

Ulteriori Conoscenze di Informatica e Statistica

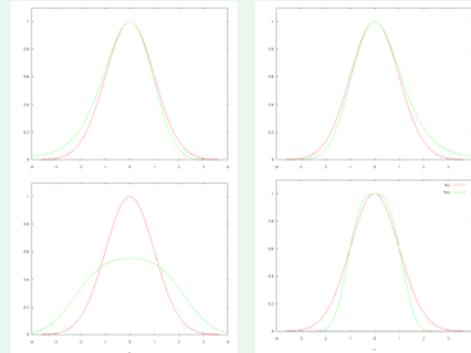
Carlo Meneghini

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

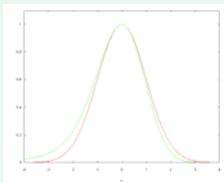
meneghini@fis.uniroma3.it

Indici di forma

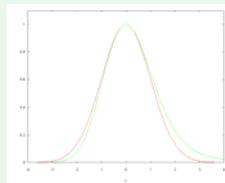
Descrivono le asimmetrie della distribuzione



Una distribuzione non simmetrica si dice obliqua



Distribuzione obliqua sx:
la media è minore della mediana



Distribuzione obliqua dx:
la media è maggiore della mediana

Momenti della distribuzione

$$m_z = \frac{1}{N} \sum_{i=1}^N (x_i)^z = \sum_{i=1}^N f_i (x_i)^z$$

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

$$\sum_{i=1}^N N_i = N$$

$$M_z = \sum_{i=1}^N f_i (x_i - m_0)^z$$

Momenti centrali
(momenti rispetto alla media)

f_i = frequenza relativa

N_i = frequenza assoluta

Media = m_0
varianza = $\sigma^2 = M_2$

$$\bar{x} = m_0 = \sum_{i=1}^N f_i x_i \quad \sigma^2 = \sum_{i=1}^N f_i (x_i - m_0)^2 = \sum_{i=1}^N f_i x_i^2 - \bar{x}^2 = \overline{x^2} - \bar{x}^2$$

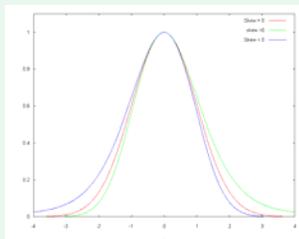
Coefficiente di
asimmetria (a_{simm})

$$m_3 = \sum_{i=1}^N f_i (x_i - m_0)^3$$

$$a_{simm} = \frac{m_3}{\sigma^3}$$

Skewness

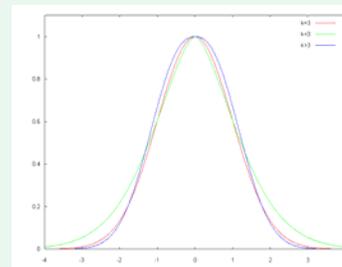
per la distribuzione
Normale (Gauss)
 $a_{simm} = 0$



Curtosi (Kurtosis) (a_4)

$$m_4 = \sum_{i=1}^N f_i (x_i - m_0)^4$$

$$a_4 = \frac{m_4}{\sigma^4}$$



per la distribuzione
Normale (Gauss)

$$a_4 = 3$$

$a_4 > 3$ indica una
distribuzione più piatta
di una Gaussiana

$a_4 < 3$ indica una
distribuzione più
piccata di una Gaussiana

Centratura e standardizzazione

Se x è una variabile aleatoria caratterizzata da valor medio μ e varianza σ^2 la variabile

