

Esercizi di Statistica per Biotecnologie

Francesca Pizzorni Ferrarese

Esercitazione IV– Statistica descrittiva II

Es. 1

Dopo il disastro, una commissione d'inchiesta del British Board of Trade ha compilato una lista di tutti i 1316 passeggeri con alcune informazioni aggiuntive riguardanti: l'esito (salvato, non salvato), la classe (I, II, III) in cui viaggiavano, il sesso, l'età,...

In questo esercizio ci limitiamo nella prima parte a considerare le informazioni sull'esito e la classe.

Esito	Classe			Totale
	I	II	III	
Salvato	203	118	178	499
Non Salvato	122	167	528	817

- Se possibile, calcolare la moda, la mediana e la media della serie
- Determinare la percentuale di passeggeri di prima classe non salvati
- Quale percentuale dei passeggeri salvati era di prima classe?

Es. 2

I seguenti dati sono i punteggi che 10 studenti hanno conseguito in due esami di Analisi

Matematica (punteggio massimo = 100).

<i>Analisi I</i>	<i>Analisi II</i>
51	74
68	70
97	93
55	67
95	99
74	73
20	33
91	91
74	80
80	86

- Rappresentare i dati mediante scatter plot
- Se possibile, calcolare la moda, la mediana e la media della serie
- Calcolare la covarianza e la correlazione.

Es. 3

Alla fine di una giornata di lavoro un intervistatore si accorge di aver perso i dati raccolti su un certo numero di famiglie relativamente al numero X di televisori

Numero televisori X	Numero componenti Y			Totale
	1	2	3	
0	0	3	1	?
1	3	?	7	16
2	1	?	?	?
Totale	?	13	11	?

posseduti e al numero Y di componenti della famiglia. Ricostruendo a memoria le interviste fatte, arriva alla seguente tabella

- Si completi la tabella e si dica se i due caratteri sono indipendenti.
- Si calcolino le medie e le varianze del numero di televisori condizionatamente al numero di componenti della famiglia.

Es. 4

Nel settore delle cucitrici a macchina in una fabbrica di capi di abbigliamento viene assegnato un punteggio a ogni capo finito sulla base della sua qualità. La paga di ciascuna cucitrice dipende in parte dal numero di capi completati. In una particolare

Cucitrice	Numero capi (X)	Punteggio (Y)
1	14	7,2
2	13	7,3
3	17	6,9
4	16	7,3
5	17	7,5
6	18	7,6
7	19	6,8
8	32	3,7
9	18	6,5
10	15	7,9
11	15	6,8
12	19	7,1

giornata si sono registrati i seguenti punteggi medi per pezzo:

- Calcolare il coefficiente di correlazione tra X e Y e commentare.
- Tracciare il diagramma di dispersione (scatter plot) dei dati
- Si costruisca la retta di regressione e sovrapporre il grafico della retta di regressione al diagramma di dispersione di X e Y .
- Calcolare l'indice della bontà di adattamento della retta ai dati.
- In base al modello di regressione si preveda il punteggio corrispondente ad un numero di capi pari a 20.
- Ricalcolare la retta di regressione dopo aver eliminato l'osservazione sull'ottava cucitrice e tracciarne il grafico.

Es. 5

Si considerino i seguenti dati relativi alle importazioni e alle esportazioni in migliaia di milioni di euro di alcune regioni italiane nel 2010:

	Piemonte	Lombardia	Veneto	Toscana	Lazio
Importazioni	33,6	118,5	35,9	21,1	22,5
Esportazioni	51,0	115,8	54,0	32,3	14,1

- a) Si calcoli la retta di regressione delle esportazioni rispetto alle importazioni.
- b) Si tracci la retta di regressione sul grafico di dispersione delle due variabili.
- c) Si calcoli la bontà di adattamento.
- d) Si prevedano le esportazioni della regione Emilia-Romagna sapendo che le importazioni sono state pari a 22,58.

Es. 6

Nella seguente tabella si riportano le misure dell'ossigeno consumato da una persona che cammina, in corrispondenza a varie velocità della persona.

<i>velocità (Km/h)</i>	<i>ossigeno (litri/h)</i>
0	19
1	20
2	20.5
3	21.5
4	22
5	23
6	23
7	23.5
8	24

Verificare che esiste dipendenza lineare del consumo di ossigeno dalla velocità, determinare la retta di regressione e verificare la bontà del modello.

