

- (HABBIAMO MEDIA CAMPIONARIA (o SEMPLICEMENTE MEDIA) DI n DATI NUMERICI $\{x_1, \dots, x_n\}$ LA QUANTITÀ

$$\bar{x} = \frac{1}{n} \sum_{j=1}^n x_j$$

ESEMPIO LE MISURE OTTENUTE SU UN CAMPIONE SONO

18 6 31 71 84 17 23 1 9 43

ABBIAMO $n=10$ E UN SEMPLICE CONTO FORNISCE $\bar{x}=30,3$.

MEDIA PER DATI RAGGRUPPATI. SE I DATI SONO STATI

GIÀ SUDDIVISI IN CLASSI E SI CONOSCE SOLO LA FREQUENZA,

ALLORA SI PUÒ DEFINIRE LA MEDIA IN QUESTO MODO:

SI ASSUME CHE I DATI NELLA SINGOLA CLASSE SIANO

DISTRIBUITI IN MODO UNIFORME E CHE QUINDI IL LORO

VALORE MEDIO COINCIDA CON IL VALORE MEDIO DELL'INTER-

VALLINO CHE REALIZZA LA CLASSE. SIA \bar{x}_k TALE VALORE

PER LA CLASSE k -ESIMA. ASSUMENDO CHE ABBIAMO N

CLASSI, \bar{x} È OTTENUTO ATTRAVERSO UNA MEDIA Ponderata

DEGLI \bar{x}_k CON PESI $f_a(k)$:

$$\bar{x} = \sum_{k=1}^N f_r(k) \bar{x}_k = \frac{1}{n} \sum_{k=1}^N f_a(k) \bar{x}_k$$

TORNANDO ALL'ESEMPIO CON IL PESO DEGLI STUDENTI, ABBIAMO

$$\bar{p} = \frac{5 \cdot 61 + 18 \cdot 64 + 42 \cdot 67 + 27 \cdot 70 + 8 \cdot 73}{100} = 67,45 \text{ kg}$$

QUALCHE SEMPLICE PROPRIETÀ

- SE APPLICHIAMO UNA TRASFORMAZIONE LINEARE AI DATI

$$y_j = ax_j + b \Rightarrow \bar{y} = a\bar{x} + b$$

CIOÈ LA MEDIA CALCOLATA CON LA TRASFORMAZIONE

- LA SOMMA DEGLI SCARTI DALLA MEDIA È NULLA

$$\sum_{j=1}^M (x_j - \bar{x}) = 0$$

- LA SOMMA DEI QUADRATI DAGLI SCARTI DELLA MEDIA È

MINIMA:

$$\sum_{j=1}^M (x_j - x)^2 \text{ ASSUME VALORE MINIMO PER } x = \bar{x}$$

OSSERVAZIONI MEDIA E MEDIANA NON COINCIDONO

NECESSARIAMENTE, MA SONO TANTO PIÙ VICINE QUANTO PIÙ

I DATI SONO DISPOSTI UNIFORMEMENTE. LA MEDIA È PIÙ

FACILE DA CALCOLARE, MA DIPENDENDO DA I VALORI DI TUTTI I

DATI RISULTA SENSIBILE ALLA PRESENZA DI ANOMALIE.

LA MEDIANA È PENO SENSIBILE, E PUÒ RISULTARE UGUALE TRA DUE SET DI DATI ANCHE SE QUESTI DIFFERISCONO NOTEVOLMENTE.

IL PRINCIPALE INDICE DI DISPERSIONE È LA VARIANZA:
DATI n DATI NUMERICI: $\{x_1, \dots, x_n\}$, È LA QUANTITÀ

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

LA RADICE QUADRATA DELLA VARIANZA SI CHIAMA DEVIAZIONE STANDARD O SCARTO QUADRATICO MEDIO

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

VARIANZA E DEVIAZIONE STANDARD SONO QUANTITÀ NON NEGATIVE, E SI ANNULLANO SOLO QUANTO $x_1 = \dots = x_n = \bar{x}$.

