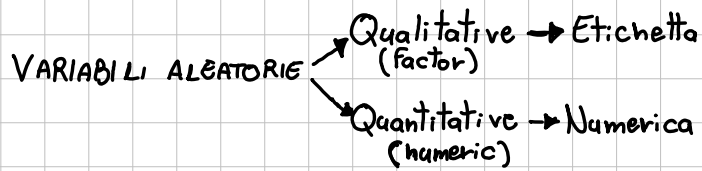


**Statistica Descrittiva:** Descrive ciò che si osserva nei DATI. Fornisce un quadro chiaro e conciso della situazione attuale. Utilizza strumenti come la MEDIA e MEDIANA.

**Statistica Inferenziale:** Va oltre la descrizione dei DATI, si cerca di trarre conclusioni, formulare ipotesi e verificare la significatività dei risultati. Permette di generalizzare i risultati da un CAMPIONE all'intera POPOLAZIONE

**VARIABILI ALEATORIE:** Una VARIABILE ALEATORIA è una funzione che associa ad ogni esperimento casuale un numero reale, viene usata per modellare fenomeni incerti o casuali



### DESCRIZIONE CAMPIONE

1) DISTRIBUZIONE UNITARIA → Elenco di tutte le osservazioni  
 $x_1, x_2, \dots, x_m$   
↳ Ampiezza Campione

2) DISTRIBUZIONE DI FREQUENZE

	FREQUENZE →	ASSOLUTA	RELATIVA	CUMULATA
$x_1$	$m_1$	$m_1$	$f_1 = \frac{m_1}{n}$	$F_1 = f_1$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$x_i$	$m_i$	$m_i$	$f_i = \frac{m_i}{n}$	$F_i = F_{i-1} + f_i$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$x_m$	$m_K$	$m_K$	$f_K = \frac{m_K}{n}$	$F_K = F_{K-1} + f_K$

**FREQUENZA ASSOLUTA:** Numero di tutte le unità statistiche che assumono un certo valore o modalità in relazione ad un carattere. Sono tutti quei dati che assumono quel determinato valore o modalità

Esempio:

25 occhi marroni	freg. assolute
14 " verdi	25
8 " azzurri	14
	8

**FREQUENZA RELATIVA:** Rapporto tra la **FREQUENZA ASSOLUTA** e la numerosità della popolazione o del campione statistico. La somma di tutte le freg. relative deve essere  $= 1$

$$\text{FREQ. RELATIVA} = \frac{\text{FREQUENZA}}{\text{N}^{\circ} \text{tot. Dati}}$$

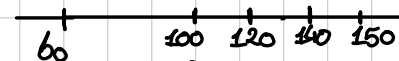
$$\text{oss: } \sum_{i=1}^K m_i = m \quad \sum_{i=1}^K f_i = \sum_{i=1}^K \frac{m_i}{m} = \frac{1}{m} \sum_{i=1}^K m_i = \frac{1}{m} \cdot m = 1$$

**FREQUENZA CUMULATA** → su R `cumsum()`

- `plot()` diagramma a linee
- `barplot()` diagramma a barre ≠ istogramma

Es. distribuzione classi:

$$\begin{aligned} C_0 &= 60 & C_1 &= 100 \\ C_2 &= 120 & C_3 &= 140 \\ C_K &= C_4 & &= 150 \end{aligned}$$



$[C_0, C_1)$

$[C_1, C_2)$

$[C_{K-1}, C_K]$

Freq. assoluta

$m_1$

$m_2$

$m_K$

Freq. Relativa

$$f_1 = \frac{m_1}{K}$$

$$f_2 = \frac{m_2}{K}$$

$$f_K = \frac{m_K}{K}$$

Densità di Frequenza

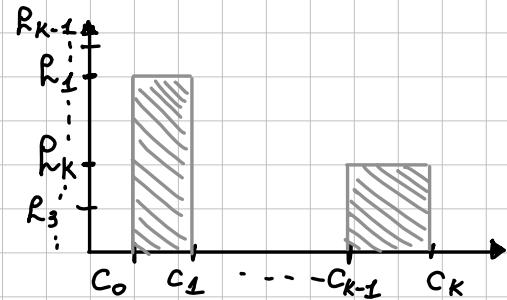
$$h_1 = f_1 / C_1 - C_0$$

$$h_2 = f_2 / C_2 - C_1$$

$$h_K = f_K / C_K - C_{K-1}$$

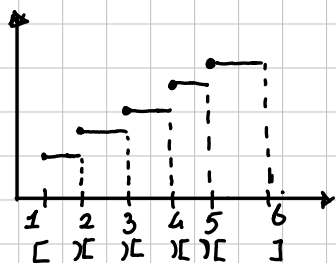
$$\text{DENSITA DI FREQUENZA} = \frac{\text{FREQ}}{\text{AMPIEZZA CLASSE}}$$

