

Popolazioni, individui e caratteri. Indicatori sintetici di campioni monovariati

La statistica descrittiva si occupa dell'analisi di dati raccolti da una popolazione, ovvero da un insieme di individui. In sintesi, dato un insieme molto grande di dati, così grande che non è *utile* guardarlo dato per dati, si cerca di estrarne delle informazioni sintetiche e tuttavia significative.

Gli oggetti con cui abbiamo a che fare sono dunque

- gli **individui** oggetto dell'indagine: ciascun individuo è un oggetto singolo dell'indagine.
- la **popolazione**, ovvero l'insieme degli individui oggetto dell'indagine.
- il **carattere** osservato o variabile, che è la quantità misurata o la qualità rilevata su ciascun individuo della popolazione.

Esempio 1.0.1. Rilevo l'altezza di ciascun abitante del Comune di Firenze. Ogni residente del Comune di Firenze è un individuo; la popolazione è l'insieme di tutti i residenti nel Comune di Firenze; il carattere in esame è l'altezza misurata, per esempio, in centimetri.

Esempio 1.0.2. Rilevo il reddito annuo di ciascun nucleo familiare del Comune di Firenze. Ogni nucleo familiare è un individuo; la popolazione è l'insieme dei nuclei familiari registrati all'Anagrafe del Comune di Firenze; il carattere osservato è il reddito annuo familiare misurato in Euro.

Esempio 1.0.3. Rilevo il numero dei componenti di ciascun nucleo familiare del Comune di Firenze. Come nell'esempio precedente ogni nucleo familiare è un individuo; la popolazione è l'insieme dei nuclei familiari registrati all'Anagrafe del Comune di Firenze. Il carattere osservato è il numero dei componenti di ciascun nucleo familiare, cioè un numero intero maggiore-uguale di 1.

Esempio 1.0.4. Per ogni studente presente in aula rilevo il colore degli occhi. Ogni studente presente in aula è un individuo. La popolazione è l'insieme degli studenti presenti ed il carattere osservato è il colore degli occhi.

In questi esempi abbiamo incontrato i due tipi fondamentali di carattere:

- **caratteri numerici o quantitativi** come l'altezza, il reddito familiare, il numero dei componenti del nucleo familiare;
- **caratteri qualitativi** come il colore degli occhi.

I caratteri numerici a loro volta si possono suddividere in

- **caratteri numerici discreti** che possono assumere solo un insieme discreto di valori, come il numero dei componenti dei nuclei familiari;
- **caratteri numerici continui** che variano con continuità ovvero con una estrema accuratezza, eccessiva rispetto ai fini dell'indagine, come l'altezza delle persone o il reddito annuo familiare.

Campione statistico, modalità e classi modali

Supponiamo di aver osservato un certo carattere su una popolazione di n individui. Abbiamo un *vettore delle osservazioni*

$$x = (x_1, \dots, x_n)$$

che chiamiamo **campione statistico** di cardinalità n .

Se il campione è relativo ad un carattere qualitativo o numerico discreto, chiamo **modalità** i valori che esso assume su un campione.

Se il campione è relativo ad un carattere numerico continuo si procede nel seguente modo: la popolazione in esame è comunque un insieme finito, quindi il carattere, per quanto continuo, nel campione assume solo un numero finito di valori. Sia $[a, b)$ un intervallo che contiene tutti i valori x_i , $i = 1, \dots, n$ assunti dal carattere sugli individui della popolazione. Suddividiamo l'intervallo $[a, b)$ in N parti uguali (N sarà suggerito dall'esperienza). Otteniamo N intervalli

$$I_j := \left[a + (j-1) \frac{b-a}{N}, a + j \frac{b-a}{N} \right), \quad j = 1, \dots, N.$$

Chiamo ciascuno di questi intervalli **classe di modalità**, se esso contiene almeno una osservazione.

Frequenza assoluta e frequenza relativa

Consideriamo un campione $x = (x_1, \dots, x_n)$ relativo ad un carattere qualitativo o numerico discreto. Nel campione, cioè nella popolazione in esame, il carattere osservato assume un certo numero di valori distinti

$$z_1, \dots, z_k, \quad 1 \leq k \leq n.$$

Per ogni $j = 1, \dots, k$ chiamo **effettivo** o **frequenza assoluta** della modalità z_j il numero

$$n_j := \# \{i \in \{1, \dots, n\} : x_i = z_j\}$$

mentre chiamo **frequenza relativa** della modalità z_j il numero

$$p_j := \frac{n_j}{n}.$$

Se il carattere osservato è numerico continuo, si considera ciascuna classe di modalità individuata

$$I_j := \left[a + (j-1) \frac{b-a}{N}, a + j \frac{b-a}{N} \right), \quad j = 1, \dots, N$$

e si chiama **frequenza assoluta o effettivo** della classe di modalità I_j il numero

$$n_j := \# \{i \in \{1, \dots, n\} : x_i \in I_j\}.$$

Come prima definiamo **frequenza relativa** della classe I_j il numero $p_j := \frac{n_j}{n}$.