INTUITIVAMENTE, PIÙ È PICCOLA LA VARIANZA PIÙ I VALORI SI CONCENTRANO VICINO A \bar{x} (PER QUESTO LA VARIANZA È UN INDICE DI DISPERSIONE).

PER IL CALCOLO DELLA VARIANZA SI PUÒ USARE LA SEGUENTE FORMULA DI FACILE VERIFICA

FORMULA, DI FACILE VERIFICA

$$\sigma^2 = \left(\frac{1}{n} \sum_{j=1}^M x_j^2 \right) - (\bar{x})^2$$

OSSERVAZIONE SE I DATI SONO DIMENSIONALI (PER ESEMPIO FORNITI IN METRI, KG ECC) ALLORA MEDIA E DEVIAZIONE STANDARD HANNO STESSA DIMENSIONE.

VARIANZA PER DATI RAGGRUPPATI. SE I DATI SONO STATI GIÀ SUDDIVISI IN CLASSI E SI CONOSCE SOLO LA FREQUENZA, IL CALCOLO ESATTO NON È POSSIBILE VISTO CHE ABBIAMO PERSO PARTE DELL'INFORMAZIONE. CALCOLIAMO ALLORA LA VARIANZA APPROSSIMANDO I DATI IN QUESTO MODO: SOSTITUIAMO A I DATI NELLA CLASSE k -ESIMA (CHE NON CONOSCIAMO) $f_a(k)$ VOLTE IL VALORE MEDIO DELL'INTERVALLINO CHE INDIVIDUA LA CLASSE, CHIAMIAMOLO \bar{x}_k . IN FORMULE

$$\sigma^2 = \frac{1}{n} \sum_{k=1}^N f_a(k) (\bar{x}_k - \bar{x})^2 = \frac{1}{n} \left(\sum_{k=1}^N f_a(k) \bar{x}_k^2 \right) - (\bar{x})^2$$

BADARE CHE È COME AVERE n DATI, QUINDI DIVIDIAMO PER n E NON PER IL NUMERO DI CLASSI N .

ESEMPIO Si consideri la seguente tabella relativa alle frequenze dei pesi in Kg di 100 individui adulti

	Peso p in Kg	f_a	
$\bar{x}_1 = 52,5$	$50 \leq p < 55$	20	$N = 100$
$\bar{x}_2 = 57,5$	$55 \leq p < 60$	15	$N = 6$
$\bar{x}_3 = 62,5$	$60 \leq p < 65$	18	
$\bar{x}_4 = 67,5$	$65 \leq p < 70$	22	
$\bar{x}_5 = 72,5$	$70 \leq p < 75$	18	
$\bar{x}_6 = 77,5$	$75 \leq p \leq 80$	7	

Abbiamo

$$\bar{p} = \frac{1}{100} (20 \cdot 52,5 + 15 \cdot 57,5 + 18 \cdot 62,5 + 22 \cdot 67,5 + 18 \cdot 72,5 + 7 \cdot 77,5)$$

$$= 63,7 \text{ kg} \quad [\text{peso medio}]$$

$$(\bar{x}_1 - \bar{p})^2 = 125,44 \quad (\bar{x}_2 - \bar{p})^2 = 38,44 \quad (\bar{x}_3 - \bar{p})^2 = 1,44$$

$$(\bar{x}_4 - \bar{p})^2 = 14,44 \quad (\bar{x}_5 - \bar{p})^2 = 77,44 \quad (\bar{x}_6 - \bar{p})^2 = 190,44$$

$$\sigma^2 = \frac{1}{100} (20 \cdot 125,44 + 15 \cdot 38,44 + 18 \cdot 1,44 + 22 \cdot 14,44 + 18 \cdot 77,44 + 7 \cdot 190,44)$$

$$\approx 61,56 \text{ kg}^2 \quad \sigma \approx 7,85 \text{ kg} \quad [\text{deviazione standard}]$$

TORNIAMO ALLA PROBABILITÀ. DEFINIREMO DOUE NUOVE DENSITÀ PER V.A. ASSOLUTAMENTE CONTINUE. SARANNO MOLTO UTILI IN STATISTICA.

