

CONCETTI di BASE

<p>Carattere X [o A]</p>	<p>caratteristica quantitativa [o <i>qualitativa</i>] rappresentativa di un “<u>fenomeno</u>” sottoposto ad <u>indagine</u></p>
<p>Popolazione</p>	<p>insieme dei <u>soggetti portatori</u> del “<u>carattere</u>” in esame (ossia TUTTI gli <i>elementi della popolazione</i>).</p>
<p>Campione</p>	<p><u>parte</u> degli “elementi della popolazione”</p>
<p>Unità Statistica di Rilevazione</p>	<p>è un “elemento” della popolazione, ed è il <u>soggetto su cui viene rilevata la caratteristica di interesse</u></p>
<p>Modalità x_i [o a_i] con $i = 1, 2, 3, \dots, n$</p>	<p><u>determinazione</u> di un “carattere”, ossia il modo in cui si manifesta il carattere stesso al momento della rilevazione</p>
<p>Frequenza y_i</p>	<p>[<i>frequenza assoluta</i>] = <u>numero di casi</u> in cui si presenta una data <i>modalità</i> di un carattere: <i>frequenza della modalità x_i di X</i> [o della <i>modalità a_i di A</i>]</p>
<p>Totale frequenze</p>	<p>$\sum_{i=1}^n y_i = N$ indica il “totale dei casi”</p>
<p>Frequenza relativa $f_i = y_i / \sum_{i=1}^n y_i$</p>	<p>“<i>frazione di casi</i>” in cui è stata rilevata la <i>modalità corrispondente x_i</i> $\sum_{i=1}^n f_i = 1$</p>
<p>Percentuale p_i</p>	<p><u>frazione di 100 soggetti</u>, su cui è stata rilevata la <i>modalità x_i</i>. $p_i = f_i \times 100$</p>

TIPOLOGIE di FREQUENZE

♣ Sulla base dei dati riportati nella tabella seguente, si chiede di calcolare le percentuali di soggetti corrispondenti a ciascuna delle classi di età considerate; si chiede inoltre di calcolare i valori accumulati delle frequenze assolute e spiegarne il significato .

Sulla base di quanto richiesto utilizziamo il seguente “prospetto di calcolo”:

Tav. 1

x_i	y_i	$f_i = \frac{y_i}{\sum_{i=1}^n y_i}$	$p_i = f_i \times 100$	$F_i = \sum_{h=1}^i y_h$
22 — 26	13	0,1857	18,57	13
26 — 30	22	0,3143	31,43	35
30 — 34	12	0,1714	17,14	47
34 — 38	10	0,1429	14,29	57
38 — 42	6	0,0857	8,57	63
42 — 46	2	0,0286	2,86	65
46 — 50	2	0,0286	2,86	67
50 — 54	3	0,0428	4,28	70
Totali	70	1,0000	100,00	//

Ricordando che le *percentuali* vanno calcolate moltiplicando le frequenze relative per 100, i valori richiesti sono riportati nella quarta colonna della tabella, dove :

$$p_i = f_i \times 100$$

$$e \quad \sum_{i=1}^n p_i = 100$$

Sul significato dei risultati ottenuti si può dire che:

- Le **frequenze assolute** (y_i) esprimono l'ammontare dei soggetti che possiedono la "caratteristica" espressa dalla corrispondente modalità (x_i).
- Le **frequenze assolute cumulate** (F_i) esprimono l'ammontare dei soggetti che possiedono la "caratteristica" fino alla modalità posta al limite superiore della classe corrispondente:
 $F_2 = 35$ indica che sono 35 gli individui che non superano i 30 anni di età .
- Le **frequenze relative** (f_i) , esprimono la frazione di soggetti che possiedono il "carattere in esame" con modalità compresa nell'intervallo corrispondente [o pari alla modalità corrispondente, se tali modalità sono espresse in "valori singoli"]:
 $f_4 = 0,1429$ individua la frazione di soggetti con "età" compresa tra 34 e 38 anni.

♣ Le frequenze relative possono essere sfruttate anche nel caso di variabili qualitative, ma è scorretto considerarne le cumulate nel caso in cui il carattere considerato è di tipo *nominale* (ossia senza un ordine). Infatti, se si utilizzano i dati dell'esempio seguente (relativo ad una variabile qualitativa nominale), è lecito calcolare le frequenze relative e le percentuali (di cui si riporta il significato); ma non si possono prendere in considerazione le frequenze cumulate (F_i):

Tab.2

<i>Sesso</i>	y_i	f_i	p_i
<i>Maschio</i>	7	0,470	47 ,0
<i>Femmina</i>	8	0,530	53,0
<i>Totali</i>	15	1,000	100,0

Considerando le *frequenze relative*, diremo che la frazione dei casi in cui è stata rilevata la modalità "femmina" del carattere *sesso* è pari a 0,53 [$f_2 = 0,53$]; ed, usando le percentuali, si dirà che è del 53% la percentuale di soggetti di sesso femminile, all'interno del gruppo considerato.