$$Z = (x - \mu) / \sigma$$

È una variabile aleatoria con media nulla e varianza unitaria

Tabelle di contingenza

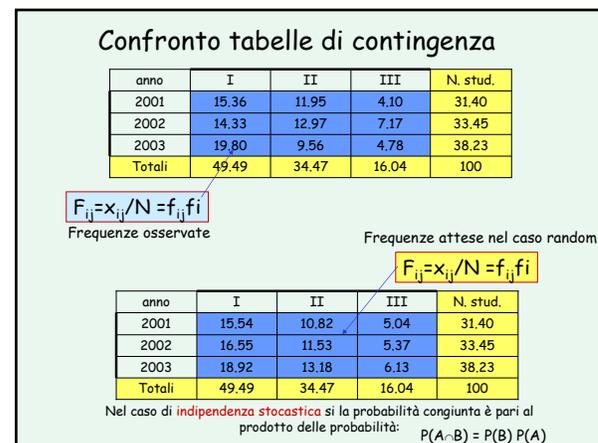
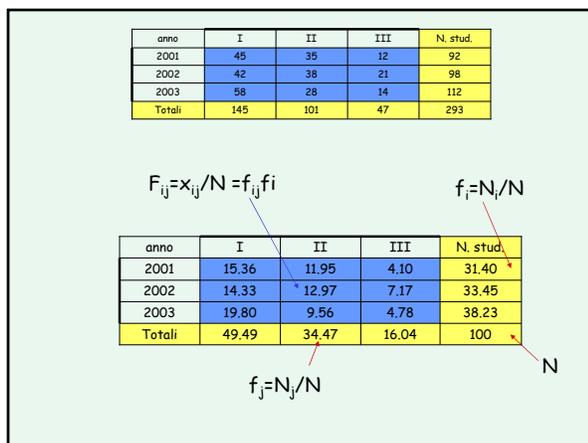
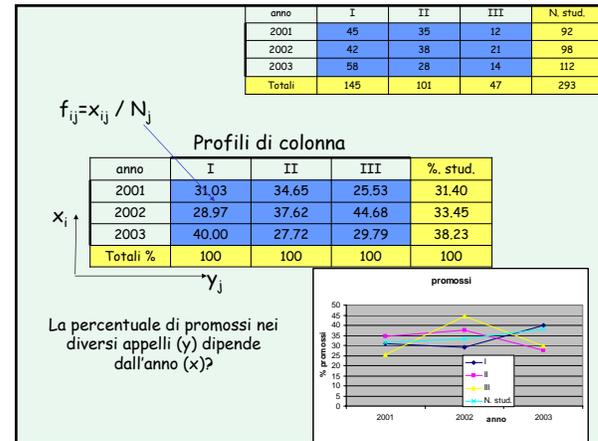
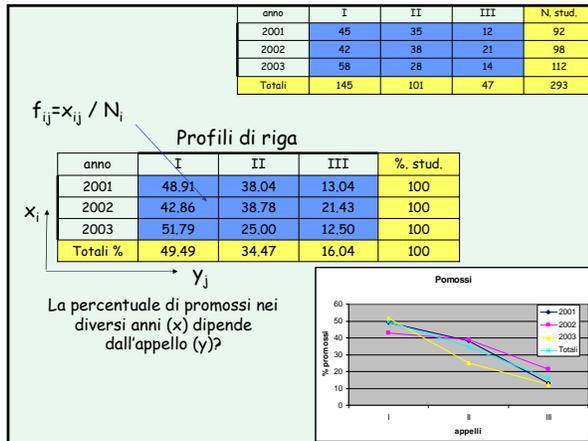
Ad ogni valore i -esimo della variabile (qualitativo o quantitativo) è associato il numero delle osservazioni

Modalità

anno	I	II	III	N. stud.
2001	45	35	12	92
2002	42	38	21	98
2003	58	28	14	112
Totale	145	101	47	293

Distribuzioni marginali

osservazioni



anno	J	K	L	N stud	Riga 1	Riga 2	Riga 3	col 1	col 2	col 3
2001	45	35	12	92	Riga 1	1		col 1	1	
2002	42	38	21	98	Riga 2	0.99	1	col 2	-0.99	1
2003	58	28	14	112	Riga 3	0.91	0.85	1	col 3	-0.47
Totale	145	101	47	293						

Strumenti
Analisi dati
Correlazione

anno	J	K	L	N stud	Riga 1	Riga 2	Riga 3	col 1	col 2	col 3
2001	31.03	34.65	25.53	31.40	Riga 1	1		col 1	1	
2002	28.97	37.62	44.65	33.45	Riga 2	0.99	1	col 2	-0.99	1
2003	40.00	27.72	29.79	38.23	Riga 3	0.91	0.85	1	col 3	-0.47
Totale	100	100	100	100						

Tabella teorica

Frequenza e Probabilità

L'insieme dei valori assunti da una variabile aleatoria come risultato di esperimenti e osservazioni costituisce una "distribuzione"

Per ogni valore la "frequenza" è il numero di volte che questo valore compare.

