Devianza totale, SST

$\Sigma(y - \bar{y})^2 = \Sigma(\hat{y} - \bar{y})^2 + \Sigma(y - \hat{y})^2$

Devianza spiegata dalla regressione, SSR

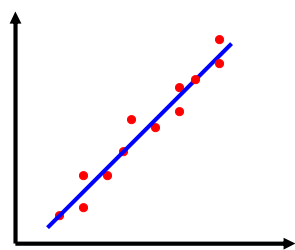
Devianza residua, SSE

Regressione lineare semplice

Si cerca di trovare la retta che meglio interpola,
che meglio si adatta alla nuvola di punti.

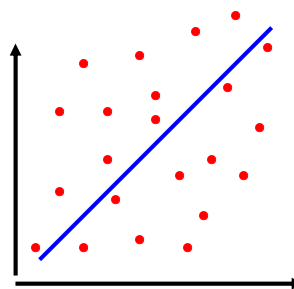
METODO DEI MINIMI QUADRATI

Si sceglie la retta che **riduce al minimo la devianza residua, SSE, $\Sigma(y - \hat{y})^2$**



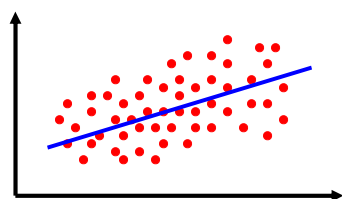
La retta di regressione
approssima bene i dati.

devianza spiegata > dev.residua

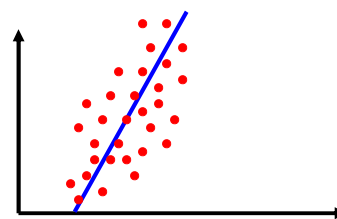


La retta di regressione
approssima male i dati.

devianza spiegata < dev.residua



b prossimo a zero



b elevato

COEFFICIENTE DI DETERMINAZIONE

$$r^2 = \frac{\sum (\hat{y} - \bar{y})^2}{\sum (y - \bar{y})^2}$$

$$r^2 = 1 - \frac{\text{SSE}}{\text{SST}} = \frac{\text{SSR}}{\text{SST}}$$

SST = Devianza (o Somma dei quadrati) totale

SSR = Devianza (o Somma dei quadrati) spiegata dalla
Regressione

**Proporzione di variazione totale della variabile dipendente
Y che è spiegata dalla variabile indipendente X**

Regressione = relazione di tipo asimmetrico:
una variabile casuale (Y) dipende da una
variabile fissa (X)

Correlazione = relazione di tipo simmetrico:
le due variabili sono sullo stesso piano

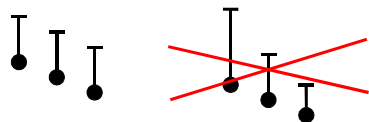
Regressione lineare semplice

$$y = \beta_0 + \beta x + \varepsilon$$

ASSUNZIONI

1) Il valore atteso degli errori $E(\varepsilon)$ deve essere pari a ZERO

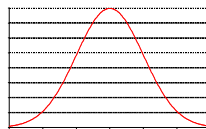
2) **OMOSCEDASTICITA'** (La varianza degli errori rimane costante)



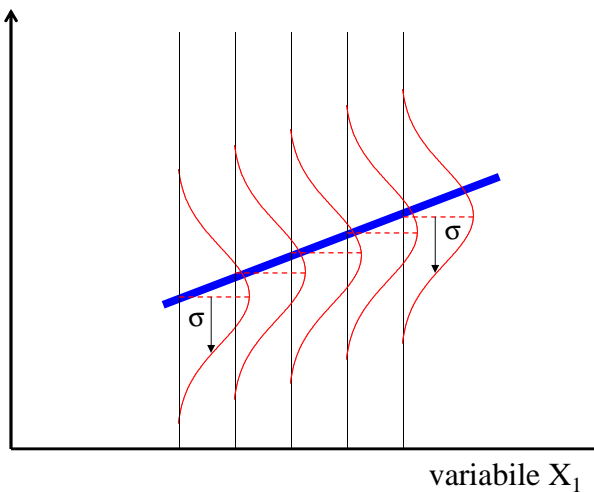
3) **INDIPENDENZA** degli errori

se le provette tra un esame e l'altro non vengono lavate adeguatamente, una determinazione risente della determinazione precedente

4) **Distribuzione NORMALE** degli errori



variabile Y



variabile X_1

Inferenza sui parametri: i dati “supportano” il modello proposto?