Moda e valori modali

Sia $x = (x_1, \dots, x_n)$ un campione statistico e siano z_1, z_2, \dots, z_k le modalità assunte (o I_1, \dots, I_k le classi di modalità assunte) e siano p_1, \dots, p_k le relative frequenze relative.

Se esiste uno ed un solo indice $\bar{j} \in \{1, 2, \dots, k\}$ tale che la modalità $z_{\bar{j}}$ (o la classe $I_{\bar{j}}$) ha frequenza massima, ovvero se esiste un unico $\bar{j} \in \{1, 2, \dots, k\}$ tale che $p_{\bar{j}} \geq p_j \forall j = 1, \dots, k$, allora la modalità $z_{\bar{j}}$ (o la classe $I_{\bar{j}}$) si dice **moda** del campione x .

Se esistono due o più indici $\bar{j}_1, \bar{j}_2, \dots, \bar{j}_s$ tali che le modalità $z_{\bar{j}_1}, z_{\bar{j}_2}, \dots, z_{\bar{j}_s}$ (o le classi $I_{\bar{j}_1}, I_{\bar{j}_2}, \dots, I_{\bar{j}_s}$) hanno frequenza massima, allora queste modalità (o classi) si dicono **valori (o classi) modali**.

Possiamo visualizzare con degli istogrammi, vedi Figura 1.3

Mediana

D'ora innanzi consideriamo solo caratteri numerici.

Sia dunque $x = (x_1, \dots, x_n)$ un campione relativo ad un carattere numerico. Ordiniamo i dati del campione in ordine crescente:

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n-1)} \leq x_{(n)}$$

e distinguiamo due casi:

- n dispari: $n = 2m + 1$

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(m)} \leq x_{(m+1)} \leq x_{(m+2)} \leq \dots \leq x_{(2m)} \leq x_{(2m+1)}$$

Il dato $x_{(m+1)}$ è maggiore-uguale di m dati e minore-uguale di altrettanti dati. Diciamo che il dato $x_{(m+1)}$ è la **mediana** del campione.

- n pari: $n = 2m$

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(m-1)} \leq x_{(m)} \leq x_{(m+1)} \leq x_{(m+2)} \leq \dots \leq x_{(2m)}$$

Il dato $x_{(m)}$ è maggiore-uguale di $m - 1$ dati e minore-uguale di m dati. Il dato $x_{(m+1)}$ è maggiore-uguale di m dati e minore-uguale di $m - 1$ dati.

Chiamiamo **mediana** del campione il numero $\frac{x_{(m)} + x_{(m+1)}}{2}$.

Media e varianza campionaria. Scarto quadratico medio (o deviazione standard)

