

# LEZIONE 3 | DISTRIBUZIONI DI FREQUENZA

• LA MATRICE di DATI  $\Rightarrow$  è una tabella che presenta i dati raccolti per ogni unità statistica della popolazione relativamente a ciascuna variabile rilevata

RIGHE: dati relativi a una particolare unità statistica

COLONNE: (eccetto la prima se si riporta il numero o l'etichetta dell'unità statistica) rappresentano una variabile

Dalla matrice si possono ottenere DIVERSE RAPPRESENTAZIONI DEI DATI, la più usata è la **TABELLA di FREQUENZA**

**NB:** Molto spesso le variabili **QUALITATIVE** vengono trasformate in un numero opportuno di **VARIABILI NUMERICHE**  $\Rightarrow$  si chiama in una variabile indicatore di ciascuna delle modalità possibili

$\Downarrow$   
es: Abbiamo una variabile qualitative con  $k$  modalità

da questa v. qualitativa si generano  $k-1$  **VARIABILI DUMMYS** (variabili indicatore)

Abbiamo così trasformato in una variabile qualitativa in variabile numerica

$\Downarrow$   
vale 1 se si osserva una certa modalità  
vale 0 altrimenti

Vedi slides per ESEMPI

## Distribuzioni di frequenza (parte 1)

- I dati raccolti rappresentano l'informazione disponibile su una variabile  $X$
- Come ottenere una visione d'insieme del fenomeno?

$\Downarrow$   
Serve aggregare le unità statistiche che presentino modalità identiche

$\Leftarrow$  **OVVERO**: quante volte ogni singola modalità è presente nella **DISTRIBUZIONE DISAGGREGATA**

$\Downarrow$   
useremo le **DISTRIBUZIONI di FREQUENZA**

$\Downarrow$   
si deve ristrutturare i dati per avere una visione d'insieme del fenomeno  $\Rightarrow$  costruire una distribuzione che non sia disaggregata ma che abbia un certo tipo di aggregazione

## OPERATORE SOMMATORIA

$a_1, a_2, a_3, \dots, a_n$

la sommatoria permette di scrivere la sommatoria:  $a_1 + a_2 + a_3 + \dots + a_n$  in modo + compatto

$\sum_{i=1}^n a_i$   $\rightarrow$  ARGOMENTO  
 da dove parte la sommatoria

### ESEMPLI

$\sum_{i=1}^n 3a_i = 3 \cdot \sum_{i=1}^n a_i \Rightarrow$  la costante è indipendente (c'è  $\cdot$ )

$\sum_{i=1}^n (2 + a_i) =$  spezzare la sommatoria due c'è  $= \sum_{i=1}^n 2 + \sum_{i=1}^n a_i = 2n + \sum_{i=1}^n a_i$   
 2 per n volte

NB:  $2n = 2+2+2+\dots+2$   
 $\neq 2^n$

$\sum_{i=1}^n \frac{4}{a_i} \neq \frac{4}{\sum_{i=1}^n a_i}$

$\sum_{i=1}^n \log a_i \neq \log \sum_{i=1}^n a_i$

### SOMMATORIA DOPPIA

$\sum_{i=1}^3 \sum_{j=1}^4 (2i - 4j) = \sum_{i=1}^3 \sum_{j=1}^4 2i - \sum_{i=1}^3 \sum_{j=1}^4 4j = 2 \sum_{i=1}^3 \sum_{j=1}^4 i - 4 \sum_{i=1}^3 \sum_{j=1}^4 j =$   
 $i \neq j \Rightarrow$  i e come se fosse una costante  
 $= 2 \sum_{i=1}^3 (i+i+i+i) - 4 \sum_{i=1}^3 (1+2+3+4) = 2 \sum_{i=1}^3 4i - 4 \sum_{i=1}^3 10 =$   
 $= 8 \sum_{i=1}^3 i - 4 \cdot 3 \cdot 10 = 8 \cdot (1+2+3) - 120 = 8 \cdot 6 - 120 = -72$

$\sum_{i=1}^4 \sum_{j=1}^3 2ij = 2 \sum_{i=1}^4 \left( \sum_{j=1}^3 ij \right) = 2 \sum_{i=1}^4 i \sum_{j=1}^3 j = 2 \sum_{i=1}^4 i (1+2+3) = 2 \sum_{i=1}^4 i \cdot 6 = 12 \sum_{i=1}^4 i =$   
 $= 12 \cdot (1+2+3+4) = 12 \cdot 10 = 120$

## - DISTRIBUZIONI di FREQUENZA - (parte 2)

- I dati raccolti rappresentano l'informazione disponibile su una variabile  $X$
- Come ottenere una visione di insieme del fenomeno?