SIANO X_1, \dots, X_M M VARIABILI ALGEBRICHE INDIPENDENTI E TUTTE CON DENSITÀ GAUSSIANA STANDARD, CIOÈ $X_k \sim N(0,1)$ $k=1, \dots, M$. PONIAMO $Y = \sum_{k=1}^M X_k^2$. LA DENSITÀ DI Y VIENE

CHIAMATA DENSITÀ DEL χ^2 QUADRATO AD M GRADI DI LIBERTÀ ED È INDICATA CON IL SIMBOLO $\chi^2(M)$.

ABBIAMO VISTO IN UNA DELLE LEZIONI SCORSE CHE SE $X \sim N(0,1)$ ALLORA $Y = X^2$ HA DENSITÀ

$$g(t) = \begin{cases} \frac{1}{\sqrt{2\pi t}} e^{-t/2} & \text{PER } t > 0 \\ 0 & \text{PER } t \leq 0 \end{cases}$$

CIOÈ $g = \Gamma(\frac{1}{2}, \frac{1}{2})$. QUESTO SIGNIFICA CHE $\chi^2(1) = \Gamma(\frac{1}{2}, \frac{1}{2})$.

ABBIAMO ANCHE VISTO CHE SOMMANDO V.A. INDIPENDENTI DI TIPO Γ SI OTTIENE ANCORA UNA V.A. DI TIPO Γ . PIÙ PRECISAMENTE, SE $Y_k \sim \Gamma(\alpha_k, \lambda)$, ALLORA $\sum_{k=1}^M Y_k \sim \Gamma(\sum_{k=1}^M \alpha_k, \lambda)$

NE DEDUCIAMO QUINDI CHE

$$\chi^2(M) = \Gamma\left(\frac{M}{2}, \frac{1}{2}\right)$$

IN PARTICOLARE, PER QUANTO VISTO SULLE V.A. DI TIPO Γ

- SE $Y \sim \chi^2(m)$ ALLORA $E[Y] = m$ E $VAR Y = 2m$ (IN QUANTO SE $Y \sim \Gamma(\alpha, \lambda)$ ALLORA $E[Y] = \alpha/\lambda$ E $VAR Y = \alpha/\lambda^2$).

IL CALCOLO DI $E[Y]$ PUO' ESSERE FATTO ANCHE DIRETTAMENTE

A MANO: TENUTO CONTO CHE $E[X_k] = 0$, ABBIAMO

$$E[X_k^2] = VAR X_k = 1, \text{ DA CUI } E[Y] = \sum_{k=1}^m E[X_k^2] = m.$$

- SE $Y_1 \sim \chi^2(m_1)$ E $Y_2 \sim \chi^2(m_2)$, E SONO INDIPENDENTI,

ALLORA

$$Y_1 + Y_2 \sim \chi^2(m_1 + m_2)$$

QUESTO E' DOVUTO AL RISULTATO SULLA SOMMA DI V.A. DI TIPO Γ .

PER m GRANDE (DICIAMO $m > 30$) POSSIAMO APPROSSIMARE

UNA V.A. DI TIPO $\chi^2(m)$ CON UNA DI TIPO GAUSSIANO. INFATTI,

SE $\{X_k\}_{k \in \mathbb{N}}$ E' UNA SUCCESSIONE DI V.A. INDIPENDENTI DI

TIPO $\chi^2(1)$, ALLORA $Y_m = \sum_{k=1}^m X_k \sim \chi^2(m)$ E PER IL TEOREMA

DEL LIMITE CENTRALE

$$Y_m \simeq Y \text{ CON } Y \sim N(m, 2m)$$

(LE X_k HANNO MEDIA 1 E VARIANZA 2), CIOE' IL COMPORTA-

MENTO DELLA V.A. Y_m CON DENSITA' $\chi^2(m)$ E' ASINTOTICA-

MENTE QUELLO DI UNA V.A. Y CON DENSITA' GAUSSIANA

$N(m, 2m)$. In particolare

$$P\{Y_m \leq x\} \simeq \Phi\left(\frac{x-m}{\sqrt{2m}}\right) \quad \forall x \in \mathbb{R}$$

Come al solito Φ indica la funzione di ripartizione associata alla Gaussiana.

