

FORMULARIO STATISTICA I

FREQUENZE ASSOLUTE

$$n_1 + n_2 + \dots + n_k = \sum_{j=1}^k n_j = n$$

$$j = 1, \dots, k \quad 0 \leq n_j \leq n$$

FREQUENZE RELATIVE

$$F_j = \frac{\text{FREQUENZE ASSOLUTE}}{\text{NUMEROSITÀ CAMPIONARIA}}$$

$$F_j = \frac{n_j}{n} \quad j = 1, \dots, k$$

FUNZIONE DI RIPARTIZIONE

EMPIRICA

$$F(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(x_i \leq x)$$

↓
FUNZIONE
INDICATRICE

PUNTAZIONE CHE
CONTA

- $x_i \leq x = 1$
- $x_i > x = 0$

• Il massimo è 1

• si sommano le
FREQUENZE RELATIVE
SI CUMULANO

MEDIA ARITMETICA

$$\bar{x} = \frac{\overset{\text{VALORI}}{x_1 + x_2 + \dots + x_n}}{\underset{\text{NUMEROSITÀ CAMPIONARIA}}{n}} = \frac{1}{n} \sum_{i=1}^n x_i \quad \leftarrow \text{PAR OGNI } x$$

$$\bar{x} = \frac{5 + 2 + 3 + 8}{4} = 4,5$$

• compresa tra il MIN e MAX

• all'interno può esserci una media
particolare dei dati → NON CAMBIA

• se $y_i = a + bx_i \quad i = 1, \dots, n$

$$\text{allora } \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = a + b\bar{x}$$

SONO 2 NUMERI QUALSIASI

MEDIA PROPRIETÀ

DEL BARICENTRO

$$\sum_{i=1}^n (x_i - \bar{x}) = (x_1 - \bar{x}) + \dots + (x_n - \bar{x}) = 0$$

• (ci saranno valori negativi sommati a positivi) = 0

• la somma degli SCARTI È = 0

MEDIA PROPRIETÀ

SCARTI QUADRATICI

• MINIMIZZA LA SOMMA DEGLI SCARTI QUADRATICI

ESERCIZIO ESAMPO

distribuzione scartici = 2, 1, 1, 4

$$\bar{x} = \frac{2 + 1 + 1 + 4}{4} = 2$$

SOMMA SCARTI $\rightarrow (2-2)^2 + (1-2)^2 + (1-2)^2 + (4-2)^2 = 6$
→ MEDIA

$$\bar{x} = \text{ARG MIN} \sum_{i=1}^n (x_i - \bar{x})^2$$

NON OTTENGONO MAI UN RISULTATO PIÙ PICCOLO

CALCOLO RICORSIVO

• SE VA AGGIUNTO ALLA MEDIA UN NUOVO DATO

$$\bar{X}_{n+1} = \frac{n}{n+1} \bar{X}_n + \frac{1}{n+1} X_{n+1}$$

ESERCIZIO

$$\bar{X}_9 = 26$$

$$X_{10} = 30$$

$$\bar{X}_{10} = ?$$

$$\bar{X}_{10} = \frac{26 \cdot 9}{9+1} + \frac{1}{9+1} \cdot 30 = 26.4$$

\uparrow \bar{X}_9
 \uparrow X_{9+1}

CONVESSITÀ / CONCAVITÀ

$$F(\lambda x + (1-\lambda)y) \leq \lambda F(x) + (1-\lambda)F(y) \quad \text{CONVESSA} \quad \forall x, y \in (a, b) \text{ e } 0 < \lambda < 1$$

$$\geq \quad \quad \quad \text{CONCAVA}$$

TEOREMA DISUGUAGLIANZA

DI JENSEN

$$\frac{1}{n} \sum_{i=1}^n f(x_i) \geq f\left(\frac{1}{n} \sum_{i=1}^n x_i\right)$$