I punti hanno distanze dalla retta di regressione che sono sensibilmente minori di quelle dalla media.

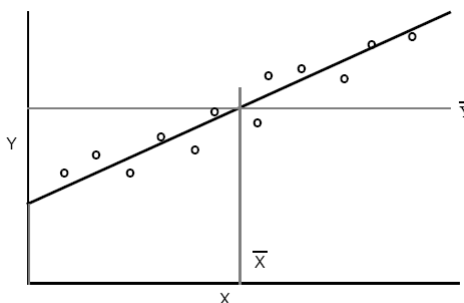


Figura A

Il valore stimato mediante la retta di regressione si avvicina molto all'osservazione reale.

Raccogliendo altri punti campionari, la retta calcolata resterebbe probabilmente immutata. La retta di regressione esprime **la relazione reale** che esiste tra X ed Y.

Inferenza sui parametri: i dati “supportano” il modello proposto?

La retta calcolata non rappresenta un miglioramento effettivo della distribuzione dei punti, rispetto alla loro media.

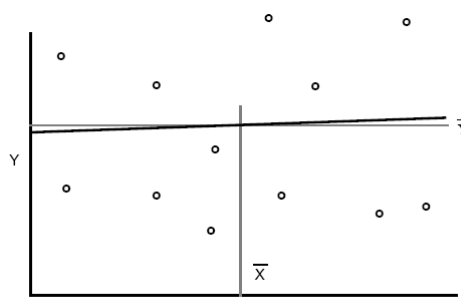


Figura B

La retta calcolata non è rappresentativa di una relazione reale tra X ed Y.

**Inferenza sui parametri:
i dati “supportano” il modello proposto?**

Test t di Student, basato su b_1 (stima di β_1)

$$\begin{cases} H_0: \beta_1 = 0 \\ H_1: \beta_1 \neq 0 \end{cases} \quad b_1 \sim N \left(\beta_1, \frac{\sigma_\varepsilon^2}{\sum(x - \bar{x})^2} \right)$$

Non conosco la quantità σ_ε^2 , ne ottengo una stima:

$$\hat{\sigma}_\varepsilon^2 = \frac{\sum(y - \hat{y})^2}{n-2}$$

$$t = \frac{b_1 - 0}{ES_b} = \frac{b_1}{\frac{\hat{\sigma}_\varepsilon^2}{\sum(x - \bar{x})^2}} = \frac{b_1}{\frac{\sum(y - \hat{y})^2}{(n-2) \sum(x - \bar{x})^2}} \sim t_{n-2 \text{ g.d.l.}}$$

ESEMPIO:

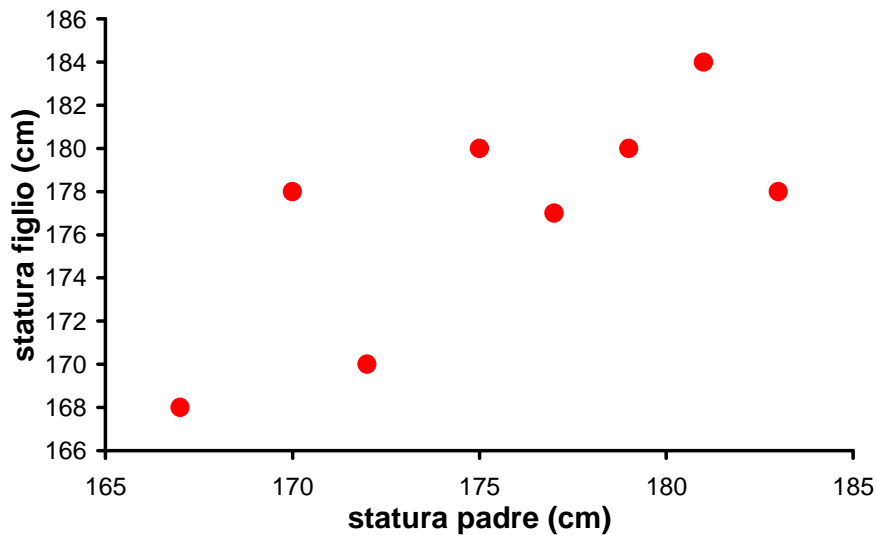
Esiste una relazione tra altezza dei padri e altezza dei figli maschi?

