

# LEZIONE\_3

2022-10-19

## STATISTICA MULTIVARIATA

PRENDE IN CONSIDERAZIONE 2 O PIÙ VARIABILI CONTEMPORANEAMENTE

### DISTRIBUZIONI DI FREQUENZA

- **CONGIUNTA**

- si prendono in considerazione entrambe le variabili durante una osservazione/misurazione
- $(x_i, y_k)$ 
  - \*  $i \in [1, |X|]$
  - \*  $k \in [1, |Y|]$

- **DISGIUNTA**

- viene presa in considerazione una variabile per volta

- **CONDIZIONATA**

- in base a una condizione dettata da una delle due variabili si cercano le osservazioni congiunte rispetto all'altra variabile
- Date X e Y due variabili
  - \*  $Y | X = \text{valore}$
  - \* seleziona i valori di Y associati ai casi in cui X = valore indicato

```
head(mtcars)
```

```
##           mpg  cyl  disp  hp  drat    wt  qsec vs  am  gear  carb
## Mazda RX4      21.0   6  160  110 3.90 2.620 16.46 0   1    4    4
## Mazda RX4 Wag  21.0   6  160  110 3.90 2.875 17.02 0   1    4    4
## Datsun 710      22.8   4  108   93 3.85 2.320 18.61 1   1    4    1
## Hornet 4 Drive  21.4   6  258  110 3.08 3.215 19.44 1   0    3    1
## Hornet Sportabout 18.7   8  360  175 3.15 3.440 17.02 0   0    3    2
## Valiant         18.1   6  225  105 2.76 3.460 20.22 1   0    3    1
```

```
hp6Cyl = mtcars[mtcars$cyl==6, "hp"]
length(hp6Cyl)
```

```
## [1] 7
```

### RAPPRESENTAZIONI GRAFICHE

Sono le stesse spiegate nel capitolo precedente, in quanto le funzioni grafiche sono in grado di accettare più variabili contemporaneamente, ognuna con il suo significato visivo

# DIPENDENZA

Date due variabili X e Y esse possono essere dipendenti tra loro.

## IPOTESI:

- **X variabile indipendente** (per convenzione matematica)
- **Y variabile dipendente = f(X)**
  - f è la funzione che regola la dipendenza
  - può essere di qualsiasi tipo

## TIPOLOGIE

### DIPENDENZA

#### 2 QUALITATIVE

### DIPENDENZA MEDIA

#### 1 QUALITATIVA 1 QUANTITATIVA

### REGRESSIONE o CORRELAZIONE

#### 2 QUANTITATIVE

# TABELLA CONTINGENZA

Riporta le distribuzioni di frequenza associate alle due variabili, perciò gli assi contengono gli elementi del supporto delle variabili, cioè i possibili valori ammessi senza ripetizioni

	$y_1$	$y_2$	$\dots$	$y_k$	
$x_1$	$n_{11}$	$n_{12}$	$\dots$	$n_{1k}$	$n_{1+}$
$x_2$	$n_{21}$	$n_{22}$	$\dots$	$n_{2k}$	$n_{2+}$
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$
$x_m$	$n_{m1}$	$n_{m2}$	$\dots$	$n_{mk}$	$n_{m+}$
	$n_{+1}$	$n_{+2}$	$\dots$	$n_{+k}$	$n$

- $X = [x_1, \dots, x_i, \dots, x_m], i \in [1, m], m = |Sx|$
- $Y = [y_1, \dots, y_i, \dots, y_k], i \in [1, k], k = |Sy|$

$$n = \sum_{i=1}^m n_{i+} = \sum_{j=1}^k n_{+j}$$

## FREQUENZA

- **CONGIUNTA:**
  - $n_{ij} = (x_i, y_j)$
- **MARGINALE:**
  - si tratta delle occorrenze dell'i-esimo valore di una delle variabili
    - \*  $n_{i+} = x_i = \sum_{j=1}^k n_{ij} \rightarrow i=\text{riga costante}$
    - \*  $n_{+j} = y_j = \sum_{i=1}^m n_{ij} \rightarrow j=\text{colonna costante}$

