

---

# Caso di studio : Terapie di riabilitazione per l'apprendimento motorio del braccio (correlazione e regressione)

**Laura Ventura**

Dipartimento di Scienze Statistiche

Università di Padova

<http://homes.stat.unipd.it/ventura/>

[ventura@stat.unipd.it](mailto:ventura@stat.unipd.it)

copyright©2015

## Caso di studio: Terapie di riabilitazione per l'apprendimento motorio del braccio

- **Dataset:** misurazioni relative ad uno studio sull'apprendimento motorio di un gruppo di pazienti, esposti al trattamento con realtà virtuale (IRCCS San Camillo, Lido di Venezia).
- **Variabile di interesse:** FIM (*Functional Independence Measure*), scala dell'autonomia del paziente con valori da 0 (non autosufficienza completa) a 130 (completa autonomia).
- Si hanno anche **due trattamenti:** 27 pazienti sono stati sottoposti ad una terapia di riabilitazione in un ambiente virtuale (casi, TRATTAMENTO=1) e 20 pazienti sono stati sottoposti ad una terapia convenzionale (controlli, TRATTAMENTO=2).
- La variabile FIM è stata misurata sia prima (FIMPRE) che dopo (FIMPOST) la terapia ricevuta, subito dopo un infarto.



# Caso di studio: Terapie di riabilitazione per l'apprendimento motorio del braccio

## □ DATI e VARIABILI:

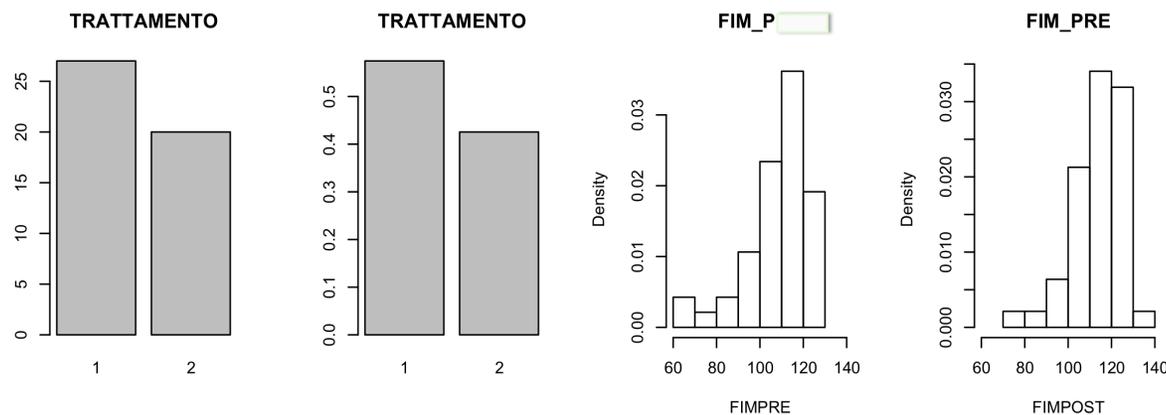
	TRATTAMENTO	FIMPRE	FIMPOST
1	2	124	124
2	2	108	110
	...		
	...		
46	1	111	113
47	1	99	108

- **Obiettivo:** Si intende verificare se c'è un miglioramento della performance motoria dell'arto a seguito della terapia, e se il gruppo trattato con la realtà virtuale ha un miglioramento superiore rispetto al gruppo di controllo.

## Riassunto delle "lezioni" precedenti: Analisi esplorativa

La prima fase di ogni analisi statistica è rappresentata dall'organizzazione e dalla sintesi dei **DATI**, le informazioni raccolte sulle **UNITÀ STATISTICHE** che compongono il **CAMPIONE**.

- **TRATTAMENTO (variabile qualitativa)**: i casi sono 27 (57.4%) e i controlli sono 20 (42.6%).
- **FIM (variabile quantitativa)**: La media di FIMPRE è di 109.3 (sd=13.8) e la media di FIMPOST è di 114.6 (sd=10.9). La media di FIMPRE è di 113.3 per i casi (sd=11.4) e 103.95 per i controlli (sd=15.14). La media di FIMPOST è 118.9 per i casi (6.81) e 108.65 per i controlli (12.6). La mediana di FIMPRE è di 116 per i casi e 107.5. La mediana di FIMPOST è 120 per i casi e 110.



## La "lezione" di oggi: Dai dati univariati ai dati bivariati

- In molte situazioni interessa **studiare** se esiste una relazione tra due variabili misurate sulle stesse unità. Esempi:
  - *“Le misurazioni della FIM prima della terapia sono in relazione con le misurazioni dopo la terapia?”*
  - *“il voto di maturità è in relazione con la performance universitaria?”*
- Oppure si desidera **prevedere** il valore di una variabile conoscendo il valore di un'altra. Esempi:
  - *“conoscendo il valore della FIMPRE, si può stimare il valore della FIMPOST?”*
  - *“conoscendo l'altezza del padre, è possibile prevedere l'altezza di un figlio?”*
- La statistica permette di rispondere a questo tipo di domande, con strumenti adatti alla natura delle variabili in esame. A tale scopo, **per variabili quantitative**, si tratteranno:
  - La **CORRELAZIONE**, che misura la dipendenza lineare tra due variabili;
  - La **REGRESSIONE**, che valuta la relazione lineare tra due variabili.

# Correlazione

- La **correlazione** misura l'associazione tra due variabili quantitative. È lo strumento che si utilizza quando si hanno a disposizione coppie di valori di variabili  $\Rightarrow$  **permette di valutare come variano i valori di una variabile al variare dell'altra e viceversa.**
- Esempi:
  - Numero di sigarette fumate in gravidanza e tasso di crescita del feto  $\Rightarrow$  all'aumentare del numero di sigarette fumate diminuisce il tasso di crescita (**correlazione negativa**).
  - Livello di colesterolo e BMI (Body Mass Index = peso (kg)/altezza<sup>2</sup> (m<sup>2</sup>))  $\Rightarrow$  tanto è maggiore il BMI quanto è maggiore il livello di colesterolo (**correlazione positiva**).
  - Il valor medio della temperatura (ambiente) e il BMI  $\Rightarrow$  non c'è motivo di pensare che la temperatura influenzi il BMI delle persone (**assenza di correlazione**).
- La relazione può essere valutata tramite:
  - Un **grafico** (**grafico di dispersione**)
  - Un **indice** che quantifica il grado di correlazione (**coefficiente di correlazione**)