Se utilizzo Freq. assoluta ho densità assoluta  
Se utilizzo Freq. relativa ho densità Relativa



# FUNZIONE DI RIPARTIZIONE

	$p_i$	$F_i$
$m_1$	1	0,1
$m_2$	2	0,3
$m_3$	3	0,5
$m_4$	5	0,7
$m_5$	8	0,9
$m_6$	11	1,0



$p_i$  = frazioni di osservazioni uguali a  $m_i$

$F_i$  = frazioni di osservazioni  $p_i$  e  $m_i$

Se abbiamo una v.a. ordinata ci possiamo chiedere quale frazione di osservazioni sia più piccola o uguale ad un valore  $x$

$F_m(x)$  = "frazioni di osservazioni del campione  $\leq x$ "

↳ FORMULA DI RIPARTIZIONE EMPIRICA:

- Costante a tratti

- Continua a destra → gli intervalli sono chiusi a sx

Cioè non DECRESCENTE e  $\lim_{x \rightarrow +\infty} F(x) = 1$

MODA: Modalità più rappresentativa (frequenza maggiore relativa) può essere univoca, biunivoca o non esistere:

$m_i$	$p_i$
3	0,1
7	0,3
11	0,2
15	0,3
18	0,1

Moda = {7, 15}

Possiamo definire la moda delle distribuzioni in classi

	$m_i$	$p_i$	$R_i$
$[1, 5)$	3	3/10	15/200
$[5, 10)$	4	4/10	16/200
$[10, 12)$	2	2/10	20/200
$[12, 17)$	1	1/10	4/200

La MODA è la classe con la DENSITÀ DI FREQUENZA maggiore (il picco dell'istogramma)

$$p_i = C_i - C_{i-1}$$

$$\frac{3}{40} + \frac{4}{50} + \frac{1}{10} + \frac{1}{50} =$$

$$\frac{15 + 16 + 20 + 4}{200}$$

**MEDIANA:** È un valore  $\bar{m}$  tale che almeno la metà delle osservazioni è minore o uguale a  $\bar{m}$  o maggiore o uguale a  $\bar{m}$

1 2 ③ 5 7

↳ Mediana

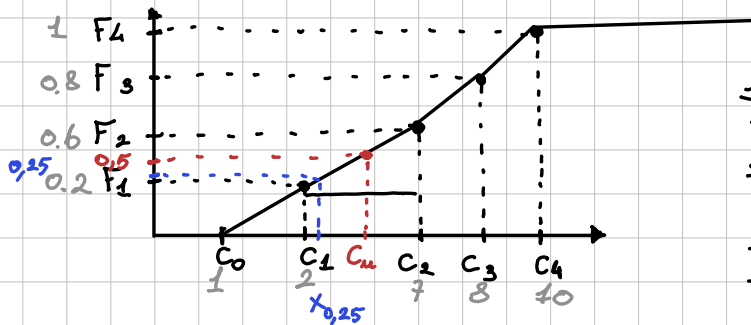
9 ③ 4 7 7

↳ Faccio la media = ⑤ → Mediana

Oss: La mediana può essere un valore che  $\notin$  al campione.

• È un indicatore ROBUSTO → non cambia se variano gli estremi

Mediana a partire dalla funzione di ripartizione della distribuzione in classi di una v.a. continua



$i^*$  = classe in cui si trova la mediana

$$\frac{1}{2} = (\bar{m} - C_{i-1}) \cdot h_i + F_{i-1}$$

$$\frac{1}{2} = (\bar{m} - C_{i^*-1}) h_{i^*} + F_{i^*-1}$$

$$\frac{\frac{1}{2} - F_{i^*-1}}{h_{i^*}} + C_{i^*-1} = \bar{m}$$

**QUANTILI** → "Generalizzazioni della mediana"

$X_p$  valore tale per cui  $p \cdot n$  osservazioni  $\leq X_p$   $(1-p) \cdot n$  osservazioni  $\geq X_p$