La frequenza relativa è il numero di volte che questo compare, normalizzato al numero di osservazioni.

Le frequenze sono **grandezze "sperimentali"**. La **probabilità** associata ad un dato risultato è il **risultato di un procedimento matematico**.

Frequenza e Probabilità

Come determinare la probabilità di un evento?

Insieme infinito di "prove"

$$p_i = \lim_{N \rightarrow \infty} \frac{n_i}{N} = \lim_{N \rightarrow \infty} f_i$$

$$\sum_{i=1}^M p_i = 1$$

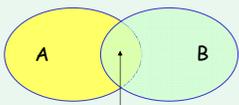
$$p_i \geq 0$$

Calcolo delle probabilità

P(A): probabilità dell'evento A

P(B): probabilità dell'evento B

$$P(A+B) = P(A \cup B) = P(A) + P(B) - P(A \cap B)$$



AB, A ∩ B,
intersezione

Probabilità composta

$$P(A \cap B) = P(B) P(A|B) = P(A) P(B|A)$$

per cui la probabilità che due eventi A e B si verifichino contemporaneamente è pari alla probabilità di uno dei due eventi moltiplicato con la probabilità dell'altro evento condizionato al verificarsi del primo.

Nel caso di **indipendenza stocastica** [$P(B|A) = P(B)$] la probabilità congiunta è pari al prodotto delle probabilità:

$$P(A \cap B) = P(B) P(A)$$

Lancio di 2 dadi



Risultati possibili

2 - 1,1	m_{ij}	$P_2 = 1/36$
3 - 1,2 2,1		$P_3 = 2/36$
4 - 1,3 2,2 3,1		$P_4 = 3/36$
5 - 1,4 2,3 3,2 4,1		$P_5 = 4/36$
6 - 1,5 2,4 3,3 4,2 5,1		$P_6 = 5/36$
7 - 1,6 2,5 3,4 4,3 5,2 6,1		$P_7 = 6/36$
8 - 2,6 3,5 4,4 5,3 6,2		$P_8 = 5/36$
9 - 3,6 4,5 5,4 6,3		$P_9 = 4/36$
10 - 4,6 5,5 6,4		$P_{10} = 3/36$
11 - 5,6 6,5		$P_{11} = 2/36$
12 - 6,6		$P_{12} = 1/36$

$p_i = \frac{1}{6}$

$p_{ij} = m_{ij} p_i p_j = \frac{m_{ij}}{36}$

Distribuzione di Bernulli

Descrive una variabile casuale che può assumere solo valori 0,1: $X \in [0,1]$

$P(1) = P(X=1) = p$ $P(0) = P(X=0) = 1-p$

$\mu = p$
 $\sigma^2 = p(1-p)$

Distribuzione Binomiale

Descrive una variabile casuale X che rappresenta il numero di successi su n prove, ognuna con probabilità di successo p (X è una sommata di variabili casuali di tipo "bernulli")

$$p(k) = \mathbb{P}(X = k) = \binom{n}{k} p^k (1-p)^{n-k}$$

k = numero di successi, n = numero di prove

$\binom{n}{k} = \frac{n!}{k!(n-k)!}$

$\mu = np$ $\sigma = \sqrt{npq}$

Distribuzione Binomiale (Bernulli)

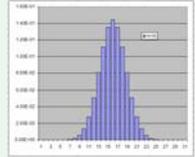
Pb.: Lancio 5 ($N=5$) volte una moneta, quale è la probabilità di avere 3 ($k=3$) teste

$$P(N, k) = \binom{N}{k} p^k q^{N-k} = \frac{N!}{k!(N-k)!} p^k (1-p)^{N-k}$$

1 ●●●○○
 2 ○●●●○
 3 ○●●●○
 4 ●●●○○
 5 ●●○○○
 6 ○●○○○
 7 ●○○○○
 8 ○○○○○
 9 ●○○○○
 10 ○○○○○

Valore atteso $\mu = np$

Dev. St. $\sigma = \sqrt{npq}$

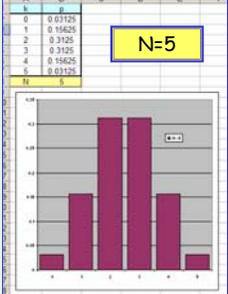


Distribuzione Binomiale (Bernulli)

