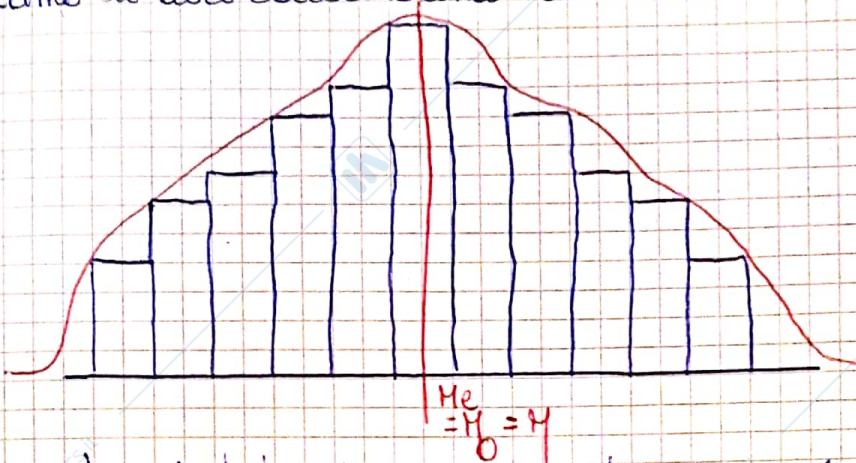


# STATISTICA

Come descrivere una variabile statistica:

- la forma della distribuzione

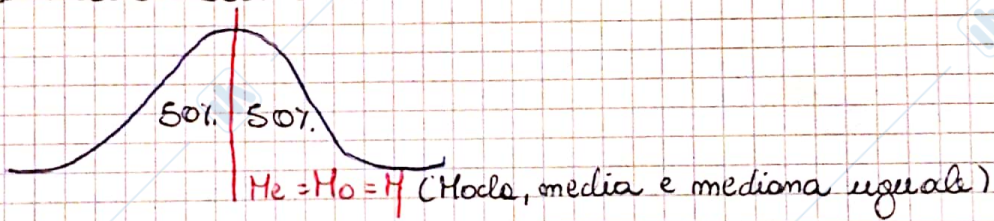
Immaginiamo di aver calcolato una variabile  $x$  e abbiamo questo istogramma:



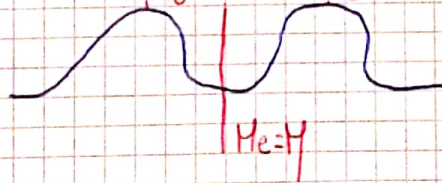
**Simmetrica:** La distribuzione è simmetrica quando dividendo la mia distribuzione a metà (in casi sproporzionali della mediana) le due metà saranno esattamente uguali.

Averemo così una variabile **simmetrica unimodale** (= ha una sola moda, perché il valore più alto è solo uno)

FORMA TIPICA:

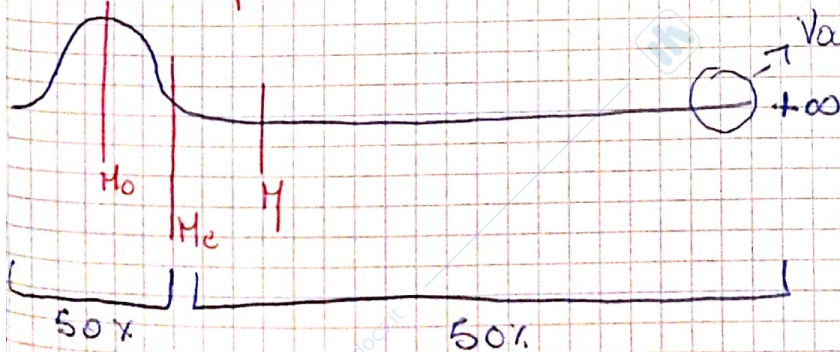


**Simmetria plurimodale:** Ha più mode



Se la nostra variabile non è simmetrica è asimmetrica:

**Asimmetria positiva:** Tende a  $+\infty$  valori, ovvero hanno valori anomali verso  $+\infty$ .



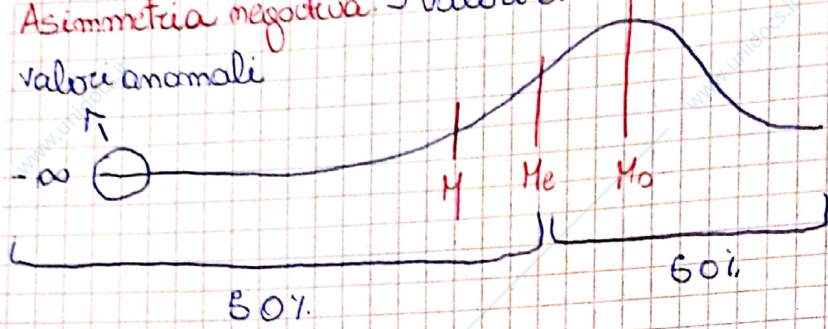
$$M > Me > Mo$$

ASIMMETRIA POSITIVA

Avevamo questo non solo capisco che c'è asimmetria positiva, ma anche che c'è mancanza di robustezza e anche che ci sono valori anomali.

Asimmetria negativa: 3 valori anomali tendono verso  $-\infty$ .

valori anomali



$$M < Me < Mo$$

ASIMMETRIA NEGATIVA

Breve:

Quando abbiamo la simmetria unimodale avremo:  $M = Me = Mo$

Quando abbiamo l'asimmetria positiva avremo:  $M > Me > Mo$

Quando abbiamo l'asimmetria negativa avremo:  $M < Me < Mo$

Quando abbiamo la simmetria plurimodale avremo:  $M = Me$ .

Variabilità:

Nel calcolo delle varianze, gli scarti più elevati acquistano molta più importanza rispetto agli scarti più piccoli.

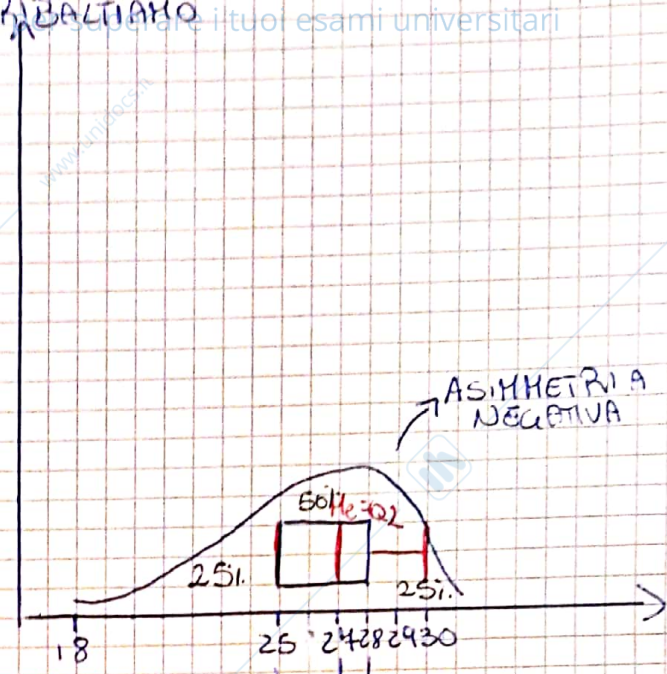
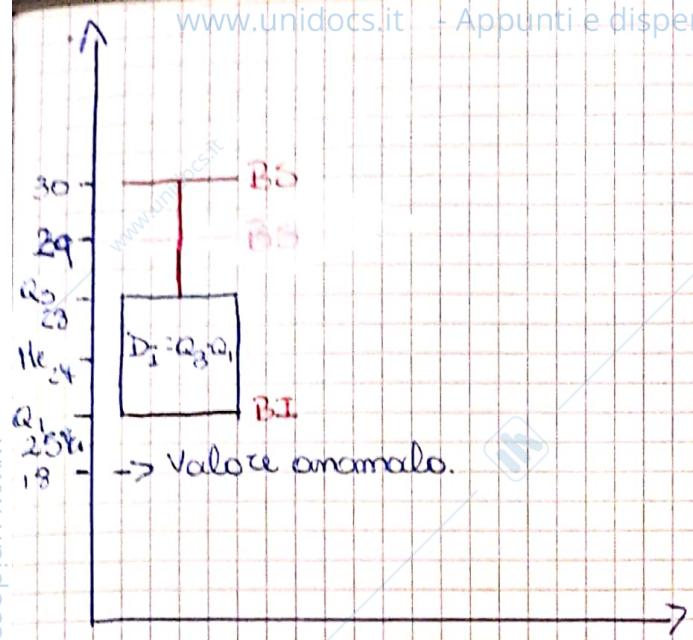
SDM = deviazione standard.