ESERCITAZIONE 1

1) Serie bi-variate: serie statistica in cui ad ogni unità statistica si riferiscono 2 caratteri
 $o_i = (x_i; y_i)$ $O = \{o_i\}$

Frequenze in valore assoluto ad ogni modalità i, j

		CLASSE C			TOTALE
		I	II	III	
Tabella a doppia entrata $m_{i,j}$: freq. assolute	SALARIO S				
	SI	203	118	178	499
	NO	122	167	528	817
	TOTALE	325	285	706	1316

$m_{i,j}$ frequenze assolute

a) moda: modalità con frequenza maggiore $mo = (NO; III)$

mediana: osservazione che bipartisce la popolazione → solo dati ordinabili
 → la mediana non si applica alle serie bivariate → NO

media: media delle osservazioni → solo dati quantitativi → NO

b) Distribuzione condizionata

$$P(S=NO | C=I) = \frac{122}{325} \cdot 100 = 37,54\%$$

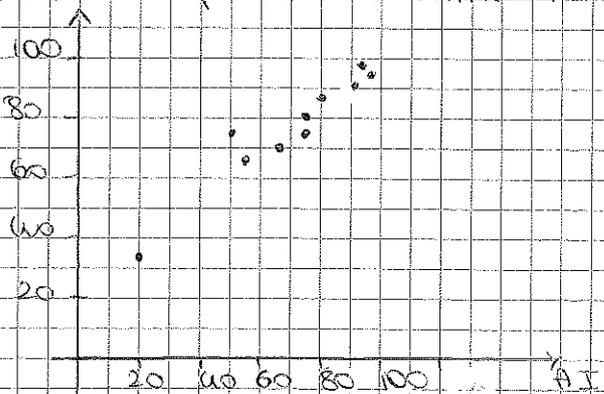
$$c) P(C=I | S=SI) = \frac{203}{499} \cdot 100 = 40,68\%$$

2) Scatter plot: rappresentazione grafica usata per serie quantitative, in cui ogni osservazione viene riportata in un piano cartesiano come se fosse un pt.

b) $mo = (97; 99)$

mediana NO

media: somma più vertici e diviso
 il valore per n^2



$$\left[(81; 76) + (68; 70) + (97; 93) + (55; 67) + (95; 99) + (76; 73) + (20; 33) + (91; 91) + (76; 80) + (80; 86) \right] / n^2 =$$

$$= (51 + 68 + 97 + 55 + 95 + 74 + 20 + 91 + 74 + 80; 74 + 70 + 93 + 67 + 99 + 73 + 33 + 91 + 80 + 86) / 10 =$$

$$= (705/10; 766/10) = (70,5; 76,6) \leftarrow \begin{array}{l} \text{Ead media) oltro) sive} \\ \text{E oltro) dalla media) \\ \text{dei carateri.} \end{array}$$

c) non è possibile adattare in gruppo ad un'unica scala di dimensioni \rightarrow esiste solo la covarianza, $S = (\bar{x}, \bar{y})$

Come misura il "grado di legame" fra 2 carateri?

Dati in associazione $(m_i; y_i)$ valutati separatamente (2 scale)

$$(m_i - \bar{m})(y_i - \bar{y}) \begin{cases} > 0 & \text{scarti concordi} \rightarrow \text{d'ordine} \\ < 0 & \text{d'inverso} \end{cases} \quad \begin{array}{l} \uparrow \\ \text{proporzionalità} \end{array}$$

\Rightarrow Dati una serie mirata alle 2 associazioni $(m_i; y_i)$ di un solo carateri

$$S_{xy} = \frac{1}{N} \sum_{i=1}^N (m_i - \bar{m})(y_i - \bar{y})$$

$$\text{Varianza } x = (51^2 + 68^2 + \dots + 80^2) / 10 - 70,5^2 = 501,45$$

$$\text{Varianza } y = (74^2 + 70^2 + \dots + 86^2) / 10 - 76,6^2 = 405$$

m_i	y_i	$m_i - \bar{m}$	$y_i - \bar{y}$	$(m_i - \bar{m}) \cdot (y_i - \bar{y})$
51	74	-19,5	-2,6	50,7
68	70	-2,5	-6,6	16,5
97	93	26,5	16,4	434,6
55	67	-15,5	-9,6	148,8
95	99	24,5	22,4	548,8
74	73	3,5	-3,6	-12,6
20	33	-50,5	-43,6	2201,8
91	91	20,5	14,4	295,2
74	80	3,5	3,4	11,9
80	86	9,5	9,4	89,3

$$\Rightarrow S_{xy} = \frac{3785}{10} = 378,5 \rightarrow \text{matrice varianze} = \Sigma = \begin{bmatrix} S_x^2 & S_{xy} \\ S_{xy} & S_y^2 \end{bmatrix} = \begin{bmatrix} 501,45 & 378,5 \\ 378,5 & 405 \end{bmatrix}$$