MEDIA ARITMETICA PONDERATA

$$\bar{X}_w = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$

\rightarrow PESI NUMERICI
 \rightarrow PESI STANDARDIZZATI

APPROSSIMAZIONE DELLA MEDIA

$$\bar{x} \approx \frac{1}{n} \sum_{i=1}^n h_i m_i$$

\rightarrow PUNTO CENTRALE DELL'INTERVALLO

MEDIANA

$$Me = \begin{cases} x_{(\frac{n+1}{2})} & \text{se } n \text{ è dispari} \\ \frac{(x_{(n/2)} + x_{(n/2+1)})}{2} & \text{se } n \text{ è pari} \end{cases}$$

$$F(Me) \approx \frac{1}{2}$$

APPROSSIMAZIONE MEDIANA

$$Me \approx z_{j-1} + (z_j - z_{j-1}) \frac{1/2 - F(z_{j-1})}{F(z_j) - F(z_{j-1})}$$

- IN CASO DI DATI RAGGRUPPATI
- z_{j-1} e z_j = ESTREMI INTERVALLO
- $F(x)$ = FUNZIONE DI RIPARTIZIONE

ESERCIZIO

CLASSI	(0, 1]	(1, 2]	(2, 3]	(3, 4]	(4, 5]
FREQ. ABS.	1	4	4	2	1

$$Me \approx z_{j-1} + (z_j - z_{j-1}) \frac{1/2 - F(z_{j-1})}{F(z_j) - F(z_{j-1})} = 2 + (3-2) \frac{1/2 - 5/12}{9/12 - 5/12} = 2.25$$

ESEMPIO VARIANZA

$k = 3$ (MODALITÀ)

MODALITÀ c_j	4	6	7
FR. ASSOL. h_j	2	8	3

• MEDIA $\bar{x} = (4 \times 2 + 6 \times 8 + 7 \times 3) / 13 = 5,9231$

• MEDIA AL QUADRATO $\bar{x}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 = \frac{1}{n} \sum_{j=1}^k h_j c_j^2 = (2 \times 4^2 + 8 \times 6^2 + 3 \times 7^2) / 13 = 35,9230$

VARIANZA $\sigma^2 = 35,9230 - 5,9231^2 = \underline{0,840}$

PROPRIETÀ

La varianza σ^2 è sempre MAGGIORE O UGUALE a 0
 $\sigma^2 \geq 0$

La varianza $\sigma^2 = 0$ se e solo se i dati sono UGUALI TRA LORO

SCARTO QUADRATICO MEDIO

• è la radice quadrata della varianza

(SCARTO QUADRATICO MEDIO) $= \sigma = \sqrt{\sigma^2}$ (coincide con l'unità di misura dei dati)

CAMPO DI VARIAZIONE

differenza tra MAX e MIN

- CAMPO DI VARIAZIONE $= X_{(n)} - X_{(1)}$

SCARTO INTERQUANTILE

S.I. $= Q_{0,75} - Q_{0,25}$

MAD median absolute deviation

MAD = Mediana $(|x_1 - M_{ex}|), \dots, (|x_n - M_{ex}|)$ $M_{ex} = \text{MEDIANA}(x_1, \dots, x_n)$

COEFFICIENTE DI VARIAZIONE

C.V. $= \frac{\sigma}{\bar{x}}$ \rightarrow SCARTO QUAD MEDIO / \rightarrow MEDIA ARITMETICA $= \sqrt{\frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\bar{x}}\right)^2}$

STANDARDIZZAZIONE DEI DATI

trasformare un insieme di dati in modo tale che abbia MEDIA NULLA

$z_i = \frac{x_i - \bar{x}}{\sigma}$ $i = 1, \dots, n$ **STANDARDIZZATI**
 VARIANZA UNITARIA $\rightarrow \text{VAR}(z) = 1$

DISUGUAGLIANZA DI CHEBYSHEV

Siano x_1, \dots, x_n con \bar{x} e σ . per qualsiasi valore di $k \geq 1$ vale

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1}(|x_i - \bar{x}| > k\sigma) \leq \frac{1}{k^2}$$