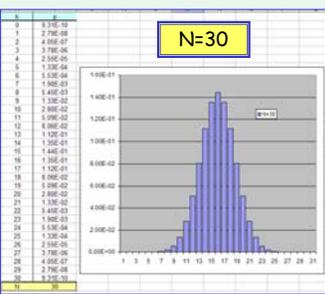
$p = q = 0.5$

$$P(N, k) = \binom{N}{k} p^k q^{N-k} =$$

N=5



N=30



Distribuzione di Poisson

La variabile stocastica X può assumere valori discreti $X \in \{0,1,2,3,4,\dots\}$

Distribuzione Binomiale con: $N \gg 1$ $p \ll 1$

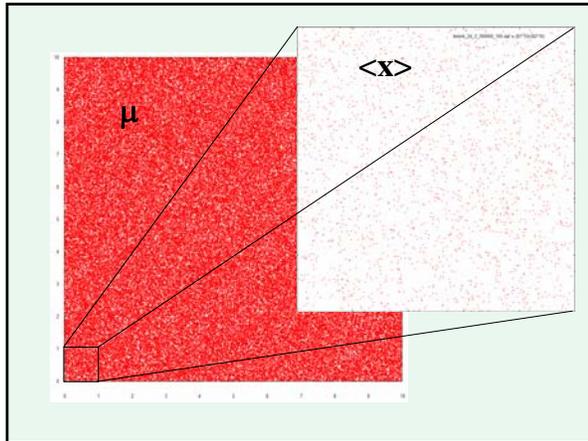
$$p(k) = \frac{\lambda^k}{k!} e^{-\lambda}$$

Valore atteso: $\mu = \lambda$
 varianza: $\sigma^2 = \lambda$

Può essere usata per descrivere il numero di cellule in una data area, il numero di errori di battitura per pagina, etc...

Descrive distribuzioni di oggetti caratterizzati da:

- densità costante (numero di oggetti proporzionale alla dimensione della regione di campionamento (superficie, volume, lunghezza etc..))
- i conteggi in regioni disgiunte sono indipendenti
- i numero di conteggi tende a zero se le dimensioni della regione tendono a zero.



Distribuzioni di probabilità notevoli

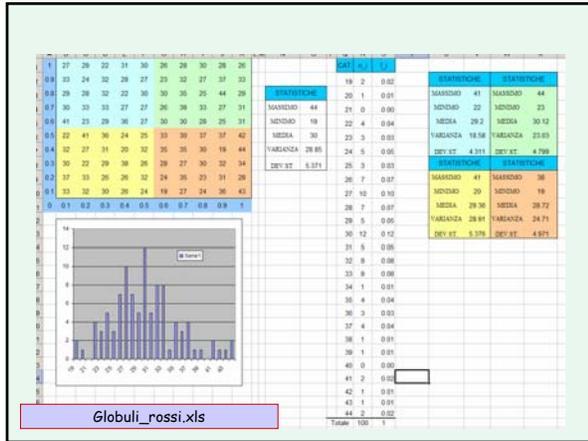
Distribuzione Poisson

Distribuzione Binomiale con: $N \gg 1$ $p \ll 1$

$$p(k) = \frac{\lambda^k}{k!} e^{-\lambda}$$

Valore atteso: $\mu = \lambda$

1	27	29	22	31	30	26	29	30	28	26
0.9	33	24	32	28	27	23	32	27	37	33
0.8	29	28	32	22	30	30	35	25	44	29
0.7	30	33	33	27	27	26	39	33	27	31
0.6	41	23	29	36	27	30	30	28	25	31
0.5	22	41	36	24	25	33	30	37	37	42
0.4	32	27	31	20	32	35	35	30	19	44
0.3	30	22	29	38	26	28	27	30	32	34
0.2	37	33	26	26	32	24	35	23	31	28
0.1	33	32	30	26	24	19	27	24	36	43
0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1