Correlazione di Pearson: data una serie bivariata $(m_i; y_i)$ di unisco collettivo a volte $R = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$

$\Rightarrow R = \frac{378,5}{\sqrt{501,45} \sqrt{405}} = 0,89$ R è meno influenzato dalla variabilità

3) a)

NUMERO TELEFONI X	NUMERO COMPONENTI Y			TOTALE
	1	2	3	
0	0	3	1	4
1	3	6	7	16
2	1	6	3	8
TOTALE	4	13	11	28

X	f(x)	f(x y ₁)	f(x y ₂)	f(x y ₃)
0	4/28	0/4	3/13	1/11
1	16/28	3/4	6/13	7/11
2	8/28	1/4	4/13	3/11

non sono indipendenti.

b) Data una serie bivariata $(m_i; y_i)$ dove $M = M_x M_y$ matriciale, con frequenze relative $f_{i,j}$ di unisco collettivo

$$\sigma_{xy} = \sum_{i=1}^M \sum_{j=1}^N (m_i - \bar{m})(y_j - \bar{y}) f_{i,j}$$

$$\bar{m} = \frac{4 \cdot 0 + 16 \cdot 1 + 8 \cdot 2}{28} = 1,28$$

$$\bar{y} = \frac{4 \cdot 1 + 13 \cdot 2 + 11 \cdot 3}{28} = 2,25$$

$$\begin{aligned} \sigma_{xy} &= \sum_{i=1}^3 \sum_{j=1}^3 (m_i - 1,28)(y_j - 2,25) f_{i,j} \\ &= \frac{(0 - 1,28)(1 - 2,25) \cdot 0}{-1,28 \cdot -1,25} + \frac{(0 - 1,28)(2 - 2,25) \cdot 3}{-1,28 \cdot -0,25} + \\ &+ \frac{(0 - 1,28)(3 - 2,25) \cdot 1}{-1,28 \cdot 0,75} + \frac{(1 - 1,28)(1 - 2,25) \cdot 3}{-0,28 \cdot -1,25} + \\ &+ \frac{(1 - 1,28)(2 - 2,25) \cdot 6}{-0,28 \cdot -0,25} + \frac{(1 - 1,28)(3 - 2,25) \cdot 7}{-0,28 \cdot 0,75} + \\ &+ \frac{(2 - 1,28)(1 - 2,25) \cdot 1}{0,72 \cdot -1,25} + \frac{(2 - 1,28)(2 - 2,25) \cdot 4}{0,72 \cdot -0,25} + \\ &+ \frac{(2 - 1,28)(3 - 2,25) \cdot 3}{0,72 \cdot 0,75} = 0 + 0,96 - 0,96 + 1,05 + 0,42 + \\ &- 1,47 - 0,9 - 0,72 + 1,62 = 0 \end{aligned}$$

9) a) $R = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}}$

$$\text{me}(x) = \frac{14+13+17+\dots+15+19}{12} = 17,75$$

$$\text{me}(y) = \frac{7,2+7,3+6,9+\dots+6,8+7,1}{12} = 6,88$$

$$\text{Var}(x) = \frac{14^2+13^2+17^2+\dots+15^2+19^2}{12} - 17,75^2 = 336,92 - 315,0625 = 21,8575$$