Consideriamo un campione relativo ad un carattere numerico

$$x = (x_1, \dots, x_n).$$

Chiamo **media aritmetica** o, più semplicemente, **media** il numero

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

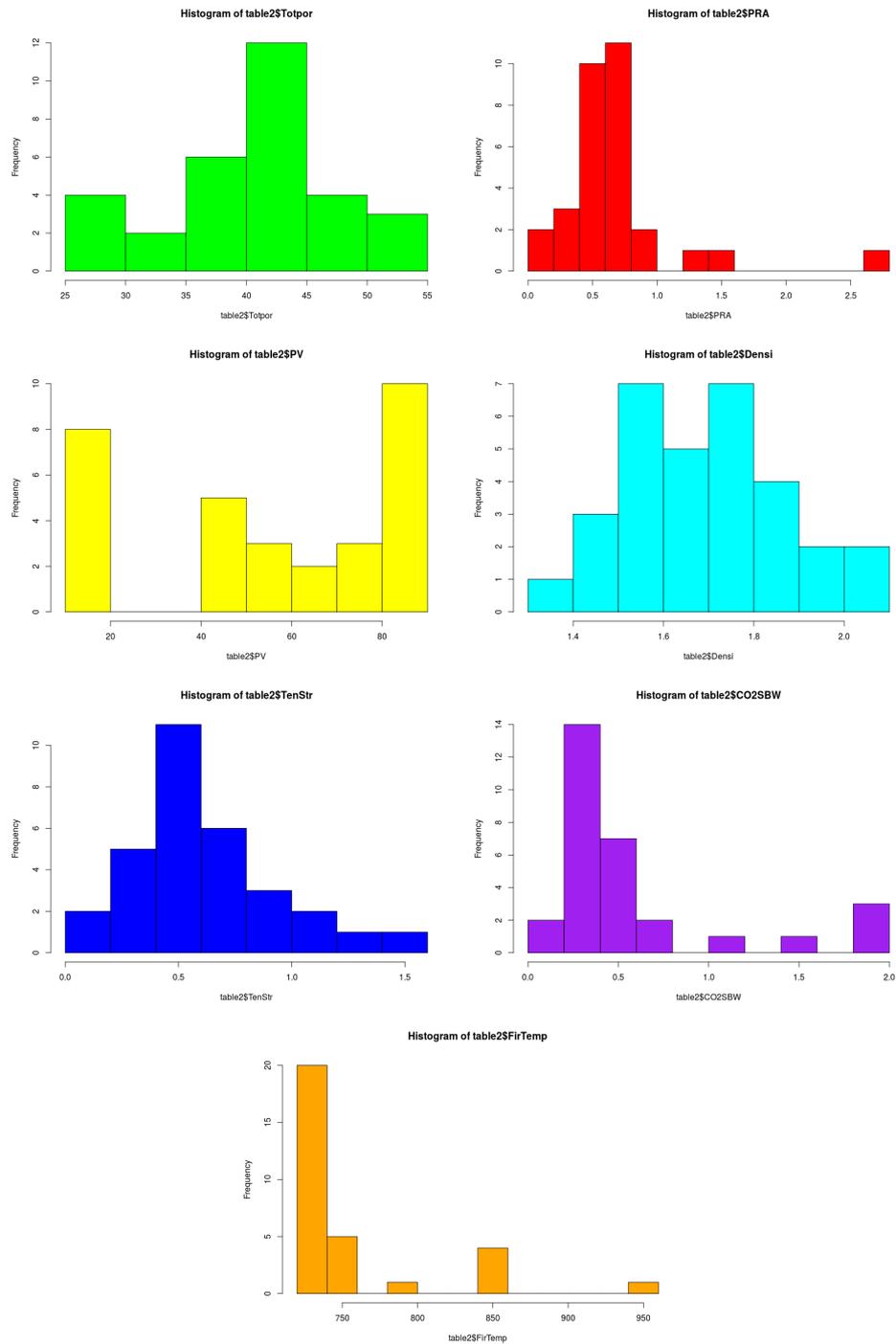


Figura 1.1: Alcuni istogrammi dall'Esempio 1.5.1

Supponiamo che nel campione siano presenti k modalità z_1, z_2, \dots, z_k con rispettive frequenze assolute N_1, N_2, \dots, N_k e frequenze relative p_1, p_2, \dots, p_k . Allora

$$\begin{aligned} \bar{x} &= \frac{1}{n} (x_1 + x_2 + \dots + x_n) = \frac{1}{n} (N_1 z_1 + N_2 z_2 + \dots + N_k z_k) = \\ &= p_1 z_1 + p_2 z_2 + \dots + p_k z_k = \sum_{j=1}^k p_j z_j. \end{aligned}$$

Chiamo **varianza campionaria** di x il numero non-negativo

$$\sigma_x^2 = \text{Var}[x] := \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2.$$

Osserviamo che la media è un valore centrale attorno al quale si dispongono i dati x_1, \dots, x_n mentre la varianza è un *indice di dispersione*: la varianza è nulla se e solo se tutti i dati del campione sono uguali (e dunque coincidono con la media). Una varianza bassa indica che comunque i dati sono *vicini* al valore medio \bar{x} mentre una varianza alta indica una maggiore dispersione dei dati.

La radice quadrata della varianza campionaria

$$\sigma_x = \text{Std}[x] := \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

si chiama **scarto quadratico medio** o **deviazione standard** del campione x .

Anche per la varianza campionaria possiamo scrivere una formula che coinvolga solo le modalità e le rispettive frequenze.

$$\begin{aligned} \sigma_x^2 &= \frac{1}{n-1} ((x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2) = \\ &= \frac{1}{n-1} (N_1(z_1 - \bar{x})^2 + N_2(z_2 - \bar{x})^2 + \dots + N_k(z_k - \bar{x})^2) = \\ &= \frac{n}{n-1} (p_1(z_1 - \bar{x})^2 + p_2(z_2 - \bar{x})^2 + \dots + p_k(z_k - \bar{x})^2) = \\ &= \frac{n}{n-1} \sum_{j=1}^k p_j(z_j - \bar{x})^2. \end{aligned}$$

Esempio 1.5.1. Nella tabella che segue, tratta da [2], riportiamo alcuni dati relativi a campioni di laterizio e che useremo per fare alcuni esempi relativi alle nozioni introdotte mediante il software R <http://cran.r-project.org/>. Per una introduzione si rimanda ai manuali [3] e [1].