## Box Plot

È un grafico che si utilizza per le variabili numeriche e serve per confrontare due variabili insieme. Nel box plot inseriamo tutti i dati studiati.

Esempio voti d'esame:

x	n	%	F
18	1	5%	5%
25	5	25%	30% = $Q_1$
27	4	20%	50% = $Me$
28	5	25%	75% = $Q_3$
29	1	5%	80%
30	1	5%	85%
	20	100%	



Box Plot:

Basso superiore =  $BS = Q_3 + 1,5 \cdot D_1$

Basso inferiore =  $BI = Q_1 - 1,5 \cdot D_1$

Hi diammo una banda, se le variabili superano questi valori avremo dei valori anomali.

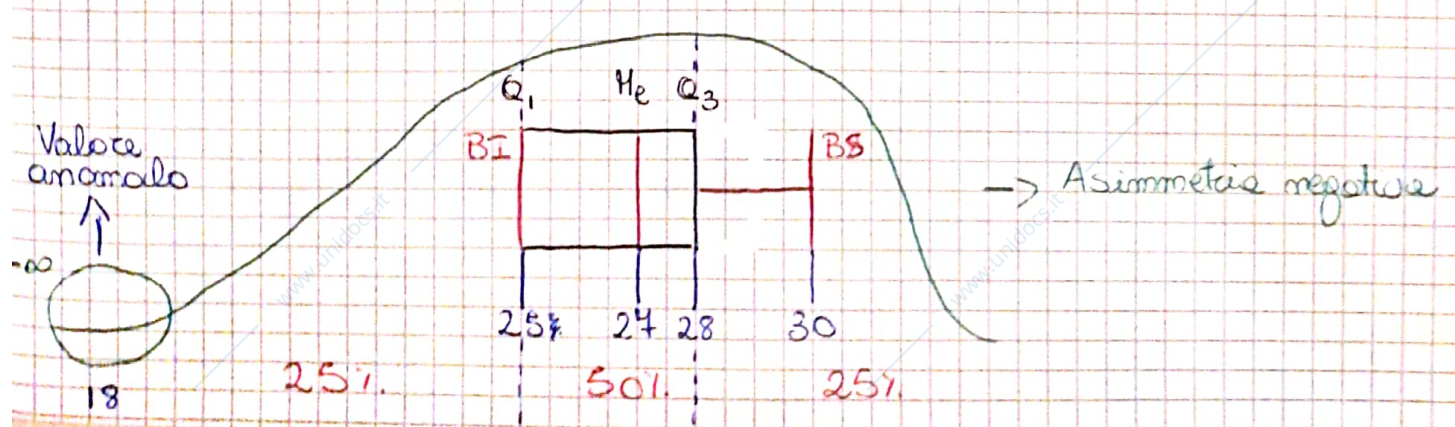
$D_1 = Q_3 - Q_1 = 28 - 25 = 3$

$BS = 28 + (1,5 \cdot 3) = 28 + 4,5 = 32,5$  → Il basso lo disegniamo a 30

$BI = 25 - (1,5 \cdot 3) = 25 - 4,5 = 20,5$  → Il valore più vicino interno sarà 25

Se come non ho valori uguali a 32,5, prendo valori interni più vicini, quindi 30 (BS) e non avendo valori uguali 20,5, prendo i valori interni, quindi 25 (BI).

Tutti i valori esterni a questi valori saranno valori anomali. Non abbiamo variabili superiori, quindi esterni, a 32,5, ma abbiamo valori inferiori, quindi esterni, a 20,5, che sarebbe 18 (È oltre la banda inferiore)



## ESERCITAZIONE

L'indice di Gini, è un indice che misura il grado di eterogeneità (omogeneità) in una distribuzione statistica di dati suddivise in  $k$  categorie (le variabili qualitative) ciascuna delle quali con frequenza relativa al quadrato.