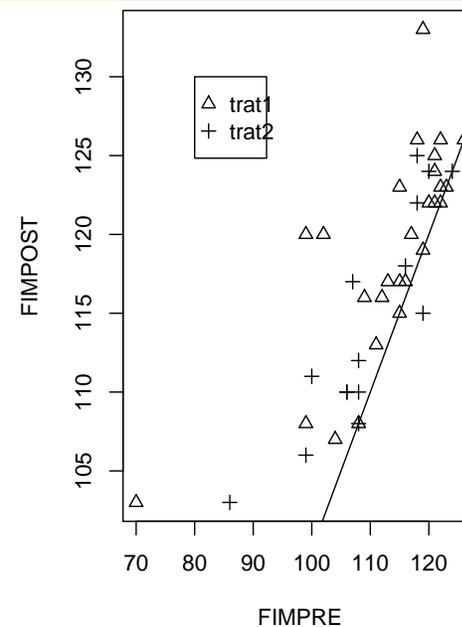
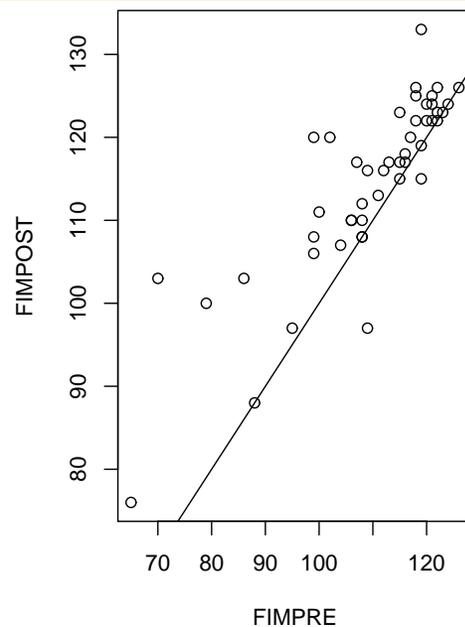
## Diagramma di dispersione

- Nello studio dell'associazione tra due variabili quantitative misurate sulle stesse unità statistiche, indicate con  $X$  e  $Y$ , è molto utile disegnare un grafico, il **diagramma di dispersione**, prima di procedere con altre analisi formali.

Nel grafico di dispersione le coppie

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

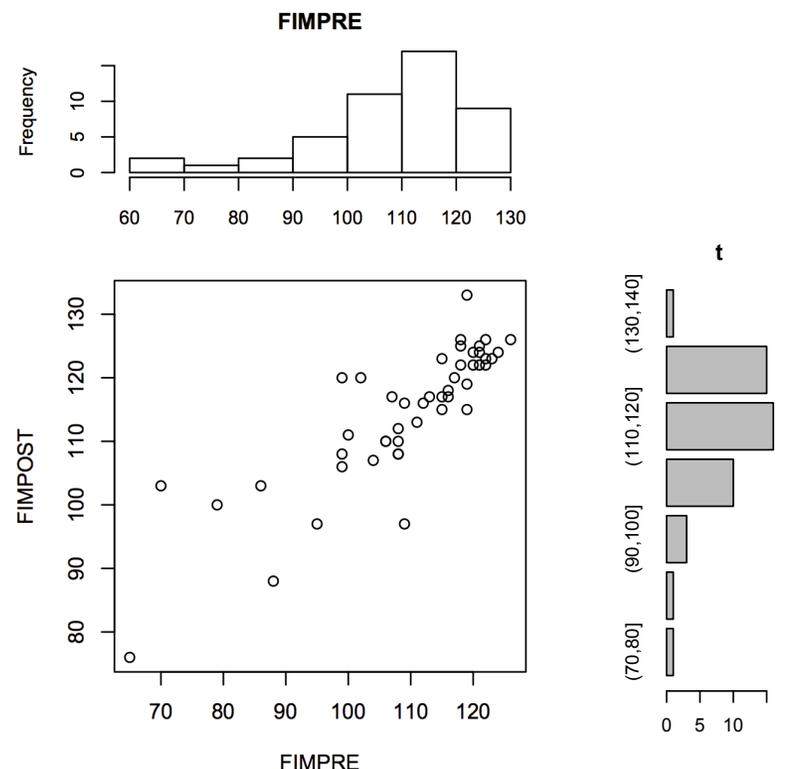
di valori di due variabili quantitative misurate sulle  $n$  unità sono rappresentati come punti di un piano cartesiano, i cui assi corrispondono alle due variabili.



## Diagramma di dispersione

- Ogni punto del grafico rappresenta una unità.
- Permette di verificare visivamente se le coppie di punti presentano una qualche forma di regolarità e per vedere come i punti si disperdono intorno a un particolare punto di riferimento: il **baricentro** della nuvola dei punti, ossia il punto di coordinate  $(m_x, m_y)$ .
- La nuvola di punti ha una forma allungata verso l'alto  $\Rightarrow$  a modalità crescenti della  $X$  corrispondono più frequentemente modalità crescenti della  $Y$ .
- Si possono considerare convenzioni grafiche per punti ripetuti.

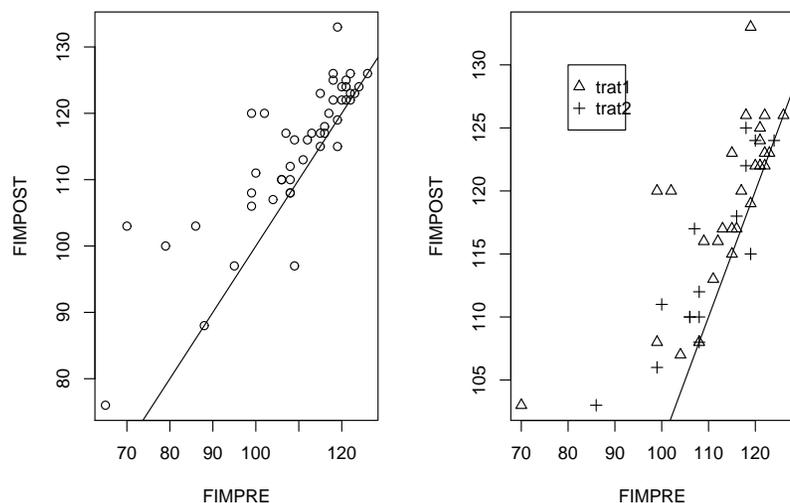
*P.s. La media aritmetica e la varianza di  $X$  sono  $m_x = \frac{x_1+x_2+\dots+x_n}{n} = \frac{1}{n} \sum_{i=1}^n x_i$  e  $S_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - m_x)^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - m_x^2$ . Analogamente, si indicano con  $m_y$  e  $S_y^2$  media e varianza di  $Y$ .*



