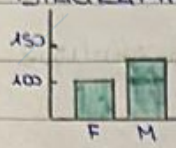


Rappresentazione immediata dei dati = GRAFICI

• **DIAGRAMMA a BARRE (BARPLOT):** Per ogni categoria si disegna un rettangolo con base costante e altezza pari alla frequenza, staccati ed equidistanti fra loro se la variabile è ordinata si rispetta l'ordine



• **DIAGRAMMA a TORTA:**

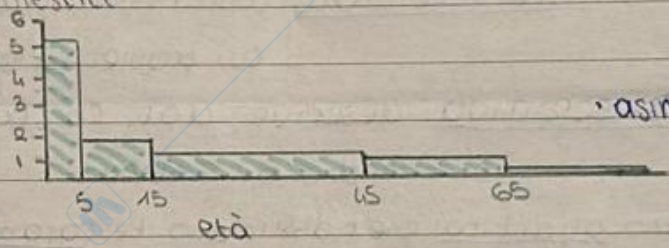


Per le variabili continue → **Densità di frequenza:** rapporto tra la frequenza $d_i = n_i / e_i$ e l'ampiezza e della classe

nell'istogramma l'area totale sarà 1 e l'altezza = densità di frequenza relativa

es: vittime incidenti domestici

età	fr. rel.	densità
0-4	25.3	5.6
5-14	18.9	1.89
15-64	30.3	1.01
45-64	13.6	0.68
65+	11.7	0.33



• asimmetrico positivamente

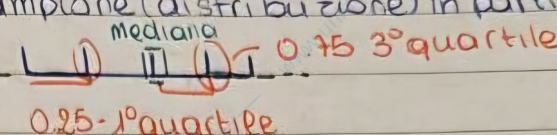
Le caratteristiche più rilevanti vengono rappresentate da numeri - **INDICI di SINTESI**
 indici di posizione

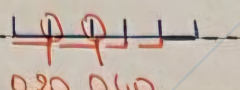
- **MODA:** è il valore più frequente per ogni tipo di dato (unimodale) quando ci sono più mode diverse = **BIMODALE** o **MULTIMODALE**
- **MEDIA:** è la somma di tutti i valori trovati divisa per n delle osservazioni
- **MEDIANA:** è il valore centrale della distribuzione che considero per cui al di sotto di essa trovo la metà delle osservazioni più piccole
- **QUANTILE:** valore che divide la distribuzione così che al di sotto di esso vengono lasciate un numero p di osservazioni

- ordino n osservazioni → la scala ha $n+1$ punti.
 - la porzione che sta sotto l' i -esima osservazione è $i/(n+1)$
 $i = p(n+1)$ trovo l'osservazione che corrisponde al quantile
 se i è intero ok, se non è intero il quantile starà tra l'intero e il n° succ.
 $x_d + (x_{d+1} - x_d) \times (i - d)$

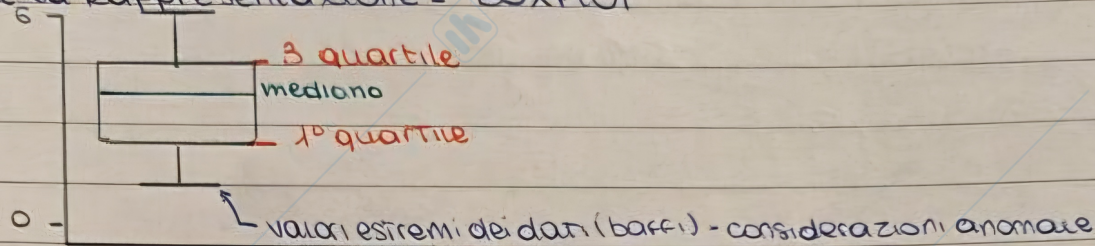
$p \max = 1$
 $n = n^{\circ} \text{osservaz.}$
 $i = \text{osservazione}$