- CONDIZIONATA

- $n_{i1}$  = vettore delle frequenze di X dato  $Y = y_1$

## ESEMPIO

```
attitudine <- rbind(cbind(rep("S",1),rep("S",1)),cbind(rep("S",3),rep("B",3)),cbind(rep("B",1),rep("S",1)))

# in data frame
colnames(attitudine) <- c("X","Y") # nomi delle colonne
attitudine$X <- ordered(attitudine$X, levels=c("S","B","O"))
attitudine$Y <- ordered(attitudine$Y, levels=c("S","B","O"))
str(attitudine)

## 'data.frame': 15 obs. of 2 variables:
## $ X: Ord.factor w/ 3 levels "S"<"B"<"O": 1 1 1 1 2 2 2 2 2 2 ...
## $ Y: Ord.factor w/ 3 levels "S"<"B"<"O": 1 2 2 2 1 2 2 2 3 3 ...

tab <- table(attitudine$X,attitudine$Y) # tabella di contingenza
#(distribuzione di frequenza assoluta congiunta)
tab

##
##      S B O
## S 1 3 0
## B 1 3 2
## O 2 1 2

# distribuzione marginale di Y
# (frequenza assoluta)
margin.table(tab,2)

##
## S B O
## 4 7 4

# distribuzione condizionata di Y|X=S (frequenza assoluta)
tab[1,]

## S B O
## 1 3 0

# distribuzione condizionata di Y|X=B (frequenza assoluta)
tab[2,]

## S B O
## 1 3 2

tab[3,] # distribuzione condizionata di Y|X=O (frequenza assoluta)

## S B O
## 2 1 2

tab[,1] # distribuzione condizionata di X|Y=S (frequenza assoluta)

## S B O
## 1 1 2

tab/sum(tab) # distribuzione di frequenza relativa congiunta
```

```
##
##           S           B           0
##  S 0.06666667 0.20000000 0.00000000
##  B 0.06666667 0.20000000 0.13333333
##  0 0.13333333 0.06666667 0.13333333

# in alternativa, prop.table(tab)
# distribuzione marginale di X (frequenza relativa)
margin.table(tab,1)/sum(margin.table(tab,1))

##
##           S           B           0
## 0.26666667 0.40000000 0.33333333

# distribuzione marginale di Y (frequenza relativa)
margin.table(tab,2)/sum(margin.table(tab,2))

##
##           S           B           0
## 0.26666667 0.46666667 0.26666667

# distribuzione condizionata di Y|X=S (frequenza relativa)
tab[1,]/sum(tab[1,])

##      S      B      0
## 0.25 0.75 0.00

# distribuzione condizionata di Y|X=B (frequenza relativa)
tab[2,]/sum(tab[2,])

##           S           B           0
## 0.16666667 0.50000000 0.33333333
```

## INDIPENDENZA STATISTICA

Viene misurata in maniera SIMMETRICA, perchè X e Y possono influenzarsi a vicenda, dipende dal punto di vista

$$\frac{n_{rc}}{n_{+c}} = \frac{n_{r+}}{n}$$

- $r \in [1; m]$
- $c \in [1; k]$

$$n_{rc} = \frac{n_{r+} * n_{+c}}{n}$$

questo valore corrisponde al valore ideale che la frequenza dovrebbe avere in caso di completa INDIPENDENZA tra le due variabili X e Y considerate

## INDICE DI CONNESSIONE

Determina la forza della dipendenza che c'è tra le due variabili considerate

$$\chi^2 = \sum_{r=1}^m \sum_{c=1}^k \frac{(n_{rs} - n_{rs}^*)^2}{n_{rs}^*}$$

$$n_{rs}^* = \frac{n_{r+} * n_{+c}}{n}$$

$n_{rs}^*$  = valore nel caso di completa indipendenza tra X e Y

## INDIPENDENZA

$$\chi^2 = 0 = (n_{rs} - n_{rs}^*)^2 = (n_{rs} - n_{rs}^*), \forall r \in [1; m], \forall s \in [1; k]$$

I valori attesi coincidono con quelli osservati, quindi vi è completa indipendenza

## DIPENDENZA

$$\chi^2 \in ]0; \min(m-1, k-1)]$$

## DIPENDENZA MEDIA

Si misura in maniera ASIMMETRICA

- X = QUALITATIVA INDIPENDENTE
- Y = QUANTITATIVA DIPENDENTE = in funzione di X

Non viene misurata la distribuzione di frequenza della variabile Y, ma solo la sua MEDIA

$$E(Y|X = x_i)$$

Media dei valori di Y associati al valore  $x_i$

## INDIPENDENZA IN MEDIA

$$E(Y) = E(Y|X = x_i) = E(Y|X = x_k), \forall i \neq k$$

## CORRELAZIONE

## REGRESSIONE