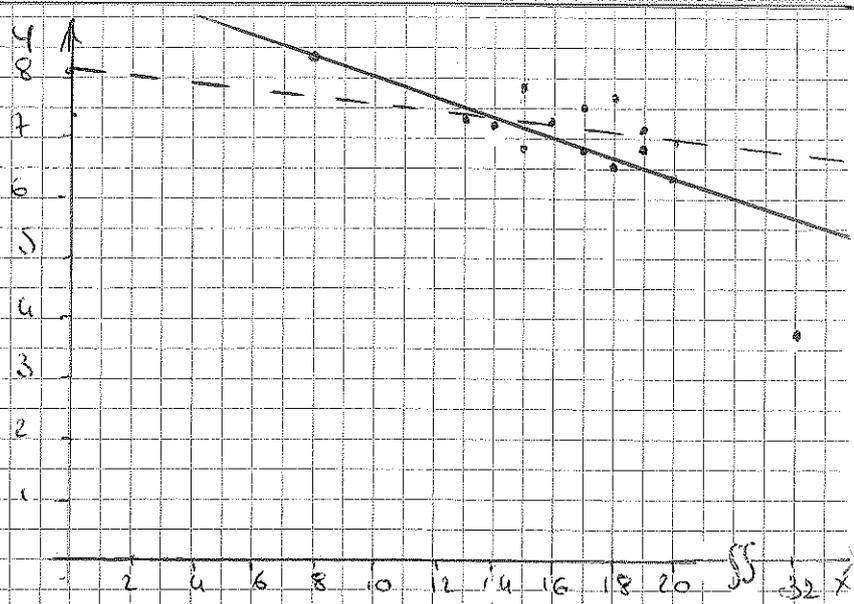
$$\text{Var}(y) = \frac{7,2^2+7,3^2+6,9^2+\dots+6,8^2+7,1^2}{12} - 6,88^2 = 48,66 - 47,3344 = 1,3256$$

x_i	y_i	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$
14	7,2	-3,75	0,32	-1,2
13	7,3	-4,75	0,42	-1,995
17	6,9	-0,75	0,02	-0,015
16	7,3	-1,75	0,42	-0,735
17	7,5	-0,75	0,62	-0,465
18	7,6	0,75	0,72	0,54
19	6,8	1,75	-0,08	-0,14
30	3,7	10,25	-3,18	-32,795
18	6,5	0,75	-0,38	-0,285
15	7,9	-2,75	1,02	-2,805
15	6,8	-2,75	-0,08	0,22
19	7,1	1,75	0,22	0,385

$$\sum xy = \frac{52,25}{12} = 4,35 \quad \text{and} \quad -52,25$$

$$\Rightarrow R = \frac{4,35}{\sqrt{21,8575 \cdot 1,3256}} = \frac{4,35}{4,92} = 0,88$$

b)



c) Y = variabile dipendente
 m^o capi X = variabile indipendente

Retta di regressione

- modello di tipo affine $m_i: \hat{y}_i = a m_i + b$

- parametri valutati facendo residuo minimo

$$SSR = \sum_{i=1}^N (\hat{y}_i - y_i)^2 = \sum_{i=1}^N (a m_i + b - y_i)^2$$

$$\arg \min_{a, b} \sum_{i=1}^N (a m_i + b - y_i)^2$$

- si dimostra che $a = \frac{\sum m_i y_i}{\sum m_i^2}$ $b = \bar{y} - a \bar{m}$

