

Ulteriori Conoscenze di Informatica e Statistica

Carlo Meneghini

Dip. di fisica - via della Vasca Navale 84,
st. 83 (I piano) tel.: 06 55 17 72 17

meneghini@fis.uniroma3.it

Comportamento in laboratorio

Durante le esercitazioni in laboratorio evitare di:

- utilizzare le macchine del laboratorio di calcolo per navigare in internet consultando pagine e siti non pertinenti al corso o alle esercitazioni
- scaricare, installare e utilizzare programmi non pertinenti al corso o alle esercitazioni
- giocare al computer
- cambiare le impostazioni del sistema (es. home page, password di accesso)

Non utilizzare il **desktop** per salvare il proprio lavoro ma creare la **propria directory** e salvare i lavori in opportune sottodirectories.

Alla fine della lezione effettuare correttamente lo shutdown del sistema e **spegnere il PC**

Note

I PC a disposizione nel centro di calcolo del Dip. di Fisica sono macchine *dual-boot*, possono cioè utilizzare più di un sistema operativo. Nel nostro caso W2K e Linux. Dopo il caricamento del BIOS il computer aspetta per alcuni secondi in attesa che voi scegliate il sistema da usare, **Scegliete windows**.

Dopo il caricamento del sistema una finestra vi chiederà il nome utente e la password. Digitando il vostro nome utente e password entrate nel sul desktop di Windows 2000.

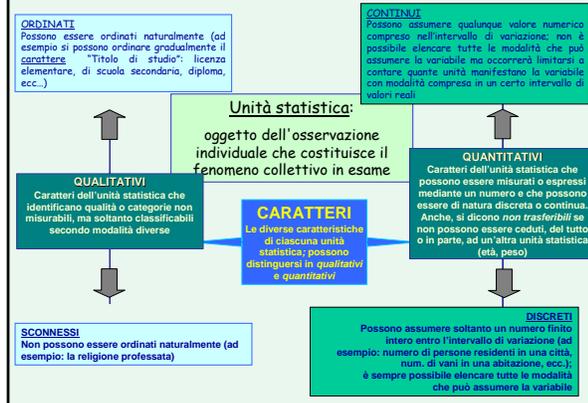
Piano della lezione

Statistica descrittiva:

- **Strumenti software: il foglio elettronico***
- **Le variabili aleatorie: caratteristiche**
- **distribuzioni di probabilità**

***Nota:** programmi freeware (come OpenOffice della Sun) sono del tutto equivalenti e quasi totalmente compatibili con il pacchetto Excel della Microsoft. La versione integrale del pacchetto OpenOffice è gratuita per studenti e istituzioni accademiche.

Statistica descrittiva



La rappresentazione dei dati statistici deve essere organizzata in modo da:

- semplificare i confronti
- sintetizzare i risultati

Esp. 1		Esp. 2	
oss.	tip.	oss.	tip.
1	M	1	M
2	F	2	F
3	F	3	F
4	F	4	F
5	M	5	F
6	M	6	M
7	F	7	F
8	M	8	M
9	M	9	M
10	F	10	F
11	F	11	F
12	F	12	M
13	M	13	M
14	M	14	M
15	F	15	F
16	F	16	F
17	M	17	M
18	M	18	M
19	M	19	M
20	F	20	F
21	M	21	M
22	F	22	F
23	F	23	F
24	M	24	M
25	M	25	M
26	F	26	F

Esp. 1		Esp. 2	
M	F	M	F
7	8	14	12

Freq. Assoluta

Freq. Relativa

Esp. 1		Esp. 2	
M %	F %	M	F
46,7	53,3	53,8	46,2

Fenomeni

deterministici:

se ripetuti nelle medesime condizioni producono gli stessi risultati

Caduta dei gravi

$$F = mg$$

$$v = mgt$$

$$x = \frac{1}{2} gt^2$$

$x, v =$ variabili deterministiche

aleatori:

pur ripetuti nelle medesime condizioni possono produrre risultati differenti

Lancio di dadi



Quale numero?

$N =$ variabile aleatoria

Testa o croce?