è più utile dividere il campione (distribuzione) in parti uguali

• in 4 parti = quartile 

• in 5 parti = quintile 

Per la rappresentazione = BOXPLOT



indici di dispersione

- **RANGE**: intervallo tra il valore minimo e il valore massimo (non è molto aff.)
meglio usare il range interquartile ovvero l'intervallo tra il 1° e il 3° q.
- **VARIANZA**: calcolo la deviazione della media = differenza tra ogni osservazione e la media
 → elevo ogni risultato al quadrato così da eliminare i problemi di segno
 e faccio la Σ = somma dei quadrati intorno alla media
 → li divido per $n-1$ (gradi di libertà)

per cui **VARIANZA** $s^2 = \frac{1}{n-1} \cdot \Sigma (x_i - \bar{x})^2$ oppure $s^2 = \frac{1}{n-1} \left(\Sigma x_i^2 - \frac{(\Sigma x_i)^2}{n} \right)$

- **DEVIAZIONE STANDARD**: per tornare alla stessa unità di misura delle osservazioni
faccio la radice della varianza: $SD = \sqrt{s^2}$

ovvero $SD = \sqrt{\frac{1}{n-1} \cdot \Sigma (x_i - \bar{x})^2}$

● STATISTICA INFERENZIALE

estende i risultati a tutta la popolazione da cui è stato preso il campione

↳ USA LE PROBABILITÀ: probabilità con cui un determinato evento accade in determinate circostanze

- insieme dei possibili risultati = spazio campionario (es dado $\{1, 2, 3, 4, 5, 6\}$)
es $A = \{1, 3, 5\}$ $B = \{4, 5, 6\}$
 - evento unione: $A \cup B = \{1, 3, 4, 5, 6\}$
 - evento intersezione: $A \cap B = \{5\}$
 - evento complementare o non $A = \{2, 4, 6\}$

▷ LA FREQUENZA RELATIVA è il rapporto tra il numero di volte in cui è possibile

il risultato e il numero dei risultati possibili $n(A)/n$ è sempre compreso tra 0 e 1

▷ Se EVENTI MUTUAMENTE ESCLUSIVI allora $A \cap B = \emptyset$ e $A \cup B = P(A) + P(B)$

es. dado $n = 6$, la probabilità di avere 1 o 3? $P(1) = 1/6$ $P(3) = 1/6 \rightarrow 1/6 + 1/6 = 1/3$

▷ Se eventi INDIPENDENTI allora $A \cup B = P(A) \cdot P(B)$

es. dado $n = 6$ ma lo lancio 2 volte, la probabilità di avere 1 e 3? $P(1) 1/6$ $P(3) 1/6 \rightarrow 1/36$

• VARIABILE ALEATORIA (CASUALE) = - DISCRETA = numero finito ≥ 0 e per cui $\sum P = 1$

CONTINUA = intervallo per cui $P(a \leq x \leq b) = \int_a^b f(x) dx$
↑ densità di probabilità

● • MEDIA = μ (valore atteso): è la \sum del valore assunto · probabilità che accada
es. 2 monete, quante volte esce testa?

Possibilità = testa-testa, testa-croce, croce-testa, croce-croce

$x = 2$

$x = 1$

$x = 1$

$x = 0$

per cui $x = 2$ $P = 1/4$

$x = 1$ $P = 2/4$

media = $2 \cdot 1/4 + 1 \cdot 1/2 + 0 \cdot 1/4 = 1/2 + 1/2 = 1$

$x = 0$ $P = 1/4$

• VARIANZA = σ^2 valore atteso del quadrato degli scarti

$x = 0$ $\sigma = x - \mu = -1 \rightarrow 1 \cdot 1 = 1$

$x = 1$ $\sigma = 1 - 1 = 0$

varianza = $(0 - 1)^2 \cdot 1/4 + (1 - 1)^2 \cdot 1/2 + (2 - 1)^2 \cdot 1/4 = 1/2$