SAMPLE CODE	POROSITÀ TOTALE (%)	RAGGIO MEDIO DEL PORO (μm)	VOLUME DEI PORI SU DIMEN- SIONE DEI PORI 0.3–0.8 μm	DENSITÀ (g/cm^3)	RESISTENZA ALLA TRA- ZIONE (MPa)	CO ₂ /SBW	TEMPERATURA DI COTTURA (DTA)
AS1	41.460	0.528	80.0	1.550	0.403	0.38	740
AS2	47.210	0.467	81.2	1.650	0.645	0.70	740
AS3	43.670	0.697	78.5	1.710	0.527	0.46	740
AS4	52.390	0.422	77.3	1.520	0.143	0.48	740
AS5	44.700	0.411	87.4	1.500	0.593	0.29	740
AS6	51.330	0.422	88.6	1.480	0.463	0.33	740
AS7	31.460	0.718	80.6	1.900	0.955	0.23	740
AS8	40.900	0.458	80.4	1.680	0.195	0.41	740
AS9	45.540	0.492	80.8	1.620	1.328	0.50	750
AS10	45.620	0.734	86.2	1.620	1.405	0.34	750
AS11	44.140	0.730	85.7	1.590	0.256	0.42	750
AS12	40.710	0.543	87.8	1.750	0.309	0.20	750
AS13	35.700	0.686	84.3	1.520	0.472	0.05	740
C1	40.290	0.306	43.5	1.760	0.520	0.43	740
C2	36.570	0.625	42.3	1.750	0.738	0.36	740
C3	42.130	0.249	63.2	1.630	0.410	0.25	740
C4	37.830	0.731	47.9	2.020	0.601	0.28	740
C5	42.180	0.407	59.4	1.580	0.376	0.34	740
C6	41.600	0.446	42.8	1.850	0.473	0.26	740
C7	32.660	0.664	64.3	1.850	0.695	0.25	740
C8	36.070	0.673	58.2	1.780	0.624	0.29	740
C9	36.040	1.397	55.6	1.730	0.582	0.38	740
C10	36.640	0.861	45.2	1.750	0.650	0.47	740
R1	42.890	0.785	10.2	1.540	0.453	1.04	850
R2	26.850	0.315	14.7	2.010	1.124	1.86	960
R3	28.550	0.158	18.6	1.920	0.937	1.96	850
R4	29.860	0.158	15.3	1.890	1.020	1.48	850
R5	45.700	0.984	12.8	1.500	0.328	–	800
R6	54.640	1.525	12.5	1.340	0.267	0.67	750
R7	27.550	2.657	14.6	1.920	0.892	0.40	730
R8	40.820	0.622	15.3	1.570	0.502	1.94	860

Inseriamo la tabella in R

```
> library(readr)
> table2 <- read_delim("~/Documents/didattica/2017-18_analisi_reale/alcuni_appunti/table2.
+ "\t", escape_double = FALSE, trim_ws = TRUE)
Parsed with column specification:
cols(
  Code = col_character(),
  Totpor = col_double(),
  PRA = col_double(),
  PV = col_double(),
  Densi = col_double(),
  TenStr = col_double(),
  CO2SBW = col_character(),
  FirTemp = col_integer()
```

```
)
> View(table2)
```

	Code	Totpor	PRA	PV	Densi	TenStr	CO2SBW	FirTemp
1	AS1	41.46	0.528	80.0	1.55	0.403	0.38	740
2	AS2	47.21	0.467	81.2	1.65	0.645	0.70	740
3	AS3	43.67	0.697	78.5	1.71	0.527	0.46	740
4	AS4	52.39	0.422	77.3	1.52	0.143	0.48	740
5	AS5	44.70	0.411	87.4	1.50	0.593	0.29	740
6	AS6	51.33	0.422	88.6	1.48	0.463	0.33	740
7	AS7	31.46	0.718	80.6	1.90	0.955	0.23	740
8	AS8	40.90	0.458	80.4	1.68	0.195	0.41	740
9	AS9	45.54	0.492	80.8	1.62	1.328	0.50	750
10	AS10	45.62	0.734	86.2	1.62	1.405	0.34	750
11	AS11	44.14	0.730	85.7	1.59	0.256	0.42	750
12	AS12	40.71	0.543	87.8	1.75	0.309	0.20	750
13	AS13	35.70	0.686	84.3	1.52	0.472	0.05	740
14	C1	40.29	0.306	43.5	1.76	0.520	0.43	740
15	C2	36.57	0.625	42.3	1.75	0.738	0.36	740
16	C3	42.13	0.249	63.2	1.63	0.410	0.25	740
17	C4	37.83	0.731	47.9	2.02	0.601	0.28	740
18	C5	42.18	0.407	59.4	1.58	0.376	0.34	740
19	C6	41.60	0.446	42.8	1.85	0.473	0.26	740
20	C7	32.66	0.664	64.3	1.85	0.695	0.25	740
21	C8	36.07	0.673	58.2	1.78	0.624	0.29	740
22	C9	36.04	1.397	55.6	1.73	0.582	0.38	740
23	C10	36.64	0.861	45.2	1.75	0.650	0.47	740
24	R1	42.89	0.785	10.2	1.54	0.453	1.04	850
25	R2	26.85	0.315	14.7	2.01	1.124	1.86	960
26	R3	28.55	0.158	18.6	1.92	0.937	1.96	850
27	R4	29.86	0.158	15.3	1.89	1.020	1.48	850
28	R5	45.70	0.984	12.8	1.50	0.328	--	800
29	R6	54.64	1.525	12.5	1.34	0.267	0.67	750
30	R7	27.55	2.657	14.6	1.92	0.892	0.40	730
31	R8	40.82	0.622	15.3	1.57	0.502	1.94	860