Distribuzione di Gauss o Normale

$$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\bar{x})^2}{2\sigma^2}}$$

\bar{x} Valor medio

σ Dev. standard

$$P(x_1 \leq x \leq x_2) = \int_{x_1}^{x_2} f(x) dx$$

Distribuzioni continue

probabilità

$$P(x_1 \leq x \leq x_2) = \int_{x_1}^{x_2} f(x) dx$$

Densità di probabilità

$$\bar{x} = \int_{x_1}^{x_2} x f(x) dx$$

$$\sigma^2 = \int_{x_1}^{x_2} (x - \bar{x})^2 f(x) dx$$

$$M_z = \int_{x_1}^{x_2} (x - \bar{x})^z f(x) dx$$

Distribuzione uniforme

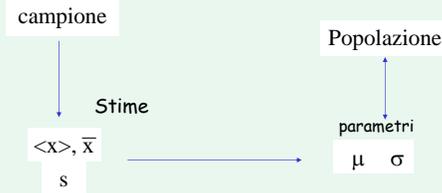
$$f(x) = \frac{1}{B - A} \quad A \leq x \leq B$$

$$\bar{x} = \frac{B + A}{2}$$

$$\sigma^2 = \frac{|B - A|^2}{12}$$

I risultati di un esperimento sono variabili aleatorie.

Un esperimento non consente di esaminare ogni elemento di una popolazione o di effettuare tutte le misure possibili.



Teorema del limite centrale

La distribuzione delle medie campionarie (\bar{x}) segue una distribuzione normale indipendentemente dalla distribuzione della popolazione d'origine

Il valor medio della distribuzione delle media campionarie è uguale alla media della popolazione d'origine

La deviazione standard dell'insieme di tutte le medie campionarie (errore standard della media $\sigma_{\bar{x}}$) è una funzione della deviazione standard della popolazione originaria e del numero di elementi del campione.

$$\frac{1}{\sqrt{2\pi\sigma_{\bar{x}}^2}} e^{-\frac{(\bar{x}_i - \mu)^2}{2\sigma_{\bar{x}}^2}}$$

nota: dev.st. della popolazione

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{m}}$$

Dato un campione n estratto da una popolazione N è possibile fornire una stima ($\langle x \rangle, s$) dei parametri reali della distribuzione (μ, σ).

I risultati ottenuti su un campione rappresentano una stima dei valori "veri"

I valori stimati sono variabili aleatorie

Quanto sono accurate queste stime?

Popolazione

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i = \sum_{k=1}^{n_c} x_k p(x_k) \quad \text{Valore atteso (media)}$$

$$\text{Varianza} \quad \sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2 = \sum_{k=1}^{n_c} p(x_k) (x_k - \mu)^2$$

Campione

$$\bar{x} = \frac{1}{m} \sum_{j=1}^m x_j = \sum_{k=1}^{n_c} x_k f(x_k) \quad \text{Media campionaria}$$

$$\text{Varianza campionaria} \quad s^2 = \frac{1}{m-1} \sum_{j=1}^m (x_j - \bar{x})^2$$

L'errore standard della media

$$\sigma_{\bar{x}}$$

indica il grado di incertezza da associare alla stima della media ottenuta utilizzando un campione dell'intera popolazione

Accuratezza delle stime

Il valor medio ottenuto da un solo campione di m elementi è una stima del valore atteso della popolazione.

L'errore standard della media rappresenta una stima dell'errore fatto nella stima del valore atteso.

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{m}} \sim \frac{s}{\sqrt{m}}$$

Risultato di un'osservazione:

$$\bar{x} \pm \sigma_{\bar{x}}$$

Accuratezza delle stime

Per migliorare la stima del valore atteso si può ripetere l'esperimento utilizzando K campioni indipendenti