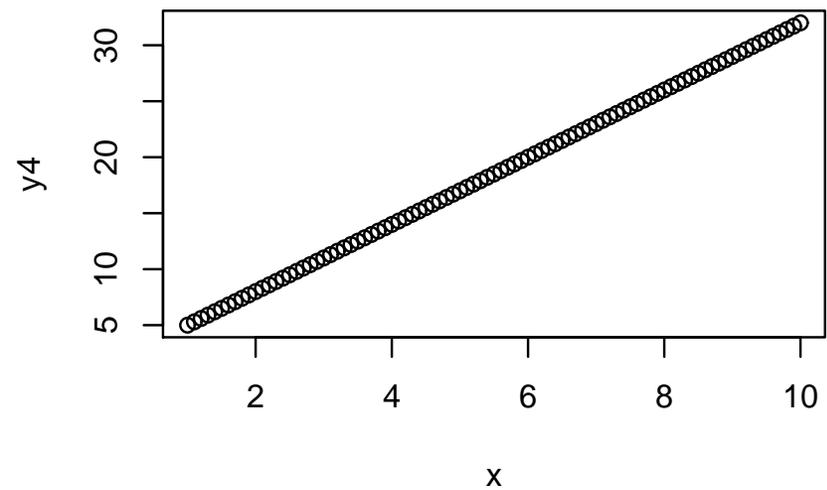
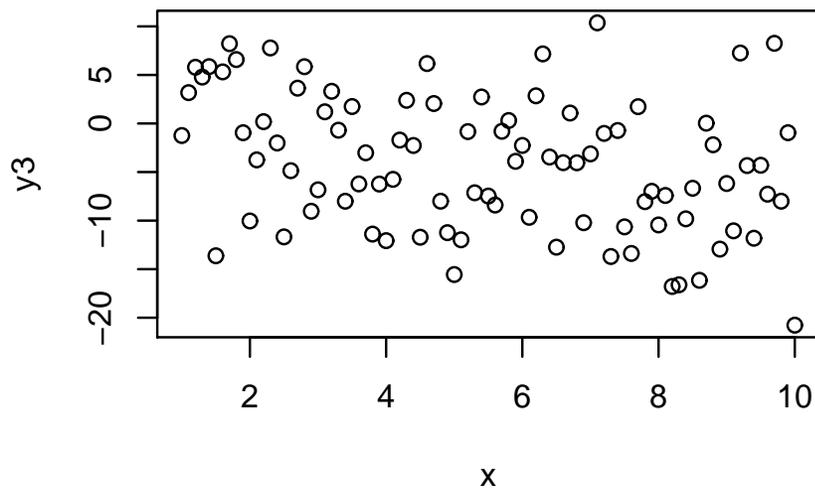
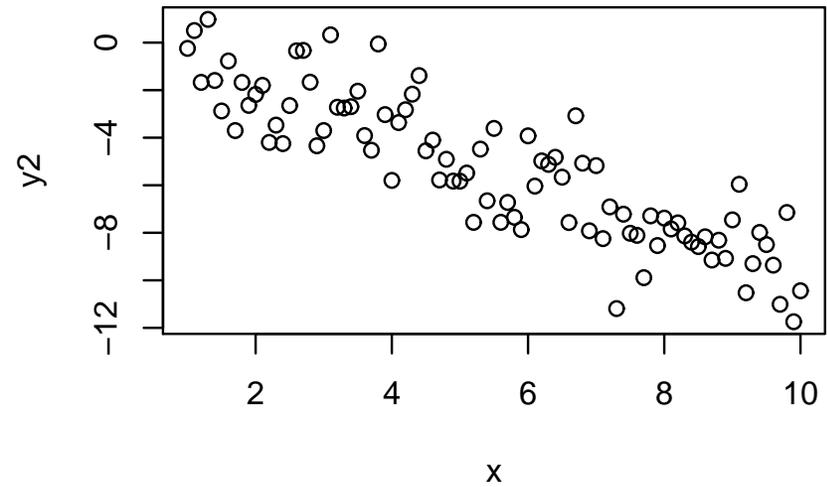
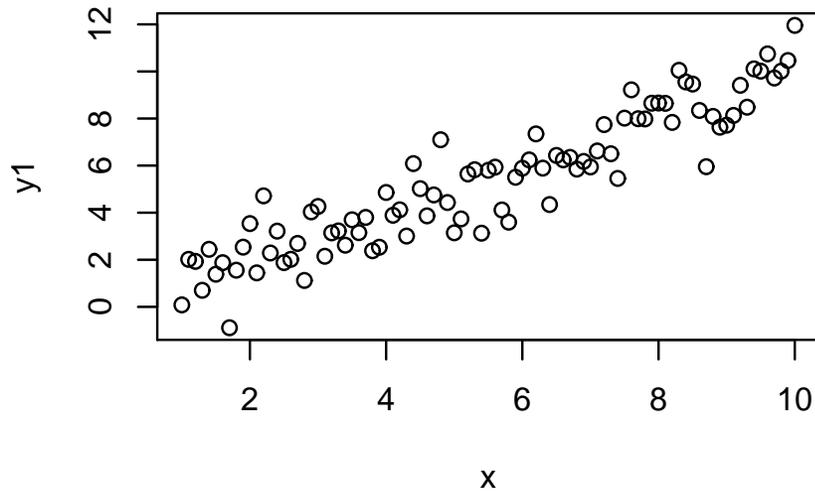
## Il ruolo delle variabili $X$ e $Y$ è simmetrico?

- A volte può essere importante spiegare una delle due variabili in funzione dell'altra. Si avrà quindi una **VARIABLE ESPLICATIVA  $X$**  e una **VARIABLE RISPOSTA  $Y$** .
- Ma a volte non ha importanza quale sia l'una e quale sia l'altra.

Nell'ESEMPIO della FIM è ragionevole voler esprimere la FIMPOST ( $Y$ ) a partire dalla FIMPRES ( $X$ ), misurabile a inizio trattamento. Dal grafico di dispersione si vede che, in generale, nei pazienti con FIMPRES elevata anche la FIMPOST è elevata  $\Rightarrow$  **correlazione positiva**.

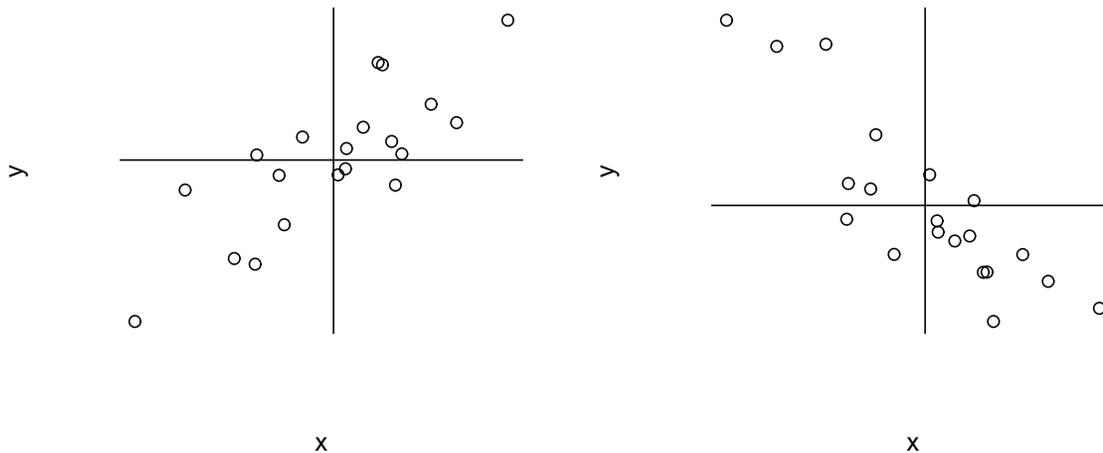


... qualche situazione tipo ...



## La covarianza

- Per avere una valutazione analitica del grado di associazione tra due variabili quantitative, esiste un indice che misura la dispersione nel piano dei punti dal proprio centro: la **COVARIANZA**.
- Il nome lascia intuire che si tratta di un'estensione al caso di due variabili della varianza. La covarianza si basa infatti sugli scarti delle  $x_i$  dalla propria media,  $(x_i - m_x)$ , e delle  $y_i$  dalla propria media,  $(y_i - m_y)$ .
- La covarianza, a differenza della varianza che è sempre positiva, misura l'eventuale direzione del legame, ovvero se le due variabili si muovono nella stessa direzione o in direzioni opposte. Il segno della covarianza riflette il senso crescente o decrescente dell'allineamento tendenziale.