DATA UNA SUCCESSIONE $\{X_k\}_{k \in \mathbb{N}}$ DI VARIABILI ALGEBRAICHE SU UN CERTO SPAZIO DI PROBABILITÀ (Ω, \mathcal{A}, P) , CHIAMIAMO VARIANZA CAMPIONARIA LA V.A.

$$S_m^2 := \frac{1}{m-1} \sum_{k=1}^m (X_k - \bar{X}_m)^2 \quad \left[\begin{array}{l} \text{BEN DEFINITA} \\ \text{SE } m \geq 2 \end{array} \right]$$

AL SOLITO $\bar{X}_m = \frac{1}{m} \sum_{k=1}^m X_k$ INDICA LA MEDIA CAMPIONARIA.

PROPOSIZIONE SUPPONIAMO CHE LE X_k SIANO INDIPENDENTI E CON DISTRIBUZIONE GAUSSIANA $N(\mu, \sigma^2)$. ALLORA

$$(i) \quad \sum_{k=1}^m \left(\frac{X_k - \mu}{\sigma} \right)^2 \sim \chi^2(m)$$

$$(ii) \quad \sum_{k=1}^m \left(\frac{X_k - \bar{X}_m}{\sigma} \right)^2 = \frac{(m-1)S_m^2}{\sigma^2} \sim \chi^2(m-1)$$

(iii) LE V.A. S_m^2 E \bar{X}_m SONO TRA LORO INDIPENDENTI

PROOF PARZIALE

IL PUNTO (i) SEGUE DAL FATTO CHE $\frac{X_n - \mu}{\sigma} \sim N(0,1)$ E DALLA DEFINIZIONE DI $\chi^2(n)$.

OSTENDIAMO LA DIMOSTRAZIONE DEGLI ALTRI DUE PUNTI, IN QUANTO DELICATA. OSSERVIAMO CHE (ii) DIFFERISCE DA (i) IN QUANTO AL POSTO DELLA MEDIA μ PRESENTA LA MEDIA CAMPIONARIA.

IL COMPORTAMENTO È DIVERSO PERCHÉ LE n V.A. $X_k - \bar{X}_n$ NON SONO INDIPENDENTI: LA LORO SOMMA È NULLA E QUINDI SONO RELAZIONATE.

ESEMPIO UNA DITTA CONFEZIONA KIWI. DALLO STORICO È NOTO CHE LA VARIANZA DELLE DIMENSIONI DEL FRUTTO È $1,26 \text{ cm}^2$. DOVENDO FORNIRE FRUTTI CON SIMILI DIMENSIONI, LA DITTA SCARTA UNA PARTITA DI FRUTTI SE LA VARIANZA CAMPIONARIA DI 40 PEZZI SCELTI A CASO SUPERA 2.

ASSUMENDO CHE LA DIMENSIONE SEGUA UNA LEGGE NORMALE CON VARIANZA $\sigma^2 = 1,26$, CHE PROBABILITÀ C'È CHE LA PARTITA VENGA SCARTATA?

INDICHIANO CON X_k , $k=1, \dots, 40$, LE DIMENSIONI DEI PEZZI

SCELTI. SIANO ASSUNTO $X_n \sim \mathcal{N}(\mu, \sigma^2)$. SAPPIAMO CHE LA VARIANZA CAMPIONARIA S_n^2 SODDISFA

$$\frac{39 S_n^2}{1,26} = \frac{(n-1) S_n^2}{\sigma^2} \sim \chi^2(n-1) = \chi^2(39)$$