↓  
 aggregare le unità statistiche in modo che vengano resi disponibili in modo + fruibile

↳ aggregare le unità statistiche che presentano modalità identiche

↓  
 si conta ogni quante volte ogni singola modalità si presenta nella **DISTRIBUZIONE DISAGGREGATA** (nella nostra grande matrice di dati che raccoglie tutte le informazioni così come sono rilevate)

↓  
 si costruiscono a questo scopo le **distribuzioni di frequenza**

↓  
 da qui si parte con le analisi

⇒ si passa da una disp. disagg. ad una distribuzione in cui c'è un'aggregazione dei dati in base alle modalità identiche

### INDICHIAMO

- $X$  è una **variabile** qualitativa / quantitativa
  - $x_1, x_2, \dots, x_k$ : le **modalità** della variabile
- ↳ supponiamo che siano  $k$  modalità distinte

### Distribuzione di frequenza assoluta (poco informativa)

- $N$  osservazioni di  $X$
- Modalità  $x_1, \dots, x_k$  che si presentano ciascuna con frequenza  $n_1, \dots, n_k$

↓  
 $n_1$  si presenta  $n_1$  volte  
 $n_2$  si presenta  $n_2$  volte  
 ...  
 $n_k$  si presenta  $n_k$  volte

**NB:** L'insieme delle coppie costituite da  $(x_1, n_1), (x_2, n_2), \dots, (x_k, n_k)$

ogni modalità il numero di volte in cui quella modalità viene rilevata

costituisce la **DISTRIBUZIONE di FREQUENZA ASSOLUTA**

• si scrive:  $\{(x_1, n_1), (x_2, n_2), \dots, (x_k, n_k)\}$

• si rappresenta:

$X$	freq. ass. $n_i$
$x_1$	$n_1$
$x_2$	$n_2$
$\vdots$	$\vdots$
$x_k$	$n_k$

↳ indice che dice che va da 1 a  $k$  modalità

↳ la somma di quello che osserviamo deve essere  $= N$

↳  $N$  dati che sono stati osservati

modalità della variabile  $X$       corrispondenti valori della freq. assoluta

**NB:**  $N = \sum_{i=1}^k n_i$       OSSERVIAMO

## Distribuzione di frequenza (relativa e percentuale)

• **FREQ. RELATIVA**: insieme di coppie e le quantità di ciascuna coppia sono  
 Permettono di fare confronti

$$(x_1, f_1) \dots (x_k, f_k)$$

modalità      frequenza relativa

$k$ : numero delle  $x$

=  
 quanto è  
 frequente quella  
 particolare modalità

$$f_i = \frac{n_i}{N} = \frac{\text{frequenza assoluta corrispondente}}{N}$$

NB: la somma delle frequenze relative:

$$\sum_{i=1}^k f_i = \sum_{i=1}^k \frac{n_i}{N} = \frac{\sum_{i=1}^k n_i}{N} = \frac{N}{N} = 1$$

costante rispetto alla sommatoria

• **FREQ. PERCENTUALE**: a volte la freq. relativa viene data in termini percentuali

$$\{(x_1, f_1\%), \dots, (x_k, f_k\%)\}$$

modalità      frequenza della modalità in percentuale

$$f_i\% = \frac{n_i}{N} \%$$

NB: la somma deve fare 100%:

$$\sum_{i=1}^k f_i\% = 100$$

tabella:

X	freq. relativa $f_i$	freq. percentuale $f_i\%$
$x_1$	$f_1$	$f_1\%$
$\vdots$	$\vdots$	$\vdots$
$x_k$	$f_k$	$f_k\%$
-	1	100

**ATTENZIONE:**

Fare attenzione agli arrotondamenti nelle freq. relativa!