## La covarianza

- La covarianza segnala una concordanza (sia che  $X$  e  $Y$  decrescono o crescono) con un segno  $+$  e una discordanza (quando  $X$  cresce e  $Y$  decresce, o viceversa) con il segno  $-$ . Formalmente, l'indice è

$$S_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - m_x)(y_i - m_y) .$$

- Una formula alternativa per il calcolo della covarianza è

$$S_{xy} = \frac{1}{n} \sum_{i=1}^n x_i y_i - m_x m_y$$

- Si noti che  $S_{xx} = S_x^2$ , ossia la covarianza tra  $X$  e  $X$  coincide con la varianza di  $X$ .

## Campo di variazione della covarianza

La covarianza può assumere valori sia positivi sia negativi. In particolare, vale

$$-S_x S_y \leq S_{xy} \leq S_x S_y$$

### Dimostrazione.

La varianza della combinazione  $aX - bY$  (Appendice), per  $a$  e  $b$  costanti, è  $a^2 S_x^2 + b^2 S_y^2 - 2ab S_{xy}$ .

Si consideri ora la variabile  $T$  definita come  $T = S_y^2 X - S_{xy} Y$ . Allora, la variabile  $T$  ha varianza

$$\begin{aligned} S_T^2 &= S_y^4 S_x^2 + S_{xy}^2 S_y^2 - 2S_y^2 S_{xy} S_{xy} \\ &= S_y^4 S_x^2 - S_{xy}^2 S_y^2 \end{aligned}$$

Ma poiché vale  $S_T^2 \geq 0$ , deve valere la diseuguaglianza

$$S_y^4 S_x^2 - S_{xy}^2 S_y^2 \geq 0$$

ossia, dividendo per  $S_y^2$ ,

$$S_{xy}^2 \leq S_y^2 S_x^2$$

da cui segue la tesi.

---

# La correlazione

## Il coefficiente di correlazione

- Dalla proprietà  $-S_x S_y \leq S_{xy} \leq S_x S_y$ , può essere costruito un indice relativo semplicemente dividendo  $S_{xy}$  per il prodotto degli scarti quadratici medi di  $X$  e  $Y$ . L'indice così ottenuto prende valori in  $[-1,1]$  e viene detto **coefficiente di correlazione**:

$$r_{xy} = \frac{S_{xy}}{S_x S_y} \quad -1 \leq r_{xy} \leq 1$$

- La formula del coefficiente di correlazione non è poi così terribile come appare!! Può solo essere noioso calcolarla a mano. In genere si usa un software opportuno.
- Un modo di procedere può essere il seguente:

- Per le due variabili si calcolano le medie  $m_x = \frac{1}{n} \sum x_i$  e  $m_y = \frac{1}{n} \sum y_i$
- Si calcola la media dei prodotti  $\frac{1}{n} \sum x_i y_i$
- Si calcolano le medie dei quadrati  $\frac{1}{n} \sum x_i^2$  e  $\frac{1}{n} \sum y_i^2$
- Si calcola la covarianza  $S_{xy} = \frac{1}{n} \sum x_i y_i - m_x m_y$
- Si calcolano  $S_x = [\frac{1}{n} \sum x_i^2 - m_x^2]^{1/2}$  e  $S_y = [\frac{1}{n} \sum y_i^2 - m_y^2]^{1/2}$
- Queste sono le grandezze che servono per calcolare  $r_{xy}$

- In sintesi: come si interpreta il valore trovato di  $r_{xy}$ ?

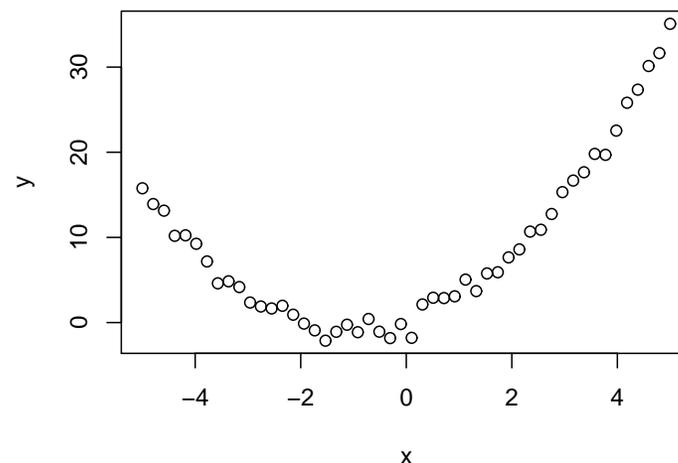
## Guida all'interpretazione di $r_{xy}$

- $-1 \leq r_{xy} \leq 1$
- $r_{xy} = +1$ : correlazione positiva perfetta (tutti i punti su una retta: concordi)
- $r_{xy} = -1$ : correlazione negativa perfetta (tutti i punti su una retta: discordi)
- $r_{xy} > 0$ : correlazione positiva
- $r_{xy} < 0$ : correlazione negativa
- $r_{xy} \cong 0$ : assenza di relazione lineare

Se  $r_{xy} = \pm 1$  le variabili sono legate da una relazione lineare perfetta (diretta o inversa, rispettivamente). Si parla di relazione lineare in quanto  $r_{xy}$  misura se le coppie di valori  $(x_i, y_i)$  sono allineate lungo una retta del tipo  $y = a + bx$ .

Quando tra  $X$  e  $Y$  non vi è una relazione lineare o essa è estremamente debole, il valore dell'indice  $r_{xy}$  è zero o circa zero, e le variabili sono dette incorrelate.

ATTENZIONE: Il coefficiente di correlazione misura una associazione lineare. Il valore  $r_{xy} = 0$  non indica tuttavia un'assenza di relazione tra le due variabili. Può esserci una relazione curvilinea.



## Esempio: $r_{xy}$ per la FIM

□ Siano  $Y = \text{FIMPOST}$  e  $X = \text{FIMPRES}$ .

□ Si ha

$$m_x = 109.3$$

$$m_y = 114.6$$

$$\sum (x_i - m_x)^2 = 8732.2$$

$$\sum (y_i - m_y)^2 = 5433.6$$

$$\sum (x_i - m_x)(y_i - m_y) = 5808.7$$

□ Allora:

$$r_{xy} = \frac{5808.7}{\sqrt{8732.2 \times 5433.6}} = 0.843$$

□ Il valore 0.843 indica una correlazione positiva elevata tra la FIMPRES e la FIMPOST (come ci si aspettava dal grafico di dispersione).

□ Con una relazione così, la FIMPOST potrebbe essere prevista in modo accurato conoscendo il valore della FIMPRES.

---

# La regressione

## La regressione

- Quando dall'analisi di un diagramma di dispersione emerge un particolare andamento della nuvola di punti di  $X$  e  $Y$ , è naturale chiedersi se esiste una qualche relazione statistica  $Y = f(X) + \text{errore}$  tra  $X$  e  $Y$ .
- Il problema è lo stesso di prima: si vuole studiare una relazione tra le variabili. La relazione non è più simmetrica!! Perché si vuole comprendere come la variabile risposta  $Y$  sia influenzata dalla variabile esplicativa  $X$ .
- Se la relazione che emerge è di tipo lineare, si può esprimere la relazione statistica tra  $X$  e  $Y$  usando un modello molto semplice: **l'equazione della retta**.