Per ciascun carattere definiamo una variabile che contenga la mediana, una per la media, una per la Varianza e una per la deviazione standard e poi stampiamo i valori (tratteremo il carattere di nome CO2SBW con attenzione perché su un individuo non è stato rilevato)

Il comando `summary` indica il numero di dati mancanti, ci dà gli indicatori di centralità ma non quelli di dispersione

```
> summary(table2)
      Code      Totpor      PRA      PV      Densi      TenStr      CO2SBW      FirTemp
Length:31  Min. :26.85  Min. :0.1580  Min. :10.20  Min. :1.340  Min. :0.1430  Min. :0.0500  Min. :730.0
Class :character  1st Qu.:36.05  1st Qu.:0.4220  1st Qu.:30.45  1st Qu.:1.560  1st Qu.:0.4065  1st Qu.:0.2900  1st Qu.:740.0
Mode  :character  Median :40.90  Median :0.6220  Median :59.40  Median :1.680  Median :0.5270  Median :0.3900  Median :740.0
      Mean :40.12  Mean :0.6733  Mean :55.33  Mean :1.693  Mean :0.6092  Mean :0.5817  Mean :764.8
      3rd Qu.:44.42  3rd Qu.:0.7305  3rd Qu.:80.70  3rd Qu.:1.815  3rd Qu.:0.7165  3rd Qu.:0.4950  3rd Qu.:750.0
      Max. :54.64  Max. :2.6570  Max. :88.60  Max. :2.020  Max. :1.4050  Max. :1.9600  Max. :960.0
      NA's :1
```

Richiediamo anche varianza campionaria e deviazione standard.

```

> medianaTotPor <- median(table2$Totpor);
> meanTotPor <- mean(table2$Totpor);
> VarTotPor <- var(table2$Totpor);
> StdTotPor <- sd(table2$Totpor)
> medianaTotPor; meanTotPor; VarTotPor; StdTotPor
[1] 40.9
[1] 40.11935
[1] 49.52185
[1] 7.037176
> medianaPRA <- median(table2$PRA);
> meanPRA <- mean(table2$PRA);
VarPRA <- var(table2$PRA);
> StdPRA <- sd(table2$PRA)
> medianaPRA; meanPRA; VarPRA; StdPRA
[1] 0.622
[1] 0.6732581
[1] 0.226613
[1] 0.4760389
> medianaPV <- median(table2$PV);
> meanPV <- mean(table2$PV);
> VarPV <- var(table2$PV);
> StdPV <- sd(table2$PV)
> medianaPV; meanPV; VarPV; StdPV
[1] 59.4
[1] 55.32903
[1] 815.0935
[1] 28.54984
> medianaDensi <- median(table2$Densi);
> meanDensi <- mean(table2$Densi);
> VarDensi <- var(table2$Densi);
> StdDensi <- sd(table2$Densi)
> medianaDensi; meanDensi; VarDensi; StdDensi
[1] 1.68
[1] 1.692903
[1] 0.02894129
[1] 0.1701214
> medianaTenStr <- median(table2$TenStr);
> meanTenStr <- mean(table2$TenStr);
> VarTenStr <- var(table2$TenStr);
> StdTenStr <- sd(table2$TenStr)
> medianaTenStr; meanTenStr; VarTenStr; StdTenStr
[1] 0.527
[1] 0.6092258
[1] 0.09882738
[1] 0.3143682

```

```
> medianaCO2SBW <- median(na.omit(table2$CO2SBW));
> meanCO2SBW <- mean(na.omit(table2$CO2SBW));
> VarCO2SBW <- var(na.omit(table2$CO2SBW));
> StdCO2SBW <- sd(na.omit(table2$CO2SBW))
> medianaCO2SBW; meanCO2SBW; VarCO2SBW; StdCO2SBW
[1] 0.39
[1] 0.5816667
[1] 0.2765868
[1] 0.5259152
> medianaFirTemp <- median(table2$FirTemp);
> meanFirTemp <- mean(table2$FirTemp);
> VarFirTemp <- var(table2$FirTemp);
> StdFirTemp <- sd(table2$FirTemp)
> medianaFirTemp; meanFirTemp; VarFirTemp; StdFirTemp
[1] 740
[1] 764.8387
[1] 2805.806
[1] 52.96986
```


Campioni bivariati: covarianza, coefficiente di correlazione e retta di regressione

Covarianza e coefficiente di correlazione

Supponiamo di avere un **campione bivariato** cioè di rilevare due caratteri sugli individui di una medesima popolazione.

Abbiamo dunque due vettori di dati

$$x = (x_1, \dots, x_n), \quad y = (y_1, \dots, y_n).$$

x_i e y_i sono le rilevazioni dei due caratteri sul medesimo individuo, l'individuo cioè che abbiamo etichettato come *individuo i*.