In questo caso la migliore stima del valore atteso è la media delle medie campionarie:

$$\bar{x} = \frac{1}{K} \sum_{i=1}^K \bar{x}_i$$

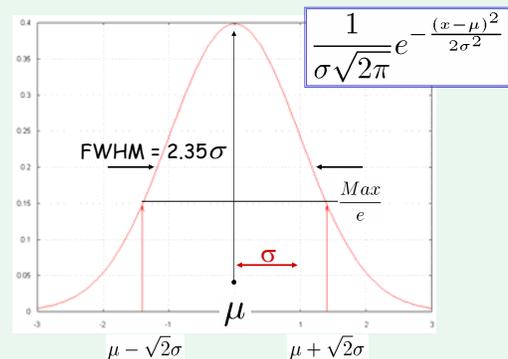
Utilizzando K campioni indipendenti l'errore standard della media è:

$$\sigma_{\bar{x}}^2 = \frac{1}{K} \sum_{i=1}^K (\bar{x}_i - \bar{x})^2$$

Risultato di un'osservazione:

$$\bar{x} \pm \sigma_{\bar{x}}$$

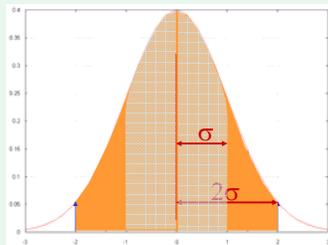
Proprietà della distribuzione di Gauss



$$\int_{\mu-\sigma}^{\mu+\sigma} g(x)dx = 0.68$$

$$\int_{\mu-2\sigma}^{\mu+2\sigma} g(x)dx = 0.95$$

$$\int_{\mu-3\sigma}^{\mu+3\sigma} g(x)dx = 0.997$$



Date due variabili aleatorie indipendenti X_a, X_b caratterizzate da $\mu_a, \sigma_a, \mu_b, \sigma_b$, la variabile $Z = X_a + X_b$ è una variabile aleatoria con:

$$\mu_z = \mu_a + \mu_b$$

$$\sigma_z = \sigma_a + \sigma_b$$

Teorema del limite centrale

La distribuzione delle medie campionarie \bar{x} su campioni di m elementi segue una distribuzione **normale** indipendentemente dalla distribuzione della popolazione d'origine

$$\frac{1}{\sqrt{2\pi}\sigma_{\bar{x}}} e^{-\frac{(\bar{x}-\mu)^2}{2\sigma_{\bar{x}}^2}}$$

La variabile: $t = \frac{\bar{x} - \mu}{\sigma_{\bar{x}}}$

è una variabile aleatoria che, per m molto grande, ha una distribuzione Normale Standard (ha media nulla e varianza unitaria):

$$g(t) = \frac{1}{\sqrt{2\pi}} e^{-t^2/2}$$

Confidenza

Pb.: un'osservazione su un campione di m elementi fornisce come risultato il valor medio \bar{x} di una variabile aleatoria.

1) costruisco una variabile aleatoria con distribuzione nota, es.:

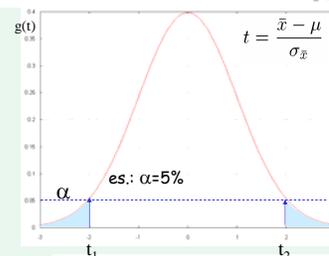
$$t = \frac{\bar{x} - \mu}{\sigma_{\bar{x}}} \quad g(t) = \frac{1}{\sqrt{2\pi}} e^{-t^2/2}$$

2) sulla base della $g(t)$ determino i valori di t che hanno una bassa probabilità di essere osservati, cioè:

- fisso un livello di confidenza α .

- determino un intervallo di valori $t_{\alpha 1} - t_{\alpha 2}$ (intervallo di confidenza) tale che la probabilità di osservare t all'esterno dell'intervallo dato sia minore di α .

$$P(t < \alpha) = P(t < t_1) + P(t < t_1) = \int_{-\infty}^{t_1} g(t)dt + \int_{t_2}^{\infty} g(t)dt = \alpha$$



Se $t_1 < t < t_2$ la probabilità di osservare il valore di t , calcolato in base ai dati, è $(1-\alpha)$

Se $t < t_1$ o $t > t_2$ la probabilità di osservare il valore di t , calcolato in base ai dati, è α

$$P(t > \alpha) = P(t_1 < t < t_2) = \int_{t_1}^{t_2} g(t)dt = 1 - \alpha$$

probabilità che t appartenga all'intervallo $t_1 - t_2$

$$t = \frac{\bar{x} - \mu}{\sigma_{\bar{x}}}$$

$$P(t_1 < t < t_1) = P(\bar{x} - \sigma_{\bar{x}}t_1 < \mu < \bar{x} + \sigma_{\bar{x}}t_1) = 1 - \alpha$$

dato il valore medio \bar{x} , osservato su un campione di m elementi, il valore aspettato della popolazione (μ) è contenuto nell'intervallo:

$$[\bar{x} - \sigma_{\bar{x}}t_1; \bar{x} + \sigma_{\bar{x}}t_1]$$

con probabilità $1-\alpha$.