Il modello è del tipo:

$$Y = a + bX + \text{errore}$$

con

$a$  = intercetta

$b$  = coefficiente angolare

errore = la deviazione dalla retta dei punti osservati

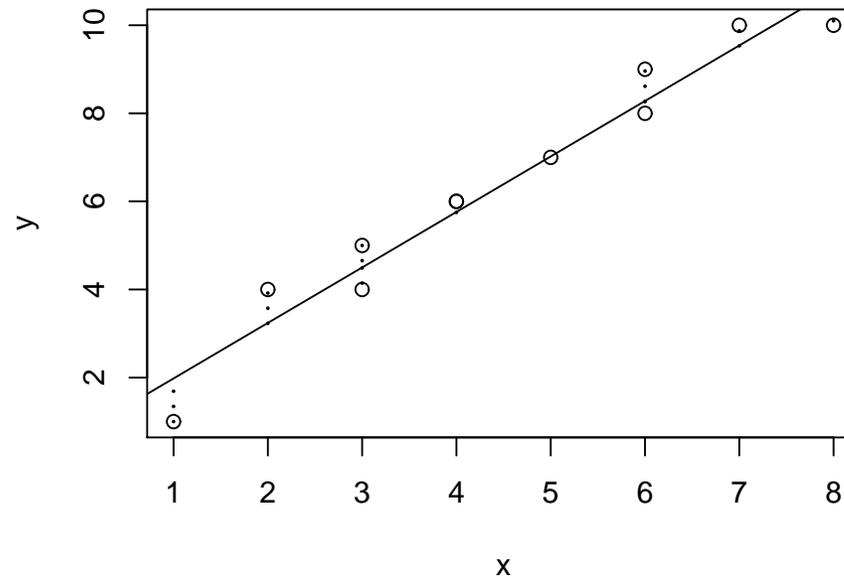
## La regressione

- Se si calcolano “opportunamente” i valori di  $a$  e  $b$ , l’equazione può essere usata per prevedere il valore della  $Y$  a partire da un qualunque valore della  $X$ .
- **PROBLEMA: come trovare la retta che si adatta nel modo migliore ai dati?**
- Si devono determinare i valori di  $a$  e  $b$  che rendono la retta la più “vicina” possibile alle coppie osservate  $(x_i, y_i)$ : la **retta interpolante**, cioè quella che passa tra i punti lasciando da essa scarti complessivamente minimi.
- I punti che stanno sulla retta sono le coppie di punti  $(x_i, \hat{y}_i) = (x_i, a + bx_i)$ , con  $\hat{y}_i$  valori **teorici** o **previsti**, cioè i valori che la variabile  $Y$  dovrebbe assumere per  $X = x_i$  se la relazione tra  $X$  e  $Y$  fosse esattamente quella ipotizzata  $Y = a + bX$ .
- $r_{xy}$  misura quanto bene i dati sono allineati lungo tale retta. Come regola empirica, valori da 0.80 a 1 (o da -1 a -0.80) rivelano una accettabile relazione lineare di tipo diretto (o inverso). Ricordiamo che quando  $r_{xy} = 0$  non è escluso che  $X$  e  $Y$  possono essere legate da altre relazioni, come  $Y = \cos(X) + \exp(X^3)$ , o altre “mostruosità” del genere.

## Minimi quadrati

- Come cerchiamo la retta **interpolante**? Si noti che le quantità  $e_i = y_i - \hat{y}_i$  misurano la **distanza** o **scarto** tra i valori di  $Y$  osservati e quelli teorici. In particolare, prendiamo la distanza **quadratica**, data da  $(y_i - \hat{y}_i)^2$ . Ne consegue che la distanza totale tra i valori osservati e teorici è

$$d(a, b) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - a - bx_i)^2 .$$



## La retta dei minimi quadrati

- La **somma dei quadrati**  $d(a, b) = \sum_{i=1}^n (y_i - a - bx_i)^2$  dipende dalle incognite  $a$  e  $b$ , mentre  $y_i$  e  $x_i$  sono numeri osservati.
- La retta interpolante è quella i cui valori di  $a$  e di  $b$  che rendono minima  $d(a, b)$ , che viene detta **retta dei minimi quadrati**.

Si mostra che i valori  $a$  e  $b$  che minimizzano  $d(a, b)$  sono dati da

$$\hat{b} = \frac{S_{xy}}{S_x^2} \quad \hat{a} = m_y - \hat{b} m_x$$

- I calcoli richiesti sono gli stessi che servono per determinare il coefficiente di correlazione ... non serve molto lavoro in più.
- Sia  $r_{xy}$  sia  $\hat{b}$  dipendono al numeratore dalla covarianza  $S_{xy}$ . Essendo le quantità al denominatore sempre positive, è evidente che i segni di  $r_{xy}$  e di  $\hat{b}$  sono concordi con il segno di  $S_{xy}$ .

## Dimostrazione

Posto  $y_i^* = y_i - bx_i$ ,  $i = 1, \dots, n$ , la somma dei quadrati  $d(a, b)$  può essere riscritta come  $\sum_{i=1}^n (y_i^* - a)^2$ . Quindi, per la proprietà dei minimi quadrati della media aritmetica, la quantità  $\sum_{i=1}^n (y_i^* - a)^2$  è minima per

$$\hat{a} = \frac{1}{n} \sum_{i=1}^n y_i^* = \frac{1}{n} \sum_{i=1}^n (y_i - bx_i) = \frac{1}{n} \sum_{i=1}^n y_i - b \frac{1}{n} \sum_{i=1}^n x_i = m_y - b m_x .$$

Sostituendo tale valore in  $d(a, b)$  si ottiene

$$\begin{aligned} \sum_{i=1}^n (y_i - m_y - bx_i + bm_x)^2 &= \sum_{i=1}^n [(y_i - m_y) - b(x_i - m_x)]^2 \\ &= \sum_{i=1}^n (y_i - m_y)^2 + b^2 \sum_{i=1}^n (x_i - m_x)^2 - 2b \sum_{i=1}^n (y_i - m_y)(x_i - m_x) \\ &= nb^2 S_x^2 - 2nb S_{xy} + nS_y^2 \end{aligned}$$