DIMOSTRAZIONE

$$\begin{aligned} \underbrace{h\sigma^2}_{\text{VAR}(x) \cdot \text{N. DI OSSERVAZ.}} &= \underbrace{\sum_{i=1}^n (x_i - \bar{x})^2}_{\text{VARIANZA}} = \sum_{i: x_i - \bar{x} < -k\sigma} (x_i - \bar{x})^2 + \sum_{i: |x_i - \bar{x}| \leq k\sigma} (x_i - \bar{x})^2 + \sum_{i: x_i - \bar{x} > k\sigma} (x_i - \bar{x})^2 \\ &\geq \sum_{i: x_i - \bar{x} < -k\sigma} (x_i - \bar{x})^2 + \sum_{i: x_i - \bar{x} > k\sigma} (x_i - \bar{x})^2 \\ &\geq \sum_{i: x_i - \bar{x} < -k\sigma} (k\sigma)^2 + \sum_{i: x_i - \bar{x} > k\sigma} (k\sigma)^2 = \\ &= \sum_{i: |x_i - \bar{x}| > k\sigma} (k\sigma)^2 = k^2 \sigma^2 \sum_{i=1}^n \mathbb{1}(|x_i - \bar{x}| > k\sigma) \end{aligned}$$

STURGES

$$n^\circ \text{ intervalli} = 1 + \log_2 n$$

DIACONIS & FREEMAN

$$n^\circ \text{ intervalli} = \frac{x_n - x_1}{2(Q_{0.75} - Q_{0.25})} n^{1/3}$$

DENSITA' DI UN INTERVALLO

$$d_j = \frac{h_j}{z_j - z_{j-1}} \quad j = 1, \dots, k$$

h_j ← freq. ASSOLUTE
↳ lunghezza intervalli

COSTRUZIONE ISTOGRAMMA

(base rettangoli) = (lunghezza intervalli)

(altre rettogod.) = $\lambda \cdot d_j$

↳ lambda > 0
MUN. QUALSIASI.
λ = 1/h

BAFFI BOXPLOT (lunghezza) OUTLIER

$$\lambda \times (Q_{0.75} - Q_{0.25}) \text{ della scatola}$$

↳ 1.5 costante arbitraria

DISTRIBUZIONE ASIMMETRICA

ASIM. POSITIVA $\bar{x} - Me > 0$

ASIM. NEGATIVA $\bar{x} - Me < 0$

DISTRIBUZIONE SIMMETRICA

media $\bar{x} = Me$ MEDIANA

INDICE DI ASIMMETRIA DI PEARSON

$$V = \frac{1}{\text{SEM}(x)^3} \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3 \ominus \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma} \right)^3$$

- se $\bar{x} = Me \rightarrow$ SIMMETRICO $V = 0$
- se $\bar{x} > Me \rightarrow$ ASIMMETRICA POSITIVA $V > 0$
- se $\bar{x} < Me \rightarrow$ ASIMMETRICA NEGATIVA $V < 0$

INDICE DI ASIMMETRIA DI BOWLEY

$$B = \frac{(Q_{0.75} - Me) + (Q_{0.25} - Me)}{Q_{0.75} - Q_{0.25}}$$

$B = 0 \rightarrow$ simmetria intorno a Me

$B > 0 \rightarrow$ asimmetria positiva

$B < 0 \rightarrow$ asimmetria negativa

INDICE DI CURTOSI DI PEARSON

$$K = \frac{1}{\text{MOM}(x)} + \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^4 \div \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma} \right)^4$$

$K \geq 0$

$K = 0$ se i dati sono COSTANTI

↳ rapporto tra due indici di variabilità

INDICE DI MUTABILITÀ DI GINI

$$G = \sum_{j=1}^k f_j(1-f_j) = 1 - \sum_{j=1}^k f_j^2$$

↳ frequenze relative

MINIMA MUTABILITÀ

$$G = 1 - \sum_{j=1}^k f_j^2 = 1 - (0^2 + \dots + 1^2 + \dots + 0^2) = 1 - 1 = 0$$

MASSIMA MUTABILITÀ

$$G = 1 - \sum_{j=1}^k \frac{1}{k^2} = 1 - \frac{k}{k^2} = 1 - \frac{1}{k} = \frac{k-1}{k}$$