Chiamiamo **covarianza di x e y** il numero

$$\text{Cov}(x, y) := \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

dove \bar{x} e \bar{y} sono le medie dei campioni x e y , rispettivamente.

Nel caso in cui né x né y siano campioni costanti (ipotesi lavorativa che sarà sempre sottintesa), definiamo **coefficiente di correlazione di x e y** il numero

$$\rho[x, y] := \frac{\text{Cov}(x, y)}{\text{Std}[x] \text{Std}[y]} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\left(\sum_{i=1}^n (x_i - \bar{x})^2\right) \left(\sum_{i=1}^n (y_i - \bar{y})^2\right)}}.$$

Osservazione 2.1.1. $\text{Cov}(x, x) = \text{Var}[x]$; $\rho[x, x] = 1$.

Osservando che $\rho[x, y]$ non è altro che il rapporto tra $\langle x - (\bar{x}, \dots, \bar{x}), y - (\bar{y}, \dots, \bar{y}) \rangle$ (prodotto scalare) e $\|x - (\bar{x}, \dots, \bar{x})\| \|y - (\bar{y}, \dots, \bar{y})\|$ (prodotto delle norme) si dimostrano facilmente le seguenti proprietà:

1. $-1 \leq \rho[x, y] \leq 1$;
2. $\rho[x, y] = 1$ se e solo se esiste $a > 0, b \in \mathbb{R}$ tale che $y_i = ax_i + b \quad \forall i = 1, \dots, n$. In tal caso i campioni x e y si dicono *positivamente correlati*;
3. $\rho[x, y] = -1$ se e solo se esiste $a < 0, b \in \mathbb{R}$ tale che $y_i = ax_i + b \quad \forall i = 1, \dots, n$. In tal caso i campioni x e y si dicono *negativamente correlati*.

Se $\rho[x, y] = 0$ i campioni x e y si dicono *scorrelati*.

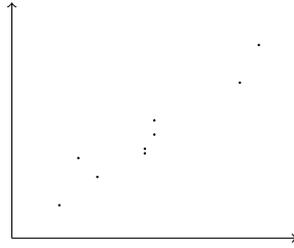


Figura 2.1: Campione bivariato *pressoché lineare*

Retta di regressione

Supponiamo di avere un campione bivariato

$$x = (x_1, \dots, x_n), \quad y = (y_1, \dots, y_n)$$

dove x_i e y_i sono i dati relativi all' i -esimo individuo. Rappresentiamo i punti (x_i, y_i) sul piano cartesiano Oxy . Capita, molto spesso, di trovarsi a disposizioni *pressoché allineate* come illustrato nella figura 2.1 Si cerca allora una retta che in qualche senso *approssimi* i punti (x_i, y_i) .

Supponiamo che $y = ax + b$ sia l'equazione della retta cercata. Per $x = x_i$ si ottiene il punto sulla retta $(x_i, ax_i + b)$. Cerchiamo la retta (ovvero i parametri a e b) che minimizza la *somma degli errori quadratici nella direzione y*

$$S(a, b) := \sum_{i=1}^n (y_i - (ax_i + b))^2.$$

Si ha

$$\begin{aligned} S(a, b) &= \sum_{i=1}^n (y_i - \bar{y} + \bar{y} - (ax_i - a\bar{x} + a\bar{x} + b))^2 = \\ &= \sum_{i=1}^n ((y_i - \bar{y}) - a(x_i - \bar{x}) + (\bar{y} - a\bar{x} - b))^2 = \\ &= \sum_{i=1}^n ((y_i - \bar{y})^2 + a^2 \sum_{i=1}^n (x_i - \bar{x})^2 + \\ &\quad + n(\bar{y} - a\bar{x} - b)^2 - 2a \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})) = \\ &= (n-1) (\text{Var}[y] + a^2 \text{Var}[x] - 2a \text{Cov}(x, y)) + n(\bar{y} - a\bar{x} - b)^2. \end{aligned}$$

L'incognita b compare solo nell'ultimo addendo, che è un quadrato. Quindi per ottenere il minimo basterà scegliere a che minimizza la funzione $f(a) := \text{Var}[y] + a^2 \text{Var}[x] - 2a \text{Cov}(x, y)$ e poi scegliere $b = \bar{y} - a\bar{x}$. Si ha

$$\begin{aligned} f'(a) &= 2a \text{Var}[x] - 2 \text{Cov}(x, y) = 0 \quad \text{se e solo se} \quad a = \frac{\text{Cov}(x, y)}{\text{Var}[x]} \\ f''(a) &= 2 \text{Var}[x] > 0 \end{aligned}$$

Il minimo della somma degli errori quadratici $S(a, b)$ si ottiene allora per

$$a = \frac{\text{Cov}(x, y)}{\text{Var}[x]}; \quad b = \bar{y} - \frac{\text{Cov}(x, y)}{\text{Var}[x]}\bar{x};$$

il minimo dell'errore S vale

$$(n - 1) \left(\text{Var}[y] - \frac{(\text{Cov}(x, y))^2}{\text{Var}[x]} \right) = (n - 1) \text{Var}[y] \left(1 - (\rho[x, y])^2 \right)$$

e la retta ha equazione

$$y = \bar{y} + \frac{\text{Cov}(x, y)}{\text{Var}[x]}(x - \bar{x}).$$

Osservazione 2.2.1. La retta così determinata si chiama **retta di regressione del campione y sul campione x** . Osserviamo infine che il punto (\bar{x}, \bar{y}) appartiene alla retta.