VALORI MEDI

Valore medio di una serie ordinata di valori è un valore compreso tra il più piccolo ed il più grande dei valori osservati.

È una definizione molto generale.

Per rispettare la natura del fenomeno in esame è preferibile fare riferimento ad una classificazione più vincolante.

Si può fare una prima distinzione tra “*valori medi algebrici*” e “*valori medi di posizione*”.

1 - I Valori Medi Algebrici

Sono sintetizzati tramite delle formule matematiche, possono essere usati solo per variabili di tipo quantitativo e prendono in considerazione tutta la distribuzione dei valori osservata.

La Media Aritmetica

Se le modalità di un fenomeno X , sono legate da una relazione di tipo additivo

$$x_1 + x_2 + \dots + x_i + \dots + x_n = \sum_{i=1}^n x_i$$

La Media Aritmetica (M) è quel valore costante che, sostituito ad ognuno dei valori osservati, ne lascia inalterata la somma:

$$\underbrace{M + M + \dots + M + \dots + M}_{n \text{ volte } M} = \sum_{i=1}^n x_i$$

Da cui si ricava la formula finale [1]:

$$M = \frac{\sum_{i=1}^n x_i}{n}$$

- Nel caso in cui i dati osservati possono essere sintetizzati in una “distribuzione di frequenza”:

con modalità : $x_1, x_2, \dots, x_i, \dots, x_n$
 e con frequenze : $y_1, y_2, \dots, y_i, \dots, y_n$

la media aritmetica diventa:

$$M = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n y_i} \quad [2]$$

ma, anche in questa formulazione « ponderata », la **media aritmetica** è sempre esprimibile come il rapporto tra l'ammontare totale del fenomeno, ossia $\sum x_i y_i$, diviso il numero totale dei casi, $\sum y_i$.

• Tra le caratteristiche della media aritmetica vi è quello di minimizzare le distorsioni o gli errori accidentali ϵ_i contenuti nei dati x_i , ottenuti come misure ripetute di un carattere X.

a) **Proprietà del Baricentro:** la somma degli scarti tra i singoli valori e la propria Media aritmetica è nulla :

$$\sum_{i=1}^n (x_i - M) = 0$$

e

$$\sum_{i=1}^n (x_i - M) y_i = 0 \quad \text{per distribuzioni di frequenza}$$

b) **Proprietà del Minimo:** la somma dei quadrati degli scarti tra i singoli valori e la propria media aritmetica è un minimo, rispetto alla somma del quadrato degli scarti dei valori da un “valore medio” qualsiasi :

$$\sum_{i=1}^n (x_i - M)^2 = \text{minimo}$$

e

$$\sum_{i=1}^n (x_i - M)^2 y_i = \text{minimo} \quad \text{per distribuzioni di frequenza}$$

2- I Valori Medi di Posizione

I “**valori medi di posizione**” sono indici che **individuano le modalità che occupano una posizione** ben precisa all’interno della serie ordinata di modalità. Rispondono all’esigenza di ripartire la distribuzione in parti di eguale numerosità e possono essere utilizzate su dati *qualitativi ordinati* e su dati *quantitativi*.

- La Mediana

Il “valore mediano” o *Mediana*, fa parte dei valori medi di posizione e pertanto può essere utilizzata in presenza di *caratteri ordinabili*, anche se qualitativi, ma solo dopo averli ordinati in ordine crescente: essa bipartisce la distribuzione in due parti di eguale numerosità.

Data una serie ordinata di modalità:

$$x_1 \leq x_2 \leq x_3 \leq \dots \leq x_i \leq \dots \leq x_n$$

La Mediana è la modalità che è preceduta e seguita dallo stesso numero di osservazioni

Per individuarla si utilizza la “posizione” occupata all’interno della serie ordinata di modalità (o valori):

- se N è dispari, il “*posto*” occupato dalla Mediana è dato da $[(N + 1) / 2]$ e la mediana è l’unica modalità che occupa quel “posto”;

- se N è pari, i “*posti*” sono due: $[N / 2]$ e $[(N / 2) + 1]$ e le mediane sono due modalità: una preceduta da metà delle osservazioni e la seconda seguita dall’altra metà delle osservazioni [possono essere due “modalità” (o valori) coincidenti].