Dati  $k=4$  (4 modalità) delle variabili soggiorno, l'eterogeneità assoluta si calcola mediante l'indice di Gini, così definito:

$$G = 1 - \sum p_i^2$$

$$G = 1 - (0,0413 + 0,0502 + 0,0182 + 0,0146 + 0,0128 + 0,0476 + 0,0028) = 0,8225$$

Poiché si ha sempre  $0 \leq G \leq \frac{k-1}{k}$  →  $0 \leq G \leq \frac{6}{4} = 0,8541$  (valore minimo → valore massimo)  
 possiamo dire che i dati sono distribuiti in modo altamente eterogeneo sulle 4 scelte (cioè le scelte relative alle tipologie di soggiorno sono molto simili tra loro).

↓  
 Questa ricerca si basa sulle tipologie di soggiorno scelte dagli italiani per trascorrere le proprie vacanze estive nel 2016 in Italia.

x	$p_i$	$p_i^2$
Albergo	0,264	0,0413
Parenti/amici	0,224	0,0502
Casa proprietà	0,135	0,0182
Villaggi turistici	0,121	0,0146
Smalaffetto	0,113	0,0128
Residence	0,087	0,0076
Campeggio	0,053	0,0028

Se vogliamo una misura relativa del grado di eterogeneità, possiamo calcolare l'indice di Gini relativo o omologizzato che ci dà più l'idea di livello di eterogeneità:

$$G_N = \frac{G}{\frac{k-1}{k}}$$

$$G_N = \frac{0,8225}{0,8541} = 0,9596$$

Totale minimo valevole anche in percentuale:  
 $Q_N \cdot 100 = Q_N \%$

$0,9596 \cdot 100 = 95,96\% \rightarrow$  (95,96% di eterogeneità = elevata eterogeneità)  
 $Q_N$  è chiaramente compreso tra 0 e 1, valori corrispondenti rispettivamente al minimo grado di eterogeneità e massimo grado. Nel nostro caso conferma l'elevata eterogeneità tra i dati della distribuzione.

**Riepilogo 2:**

Da un auto collettivo, dove  $m=5$ , la variabile  $X$  rappresenta il numero di componenti delle famiglie che presentano i seguenti valori  $\{3, 1, 2, 4, 5\}$ . Calcolare la devianza, la varianza e lo SM e il CV.

<del>xi</del>	$x_i$	$(x_i)^2$	$(x_i - \mu)$	$(x_i - \mu)^2$
<del>3</del>	3	9	0,6	0,36
<del>1</del>	1	1	-1,4	2,96
<del>2</del>	2	4	-0,4	0,16
<del>4</del>	4	16	1,6	2,56
<del>5</del>	2	4	-0,4	0,16
	12	34		5,20

$\mu = \frac{12}{5} = 2,4$

$Dev(X) = 5,20$

$Var(X) = \sigma^2 = \frac{1}{N} \sum (x_i - \mu)^2 = \frac{Dev}{N} = \frac{5,20}{5} = 1,04$

$\sigma = \sqrt{VAR(X)} = \sqrt{1,04} = 1,02$

$CV = \frac{1,02}{2,4} = 0,4249 \cdot 100 = 42,49$

**Riepilogo 2:**

La distribuzione del reddito annuo in euro dei 1.000 abitanti di un Comune è la seguente. Calcolare  $dev(x)$ ,  $var(x)$ , SM, CV.

CLASSI DI REDDITO	VALORE CENTRALE	$m_i$	$c_i \cdot m_i$	$(c_i - \mu)$	$(c_i - \mu)^2$	$(c_i - \mu)^2 \cdot m_i$
1.000 + 5.000	3.000	100	300.000	19800	392040000	39204000000
5.000 + 15.000	10.000	400	4000000	12800	163840000	65536000000
15.000 + 35.000	25.000	300	7500000	2200	4840000	1452000000
35.000 + 75.000	55.000	200	11000000	82200	1036840000	207368000000
TOTALE		1000	22800000			313560000000

$\mu = \frac{22800000}{1000} = 22.800$

$Dev(X) = 313560000000$

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

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

$\sigma = \sqrt{VAR(X)}$

$VAR(X) = \sigma^2 = \frac{1}{N} \sum (c_i - \bar{y})^2 \cdot m_i = \frac{Dev(X)}{N} = \frac{313560000}{1000} = 313560000$

$\sigma = \sqrt{VAR(X)} = \sqrt{313560000} = 17704,625$

$CV = \frac{\sigma}{\bar{y}} \cdot 100 = \frac{17704,625}{22800} = 0,7764 \cdot 100 = 77,64$