Come funzione di  $b$ , si tratta di una funzione quadratica, il cui grafico è una parabola con concavità rivolta verso l'alto. Il minimo si ha in corrispondenza del vertice, ossia per

$$\hat{b} = \frac{-(-2nS_{xy})}{2nS_x^2} = \frac{S_{xy}}{S_x^2}$$

## Esempio: FIM

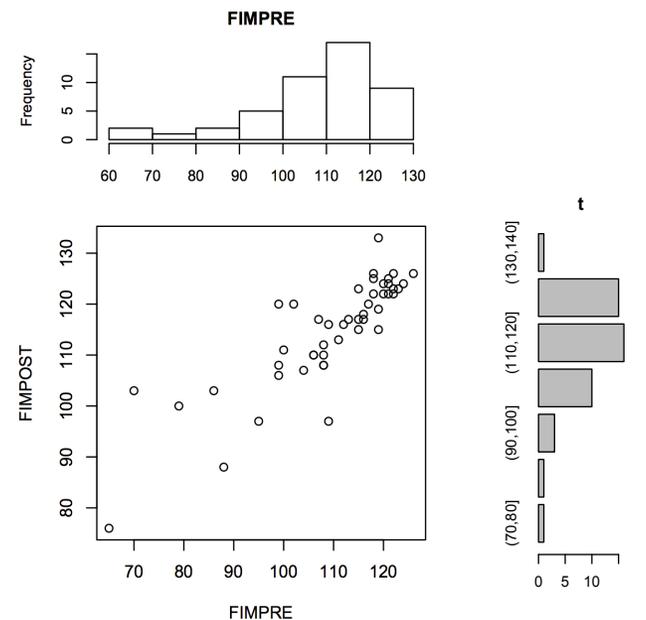
Nell'esempio dalla FIMPOST ( $Y$ ) e FIMPRE ( $X$ ) si trovano i seguenti valori di  $\hat{a}$  e  $\hat{b}$ :

$$\hat{b} = 5808.702/8732.213 = 0.67 \quad \hat{a} = 114.6 - 0.67 \times 109.3 = 41.37$$

La retta di regressione per questi dati è:

$$\hat{Y} = 41.37 + 0.67 X = 41.37 + 0.67 \text{ FIMPRE}$$

Abbiamo il risultato: ma come interpretarlo e usarlo?? La retta è UTILE per fare previsioni sulla variabile risposta. Ad esempio per  $X = 90$ , si trova  $Y = 41.37 + 0.67 \times 90 = 101.67$ .



## Bontà dell'adattamento della retta ai dati

- Come possiamo valutare se la retta si adatta bene ai dati? Abbiamo bisogno di un indice capace di riassumere l'adattamento globale e la capacità esplicativa complessiva del modello in rapporto ai dati osservati.
- Si può utilizzare ancora il coefficiente di correlazione  $r_{xy}$ . E poiché non ha importanza se la correlazione è positiva o negativa, si eleva  $r_{xy}$  al quadrato  $\Rightarrow$  **COEFFICIENTE DI DETERMINAZIONE:**

$$R^2 = r_{xy}^2$$

NOTA:

Se  $R^2 = 1$ : adattamento perfetto (tutti i punti sulla retta)

Se  $R^2 = 0$ : la retta non ha nulla da vedere con i dati

Se  $R^2 = 0.8$ : “buon livello” di adattamento

- ESEMPIO:  $r_{xy} = 0.84^2 \Rightarrow R^2 = 0.71$ , ossia la retta di regressione spiega discretamente la relazione.

## Interpretazione di $R^2$ come proporzione di varianza spiegata

- Siano  $\hat{y}_i = \hat{a} + \hat{b}x_i$ ,  $i = 1, \dots, n$ , i valori calcolati sulla retta dei minimi quadrati.
- La somma dei residui  $y_i - \hat{y}_i$  vale zero.

Infatti,  $\sum_{i=1}^n (y_i - \hat{y}_i) = \sum_{i=1}^n (y_i - \hat{a} - \hat{b}x_i) = \sum_{i=1}^n (y_i - m_y + \hat{b}m_x - \hat{b}x_i) = \sum_{i=1}^n (y_i - m_y) - \hat{b} \sum_{i=1}^n (x_i - m_x) = 0$  (proprietà di baricentro).

- Inoltre,  $\sum_{i=1}^n (y_i - \hat{y}_i)x_i = \sum_{i=1}^n (y_i - \hat{y}_i)(x_i - m_x) = \sum_{i=1}^n (y_i - m_y + \hat{b}m_x - \hat{b}x_i)(x_i - m_x) = nS_{xy} - \hat{b}nS_x^2 = 0$ .

- Allora, dall'identità  $\sum_{i=1}^n (y_i - m_y)^2 = \sum_{i=1}^n (y_i \pm \hat{y}_i - m_y)^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - m_y)^2 + 2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - m_y)$ , usando le due relazioni precedenti, si vede facilmente che l'ultima sommatoria vale zero. Dunque  $\frac{1}{n} \sum_{i=1}^n (y_i - m_y)^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - m_y)^2$  ossia

**VARIANZA TOTALE = VARIANZA RESIDUA + VARIANZA SPIEGATA**

- Si vede infine che  $R^2 = \text{VARIANZA SPIEGATA} / \text{VARIANZA TOTALE}$ .

Infatti,  $\sum_{i=1}^n (\hat{y}_i - m_y)^2 = \sum_{i=1}^n (m_y - \hat{b}m_x + \hat{b}x_i - m_y)^2 = n\hat{b}^2 S_x^2 = nS_{xy}^2 / S_x^2$ . E quindi

$$\frac{\sum_{i=1}^n (\hat{y}_i - m_y)^2}{\sum_{i=1}^n (y_i - m_y)^2} = \frac{nS_{xy}^2}{S_x^2 n S_y^2} = R^2.$$

## Esempio: Tensione, corrente e resistenza

I seguenti dati riportano  $n = 12$  misurazioni della tensione ( $V$ ) e della corrente ( $I$ ):

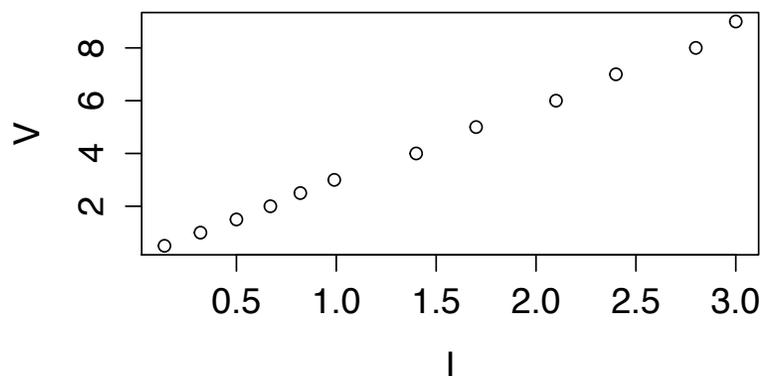
$V = (0.5, 1, 1.5, 2, 2.5, 3, 4, 5, 6, 7, 8, 9)$  in *volt*

$I = (0.14, 0.32, 0.50, 0.67, 0.82, 0.99, 1.4, 1.7, 2.1, 2.4, 2.8, 3)$  in *ampere*

La relazione lineare tra le due variabili è esprimibile come

$$V = a + bI + \text{errore}$$

e ci si attende dal modello teorico  $a \doteq 0$  *volt*,  $b = Res$  *volt/ampere*, dove  $Res$  è una costante di proporzionalità che misura la resistenza, e un valore di  $R^2$  estremamente elevato.