QUINDI

$$P\{S_n^2 > 2\} = P\left\{\frac{39 S_n^2}{1,26} > 61,9\right\} = 1 - P\left\{\frac{39 S_n^2}{1,26} \leq 61,9\right\} \approx 0,01123$$

CON L'ULTIMO STEP OTTENUTO PER COMPUTAZIONE NUMERICA DELLA FUNZIONE DI RIPARTIZIONE

$$F(x) = \int_{-\infty}^x \chi^2(39) dt \quad \text{in } x = 61,9$$

VEDIAMO UN'ULTERIORE TIPO DI V.A. ANCHE QUESTA, COME LA $\chi^2(n)$, TORNA UTILE IN STATISTICA.

SIANO $X \sim \mathcal{N}(0,1)$ E $Y \sim \chi^2(n)$ DUE VARIABILI ALGEBRAICHE INDIPENDENTI. PONIAMO $T = \frac{X}{\sqrt{Y}}$. OSSERVIAMO CHE

ESSENDO $\chi^2(n)$ NULLA SULLA SEMIRETTA $(-\infty, 0]$, ABBIAMO

$P\{Y \leq 0\} = 0$ E QUINDI L'ESPRESSIONE PER T È BEN

POSTA A NENNO DI INSIEMI DI PROBABILITÀ NULLA

LA DENSITÀ DI T VIENE CHIAMATA DENSITÀ DI STUDENT

AD N GRADI DI LIBERTÀ ED INDICATA CON IL SIMBOLO $t(n)$

È POSSIBILE SCRIVERE ESPLICITAMENTE LA DENSITÀ $t(m)$:

$$t(m)(s) = c_m \left(1 + \frac{s^2}{m}\right)^{-\frac{(m+1)}{2}} \quad \text{SE } \mathbb{R}$$

CON L'ESPRESSIONE DELLA COSTANTE c_m CHE TRALASCIAMO, ESSENDO COMPLICATA E NON RILEVANTE PER I NOSTRI SCOPI.

COMUNQUE $\lim_{m \rightarrow \infty} \left(1 + \frac{s^2}{m}\right)^{-\frac{(m+1)}{2}} = e^{-s^2/2}$ (PER NOTO LIMITE NOTEVOLE), MENTRE $\lim_{m \rightarrow \infty} c_m = 1/\sqrt{2\pi}$. ABBIAMO QUINDI

IL SEGUENTE RISULTATO

PROPOSIZIONE PER $m \rightarrow +\infty$, LA DENSITÀ $t(m)$ CONVERGE ALLA DENSITÀ $N(0,1)$

PROPOSIZIONE SUPPONIAMO DI AVERE m VARIABILI ALEATORIE X_1, \dots, X_m INDIPENDENTI E CON DENSITÀ GAUSSIANA $N(\mu, \sigma^2)$. ALLORA

$$\sqrt{m} \frac{\bar{X}_m - \mu}{\sqrt{S_m^2}} \sim t(m-1)$$

PROOF

SAPPIAMO CHE $\bar{X}_m \sim N(\mu, \frac{\sigma^2}{m})$, QUINDI $\frac{\bar{X}_m - \mu}{\sigma/\sqrt{m}} \sim N(0,1)$.

INOLTRE, PER QUANTO VISTO SULLA VARIANZA CAMPIONARIA

$\frac{(m-1)S_m^2}{\sigma^2} \sim \chi^2(m-1)$, MENTRE S_m^2 ED \bar{X}_m SONO INDIPENDENTI

Quindi:

$$\sqrt{n} \frac{\bar{X}_n - \mu}{\sqrt{S_n^2}} = \sqrt{n-1} \frac{\frac{\bar{X}_n - \mu}{S/\sqrt{n}}}{\sqrt{\frac{(n-1)S_n^2}{S^2}}} \sim t(n-1)$$



NOTA STORICA

SEALY GOSSET INTRODUSSE LA DENSITÀ DI STUDENT
AI PRIMI DEL 900. IL NOME È LO PSEUDONIMO USATO DA
GOSSET QUANDO PUBBLICÒ L'ARTICOLO.