**Riepilogo 3.**

La seguente tabella fornisce la distribuzione delle 100 famiglie di un quartiere secondo il carattere  $x =$  "numero dei figli"

$x_i$	$m_i$	$N_i$
0	30	30
1	15	45
2	20	65
3	12	77
4	10	87
5	5	92
6	4	100
TOTALE	100	

$Range = x_{max} - x_{min} = 6 - 0 = 6$

$IQR = Q_3 - Q_1 = x\left(3 \cdot \frac{N+1}{4}\right) - x\left(\frac{N+1}{4}\right) = x(75,75) - x(25) = 3 - 0 = 3$

*Differenza interquartile*

Tale valore sta ad indicare che il 50% delle famiglie analizzate hanno un numero di figli in un intervallo di ampiezza 3.

Gli indici che calcolano la dispersione tra le unità sono: Range e differenza interquartile.

Gli indici che calcolano la dispersione delle unità da un valore centrale sono: SM, devianza, varianza, coefficiente di variazione

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

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

Ripetuto 4 (quinto al 8)

$x_i$	$n_i m_i$	$(x_i - H)^1$	$(x_i - H)^2$	$(x_i - H)^2 m_i$
0	0	-2	4	0
1	15	-1	1	15
2	40	0	0	0
3	36	1	1	36
4	40	2	4	160
5	45	3	9	405
6	24	4	16	384
TOTALE	200			990

M  
S  
V  
C

# STATISTICA

## BIVARIABILE

X (TITOLO DI STUDIO)	Y (Reddito)
A	H
B	B
H	A
A	A
A	A
H	H
B	H

A=Alto  
B=Basso  
H=Medio

Nella distribuzione univariata avremo due variabili statistiche, due informazioni sul nostro individuo (Titolo di studio). Oltre a studiare due variabili singolarmente, possiamo studiare le dipendenze di una o l'altra variabile tra le due variabili.

Per farlo useremo, prima di tutto, la distribuzione in frequenza

DISTRIBUZIONE IN FREQUENZA DOPIA:

		$Y_1$ B	$Y_2$ H	$Y_3$ A
$X_1$	B	$m_{11}$	$m_{12}$	$m_{13}$
$X_2$	H	$m_{21}$	$m_{22}$	$m_{23}$
$X_3$	A	$m_{31}$	$m_{32}$	$m_{33}$
	TOTALE			

-> Andrà a vedere, per esempio, quante persone che hanno un titolo di studio hanno accanto anche un reddito basso.

-> E per colonne e per riga (un'altra tabella dove facciano le dei marginali di riga)

$m_{ij}$  = frequenza assoluta congiunta (Misura congiuntamente sia le modalità di  $x$  e le modalità di  $y$ )  
 $m_{i.}$  = TOTALE FREQUENZE ASSOLUTE  
 $N$  = Numero popolazione  
 $i = X$   
 $j = Y$

ESEMPIO:

		REDDITO			
		BASSO	MEDIO	ALTO	
TITOLO DI STUDIO (X)	BASSO	5	3	1	9
	MEDIO	1	18	1	20
	ALTO	1	4	16	21
		7	25	18	50

TOTALE MARGINALI X (E per riga)

TOTALE MARGINALI Y (E per colonna)

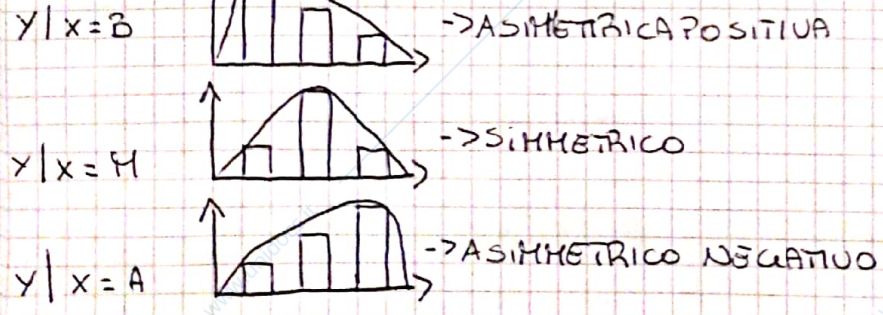
Il totale dei marginali di  $y$  e di  $x$  sarà sempre uguale, avremo cioè il numero della popolazione.