Esempio 2.2.1. Riconsideriamo l'esempio 1.5.1. Carichiamo in R la tabella dei dati.

```
> library(readr)
> table2 <- read_delim("~/Documents/didattica/2017-18_analisi_reale/alcuni_appunti/table2.csv",
+   "\t", escape_double = FALSE, trim_ws = TRUE)
Parsed with column specification:
cols(
  Code = col_character(),
  Totpor = col_double(),
  PRA = col_double(),
  PV = col_double(),
  Densi = col_double(),
  TenStr = col_double(),
  CO2SBW = col_character(),
  FirTemp = col_integer()
)
```

Tracciamo sul piano cartesiano i dati relativi ai caratteri porosità totale (in ascissa) e densità (in ordinata) e salviamo la figura in un file.

```
> library(car)
> scatterplot(Densi~Totpor, lm=TRUE, smooth=FALSE, spread=FALSE, boxplots=TRUE, span=0.5, data= table2)
```

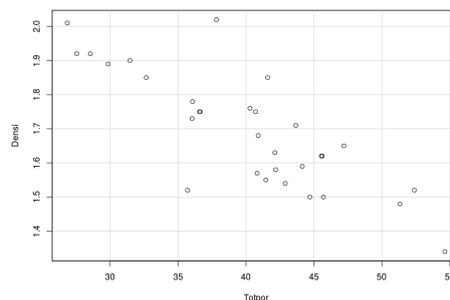


Figura 2.2: Porosità totale versus Densità

Sembrano *ragionevolmente allineati*. Calcoliamo il loro coefficiente di correlazione

```
> CorTotporDensi<- cor(table2$Totpor, table2$Densi)
> CorTotporDensi
[1] -0.8187597
```

Calcoliamo la retta di regressione del carattere Densità sul carattere Porosità Totale

```
> RegModel.Densi.Totpor <- lm(Densi~Totpor, data=table2)
> summary(RegModel.Densi.Totpor)
```

Call:

```
lm(formula = Densi ~ Totpor, data = table2)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-0.260377	-0.054570	-0.001898	0.045213	0.281783

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	2.486995	0.104930	23.70	< 2e-16 ***
Totpor	-0.019793	0.002577	-7.68	1.81e-08 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.09934 on 29 degrees of freedom

Multiple R-squared: 0.6704, Adjusted R-squared: 0.659

F-statistic: 58.98 on 1 and 29 DF, p-value: 1.814e-08

Intercept dice che l'ordinata all'origine (il coefficiente b) della retta di regressione è 2.486995 mentre il coefficiente angolare (cioè a) è -0.019793 . Ridisegniamo i punti sul piano cartesiano, aggiungendo la retta di regressione (e salviamo l'immagine in un file).