Intuitivamente  $X_{0.25}$  valore tale che il 25% delle osservazioni sono  $\leq X_{0.25}$

**MEDIA CAMPIONARIA**

Es: Tempi di attesa ad un servizio

1 3 2 11 8 8 5 6 6 8

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \Rightarrow \bar{x} = \frac{1}{10} \sum_{i=1}^{10} x_i = \frac{1}{10} \cdot 64 = 6.4$$

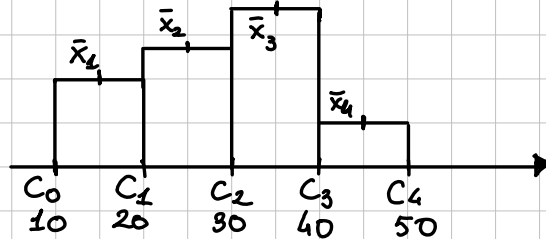
Es: calcoliamo le medie a partire dalla distribuzione delle frequenze relative

$m_i$	$m_i$	$p_i$
1	10	10/38
2	4	4/38
3	8	8/38
4	16	16/38

$$\bar{x} = \frac{1}{n} \sum_{i=1}^k m_i \cdot m_i = \sum_{i=1}^k \frac{m_i}{n} \cdot m_i = \sum_{i=1}^k p_i \cdot m_i$$

Oss:  $\sum_{i=1}^k p_i = 1$

Supponiamo di avere una distribuzione in classi



$$\bar{x} = \frac{\sum_{i=1}^k \bar{x}_i \cdot m_i}{n} = \sum_{i=1}^k p_i \cdot \bar{x}_i \quad \bar{x}_i = \frac{C_i - C_{i-1}}{2}$$

Caratterizzare la dispersione dei dati

- CAMPO DI VARIAZIONE (range)

$$x_{(n)} - x_{(1)}$$

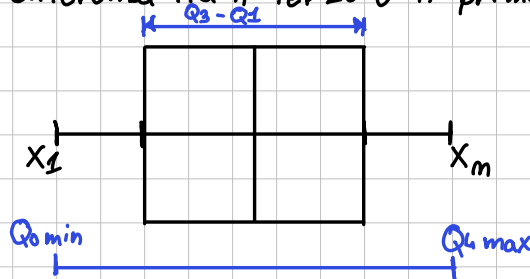
↑  
più grande

↑  
il più piccolo

- SCARTO INTERQUARTILE: differenza tra il terzo e il primo quartile

$Q_3 - Q_1$   
↑  
3° quartile  
75%

↑  
1° quartile  
(25%)



Vogliamo descrivere le distanze tipica dei valori del campione da  $\bar{x}$

$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}) = \frac{1}{n} \sum_{i=1}^n x_i - \frac{1}{n} \sum_{i=1}^n \bar{x} = \bar{x} - \frac{1}{n} \cdot n \bar{x} = 0$$

Scarti

$$\sigma_x^2 = \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 : \text{Varianza della popolazione}$$

$$S_m^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 : \text{Varianza Campionaria}$$

$S_m^2 = \frac{n}{n-1} \cdot \sigma^2$

Oss:

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

$$= \frac{1}{n} \sum_{i=1}^n x_i^2 - \frac{2}{n} \sum_{i=1}^n x_i \bar{x} + \frac{1}{n} \sum_{i=1}^n \bar{x}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \frac{2}{n} \bar{x} \underbrace{\sum_{i=1}^n x_i}_{n\bar{x}} + \frac{1}{n} n \bar{x}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - 2\bar{x}^2 + \bar{x}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2 > 0$$

( $\geq 0$ ,  $= 0 \Leftrightarrow$  tutti i  $x_i = \bar{x}$ )

Esempio: Calcolare la varianza a partire dalla distr. di frequenza

$m_1$	$m_2$	$\vdots$	$m_K$
$\frac{p_1}{m_1}$	$\frac{p_2}{m_2}$	$\vdots$	$\frac{p_K}{m_K}$

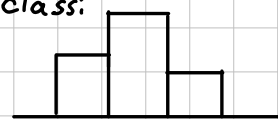
$$S_m^2 = \frac{p}{n-1} \sum_{i=1}^K (m_i - \bar{x})^2 \cdot m_i$$

$$S_m^2 = \frac{n}{n-1} \sigma^2$$

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^K (m_i - \bar{x})^2 \cdot m_i = \sum_{i=1}^K (m_i - \bar{x})^2 \cdot p_i$$