Se andiamo a prendere una singola riga prenderemo una sotto popolazione e ognuna delle righe sarà distribuzione condizionata, per esempio se ci soffermiamo sulle prime righe analizzeremo come varia la variabile  $y$  condizionata a  $x = \text{BASSO}$  ( $Y|X=B$ ). Lo stesso discorso vale per le colonne. Quindi avremo:

- Distribuzione condizionata per  $x$ ;
- Distribuzione condizionata per  $y$ .

Più le distribuzioni condizionate sono diverse tra loro, più c'è dipendenze, più sono uguali tra loro, più non c'è dipendenze.

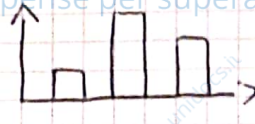
GRAFICA:



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

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

DISTRIBUZIONE CONDIZIONATA PER Y:



Dai grafici vediamo, quindi, che le distribuzioni sono diverse tra loro, quindi, hanno un rapporto di dipendenza, se fossero stati uguali non c'era dipendenza. PER COSTRUIRE IL GRAFICO DOBBIAMO FARE LE FREQUENZE RELATIVE, NON ASSOLUTA

$(m_{i,j} - m)$  -> FREQUENZA RELATIVA.

FREQUENZE RELATIVE ESERCIZIO:

	B	M	A	
B	0,56	0,33	0,11	1
M	0,05	0,30	0,05	1
A	0,05	0,19	0,46	1
	4/150	25/150	18/150	-> PROFILO MEDIO RIGA

TOTALE FREQUENZE RELATIVE.

Le frequenze relative per riga se chiamiamo **PROFILI DI RIGA**

La frequenza relativa la puoi calcolare sia per righe che per colonne, è uguale. (**PROFILO DI COLONNA**)

Se i profili di riga (o profili di colonna) sono uguali al profilo medio, abbiamo un rapporto di **indipendenza** tra x e y.

$$\frac{m_{i,j}}{m_i} = \frac{m_{.j}}{N}$$

$$\rightarrow m_{i,j} = \frac{m_i \cdot m_{.j}}{N}$$

-> FREQUENZA CONGIUNTA OSSERVATA.

TABELLA INIZIALE (Distribuzione in frequenza)

PROFILO MEDIO RIGA ESERCIZIO

0,14	0,50	0,36	1
------	------	------	---

FREQUENZA CONGIUNTA TEORICA ( $m_{i,j}^* =$  FREQUENZA TEORICA)

	B	M	A
B	$\frac{4,9}{50} = 1,26$	4,50	3,24
M	2,80	10,00	4,20
A	2,94	10,50	4,56

Ci sarà dipendenza quando:  
 $m_{i,j} \neq m_{i,j}^*$

La frequenza congiunta osservata è diversa alla frequenza teorica

PIÙ C'È DIFFERENZA PIÙ C'È DIPENDENZA

Ci sarà indipendenza quando  
 $m_{i,j} = m_{i,j}^*$

Per vedere quanta differenza abbiamo, calcoliamo la **contingenza**, è come se fosse uno scudo. Faremo la differenza tra  $m_{i,j}$  e  $m_{i,j}^*$  ( $m_{i,j} - m_{i,j}^*$ )

Se la differenza è uguale a 0, avremo indipendenza  
la contingenza la indichiamo con  $C_{ij}$

**CONTINGENZA ESERCIZIO:**

	B	H	A	
B	3,74	-1,50	-2,24	0
H	-1,80	8	-6,20	0
A	-1,34	-6,50	8,44	0
	0	0	0	

La somma delle contingenze (per riga / per colonna) è 0.  
Dobbiamo stabilire quanto dipendenza c'è.

$$\chi^2_{MAX} = N \cdot \min(I - 1 \text{ e } 5 - 1)$$

$$465 \cdot 1 = 465$$

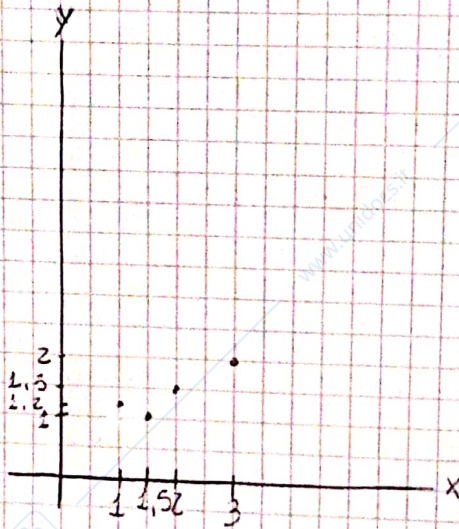
$$V = \frac{4,08}{465} = 0,01 \text{ } \text{C'è pochissima dipendenza}$$