DISTANZA DI HAMMING

(fra x_i e x_j)

$$\mathbb{1}(x_i \neq x_j) = \begin{cases} 0 & \text{se } x_i = x_j \\ 1 & \text{se } x_i \neq x_j \end{cases}$$

DEMOSTRAZIONE INDICE DI GINI

$$\begin{aligned} G &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \mathbb{1}(x_i \neq x_j) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^k h_j \mathbb{1}(x_i \neq c_j) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^k n_i h_j \mathbb{1}(c_i \neq c_j) \\ &= \frac{1}{n^2} (0 \times n_1 h_1 + h_1 h_2 + \dots + 1 \times n_1 h_k + h_2 h_1 + 0 \times h_2 h_2 + \dots + h_2 h_k + \dots + 0 \times h_k h_k) = \\ &= \frac{1}{n^2} [n_1(n - n_1) + n_2(n - n_2) + \dots + n_k(n - n_k)] = \\ &= \frac{1}{n^2} \sum_{j=1}^k n_j(n - n_j) = 1 - \frac{1}{n^2} \sum_{j=1}^k n_j^2 = 1 - \sum_{j=1}^k f_j^2 \end{aligned}$$

GINI ALTERNATIVO

$$G = 1 - \frac{1}{n^2} \sum_{j=1}^k n_j^2$$

GINI NORMALIZZATO

$$G_{\text{norm}} = \frac{k}{k-1} G$$

ENTROPIA DI SHANNON

$$H = - \sum_{j=1}^K f_j \log f_j$$

↳ freq. relative

SHANNON NORMALIZZATA

$$H_{norm} = \frac{H}{\log K}$$

COVARIANZA

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

TRASFORMAZIONE LINEARE COV

$$\text{cov}(v, w) = b_x b_y (x, y)$$

VARIANZA DI UNA SOMMA DI DUE VARIABILI

$$\text{var}(w) = \text{var}(x) + \text{var}(y) + 2 \text{cov}(x, y)$$

VARIANZA CON COVARIANZA

$$\text{var}\left(\frac{x}{\sigma_x} + \frac{y}{\sigma_y}\right) = 2 \left(1 + \frac{1}{\sigma_x \sigma_y} \text{cov}(x, y)\right)$$

COEFFICIENTE DI CORRELAZIONE

$$r_{xy} = \rho = \frac{\text{cov}(x, y)}{\text{sqm}(x) \text{sqm}(y)} = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma_x}\right) \left(\frac{y_i - \bar{y}}{\sigma_y}\right)$$

E' LA COVARIANZA STANDARDIZZATA

FUNZIONE CHI MINIMIZZA LA PERDITA

$$l(\alpha, \beta) = \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2$$

REGRESSIONE LINEARE

$$y \approx \alpha + \beta x + e_i$$

$(\hat{\alpha}, \hat{\beta})$ = valori che minimizzano la funzione d. per

le cose prese sono il MINO POSSIBILE

$$(\hat{\alpha}, \hat{\beta}) = \underset{\alpha, \beta}{\text{argmin}} \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2 = \begin{cases} \hat{\alpha} = \bar{y} - \hat{\beta} \bar{x} \\ \hat{\beta} = \frac{\text{cov}(x, y)}{\text{var}(x)} \end{cases}$$

valore = $\alpha + \beta(\text{diametro}) + (\text{errore})$

RESIDUI

$$r_i = y_i - (\hat{\alpha} + \hat{\beta} x_i) \quad i = 1, \dots, n$$

- VARIANZA DEI RESIDUI UTILIZZATA A VALUTARE LA BONTA' DI ADATTAMENTO DEL MODELLO AI DATI
- LA MEDIA DEI RESIDUI E' NULLA

COEFFICIENTE DI DETERMINAZIONE

R^2

$$R^2 = 1 - \frac{\text{var}(x)}{\text{var}(y)} = \text{cov}(x, y)^2$$

$$0 \leq R^2 \leq 1$$

COEFFICIENTE DI DETERMINAZIONE 2

$$R^2 = \frac{\text{cov}(x, y)^2}{\text{var}(x) \text{var}(y)} = \text{cov}(x, y)^2$$