$x = 2$ $\sigma = 2 - 1 = 1$

con le variabili discrete "piccole" uso la DISTRIBUZIONE BINOMIALE

es. soggetti guariti dopo aver trattato un campione di n pazienti

$$0! = 1$$

n = numero di risultati possibili

$$n! = n(n-1) \cdot (n-2) \cdot \dots \cdot 1$$

r = numero di successi

$$Prob(r) = \frac{n!}{r!(n-r)!} \cdot p^r \cdot (1-p)^{n-r}$$

p = Probabilità di successo

es. ~~con due successi~~ la probabilità di avere 6

~~11/12~~

in una popolazione in cui il 20% è fumatore, prendendo un campione di 7 uomini, quale è la probabilità che 6 fumino?

$$r = 6$$

$$p = 20\% = 0.2$$

$$n = 7$$

$$\frac{7!}{6!(7-6)!} \cdot 0.2^6 \cdot (1-0.2)^1 = 0.29$$

con le Variabili continue uso la DISTRIBUZIONE NORMALE

• è simmetrica rispetto alla media e μ = valore massimo e moda

• $\mu - \sigma$ e $\mu + \sigma$ punti di flesso

• ha dei valori standard Z

- per cui $P(X > x_1) = P\left(Z > \frac{x_1 - \mu}{\sigma}\right)$ -> per cui = $1 - P\left(Z \leq \frac{x_1 - \mu}{\sigma}\right)$

es. $\mu = 120$ $\sigma = 13$ $P(X > 140)$

$$P(X > 140) = P\left(Z > \frac{140 - 120}{13}\right) \rightarrow P(Z > 1.54) = 1 - P(Z \leq 1.54) = 1 - 0.9382 = 0.0618$$

↳ guardo sulla tabella

- se in un range:

$$P(x_1 \leq X \leq x_2) = P\left(\frac{x_1 - \mu}{\sigma} \leq Z \leq \frac{x_2 - \mu}{\sigma}\right) \text{ per cui } P(Z \leq z_2) - P(Z < z_1)$$

es di prima ma $100 \leq X \leq 130$: $P\left(Z \leq \frac{130 - 120}{13}\right) - P\left(Z < \frac{100 - 120}{13}\right) = P(Z \leq 0.77) - P(Z < -1.54)$

inverso, do P voglio X

$$\mu = 120 \quad \sigma = 13 \quad X > x_1 \sim P(X > x_1) = 5\% = 0.05$$

$$P\left(Z > \frac{x_1 - 120}{13}\right) = 0.05 \rightarrow 0.05 = 1 - P\left(Z < \frac{x_1 - 120}{13}\right)$$

$$+ 0.95 = + P\left(Z < \frac{x_1 - 120}{13}\right) = Z = 1.654$$

$$1.654 = \frac{X - 120}{13}$$

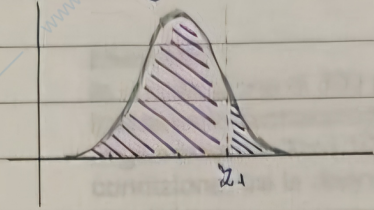
NB

quando devo trovare la densità di frequenza \rightarrow Paragono la curva allo NORMALE STANDARD

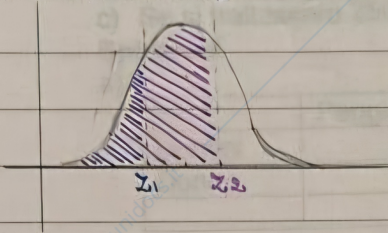
Per cui ad ogni valore di Z corrisponde l'area della parte $Z < z_1$

per cui $P(X > x_1) = P(Z > z_1)$ $z_1 = \frac{x_1 - \mu}{\sigma}$

trovo z_1 e guardo sulla tabella a quale area mi corrisponde $P(Z > z_1) = 1 - P(Z < z_1)$