Se si arrotonda bisogna cercare di arrotondare in modo che la somma delle freq. relative faccia 1

- **PERCHÉ** si ricorre a queste frequenze?
  - x creare una distribuzione aggregata (moda + compatto per partire con i dati)
  - effettuare confronti circa il comportamento della variabile X
- **COME** scegliere quale frequenza usare? (ASSOLUTA, **RELATIVA** o PERCENTUALE?)

VANTAGGI e SVANTAGGI

Solitamente si usa quella relativa

**ATTENZIONE:**

- si perde la percezione dell'ordine di grandezza di X

- permette di valutare la composizione interna delle popolazioni velocemente
- permette di effettuare confronti (riguardo alla distribuzione della variabile X o il comportamento in  $\neq$  gruppi che hanno numerosità  $\neq$ )

Distribuzione di frequenza assoluta cumulata

- è costruita accumulando le osservazioni si parte dalle modalità distinte  $u_i$ 
  - per ciascuna modalità  $u_i$  si costruisce  $N_i$  (somma di tutte le unità statistiche che ci sono fino a quella particolare modalità)
- ci deve essere un certo ordine  $\leftarrow$ 
  - tutte quelle modalità statistiche che hanno un valore osservato  $\leq$  minore o uguale rispetto a quella particolare modalità  $u_i$
- infatti: la distribuzione di frequenza di tipo cumulato si può costruire per:
  - le variabili **QUANTITATIVE** (quando le modalità sono dei numeri)
  - le variabili **ORDINALI** (qualitative ordinali)

**ATTENZIONE!** se è **SCALA SCONNESSA**  $\textcircled{\text{non}}$  è ordinabile!

• la freq. assoluta cumulata si costruisce così:

quando si osserva la prima modalità ( $u_1$ ), si osserva la sua distribuzione di frequenza assoluta ( $n_1$ ) che è uguale alla frequenza assoluta cumulata ( $N_1$ )

$$u_1 \rightarrow N_1 = n_1$$

$$u_2 \rightarrow N_2 = n_1 + n_2$$

$$u_k \rightarrow N_k = n_1 + n_2 + \dots + n_k = N \text{ totale delle freq. cumulata} = \text{tot delle osservazioni } N$$

www.unidocs.it - Appunti e dispense per superare i tuoi esami universitari

www.unidocs.it - Appunti e dispense per superare i tuoi esami universitari

## Ripetizione di frequenza (relativa cumulata)

$$F_i = \frac{N_i}{N}$$

Operare alla freq. assoluta cumulata, basta dividere per N

$$F_i = f_1 + f_2 + \dots + f_k \quad \rightarrow \text{anche } \%$$

- $n_1 \rightarrow F_1 = f_1$
- $n_2 \rightarrow F_2 = f_1 + f_2$
- $n_k \rightarrow F_k = f_1 + f_2 + \dots + f_k = 1$
- $n_1 \rightarrow F_1\% = f_1\%$
- $n_2 \rightarrow F_2\% = f_1\% + f_2\%$
- $n_k \rightarrow F_k\% = f_k\% = 100\%$

Da tenere a mente • la dist. di freq. accumulata si calcola SOLO per variabili almeno ordinali (le qualitative su scale ordinarie e per le quantitative)

• l'ultima modalità ha associata la distribuzione (come dist. di freq. relativa cumulata) l'ultimo valore che è poi a 1 se relative cumulata e N se la distr. è assoluta cumulata

→ le X devono essere ordinate nella tabella in modo crescente

• l'ultima riga che costruiamo nella tabella la costruiamo per avere un'idea dell'ammontare delle unità statistiche

Nelle colonne di distribuzione cumulata NON si sommano (come la prima colonna)

Rilevazione delle variabili in classi è molto frequente x variabili come età, Pil, gettato, investimenti, reddito, tempo, ...