CORRELAZIONE PARZIALE

$$\tilde{x} = \sqrt{x, z} = \frac{cov(x, z)}{var(z)} (z_i - \bar{z})$$

$$\tilde{y} = y, z \rightarrow y_i - \bar{y} - \frac{cov(y, z)}{var(z)} (z_i - \bar{z})$$

MEDE DEI GRUPPI

$$\bar{x}_j = \frac{1}{n_j} \sum_{i=1}^{n_j} x_{ij} \quad j = 1, \dots, k$$

MEDE TOTALE

$$\bar{x} = \frac{1}{n} \sum_{j=1}^k n_j \bar{x}_j = \frac{1}{n} \sum_{j=1}^k \sum_{i=1}^{n_j} x_{ij}$$

1° GRUPPO 2° GRUPPO

DEVIANZA

$$D_{tr}^2 = \sum_{j=1}^k n_j (\bar{x}_j - \bar{x})^2$$

SE = 0 ALGUNA DEVIANZA ZERO

DEVIANZA INTRO I GRUPPI

$$d_j^2 = \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)^2 \quad j = 1, \dots, k$$

$$D_{en}^2 = \sum_{j=1}^k d_j^2 = \sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)^2$$

DEVIANZA TOTALE

$$D^2 = \sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ij} - \bar{x})^2$$

NON C'È

• SCARTI DI TUTTE LE OSSERVAZIONI DALLA MEDIA TOTALE

TROVARE SCOMPOSIZIONE DEVIANZA

$$D^2 = D_{tr}^2 + D_{en}^2$$

FRA I GRUPPI INTRO I GRUPPI

$$0 \leq D_{tr}^2 \leq D^2$$

MASSIMO TRONCO

RAPPORTO DI CORRELAZIONE

$$\eta^2 = \frac{\text{(DEVIANZA FRA I GRUPPI)}}{\text{(DEVIANZA TOTALE)}} = 1 - \frac{\text{(DEVIANZA INTRO I GRUPPI)}}{\text{(DEVIANZA TOTALE)}}$$

$$\eta^2 = \frac{D_{tr}^2}{D^2} = 1 - \frac{D_{en}^2}{D^2}$$

INDICE NORMALIZZATO POIKHAR
 $0 \leq \eta^2 \leq 1$

• MISURA LA FORZA DELLA DIPENDENZA IN MEDIA

INDICE χ^2 DI PEARSON

$$\chi^2 = \sum_{i=1}^h \sum_{j=1}^k \frac{(n_{ij} - \hat{n}_{ij})^2}{\hat{n}_{ij}} = n \left(\sum_{i=1}^h \sum_{j=1}^k \frac{F_{ij}^2}{E_i \cdot F_{+j}} - 1 \right)$$

MODI CONDIZIONATE

$$\bar{y}_i = \frac{1}{n_{i\cdot}} \sum_{j=1}^k n_{ij} d_j \quad i=1..h$$

DIFFERENZA SEMPLICE MEDIA

$$\Delta = \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j=1}^n |x_i - x_j|$$

$$\text{OPPURE } \Delta = \frac{2}{n(n-1)} \left(\sum_{i=1}^n i \cdot x_{(i)} \right) - 2\bar{x} \frac{n+1}{n-1}$$

TRASFORMAZIONE LINEARE Δ

$$\Delta_y = |b| \Delta_x$$

 Δ_{MAX}

$$R = \frac{\Delta}{2\bar{x}}$$

$$0 \leq R \leq 1$$

MAX e MIN Δ

$$\Delta_{\text{MAX}} = 2\bar{x}$$

$$\Delta_{\text{MIN}} = 0$$

TRASFORMAZIONE LINEARE R

$$R_x = R_y$$

CURVA \Rightarrow LORENZ

R con AREA

$$q_i = \frac{1}{S} \sum_{j=1}^i x_j$$

$$R = \frac{n-1}{n}$$