(A padri più bassi corrispondono figli più bassi?)

(A padri più alti corrispondono figli più alti?)

Padre	Figlio
167 cm	168 cm
175 cm	180 cm
183 cm	178 cm
170 cm	178 cm
181 cm	184 cm
172 cm	170 cm
177 cm	177 cm
179 cm	180 cm

I passo: rappresentazione grafica mediante diagramma di dispersione (scatterplot)



II passo: si ipotizza un modello statistico, che possa essere utile ad interpretare i dati

Ipotizziamo un modello lineare del tipo: $y = \beta_0 + \beta x + \varepsilon$

$$(\text{altezza figli}) = \beta_0 + \beta (\text{altezza padri}) + \varepsilon$$

(i figli di uno stesso padre hanno statura abbastanza simile, ma non necessariamente uguale, anche se ομοιομετρός, cioè figli della stessa madre)

IV passo: Stima dei parametri del modello con il metodo dei minimi quadrati

$$b_1 = \frac{\text{codevianza}_{xy}}{\text{devianza}_x} = 150,5 / 216 = 0,697 \text{ cm/cm}$$

$$b_0 = \bar{y} - b_1 \bar{x} = 176,875 - 0,697 * 175,5 = 54,59 \text{ cm}$$

Retta di regressione:

$$\text{altezza figlio (cm)} = 54,6 \text{ cm} + 0,697 \text{ cm/cm} * \text{altezza padre (cm)}$$

Quando la statura del padre cresce di 1 cm, la statura del figlio cresce in media di 7 mm.

VI passo: Inferenza sui parametri: i dati “supportano” il modello proposto?

Test t di Student, basato su b_1 (stima di β_1)

$$\begin{cases} H_0: \beta_1 = 0 \\ H_1: \beta_1 \neq 0 \end{cases} \quad \begin{array}{l} \text{Livello di significatività} = 5\% \\ \text{Gradi di libertà} = n - 2 = 8 - 2 = 6 \\ \text{Soglia critica} = t_{6, 0,025} = 2,447 \end{array}$$

test a due code

$$t = \frac{b - 0}{ES_b} = \frac{b}{\sqrt{\text{var}_{\text{res}} / \text{dev}_x}} = \frac{0,697}{\sqrt{15,67 / 216}} = 2,588$$

$$\text{dev}_{\text{res}} = \text{dev}_y - \text{codev}_{xy}^2 / \text{dev}_x = 198,9 - 150,5^2 / 216 = 94,01$$

$$\text{var}_{\text{res}} = \text{dev}_{\text{res}} / (n-2) = 94,01 / 6 = 15,67$$

VII passo: Previsione sui valori della variabile Y

Per $x = 185$ cm, qual è il valore atteso di Y?

Retta di regressione:

$$\hat{y} = 54,6 \text{ cm} + 0,697 \text{ cm/cm} * 185 \text{ cm} = 183,49 \text{ cm}$$

$$\begin{aligned} ES_{\hat{y}} &= \sqrt{\text{var}_{\text{res}} [1/n + (x-\bar{x})^2 / \text{dev}_x]} = \\ &= \sqrt{15,67 [1/8 + (185-175,5)^2 / 216]} = 2,916 \end{aligned}$$

$$IC_{95\%} = \hat{y} \pm t_{v,\alpha/2} * ES_{\hat{y}} = 183,49 \pm 2,447 * 2,916 = \begin{bmatrix} 190,63 \\ 176,36 \end{bmatrix}$$

Tabella per l'ANALISI della REGRESSIONE

Fonte di variabilità	Gradi di libertà	Devianza	Varianza	Test-F
Spiegata	1	$\Sigma (\hat{y} - \bar{y})^2$	$DEV_{\text{spieg}} / 1$	$\frac{VAR_{\text{spieg}}}{VAR_{\text{resid}}}$
Residua (errore)	N-2	$\Sigma (y - \hat{y})^2$	$DEV_{\text{resid}} / (N-2)$	
Totale	N-1	$\Sigma (y - \bar{y})^2$		

1

40