Il diagramma a dispersione o scatter plot

è un grafico che permette di rappresentare nel grafico due variabili

ESEMPLO:

X	Y
1	1,2
1,5	1
2	1,5
3	2



Se la differenza è uguale a 0, avremo indipendenza. ( $C_{13} = 0$ )

La contingenza la indichiamo con  $C_{13}$

### CONTINGENZA ESERCIZIO:

	B	H	A	
B	3,44	-1,50	-2,24	0
H	-1,80	8	-6,20	0
A	-1,34	-6,50	8,44	0
	0	0	0	

Il valore di una variabile non cambia  
al variare dell'altra variabile



La somma delle contingenze (per riga / per colonna) ci darà sempre 0.

Dobbiamo stabilire quanto dipendenza c'è.

Usiamo il quadrato della contingenza, in quanto elimina i valori negativi e amplifica la contingenza.

Aviamo quindi  $(m_{13} - m_{13}^e)^2 = C_{13}^2$

Dopo ricaviamo lo scostamento relativo, in questo modo  $\left(\frac{C_{13}^2}{m_{13}^e}\right)$  Rapporto tra la contingenza al quadrato e la frequenza attesa.

### SCOSTAMENTO ESERCIZIO (INDICE DI PEARSON)

	B	H	A	
B	1,10	0,60	1,55	13,15
H	1,16	6,40	5,84	12,90 → PEARSON
A	1,28	4,02	9,42	14,72
	13,54	10,92	16,31	40,7 → TOTALE

↓  
PEARSON

A questo punto misuriamo la dipendenza in una tabella e dopo è indicato con l'indice di associazione di Pearson.

$$\chi^2 = \sum_i \sum_j \frac{(x_{ij} - x_{ij}^e)^2}{m_{ij}^e}$$

= 0 → INDIPENDENZA ( $C_{13} = 0 \forall_{11}$ ) ( $m_{13} = m_{13}^e \forall_{11}$ )

> 0 → DIPENDENZA (quanto detto INDIPENDENZA)

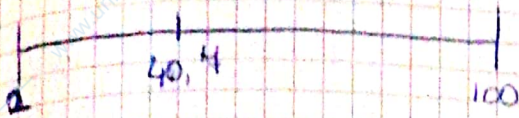
Dobbiamo quindi stabilire il valore massimo, con questa formula:

$$\chi_{MAX}^2 = N \cdot \min(I-1 \text{ e } J-1)$$

↓                      ↓  
 Numero            Numero  
 di righe            di colonne

$\chi^2_{max} = 50 \cdot 2 = 100$  (3-1) e (2-1) = (2 e 2)

Quindi:



... una sola coppia  
... valori diversi per i due  
... e valori per i due

Per capire con valori relativi si mostra un momento dipendenza o siamo la

V di Cramer

$$V = \frac{\chi^2}{\chi^2_{MAX}}$$

$$V = \frac{40,4}{100} = 0,41 = 41\%$$

Indice di associazione di Pearson.

L'indice di associazione di Pearson è un indice simmetrico, cambiando la disposizione delle variabili, l'indice non cambia. X ed Y si influenzano tra loro.

Indipendenza

Avremo un rapporto di indipendenza quando i rapporti tra i profili riga è uguale al rapporto tra i profili medi di riga.

Indipendenza

È il valore più forte di dipendenza tra due variabili.

Esercizio

	L	M	N	
H	40	110	180	} m <sub>j</sub>
F	85	200	285	
	155	310	465	
	} m <sub>i</sub>			

	L	M	N
H	60	120	180
F	95	190	285
	155	310	465

$$m_{ij}^* = \frac{m_i \cdot m_j}{N}$$

$\chi^2 = C_{ij}^2$

	L	M
H	1,64	0,89
F	1,05	0,05

$$C_{ij}^2 = \frac{(m_{ij} - m_{ij}^*)^2}{m_{ij}^*}$$

$$\chi^2 = \sum_i \sum_j \frac{(x_{ij} - x_{ij}^*)^2}{m_{ij}^*} = 4,08$$