Se $x_1 < X < x_2$ $P(x_1 < X < x_2) = P(z_1 < Z < z_2) \rightarrow P(Z < z_2) - P(Z < z_1)$



a) $P(\text{nona e peggioramento}) = \frac{20}{100} = 20\%$
 b) $250 \cdot P(\text{miglioramento}) = 250 \cdot \frac{20}{100} = 50$

esercizio 2
 100 pazienti, di cui 21 presentano disturbi alimentari e 79 sono in terapia presso due psicologi. Il psicologo che utilizza un approccio psicodinamico (P) segue 48 dei pazienti e l'altro psicologo, che utilizza un approccio cognitivo-comportamentale (CC), segue 10 dei pazienti che presentano disturbi alimentari.
 a) Costruire la tabella di contingenza delle frequenze assolute.
 b) Qual è la probabilità che un paziente del gruppo, preso a caso, sia seguito dal psicologo con approccio CC?
 c) Qual è la probabilità che un paziente soffre di disturbi alimentari e venga seguito dal psicologo con approccio CC?

	Disturbi alimentari	Fobie	Altri disturbi	Totale
CC	10	10	20	40
P	11	20	17	48
Totale	21	30	37	100

b) $P(\text{CC} | \text{Disturbi}) = \frac{10}{21} = 0.476$
 c) $P(\text{CC} \cap \text{Disturbi}) = \frac{10}{100} = 0.1$

• diff. tra Proporzioni

$$z_{obs} = \frac{P_1 - P_2}{SE}$$

$$SE = \sqrt{p(1-p) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$$

$$\text{dici } p = \frac{n_1 P_1 + n_2 P_2}{n_1 + n_2}$$

campioni indipendenti, con σ ignota

$$z_{obs} = \frac{\bar{X}_1 - \bar{X}_2}{SE}$$

$$SE = sp \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

$$e \quad sp = \sqrt{\frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{(n_1 - 1) + (n_2 - 1)}}$$

Se i campioni sono appaiati (es stesso campione prima e dopo)

$$z_{obs} = \frac{\bar{d}}{SE}$$

dove $\bar{d} = (\sum \text{prima-dopo}) / n \rightarrow$ calcolo σ e SE

Intervallo di confidenza:

se devo stimare la media di una popolazione μ

la miglior stima è la media del campione che esamino \bar{x} (più il campione è grande, meglio)

→ posso misurare l'incertezza: ERRORE STANDARD = $SE = \sqrt{\frac{P(\text{probabilità successo}) \cdot (1-P)}{n(\text{campione})}}$

$$\bar{x} = \bar{x}_1 - \bar{x}_2 \quad \text{diff. tra 2 medie} \quad SE = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

$$SE = \sqrt{\frac{P(1-P)}{n}}$$

$$P = P_1 - P_2$$

$$\text{tra 2 proporz.} \quad SE = \sqrt{\frac{P_1(1-P_1)}{n_1} + \frac{P_2(1-P_2)}{n_2}}$$

intervallo di confidenza = livello di fiducia = $(1-\alpha) \cdot 100\%$

NB. nelle distribuzioni campionarie $\sigma = \sigma/\sqrt{n} = SE$ e $z = \frac{\bar{x} - \mu}{SE}$

• se lv confidenza = 95% → $\alpha = 0,05$ $\alpha/2 = 0,0025$ $z = \pm 1,96$

$$IC = \mu \pm z_{\alpha/2} \cdot SE \quad \text{se } \bar{x} \text{ tale che } \left| \frac{\bar{x} - \mu}{SE} \right| > z_{1-\alpha/2}$$

~~Probabilità di errore~~

$$IC = 95\% \rightarrow \alpha = 0,05 \quad z = \pm 1,96$$

$$IC = 90\% \rightarrow \alpha = 0,1 \quad z = \pm 1,64$$

$$IC = 99\% \rightarrow \alpha = 0,01 \quad z = \pm 2,58$$

quando il campione è piccolo uso t student → colonna = $n-1$ e riga = $\alpha/2$

$$\text{poi } IC = \bar{x} \pm t_{\text{obs}} \cdot SE$$

Per trarre dei risultati su tutta la popolazione, estendo ad essa il risultato del campione perciò

$\mu_{\text{campione}} = \mu_{\text{popolazione}}$

$$SE = \frac{\sigma_{\text{campione}}}{\sqrt{n}} \quad \text{se } n > 30$$

$$\text{lo standardizzo per cui } z = \frac{\bar{x} - \mu_{\text{campione}}}{SE}$$