CONSIDERIAMO variabili quantitative di tipo continuo

chiamiamo CAMPO di VARIAZIONE l'intervallo iniziale di rilevazione

SUBDIVIDIAMO  
in  
SOTTOINTERVALLI

= l'intervallo costituito dai valori che vanno dal più piccolo valore osservato per le modalità della variabile fino al più grande

e assegnamo le unità statistiche a ciascun intervallo e a cui appartengono ⇒ si ottiene la rilevazione / sua suddivisione in classi della variabile che si chiama SERIAZIONE in INTERVALLI

### VANTAGGI

- utile quando la variabile si presenta con molte modalità (+pratico)
- la sua rappresentazione permette di presentare i dati in modo più compatto (x le rappresentazioni grafiche)

### SVANTAGGI

- perdita di informazioni che si ha in quanto i dati non vengono ripartiti esattamente come sono rilevati, ma si rileva solo la classe a cui appartiene
- "confusione" nella chiusura delle classi (la costruzione delle classi in modo che si eviti il passaggio da una classe all'altra)

Quando si decide di rilevare una variabile in classi



Bisogna scegliere il numero  $k$  di classi (è una scelta di tipo pratico)

**NB** Non dobbiamo dividere le osservazioni in tante classi, perché c'è il rischio di avere delle classi piccole con frequenza nulla oppure con frequenze estremamente basse

**MA** non si devono raggruppare le modalità in poche classi perché significa creare delle classi decisamente prive di osservazioni e quindi avere una grande perdita di informazioni



deve essere

- perdita di informazioni ridotta al minimo
- serve per prima del numero di classi ed efficace prestazione dei dati (poche classi)



regoletta pratica:  $k = \log_2(N) + 1$  per la scelta del numero di classi

• **Costituzione della distribuzione a classi**

- Dove serve, le classi possono essere di ampiezza diversa
- le classi non devono sovrapporsi  $\Rightarrow$  sono una di seguito all'altra
- Ogni unità statistica entra in una sola classe

• **ATTENZIONE:** alla CHIUSURA DELLE CLASSI

tipica scelta:  $[c_{i-1}, c_i]$ ,  $c_{i-1} - c_i$  classe aperta a sinistra

$[c_{i-1}, c_i)$ ,  $c_{i-1} - c_i$  classe aperta a destra

$[c_{i-1}, c_i]$

~~$c_i - c_i$~~  //  ~~$(c_{i-1}, c_i)$~~  non sono mai compresi!

Quantità che si costruiscono per poter continuare le analisi (utili)

• **VALORE CENTRALE DELLA CLASSE:** valore definito come la semi-somma degli estremi

$\frac{\text{estremo superiore} + \text{estremo inferiore}}{2}$

• è un valore rappresentativo di quella classe



La nostra variabile  $X$  viene rappresentata tramite delle modalità che sono adesso delle classi

$$\bar{m}_i = \frac{c_{i-1} + c_i}{2}$$

• **AMPIEZZA DELLA CLASSE:** differenza tra gli estremi

$$d_i = c_i - (c_{i-1})$$

• ci permette di avere un'idea della differenza tra le classi in base a come sono state costruite

- **FREQUENZA ASSOLUTA**: stabilite le unità statistiche all'interno di ciascuna classe

$n_i$

**(NB)** Questo va PRIMA di tutte le altre quantità

PRECISAZIONE: Nel momento in cui noi decidiamo di scegliere un elemento rappresentativo della classe (es. valore centrale della classe), molto spesso facciamo un'assunzione

uniforme distribuzionale delle classi

**(SE)** Non si sa nulla sulla distribuzione di  $X$  dentro la classe, allora non c'è altra scelta

assembiamo che le unità statistiche siano distribuite **EQUAMENTE**

Il valore centrale di classe sarà un buon rappresentante della classe

**Assunzione** - \* Supponiamo di dividere la classe  $(c_{i-1}, c_i]$  di freq. ass.  $n_i$  e ampiezza di  $n_i$  intervallini uguali di ampiezza  $\frac{d_i}{n_i}$

- \* Posizioniamo una unità statistica al centro di ogni intervallino
- \* Le rimanenti unità statistiche saranno collocate nell'intervallo tutte alla stessa distanza l'una dall'altra

**(NB)** Come confrontare le classi di  $\neq$  ampiezza?

**NON** sarà possibile lavorare con le freq. assolute o relative

i risultati sarebbero sbagliati

**BISOGNA**: considerare insieme alla freq. e l'ampiezza di classe su cui le unità statistiche insistono

Devono di frequenza di ogni classe  $(c_{i-1}, c_i]$

$$h_i = \frac{n_i}{d_i}$$

freq. assoluta corretta per l'ampiezza di classe

oppure

$$\frac{\frac{n_i}{N}}{d_i} = \frac{f_i}{d_i} = \frac{h_i}{N}$$

(Vedi esercizi ultima slide)