Posto  $X = I$  e  $Y = V$ , si ha:

$$m_x = 1.403 \text{ e } m_y = 4.125$$

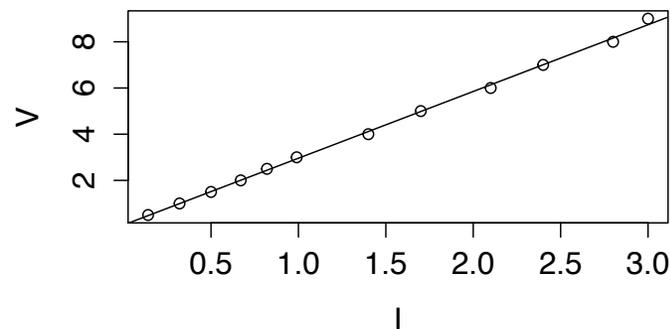
$$S_x^2 = 0.892, S_y^2 = 7.463 \text{ e } S_{xy} = 2.578$$

$$\rightarrow \hat{b} = 2.578/0.892 = 2.89 \text{ volt/ampere e } \hat{a} = 4.125 - 2.89 \times 1.403 = 0.07 \text{ volt.}$$

La retta di regressione per questi dati è:

$$\hat{Y} = 0.07 + 2.89 X$$

Con correlazione  $r_{xy} = 0.999$  ( $R^2 = 0.9985$ ), tale modello evidenzia una relazione lineare tra le due variabili. Inoltre,  $a \doteq 0$  volt come ci si aspettava dal modello teorico, mentre  $Res = 2.89$  volt/ampere.



## Esempio: Intensità luminosa e inverso del quadrato della distanza

I seguenti dati riportano  $n = 8$  misurazioni dell'intensità luminosa della luce di una lampadina ( $Y$ ) raccolta da un sensore a distanza  $d$  e la grandezza  $X = 1/d^2$ :

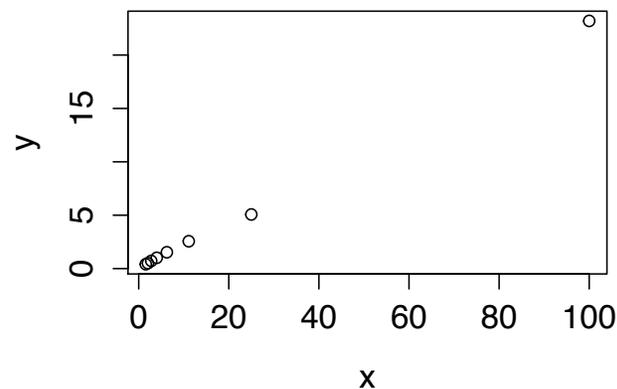
$x = (100, 25, 11.11, 6.25, 4, 2.778, 2.041, 1.563)$

$y = (23.2, 5.07, 2.56, 1.53, 1.01, 0.72, 0.51, 0.41)$

La relazione lineare tra le due variabili è esprimibile come:

$$Y = a + bX + \text{errore}$$

e ci si attende dal modello teorico  $a \doteq 0$ ,  $b = k$ , dove  $k$  è una costante di proporzionalità tale che  $Y = kX$ , e un valore di  $R^2$  estremamente elevato.



Si ha:

$$m_x = 19.09 \text{ e } m_y = 4.38$$

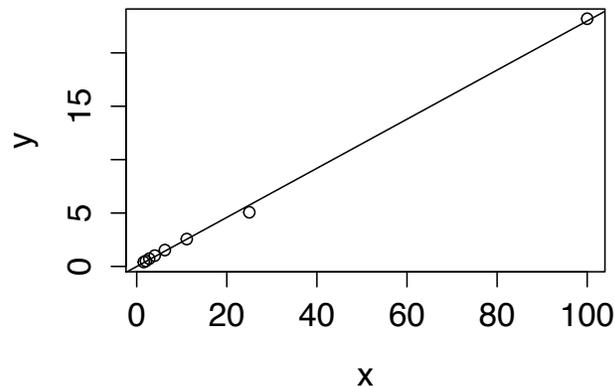
$$S_x^2 = 987.69, S_y^2 = 52.69 \text{ e } S_{xy} = 228.004$$

$$\rightarrow \hat{b} = 228.004/987.69 = 0.23 \text{ e } \hat{a} = 4.38 - 0.23 \times 19.09 = -0.01.$$

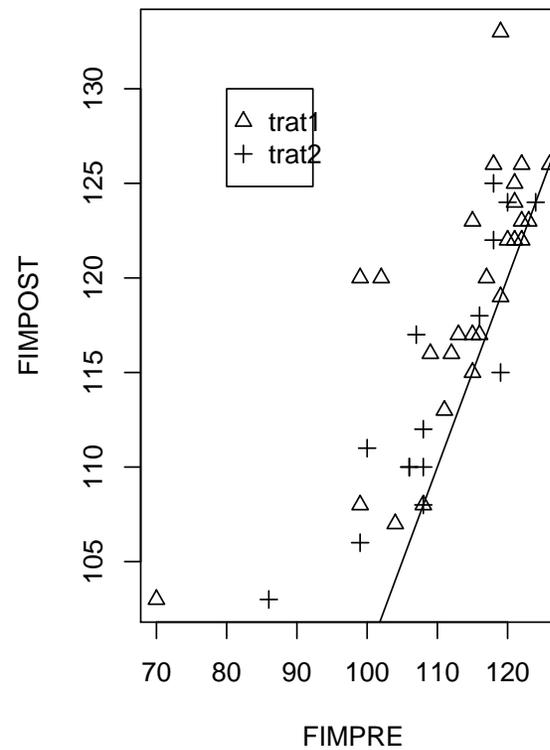
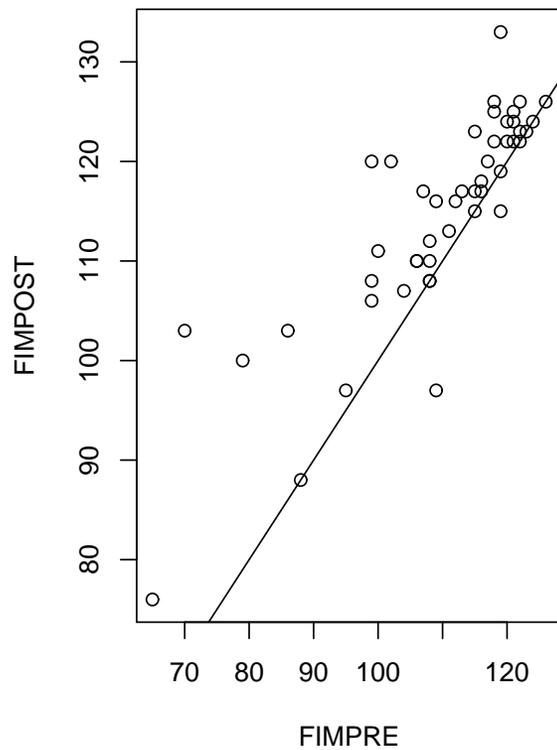
La retta di regressione per questi dati è:

$$\hat{Y} = -0.01 + 0.23 X$$

Con correlazione  $r_{xy} = 0.999$  ( $R^2 = 0.9988$ ), tale modello evidenzia una relazione lineare tra le due variabili. Inoltre,  $a \doteq 0$  come ci si aspettava dal modello teorico, mentre  $k = 0.23$ .