$p_i = \frac{m_i}{n}$

Varianza a partire in classi:



dalla distribuzione

$$\sum_{i=1}^K (\bar{x}_i - \bar{x}) \cdot m_i = S_m^2$$

"medio" dell'i-esima class

medio della distribuzione in classi

(è una scelta convenzionale)

Esempio

Tempi di attesa in secondi: 1s, 2s, 3s

Calcoliamo  $S_m^2$      $n=3$      $\bar{x}=2$      $S_m^2 = \frac{(1-\bar{x})^2 + (2-\bar{x})^2 + (3-\bar{x})^2}{3-1} = 1s^2$

Deviazione Standard (scarto quadratico medio)

$$\sigma = \sqrt{\sigma^2}$$

$$s_m = \sqrt{s_m^2}$$

$$\frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}| \text{ vs } \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}|^2 \geq \left| \frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}| \right|^2 \quad |x_i - \bar{x}|^2 = (x_i - \bar{x})^2$$

$$s_m \geq \frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}|$$

Coefficiente di Variazione

È un indice di dispersione, permette di confrontare misure di fenomeni riferite a unità di misura differenti, in quanto si tratta di una grandezza adimensionale

$$Cv(x) = \frac{\sigma_x}{|\bar{x}|} \quad Cv = \frac{\sigma}{\mu} \rightarrow \begin{array}{l} \text{Deviazione standard} \\ \text{Media Aritmetica} \end{array}$$

1) Non dipende da unità di misura

$$\begin{array}{l} x = \text{secondi} \\ y = \text{minuti} \end{array} \quad y = ax \quad \frac{\sigma_y}{|\bar{y}|} = \frac{|a| \sigma_x}{|a \bar{x}|} = \frac{\sigma_x}{|\bar{x}|}$$

RELAZIONI FRA PIU VARIABILI ALEATORIE (Cerchiamo di descrivere le connessioni/relazioni fra i fenomeni aleatorici)  
"m" osservazioni di coppie di variabili aleatorie

$$\begin{array}{l} (\bar{x}_1, \bar{y}_1) \\ (\bar{x}_2, \bar{y}_2) \\ \vdots \\ (\bar{x}_m, \bar{y}_m) \end{array}$$

x \ y	modalità delle y				
	y <sub>1</sub>	y <sub>2</sub>	y <sub>3</sub>	... y <sub>k</sub>	
x <sub>1</sub>	m <sub>1,1</sub>	m <sub>1,2</sub>	m <sub>1,3</sub>	... m <sub>1,k</sub>	m <sub>1.</sub>
x <sub>2</sub>	m <sub>2,1</sub>	m <sub>2,2</sub>	m <sub>2,3</sub>	... m <sub>2,k</sub>	m <sub>2.</sub>
...	...	...	...	...	...
x <sub>i</sub>	m <sub>i,1</sub>	m <sub>i,2</sub>	m <sub>i,3</sub>	... m <sub>i,k</sub>	m <sub>i.</sub>
...	...	...	...	...	...
x <sub>p</sub>	m <sub>p,1</sub>	m <sub>p,2</sub>	m <sub>p,3</sub>	... m <sub>p,k</sub>	m <sub>p.</sub>
...	...	...	...	...	...
x <sub>m</sub>	m <sub>m,1</sub>	m <sub>m,2</sub>	m <sub>m,3</sub>	... m <sub>m,k</sub>	m <sub>m.</sub>
	m <sub>.1</sub>	m <sub>.2</sub>	m <sub>.3</sub>	... m <sub>.k</sub>	

modalità delle x

Frequenze congiunte assolute

Frequenze assolute

MARGINALI

$m_{i.} = \sum_{j=1}^k m_{i,j}$  (delle x)

$m_{.j} = \sum_{i=1}^m m_{i,j}$  (delle y)

Esempio: x: reddito y: quartiere di residenza (A, B, C)

	A	B	C	
ALTO	7	4	1	12
MEDIO	3	5	11	19
BASSO	2	21	3	26
	12	30	15	

$$p_{i.} = \sum_{j=1}^k p_{i,j} = \frac{m_{i.}}{n}$$

$$p_{.j} = \sum_{i=1}^m p_{i,j} = \frac{m_{.j}}{n}$$

Qual è la proporzione di persone del quartiere A se sappiamo che le persone reddito alto?

