

Nel database ci sono due dimensioni:

- n -> numero soggetti, da 1 a n significa che si sta facendo ancora l'indice sui soggetti
- k -> numero categorie, significa che quel "i" va da 1 a k

Le scale quantitative dove X_i è associata al soggetto, quindi si dice che i va da 1 a n ($i \leq n$)

La variabile indicatore è una variabile che di solito indica uno stato e quindi devi dare come nome il nome della variabile (es: se la chiami "malato" una persona lo deve essere).

Quando si descrivono le caratteristiche di un dato, quelle caratteristiche le esaminiamo per due aspetti:

- tendenza centrale -> il dato come si colloca
- variabilità -> il dato quanto è variabile

La distribuzione di una variabile statistica è una relazione che associa alla variabile una frequenza, è una relazione in cui ho la variabile X_i con diversi valori ($X_1 X_2 X_3$ ec) che posso e essere categoriali o quantitativi. La frequenza si chiama n_i se è una frequenza assoluta oppure la chiamerò p_i se è relativa (espressa in percentuale). 1.13

Abbiamo tre misure di tendenza centrale:

- Media aritmetica -> si calcola solo per le variabili quantitative. La media è influenzata molto dai valori estremi
Formula -> sommatoria : n
- Mediana -> è il valore che mi taglia al centro la distribuzione, mi lascia a destra e a sinistra il 50% e non è influenzata da valori estremi. La mediana la posso calcolare anche per le variabili ordinali (dal più piccolo al più grande)
Formula: $\frac{n+1}{2}$
- Moda -> è il valore della massima frequenza (indica il picco), si può calcolare su tutte le variabili

Non c'è una tendenza centrale migliore, la soluzione più giusta è tenere in considerazione tutte e tre. Ad esempio quando ho la media uguale alla mediana ho un'informazione di dati distribuiti in modo simmetrico, simmetrico vuol dire che il valore che taglia a metà lascia esattamente 50% più alto e il 50% più basso a destra e sinistra ed è anche il valore medio. Invece, quando una media è più piccola della mediana significa che si hanno dei valori estremi a sinistra che spostano la media verso sinistra, quando è il contrario la curva si sposta a destra.

La variabilità è più importante della centralità.