Se le osservazioni sono distribuite con:

valor medio \bar{x}
 dev.st. $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{m}}$

funzione EXCEL:

CONFIDENZA(α , dev.st, m)

Intervalli di confidenza: varianza nota

Se le osservazioni sono distribuite con:

valor medio \bar{x}
 dev.st. $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{m}}$

funzione EXCEL: **CONFIDENZA(α , dev.st, m)**

La resistenza elettrica di un cavo viene misurata con uno strumento che ha un'incertezza $\sigma=0.5 \Omega$. Vengono effettuate 5 misure, ne risulta un valor medio $\bar{R}=4.52 \Omega$

CONFIDENZA(0.05, 0.5, 5)=0.438

La resistenza vera del cavo è nell'intervallo :

$R = 4.52 \pm 0.44 \Omega$ oppure: $R = [4.08, 4.96]$

Nota: α rappresenta il rischio di sbagliare, cioè la probabilità che il valore vero della resistenza sia esterno all'intervallo dato

Intervalli di confidenza: varianza campionaria

Molto più spesso non conosco la varianza della distribuzione. La migliore stima della varianza in un campione di m elementi è:

$$\sigma_{\bar{x}} = \frac{s}{\sqrt{m}} \quad \text{da cui:} \quad t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{m}}}$$

$$P(\bar{x} - t_1 \frac{s}{\sqrt{m}} < \mu < \bar{x} + t_1 \frac{s}{\sqrt{m}}) = 1 - \alpha$$

probabilità che il valore vero (μ) sia nell'intervallo:

$$[\bar{x} - t_1 \frac{s}{\sqrt{m}}; \bar{x} + t_1 \frac{s}{\sqrt{m}}]$$

funzione EXCEL: **INV.T(α ,m-1) = $t_1 \frac{s}{\sqrt{m}}$**

Nota: la variabile t così definita ha una distribuzione nota (t -Student) con $v = m-1$ gradi di libertà. La t -Student approssima una distribuzione Gaussiana per v che tende a infinito

Una misura dell'altezza di un gruppo di 20 studenti fornisce il valore medio: $H = 1.68$ m con la deviazione standard stimata $s = 9$ cm.

Determinare l'intervallo di confidenza dell'1%

Dati	
numero di osservazioni: m	20
valor medio	1.68
dev. standard: s	0.09
Confidenza	α 0.01
t_a	2.86
μ	1.68 ± 0.08

inv.T(α ; m-1)

$t_{\alpha} \frac{s}{\sqrt{m}}$

valore medio

Test di ipotesi, test statistici

- ipotesi da verificare Es.: il valore misurato è compatibile con il valore vero?
- costruisco una variabile aleatoria con distribuzione nota, es.:

$$t = \frac{\bar{x} - \mu}{\sigma_{\bar{x}}} \quad g(t) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$
- sulla base della $g(t)$ determino i valori di t che hanno una bassa probabilità di essere osservati, fissando il livello di confidenza α . Se, in base alla distribuzione scelta, il valore osservato fornisce un valore di t con bassa probabilità di essere osservato, l'ipotesi deve essere rifiutata, altrimenti può essere accettata.

$$P(t < \alpha) = P(t < t_1) + P(t < t_1) = \int_{-\infty}^{t_1} g(t)dt + \int_{t_2}^{\infty} g(t)dt = \alpha$$

$$P(t > \alpha) = P(t_1 < t < t_1) = \int_{t_1}^{t_2} g(t)dt = 1 - \alpha$$

Se $t_1 < t < t_2$ il risultato (cui è associato il valore t) è compatibile l'ipotesi fatta con una probabilità del $P = (1-\alpha)$

Se $t < t_1$ o $t > t_2$ il risultato (cui è associato il valore t) non è compatibile l'ipotesi fatta con una probabilità del $P = (1-\alpha)$

Nota: α rappresenta la probabilità di sbagliare e scartare un'ipotesi corretta.