$\frac{7}{12}$

Frequenze Condizionate

$$p_{x=i|y=j} = \frac{m_{i,j}}{m_{.j}}$$

$$p_{y=j|x=i} = \frac{m_{i,j}}{m_{i.}}$$

	BLU	VERDE	MARRONE			BLU	VERDE	MARRONE		
ALTO	2	4	14	20	→	ALTO	2/100	4/100	14/100	2/10
MEDIO	3	6	21	30		MEDIO	3/100	6/100	21/100	3/10
BASSO	5	10	35	50		BASSO	5/100	10/100	35/100	5/10
	10	20	70			$\frac{1}{10}$	$\frac{2}{10}$	$\frac{7}{10}$		

CALCOLI LE CONDIZIONATE RISPETTO AL COLORE DEGLI OCCHI

	BLU	VERDE	MARRONE	OSSI Le condizionate coincidono con le marginali;
ALTO	2/70	4/70	14/70	
MEDIO	3/70	6/70	21/70	
BASSO	5/70	10/70	35/70	

PROVIAMO A CALCOLARE LE CONDIZIONATE RISPETTO AL REDDITO

	BLU	VERDE	MARRONE	Anche in questo caso le frequenze condizionate coincidono con le frequenze marginali;
ALTO	2/20	4/20	14/20	
MEDIO	3/30	6/30	21/30	
BASSO	5/50	10/50	35/50	⇒ Le variabili sono dette INDIPENDENTI

Consideriamo due variabili indipendenti:

$$\begin{aligned}
 P_{x=i|y=j} &= \frac{m_{ij}}{m_{\cdot j}} = \frac{m_{ij}}{m} = P_{i\cdot} \\
 P_{y=j|x=i} &= \frac{m_{ij}}{m_{i\cdot}} = \frac{m_{ij}}{m} = P_{\cdot j}
 \end{aligned}
 \left\{ \begin{array}{l} \text{in entrambi i casi abbiamo} \\ m_{ij} = \frac{m_{i\cdot} \times m_{\cdot j}}{m} \end{array} \right.$$

$$\boxed{P_{ij} = P_{i\cdot} \times P_{\cdot j}}$$

Prendiamo questo  
come DEFINIZIONE  
DI INDIPENDENZA



Prendiamo due V.A.  $X, Y$  con frequenze marginali  $p_{i\cdot}, p_{\cdot j}$  con  $i=1, \dots, R$   $j=1, \dots, K$

misuriamo il "grado di dipendenza"  
 Freq. Effettiva  
 $C_{i,j} = m_{i,j} - m_{i,j}^*$   
 Contingenza

				$p_{\cdot 1}$
				$p_{\cdot 2}$
				$p_{\cdot 3}$
$p_{1\cdot}$	$p_{2\cdot}$	$p_{3\cdot}$	$p_{4\cdot}$	

Se fossero indipendenti la distribuzione congiunta  
 $p_{i,j}^* = p_{i\cdot} \times p_{\cdot j} \Leftrightarrow m_{i,j}^* = \frac{m_{i\cdot} \times m_{\cdot j}}{n}$

$$\chi^2 = \sum_{i=1}^K \sum_{j=1}^K \frac{(m_{i,j} - m_{i,j}^*)^2}{m_{i,j}^*} \quad m_{i,j}^* = \frac{m_{i\cdot} \times m_{\cdot j}}{n}$$

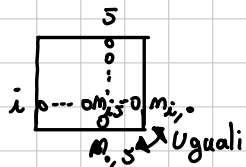
STATISTICA CHI-QUADRO

$$\begin{aligned} \chi^2 &= \sum_{i=1}^K \sum_{j=1}^K \frac{(m_{i,j} - m_{i,j}^*)^2}{m_{i,j}^*} = \sum_{i=1}^K \sum_{j=1}^K \frac{m_{i,j}^2 - 2m_{i,j} \times m_{i,j}^* + m_{i,j}^{*2}}{m_{i,j}^*} = \sum_{i=1}^K \sum_{j=1}^K \frac{m_{i,j}^2}{m_{i,j}^*} - 2 \sum_{i=1}^K \sum_{j=1}^K \frac{m_{i,j} \cdot m_{i,j}^*}{m_{i,j}^*} + \sum_{i=1}^K \sum_{j=1}^K \frac{m_{i,j}^*}{m_{i,j}^*} \\ &= m \sum_{i=1}^K \sum_{j=1}^K \frac{m_{i,j}^2}{m_{i\cdot} \times m_{\cdot j}} - 2m + m = m \left( \sum_{i=1}^K \sum_{j=1}^K \frac{m_{i,j}^2}{m_{i\cdot} \times m_{\cdot j}} - 1 \right) \geq 0 \end{aligned}$$