$$\Rightarrow a = \frac{-4,35}{21,86} = -0,199$$

$$b = 6,88 + 0,199 \cdot 17,75 = 10,61$$

\Rightarrow l'equazione di regressione stimata è $\hat{y} = -0,199 m + 10,61$

a) Valuto il coefficiente di Pearson

IRI 0,3 scarsa prob. di legame

0,3 < IRI < 0,7 moderata prob. di legame lineare

0,7 < IRI all'alta prob. " " \leftarrow IRI = 0,88

$$e) \hat{y} = am + b = -0,199 \cdot 20 + 10,41 = 6,43$$

↑
m° copri

$$f) me(x) = \frac{17,75 \cdot 12 - 32}{11} = 16,45$$

$$me(y) = \frac{6,88 \cdot 12 - 3,7}{11} = 7,17$$

$$\sigma_{xy} = \frac{-2,5635}{11} = -0,2330$$

$$\sigma_x^2 = \frac{3019}{11} - 16,45^2 = 270,45 - 270,6025 = 3,8520$$

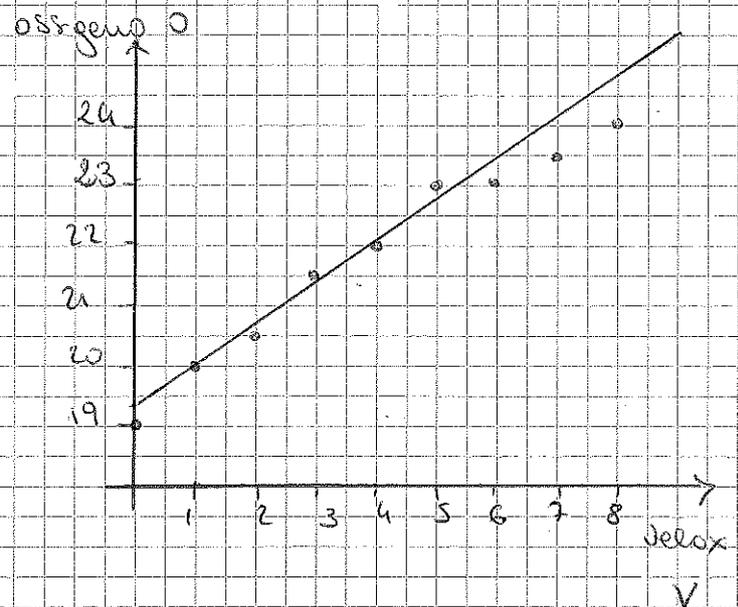
$$\sigma_y^2 = \frac{567,59}{11} - 7,17^2 = 51,6 - 51,4089 = 0,19$$

$$a = \frac{\sigma_{xy}}{\sigma_x^2} = \frac{-0,2330}{3,8520} = -0,0605$$

$$b = 7,17 + 0,0605 \cdot 16,45 = 8,1652$$

$$\hat{y}_i = -0,0605x + 8,1652$$

6) Il diagramma di dispersione mostra che c'è una relazione lineare all'incirca tra la velocità e l'altezza del pallone.



$$\bar{x} = \frac{0 + 1 + \dots + 7 + 8}{9} = \frac{36}{9} = 4$$

$$\bar{y} = \frac{19 + 20 + \dots + 23,5 + 24}{9} = \frac{196,5}{9} = 21,8\bar{3}$$

$$\sigma_x^2 = \frac{0^2 + 1^2 + \dots + 7^2 + 8^2}{9} - 4^2 = 22,67 - 16 = 6,67$$

x_i	y_i	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$
0	19	-4	-2,83	11,32
1	20	-3	-1,83	5,49
2	20,5	-2	-1,33	2,66
3	21,5	-1	-0,33	0,33
4	22	0	0,17	0
5	23	1	1,17	1,17
6	23	2	1,17	2,34
7	23,5	3	1,67	5,01
8	24	4	2,17	8,68
				<hr/> 37

$$S_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) =$$

$$= \frac{37}{9} = 4,11$$

Resta da determinare

$$a = \frac{S_{xy}}{S_x^2} = \frac{4,11}{6,67} = 0,62$$

$$b = \bar{y} - a\bar{x} = 21,83 - 0,62 \cdot 4 = 49,35$$

$$\Rightarrow \hat{y}_i = a x_i + b \quad \text{modello di regressione}$$

$$= 0,62 x_i + 49,35$$

$$R = \frac{S_{xy}}{\sqrt{S_x^2 S_y^2}}$$

$$S_y^2 = \frac{49^2 + 20^2 + \dots + 23,5^2 + 24^2}{9} \quad 21,83^2 = \frac{4313,75}{9} = 476,55 =$$

$$= 479,30 - 476,55 = 2,75$$

$$\Rightarrow R = \frac{4,11}{\sqrt{6,67 \cdot 2,75}} = 0,96 > 0,7 \quad \text{adatta per il regime lineare}$$

3) cost. matrice varianza/covarianza)

$$\Sigma = \begin{bmatrix} \sigma^2_x & \sigma_{xy} \\ \sigma_{xy} & \sigma^2_y \end{bmatrix} = \begin{bmatrix} 0,07 & 0 \\ 0 & 0,47 \end{bmatrix}$$

$$\sigma^2_x = \frac{(0^2) \cdot 4 + (1^2) \cdot 16 + (2^2) \cdot 8}{28} \quad 1,28^2 = 1,64 - 1,64 = 0,07$$

$$\sigma^2_y = \frac{(1^2) \cdot 4 + (2^2) \cdot 13 + (3^2) \cdot 11}{28} \quad 2,25^2 = 5,06 - 5,06 = 0,47$$

1) B. Proprietà R)

a) $-1 \leq R \leq 1$

b) Se per ogni osservazione $y_i = a \cdot x_i + b$ allora $|R| = 1$

c) se $(x; y)$ sono ottenute da realizzazioni di v.v. c.c. X ed Y indipendenti \Rightarrow
 $R = \text{cov}(X, Y) = 0$
 $E[R] = 0$