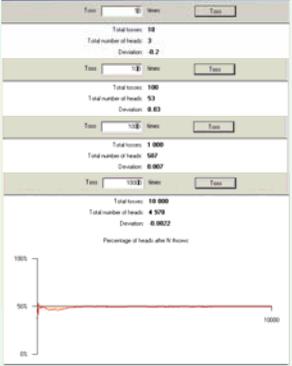


Lancia la moneta
(moneta ex)

$$p_T = \lim_{N \rightarrow \infty} \frac{n_T}{N}$$

$$p_+ = \lim_{N \rightarrow \infty} \frac{n_+}{N}$$

$p_T = 0.5 = p_+$

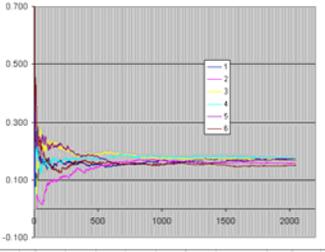


Lancio di un dado



Lancio di un dado

Lancio	1	2	3	4	5	6
1	0.000	0.000	0.000	0.000	0.000	1.000
2	0.000	0.000	0.000	0.000	0.000	1.000
3	0.000	0.000	0.000	0.000	0.333	0.667
4	0.250	0.000	0.000	0.000	0.250	0.500
5	0.200	0.000	0.000	0.000	0.400	0.400
6	0.200	0.000	0.000	0.000	0.400	0.400
7	0.5	0.5	0.5	0.5	0.5	0.5
8	4	4	4	4	4	4
9	4	4	4	4	4	4
10	6	6	6	6	6	6
11	6	6	6	6	6	6
12	4	4	4	4	4	4
13	3	3	3	3	3	3
14	6	6	6	6	6	6
15	5	5	5	5	5	5
16	1	1	1	1	1	1
17	4	4	4	4	4	4
18	4	4	4	4	4	4
19	3	3	3	3	3	3
20	2	2	2	2	2	2
21	5	5	5	5	5	5
22	1	1	1	1	1	1
23	3	3	3	3	3	3
24	5	5	5	5	5	5
25	3	3	3	3	3	3
26	6	6	6	6	6	6
27	5	5	5	5	5	5
28	4	4	4	4	4	4
29	6	6	6	6	6	6
30	5	5	5	5	5	5



Frequenza: Variabili discrete

$X =$ variabili aleatoria, $N =$ numero di osservazioni

x_1, x_2, \dots, x_i : valori assunti dalla variabile X

n_1, n_2, \dots, n_i : numero di volte che si osserva il valore i -esimo x_i

n_i : frequenza assoluta della variabile x_i con: $\sum_{i=1}^V n_i = n$

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

frequenza relativa della variabile x_i

$\sum_{i=1}^V f_i = 1$

$f_i \geq 0$

Frequenza: Variabili continue

$X =$ variabili aleatoria

x : valori assunti dalla variabile X ,

$f(x)$: densità di frequenza della variabile aleatoria X

$f(x) \geq 0$

$\int_{-\infty}^{\infty} f(x) dx = 1$

$f(x_i)$: frequenza relativa nell'intervallo $x_i < x < x_{i+1}$

$f(x_i) = \int_{x_i}^{x_{i+1}} f(x) dx$

$f(x_i) = \frac{\text{numero di osservazioni tra } x_i \text{ e } x_{i+1}}{N}$

La statistica descrittiva sintetizza l'informazione contenuta nell'insieme dei valori assunti da una variabile aleatoria (distribuzione) utilizzando:

- indici di posizione
- indici di dispersione (variabilità)
- indici di forma
- istogrammi di frequenza
- box plot

Indici di "posizione" (indici di tendenza)

indice	definizione	funzione EXCEL
Media	$\frac{1}{N} \sum_{i=1}^N x_i = \bar{x} = \langle x \rangle$ $\frac{\int_A^B f(x) dx}{\int_A^B dx} = \bar{x} = \langle x \rangle$	MEDIA(<i>dati</i>)
Moda	Valore della variabile cui corrisponde la massima frequenza	MODA(<i>dati</i>)
Mediana	Valore della variabile che permette di dividere la distribuzione delle osservazioni in due parti uguali	MEDIANA(<i>dati</i>)
Quantili		QUARTILE(<i>dati</i> ;q)

Indici di "dispersione"

indice	definizione	funzione EXCEL
Varianza	$\frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N - 1} = \sigma^2$	VAR(<i>dati</i>)
Deviazione Standard	$\sqrt{\sigma^2} = \sigma$	DEV.ST(<i>dati</i>)
Interquartile	Q3-Q1	

Media

La **media** si determina attraverso la funzione **MEDIA** [AVERAGE].