TEST D'IPOTESI

genero un'affermazione circa 1 o più parametri di una popolazione circa la probabilità della distribuzione di un parametro

- H_0 = ipotesi nulla \rightarrow non ci sono differenze (da confutare)
- H_A = ipotesi alternativa \rightarrow ci sono differenze (da dimostrare)

in genere pongo:

$$\begin{cases} H_0: \mu_{\text{campione}} = \mu_{\text{popolazione}} & \rightarrow \text{per cui } \bar{x} \approx N(\mu_0, SE) \\ H_A: \mu \neq \mu_0 \end{cases}$$

mi ricavo un IC in genere 95% per cui $\alpha = 0.05$ e controllo ~~che~~ \bar{x} dove cade

\rightarrow se cade nella regione critica = RIFIUTO H_0 (5%)

\rightarrow se cade nella regione di non rifiuto = NON RIFIUTO H_0 (95%)

(NB) POSSO STANDARDIZZARE per cui $Z_{\text{OBS}} = \frac{\bar{x} - \mu_0}{SE}$ e $|Z_{\text{OBS}}| > Z_{1-\alpha/2}$ RIFIUTO H_0
 $|Z_{\text{OBS}}| < Z_{1-\alpha/2}$ NON RIFIUTO H_0

Posso usare anche P-VALUE

$$P\text{-VALUE} = 2 \cdot P(Z > |Z_{\text{OBS}}|) = 2 \cdot (1 - P(Z < Z_{\text{OBS}}))$$

\rightarrow se $2 \cdot (1 - P(Z < Z_{\text{OBS}})) < \alpha$ RIFIUTO H_0

\rightarrow se $2 \cdot (1 - P(Z < Z_{\text{OBS}})) > \alpha$ NON RIFIUTO H_0

Se ho μ_0 o $SE(\sigma)$ o entrambe sconosciute con un campione $n < 30$

USO T STUDENT:

$$t_{\text{OBS}} = \frac{\bar{x} - \mu_0}{SE} \rightarrow \text{poi cerco il grado di libertà sulla tavola T}$$

colonna = $n-1$ (o $\geq n-1$) riga = $\alpha/2$

\rightarrow se $t_{\text{OBS}} > t_{1-\alpha/2}$ RIFIUTO H_0 con Lv significatività α

\rightarrow se $t_{\text{OBS}} < t_{1-\alpha/2}$ non RIFIUTO H_0 con Lv significatività α

(NB) se ho \bar{p}_0 successi $Z_{\text{OBS}} = \frac{p - \bar{p}_0}{\sqrt{\frac{\bar{p}_0(1-\bar{p}_0)}{n}}}$ e $|Z_{\text{OBS}}| > Z_{1-\alpha/2}$ RIFIUTO H_0 Lv signif. α
 $|Z_{\text{OBS}}| < Z_{1-\alpha/2}$ NON RIFIUTO H_0 Lvs α

(NB) se i campioni sono appaiati
 es stesso campione prima e dopo faccio la diff. prima e dopo, le sommo tutte
 e divido per n , poi mi calcolo σ e SE e $t_{\text{OBS}} = \frac{\bar{d}}{SE}$
 (d)

- Se ho 2 campioni indipendenti con μ_2 e μ_1 e $\sigma_1^2 = \sigma_2^2$ igno

$$t_{OBS} = \frac{\bar{X}_1 - \bar{X}_2}{SE}$$

$$SE = sp \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

$$sp = \sqrt{\frac{(n_1-1)S_1^2 + (n_2-1)S_2^2}{(n_1-1) + (n_2-1)}}$$

dif. tra medie

se n_1 e $n_2 > 30$

$$z_{OBS} = \frac{\bar{X}_1 - \bar{X}_2}{SE}$$

$$SE = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

- se proporzioni dif. proporzioni

$$z_{OBS} = \frac{p_1 - p_2}{SE}$$

$$SE = \sqrt{p(1-p) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$$

e $p = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2}$