# Esempio: FIM per TRATTAMENTO



- Posto  $Y_R = \text{FIMPOST}$  con realtà virtuale e  $X_R = \text{FIMPRES}$  con realtà virtuale, si ha:

$$m_{x_R} = 113.29 \text{ e } m_{y_R} = 118.93$$
$$S_{x_R}^2 = 129.75, S_{y_R}^2 = 46.38 \text{ e } S_{xy_R} = 58.09$$

La retta di regressione per questi dati è:

$$\hat{Y}_R = 68.18 + 0.45 X_R$$

La correlazione è  $r_{xy_R} \doteq 0.75$ .

- Posto  $Y_F = \text{FIMPOST}$  con fisioterapia e  $X_F = \text{FIMPRES}$  con fisioterapia, si ha:

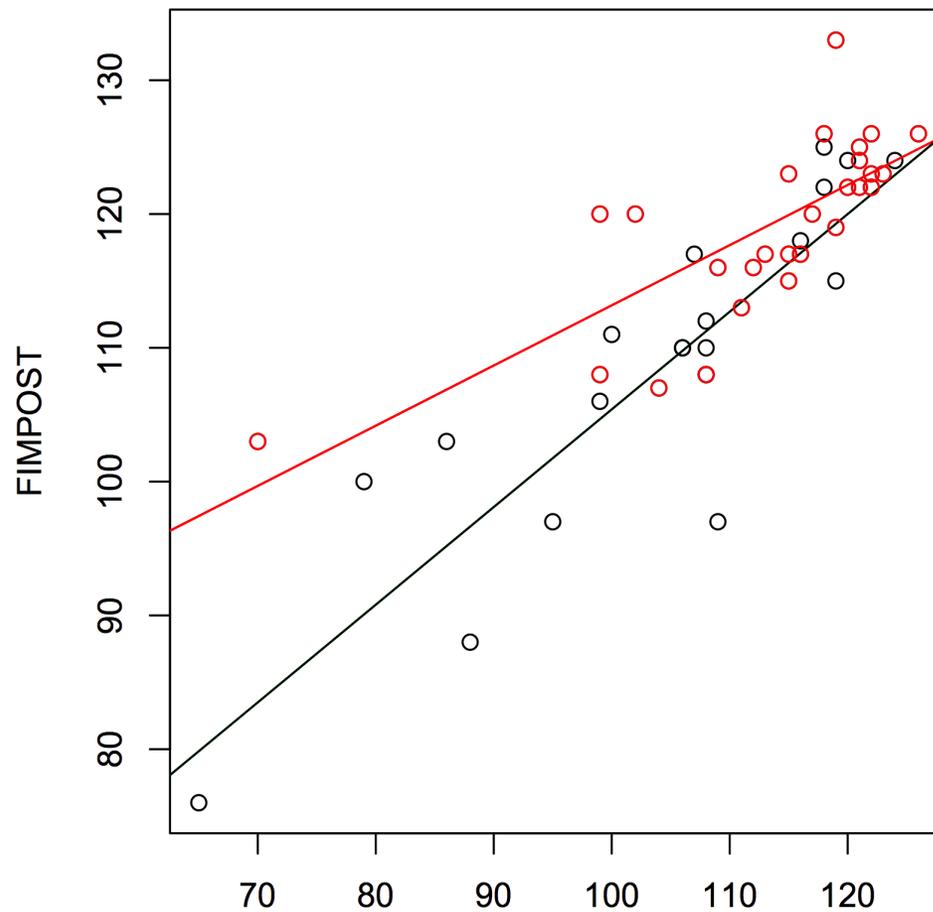
$$m_{x_F} = 103.95 \text{ e } m_{y_F} = 108.65$$
$$S_{x_F}^2 = 229.21, S_{y_F}^2 = 158.66 \text{ e } S_{xy_F} = 168.14$$

La retta di regressione per questi dati è:

$$\hat{Y}_F = 32.40 + 0.73 X_F$$

La correlazione è  $r_{xy_F} \doteq 0.88$ .

## Esempio: FIM per trattamento



## Appendice: proprietà della media e della varianza

### Media

- Linearità:  $m_{a+bx} = a + bm_x$ , con  $a, b \in \mathbb{R}$
- Combinazione lineare:  $m_{ax+by} = am_x + bm_y$ , con  $a, b \in \mathbb{R}$

### Varianza

- Invarianza rispetto a traslazioni:  $S_{a+x}^2 = S_x^2$ , con  $a \in \mathbb{R}$
- Omogeneità (di secondo grado):  $S_{bx}^2 = b^2 S_x^2$ , con  $b \in \mathbb{R}$   
 $\rightarrow S_{a+bx}^2 = b^2 S_x^2$ , con  $a, b \in \mathbb{R}$
- Combinazione lineare:  $S_{ax+by}^2 = a^2 S_x^2 + b^2 m_y^2 + 2ab S_{xy}$ , con  $a, b \in \mathbb{R}$  e  
 $S_{ax-by}^2 = a^2 S_x^2 + b^2 m_y^2 - 2ab S_{xy}$ , con  $a, b \in \mathbb{R}$

## Esercizi

- (1) La gascromatografia è una tecnica per analizzare miscele di gas. I dati che seguono mostrano la quantità di una certa sostanza ( $Y$ ) e la corrispondente misura ottenuta da un gascromatografo ( $X$ ):

quantità	0.25	0.25	0.25	1	1	1	5	5	5	20	20	20
misura	6.55	7.98	6.54	29.7	30	30.1	211	204	212	929	905	922

- 1) Disegnare il diagramma di dispersione dei dati
  - 2) Calcolare la quantità media di sostanza
  - 3) Calcolare la retta di regressione che permette di prevedere la quantità di sostanza come funzione della misura ottenuta dal gascromatografo
  - 4) Calcolare l'indice di correlazione
  - 5) Per una quantità di sostanza pari a 2, il gascromatografo ha fornito una misura pari a?
- (2) La seguente tabella mostra per vari anni il numero di incidenti stradali in una certa regione:

Anno	1991	1992	1993	1994	1995	1996	1997	1998	1999	2000
Incidenti	5413	6122	6705	6824	7790	7698	8571	8688	9422	9904

- 1) Si calcoli il numero medio di incidenti in un anno.
- 2) Si fornisca una rappresentazione grafica dei dati opportuna.
- 3) Si calcoli la retta di regressione che permette di prevedere il numero di incidenti come funzione dell'anno.
- 4) Si calcoli il coefficiente di correlazione.
- 5) Si fornisca una previsione per il numero di incidenti per il 2001.

## Alcuni riferimenti bibliografici

---

- Agresti, A., Finlay, B. (2009). *Statistica per le scienze sociali*. Pearson.
- Agresti, A., Franklin, C. (2013). *Statistics. The Art and Science of Learning from Data*. Pearson.
- Bernstein, S., Bernstein, R. (2003). *Statistica Descrittiva*. McGraw-Hill.
- Bradstreet, T.E. (1996). Teaching introductory statistics courses so nonstatisticians experience statistical reasoning. *The American Statistician*, Vol. 50, 69 – 78.
- Diamond, I., Jefferies, J. (2001). *Introduzione alla statistica per le scienze sociali*. McGraw-Hill.
- Pace, L., Salvan, A. (1996). *Introduzione alla Statistica. I Statistica Descrittiva*. Cedam.
- Rosenthal, J.S. (2005). *Le Regole del Caso: Istruzioni per l'Uso*. Longanesi.

Oppure...