Il risultato di questa funzione è la media aritmetica

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

Studenti	Altezza (cm)	Media (cm)
marco	171	173,2
antonella	170	
luca	186	
marina	154	
gianna	166	
luigi	176	
francesco	182	
michele	178	
stefania	176	
claudia	173	

Come si fa:

=MEDIA(B2:B11)

Moda

La **moda** di un collettivo, distribuito secondo un carattere, è la modalità prevalente del carattere ossia quella a cui è associata la massima frequenza.

Si determina mediante la funzione **MODA** [MODE].

Studenti	Altezza (cm)	Moda (cm)
marco	171	170
antonella	169	
luca	186	
marina	170	
gianna	166	
luigi	190	
francesco	180	
michele	172	
stefania	174	
claudia	173	

Come si fa:

=MODA (B2:B11)

Mediana

La **mediana** suddivide ogni distribuzione ordinata in due distribuzioni aventi ciascuna una numerosità (o una quantità) che è il 50% della numerosità (o della quantità) della distribuzione totale.

Si determina mediante la funzione **MEDIANA** [MEDIAN].

Studenti	Altezza (cm)	Mediana (cm)
marco	171	173,2
antonella	169	
luca	186	
marina	170	
gianna	166	
luigi	190	
francesco	180	
michele	172	
stefania	174	
claudia	173	

Come si fa:

=MEDIANA (B2:B11)

Quantili

Si può dividere la distribuzione parti (percentili) contenenti ognuna la q -esima parte della quantità della distribuzione totale.

I quantili sono le n parti in cui è stata suddivisa una distribuzione.

per $q = 4$ (più usati) si parla di **quartili**

I quartili dividono la distribuzione in quattro parti aventi ognuna il 1/4 (25%) della quantità totale;

Il I quartile (Q1) è il limite superiore della distribuzione che ha il 25% della quantità totale;

Il II quartile (Q2) è il limite superiore della seconda distribuzione e quindi da solo separa nella distribuzione totale due distribuzioni che hanno ciascuna il 50% della quantità totale, il Q2 coincide con la mediana;

Il III quartile (Q3) è il limite superiore della distribuzione che ha il 75% dell'ammontare della distribuzione totale.

I quartili si determinano mediante le funzioni **QUARTILE [QUARTILE]** e **PERCENTILE [PERCENTILE]**.

QUARTILE (sequenza di numeri o indirizzo di cella; 0 o 1 o 2 o 3 o 4) (0 = minimo; 1 = 1° quartile; 2 = mediana; 3 = 3° quartile; 4 = massimo)

A	B	C	D	E	F	G
1	Studenti	Altezza (cm)				
2	marco	171				
3	antonella	169	1° quartile (cm)	170,25		
4	luca	186	2° quartile (cm)	183		
5	marina	170	3° quartile (cm)	183		
6	gianna	165				
7	lugi	190				
8	francesco	190				
9	michele	172				
10	stefania	174				
11	claudia	173				

Come si fa:

- Q1: =QUARTILE (B2:B11,1)
- Q2: =QUARTILE (B2:B11,2)
- Q3: =QUARTILE (B2:B11,3)

PERCENTILE (sequenza di numeri o indirizzo di cella; numero compreso tra 0 ed 1) (percentile p%: inserire il numero p)

A	B	C	D	E	F
1	Studenti	Altezza (cm)			
2	marco	171			
3	antonella	169	1° quartile (cm)	170,25	
4	luca	186	2° quartile (cm)	183	
5	marina	170	Percentile 85%	180,25	
6	gianna	165			
7	lugi	190			
8	francesco	190			
9	michele	172			
10	stefania	174			
11	claudia	173			