Def: Una statistica è una qualsiasi funzione del campione

Vogliamo confrontare il valore di  $\chi^2$  con un valore di riferimento. Determiniamo il valore massimo  $\chi_{\max}^2$  che la statistica  $\chi^2$  può assumere.

Si ha dipendenza massima se conoscere  $x$  ci dice anche il valore di  $y$  e viceversa



Calcoliamo  $\chi_{\max}^2$  in un caso di questo tipo

$$m \left( \sum_{i=1}^K \sum_{j=1}^K \frac{m_{i,j}^2}{m_{i\cdot} \times m_{\cdot j}} - 1 \right) \text{ poichè } m_{\cdot j} = m_{i\cdot} \Rightarrow m \left( \sum_{i=1}^K \sum_{j=1}^K \frac{m_{i,j}^2}{m_{i\cdot}^2} - 1 \right)$$

Al massimo ci sono  $\min\{R, K\}$  elementi non nulli:  $\chi^2 = m(\min\{R, K\} - 1)$

$\frac{\chi^2}{\chi_{\max}^2} \leftarrow$  Chi-quadro RELATIVO  $[0, 1]$

Più vicino a 1 e più le variabili sono dipendenti

Vogliamo misurare il grado di dipendenza lineare di due variabili aleatorie



$$\text{Cov}(x, y) \quad \sigma_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y}) = \frac{1}{n} \sum_{i=1}^n (x_i y_i - x_i \bar{y} - \bar{x} y_i + \bar{x} \bar{y}) =$$

$$= \frac{1}{n} \sum_{i=1}^n x_i y_i - \frac{1}{n} \sum_{i=1}^n x_i \bar{y} - \frac{1}{n} \sum_{i=1}^n \bar{x} y_i + \frac{1}{n} \sum_{i=1}^n \bar{x} \bar{y} = \sum_{i=1}^n \frac{x_i y_i}{n} - \bar{x} \bar{y}$$

Ci aspettiamo  $\text{Cov}(x, y) > 0$

$$X \leadsto Z = aX + b \quad z_i = ax_i + b$$

$$Y \leadsto W = cY + d \quad w_i = cy_i + d \quad \text{Cov}(Z, W) = \frac{1}{n} \sum_{i=1}^n (z_i - \bar{z})(w_i - \bar{w}) = \frac{1}{n} \sum_{i=1}^n (ax_i + b - a\bar{x} - b)(cy_i + d - c\bar{y} - d)$$

$$= \frac{1}{n} \sum_{i=1}^n a(x_i - \bar{x})c(y_i - \bar{y}) = a \cdot c \cdot \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = a \cdot c \cdot \text{Cov}(x, y)$$

Esercizio: Siano  $X$  e  $Y$  t.c.  $\text{Cov}(X, Y) = 48$  sia  $\sigma_x = 6$  e sia  $\sigma_y = 8$  quanto vale  $\text{Cov}(\frac{1}{6}X, -\frac{1}{8}Y)$ !

$$\text{Cov}(\frac{1}{6}X, -\frac{1}{8}Y) = \frac{1}{6} \cdot (-\frac{1}{8}) \cdot \text{Cov}(X, Y) = -1 \quad \text{Cov}(X, Y) \leadsto \rho_{xy} = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y}$$

$$\text{oss: Var}(X) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x}) = \text{Cov}(X, X)$$

↑ COEFFICIENTE DI CORRELAZIONE

$$\text{oss: In generale } \text{Var}(\frac{x}{\sigma_x}) = \frac{1}{\sigma_x^2} \cdot \sigma_x^2 = 1 \quad [\text{Var}(aX) = a^2 \text{Var}(X)]$$

Siano  $Z = aX + b$  e  $W = cY + d$

calcoliamo il coeff. di correlazione di  $Z, W$

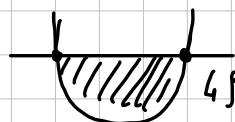
CLAIM:  $-1 \leq \rho_{xy} \leq 1$

PROOF:

Siano  $X, Y$  t.c.  $\bar{x} = 0 \quad \bar{y} = 0, \sigma_x^2 = 1 \quad \sigma_y^2 = 1$

$$0 \leq \frac{1}{n} \sum_{i=1}^n (x_i - t y_i)^2 = \underbrace{\frac{1}{n} \sum_{i=1}^n x_i^2}_{\text{Var}(x)=1} - 2t \underbrace{\frac{1}{n} \sum_{i=1}^n x_i y_i}_{\rho_{xy}} + \underbrace{\frac{1}{n} \sum_{i=1}^n t^2 y_i^2}_{t^2} = 1 - 2t \rho_{xy} + t^2$$

$$\Delta = 4 \rho_{xy}^2 - 4$$



$$| \rho_{xy} | = 1 \Leftrightarrow y = aX + b \text{ e } x = \frac{y}{a} - \frac{b}{a}$$

il valore assoluto del coeff. di corr. non cambia se le variabili aleatorie sono sottoposte ad una tras. affine

$$-1 \leq \rho_{xy} \leq 1$$

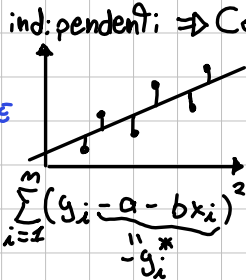
$$\rho_{xy}^2 \leq 1$$

Oss: Se due variabili aleatorie sono indipendenti  $\Rightarrow \text{Cov}(X, Y) = 0 \quad X \perp Y \Rightarrow \text{Cov}(X, Y) = 0$

## REGRESSIONE LINEARE

$$Y = aX + b + \textcircled{Z} \rightarrow \text{TERMINE DI ERRORE}$$

$$y_i = a x_i + b + z_i$$



$\bar{z} = 0$   
omoschedasticità la  
varianza della  $z$  non dipende  
dalla  $x$

Trovare  $a$  e  $b$  in modo tale che sia minimo  $\sum_{i=1}^n (y_i - \underbrace{a - bx_i}_{\hat{y}_i})^2$ : per trovarli poniamo le derivate parziali pari a 0

$$\begin{cases} \frac{d}{da} \sum_{i=1}^n (y_i - a - bx_i)^2 = 0 \\ \frac{d}{db} \sum_{i=1}^n (y_i - a - bx_i)^2 = 0 \end{cases} \Rightarrow \begin{cases} 2 \sum_{i=1}^n y_i - a - bx_i = 0 \\ -2 \sum_{i=1}^n (y_i - a - bx_i) x_i = 0 \end{cases}$$

$$a) \sum_{i=1}^n y_i - a - bx_i = 0 \Leftrightarrow n\bar{y} - na - bn\bar{x} = 0 \Leftrightarrow \bar{y} - a - b\bar{x} = 0 \Rightarrow \hat{a} = \bar{y} - b\bar{x}$$

$$b) \sum_{i=1}^n y_i x_i - a x_i - b x_i^2 = \sum_{i=1}^n y_i x_i - \bar{y} x_i + b x_i \bar{x} - b x_i^2 = \sum_{i=1}^n y_i x_i - n\bar{x}\bar{y} + n b \bar{x}^2 - b \sum_{i=1}^n x_i^2 = 0 \Rightarrow \hat{b} = \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^n x_i^2 - n\bar{x}^2}$$

dividiamo sopra e sotto per  $n$

$$\hat{b} = \frac{\frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x}\bar{y}}{\frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2} = \frac{\sigma_{xy}}{\sigma_x^2} = \frac{\rho_{xy} \cdot \sigma_x \sigma_y}{\sigma_x^2} = \boxed{\rho_{xy} \cdot \frac{\sigma_y}{\sigma_x}}$$

sostituisco  
con  $\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} \Rightarrow \sigma_{xy} = \rho_{xy} \cdot \sigma_x \sigma_y$

$$[\rho_{xy}]^2 > 0.7 \Rightarrow \text{modello "buono"}$$

CAMBIAMENTI DI SCALA  $Y = e^{(a+bX)} \Rightarrow \ln Y = a + bX$

$$Y = (a + bX)^K \Rightarrow (Y)^{1/K} = a + bX$$