• Nel caso di distribuzione di frequenza:

$$\begin{array}{cccccccc} x_1 & \leq & x_2 & \leq & x_3 & \leq & \dots & \leq & x_i & \leq & \dots & \leq & x_n \\ \text{con frequenze: } & & y_1 & & y_2 & & y_3 & & & & y_i & & & & & & y_n \end{array}$$

valgono le stesse regole per individuare il “posto” occupato dalla Mediana, ma occorre utilizzare le “frequenze cumulate” F_i , ossia:

la F_i maggiore o uguale al *posto* occupato dalla mediana, la individua

$$\text{Se } F_k \geq \text{“Posto”}, \text{ allora } M_e = x_k$$

Esempi

Es. Calcolare la media aritmetica e la mediana dei dati riportati nella seguente distribuzione di frequenza del “consumo di gasolio” sostenuto in un anno da 155 “Hotel”.

Tab.3

<i>Consumo Gasolio (in quintali)</i>	<i>Numero di Hotel</i>
148	14
152	26
156	31
160	37
164	38
168	9
<i>Totale</i>	<i>155</i>

In presenza di una distribuzione di frequenza, per il calcolo della media aritmetica si utilizza la formula ponderata.

Prospetto di calcolo

<i>Consumo Gasolio x_i</i>	<i>y_i</i>	<i>$x_i y_i$</i>	<i>F_i</i>
148	14	2072	14
152	26	3952	40
156	31	4836	71
160	37	5920	108
164	38	6232	146
168	9	1512	155
<i>Totale</i>	<i>155</i>	<i>24524</i>	<i>//</i>

- Utilizzando la formula di seguito riportata ed i calcoli della terza colonna della tabella precedente, si perviene al valore finale:

$$M = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n y_i}$$

$$M = 24524 / 155 = 158,2 \quad \text{quintali}$$

Da notare che:

- la media è espressa nella stessa unità di misura del carattere;
- il valore ottenuto è compreso tra il più piccolo (148) ed il più grande (168) delle modalità riportate in tabella.

- Per il calcolo della **mediana**, essendo in presenza di una distribuzione di frequenza, si utilizzano le frequenze cumulate, riportate nell'ultima colonna della seconda tabella.

La mediana è la modalità preceduta e seguita dallo stesso numero di casi; tale modalità, nel caso in cui il “totale frequenze” (N) è dispari, occupa il posto:

$$(N + 1) / 2 = (155 + 1) / 2 = 78^{\circ} \text{ posto}$$

La prima frequenza cumulata che supera questo valore è $N_4 = 108$ a cui corrisponde la modalità 160, che è la mediana della distribuzione:

$$M_e = 160 \text{ quintali}$$

Da notare che, anche in questo caso:

- la mediana è espressa nella stessa unità di misura del carattere;
- anche la mediana è compresa tra il più piccolo ed il più grande dei valori osservati.

Es. Sui dati della Tab.1, calcolare la mediana della distribuzione.

• Poiché il totale delle frequenze è un numero pari, vi sono due Mediane: la 1° mediana occupa il **posto** $(N / 2) = (70 / 2) = 35$ che individua **la classe mediana** [26 —| 30] , corrispondente alla frequenza assoluta cumulata $F_2 = 35$; mentre la 2° mediana occupa il **posto** $(N / 2) + 1 = 35 + 1 = 36$, che individua la classe mediana adiacente [30 —| 34], corrispondente alla frequenza assoluta cumulata $F_3 = 47$.

Per individuare **una modalità mediana all'interno della classe** si utilizza la nota proporzione:

$$[L - l] : [M_e - l] = [F_s - F_{s-1}] : [Posto^\circ - F_{s-1}]$$

- Nell'esempio considerato $M_e^{(I)}$:

$$L = 30 \quad l = 26 \quad F_s = 35 \quad F_{s-1} = 13 \quad Posto = 35^\circ$$

$$(30 - 26) : (M_e^{(I)} - 26) = (35 - 13) : (35 - 13)$$

$$M_e^{(I)} = 26 + [(30 - 26) \times (35 - 13)] / (35 - 13)$$

$M_e^{(I)} = 26 + 4 = 30$ **anni** (età massima, raggiunta dalla prima metà dei soggetti)

- Analogamente si calcola il 2° valore mediano (in questo caso non coincidente con il 1°): $M_e^{(II)}$

$$L = 34 \quad l = 30 \quad F_s = 47 \quad F_{s-1} = 35 \quad Posto = 36^\circ$$

$$(34 - 30) : (M_e^{(II)} - 30) = (47 - 35) : (36 - 35)$$

$$M_e^{(II)} = 30 + [(34 - 30) \times (36 - 35)] / (36 - 35)$$

$$M_e^{(II)} = 30 + (4 / 12) = 33,3 \bar{3} \cong 33 \text{ anni}$$

Ossia metà degli individui hanno un'età inferiore a 30 anni e l'altra metà un'età superiore a 33 anni .