Come si fa:

- Spostare il cursore nella cella C5 e digitare: Percentile 85% (cm)
- Spostare il cursore nella cella D5 e inserire la funzione: =PERCENTILE (B2:B11,0.85)

Indicatori di variabilità (dispersione)

Misurano la dispersione dei valori di una distribuzione

- Varianza
- Deviazione standard
- Ampiezza
- Interquartile

Varianza

La varianza si determina attraverso la funzione **VAR [VAR]**

Il risultato di questa funzione è la varianza campionaria (s^2) dei valori introdotti come argomento

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$$

A	B	C	D	E
1	Studenti	Altezza (cm)		
2	marco	171		
3	antonella	169		
4	luca	186		
5	marina	170		
6	gianna	165		
7	lugi	190		
8	francesco	190	varianza (cm²)	81,66666667
9	michele	172		
10	stefania	174		
11	claudia	173		

Come si fa:

=VAR(B2:B11)

Deviazione standard

La **deviazione standard** si determina attraverso la funzione **DEV.ST [STDEV]**

$$s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}}$$

A	B	C	D	E
1	Studenti	Altezza (cm)		
2	marco	171		
3	antonella	169		
4	luca	186		
5	marina	170		
6	gianna	165		
7	lugi	190		
8	francesco	190	varianza (cm²)	81,66666667
9	michele	172	devianza standard	9,03747286
10	stefania	174		
11	claudia	173		

=DEV.ST(B2:B11)

Ampiezza del campione

si ottiene come differenza tra l'estremo superiore e quello inferiore di valori osservati del campione.

= MAX(dati) - MIN(dati).

Ampiezza interquartile

si ottiene come differenza tra il terzo e il primo quartile

= quartile(dati,3) - quartile(dati,1).

Componenti Aggiuntivi di Excel

Strumenti di analisi
Fornisce interfacce e funzioni per analisi di dati scientifico e finanziari

Contestat	Valore
Media	96,2
Errore standard	$=dev.st/n^{0,5}$ 1,75
Mediana	96,5
Moda	96
Deviazione standard	3,59
Varianza campionaria	91,89
Curtosi	-0,0662
Asimmetria	0,13
Intervallo	39
Minimo	79
Massimo	110
Somma	2006
Conteggio	$=n$ 30
Livello di confidenza(95%)	3,59

Istogrammi di frequenza e indici statistici

Tabella di frequenze

Istogramma

Indici di posizione e dispersione

Foglio_istogramma.xls

Dati bivariati

Creazione guidata Grafico - Passaggio 1 di 4 - Tipo di grafico

Creazione guidata Grafico - Passaggio 3 di 4 - Layout del grafico

Creazione guidata Grafico - Passaggio 4 di 4 - Posizione grafica

Right-click

Right-click

Right-click

Right-click

Right-click

Coefficiente di correlazione lineare

Misura la correlazione tra due variabili. In Excel si usa la funzione **CORRELAZIONE** [CORREL]
 = **correlazione (dati_x, dati_y)**

Il risultato di questa funzione è il coefficiente di correlazione (r) tra i due insiemi di valori:

$$r = \frac{s_{xy}}{\sqrt{s_{xx}} \sqrt{s_{yy}}}$$

dove:

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

è n volte la covarianza fra X e Y.

	Colonna 1	Colonna 2
Colonna 1	1	
Colonna 2	0.99918	1

Box plot

E' una "scatola" in cui

- I bordi corrispondono a Q1 e Q3
- Una linea fra di essi indica il valore di Q2 (mediana)
- All'esterno vengono aggiunti:
 - Un "baffo superiore" = distanza da Q3 del più grande valore inferiore a Q3+1.5(Q3-Q1)
 - Un "baffo inferiore" = distanza da Q1 del più piccolo valore minore di Q1-1.5(Q3-Q1)
- I valori esterni all'intervallo compreso tra i due "baffi", detti "outliers", vengono rappresentati individualmente