```
> abline(lm(Densi ~ Totpor, data=table2), col="red")
```

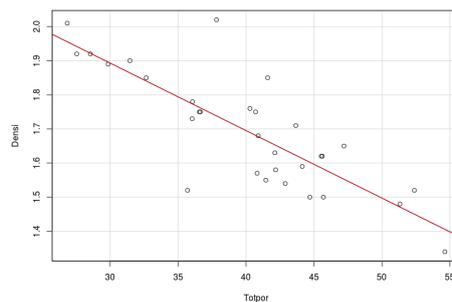


Figura 2.3: Retta di regressione lineare

Parte II

Statistica inferenziale

Campioni statistici

Introduzione

Scopo della statistica inferenziale è lo stabilire metodi rigorosi per ottenere – con un calcolabile *grado di certezza* proprietà generali di una popolazione a partire da una raccolta di dati sulla popolazione stessa.

Possiamo sintetizzare il modello matematico che applichiamo come segue

- Se rileviamo un carattere su una popolazione di n individui, consideriamo ciascun dato rilevato come il valore assunto da X_1, X_2, \dots, X_n variabili aleatorie aventi tutte la stessa distribuzione μ e che (molto spesso) si possono supporre indipendenti.
- La distribuzione μ è (parzialmente) incognita; si cercano informazioni su μ a partire dai dati rilevati. Le informazioni ricavate sulla distribuzione μ sono di natura probabilistica. Per esempio, non riusciremo ad ottenere informazioni del tipo *la media della distribuzione μ è 50* ma informazioni del tipo *la media della distribuzione μ è compresa tra 49.8 e 50.2 con probabilità del 90%*.

Comunemente si suppone di conoscere il *tipo* della distribuzione μ , ovvero si suppone di sapere se è gaussiana, esponenziale o binomiale o altro, ma di non conoscere i parametri che la caratterizzano.

Definizione 3.1.1 (Campione statistico). Una famiglia di variabili aleatorie

$$X_1, \dots, X_n$$

si dice un *campione statistico di numerosità n* se le v.a. X_1, \dots, X_n sono indipendenti ed identicamente distribuite.

Se f è la comune densità delle v.a. X_1, \dots, X_n , allora la v.a. vettoriale $X := (X_1, \dots, X_n)$ ha densità congiunta

$$g_{X_1, \dots, X_n}(x_1, \dots, x_n) = f(x_1)f(x_2)\dots f(x_n).$$

La comune distribuzione delle X_i si dice *distribuzione campionaria di X_1, \dots, X_n* .

Osservazione 3.1.1. Poiché le v.a. X_1, \dots, X_n seguono la stessa distribuzione, esse hanno anche la stessa media e la stessa varianza (se queste quantità esistono).

Definizione 3.1.2 (Statistica). Sia X_1, \dots, X_n un campione statistico. Una funzione (non dipendente da parametri) di X_1, \dots, X_n si dice una statistica.

Osservazione 3.1.2. Chiariamo cosa si intende per statistica: $3X_1 - 2X_2$ è una statistica; $\max\{X_1, \dots, X_n\}$ è una statistica. $X_1 - \mu$ $\mu \in \mathbb{R}$ non è una statistica.

Media campionaria e varianza campionaria

Definizione 3.2.1. Sia X_1, \dots, X_n un campione statistico. Chiamiamo **media campionaria** di X_1, \dots, X_n la statistica

$$\bar{X} := \frac{1}{n} \sum_{i=1}^n X_i,$$

chiamiamo **varianza campionaria** di X_1, \dots, X_n la statistica

$$S^2 := \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

Proposizione 3.2.1. Sia X_1, \dots, X_n un campione statistico di numerosità n con media μ e varianza σ^2 finite. Siano \bar{X} e S^2 la media campionaria e la varianza campionaria. Allora

$$\mathbb{E}[\bar{X}] = \mu, \quad \text{Var}[\bar{X}] = \frac{\sigma^2}{n}, \quad \mathbb{E}[S^2] = \sigma^2.$$

Dimostrazione.

$$\begin{aligned} \mathbb{E}[\bar{X}] &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] = \frac{1}{n} n\mu = \mu \\ \text{Var}[\bar{X}] &= \text{Var}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n^2} \text{Var}\left[\sum_{i=1}^n X_i\right] = \frac{1}{n^2} \sum_{i=1}^n \text{Var}[X_i] = \frac{1}{n^2} n\sigma^2 = \frac{\sigma^2}{n}. \end{aligned}$$

Per calcolare la media di S^2 osserviamo preliminarmente che

$$\begin{aligned} S^2 &= \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - 2\bar{X} \sum_{i=1}^n X_i + \sum_{i=1}^n \bar{X}^2 \right) \\ &= \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - 2\bar{X}n\bar{X} + n\bar{X}^2 \right) = \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - n\bar{X}^2 \right). \end{aligned}$$

Dunque

$$\begin{aligned} (n-1)\mathbb{E}[S^2] &= \mathbb{E}\left[\sum_{i=1}^n X_i^2 - n\bar{X}^2\right] = \mathbb{E}\left[\sum_{i=1}^n (X_i - \mu + \mu)^2 - n(\bar{X} - \mu + \mu)^2\right] \\ &= \sum_{i=1}^n \mathbb{E}\left[(X_i - \mu + \mu)^2\right] - n\mathbb{E}\left[(\bar{X} - \mu + \mu)^2\right] \\ &= \sum_{i=1}^n \mathbb{E}\left[(X_i - \mu)^2 + \mu^2 + 2\mu(X_i - \mu)\right] \\ &\quad - n\left(\mathbb{E}\left[(\bar{X} - \mu)^2\right] + \mu^2 - 2\mu\mathbb{E}[\bar{X} - \mu]\right) \\ &= \sum_{i=1}^n (\sigma^2 + \mu^2) - n\left(\frac{\sigma^2}{n} + \mu^2\right) = (n-1)\sigma^2 \end{aligned}$$

e quindi $\mathbb{E}[S^2] = \sigma^2$. □

La disuguaglianza di Chebyshev e la legge (debole) dei grandi numeri

Enunciamo alcuni importanti risultati asintotici che giustificano l'uso della media campionaria \bar{X} come stima della media μ del campione.

Teorema 3.2.1 (Disuguaglianza di Chebyshev). *Se X è una variabile aleatoria con media μ e varianza non superiore a σ^2 , allora*

$$\mathbb{P}(|X - \mu| \geq t) \leq \frac{\sigma^2}{t^2} \quad \forall t > 0.$$

Osservazione 3.2.1. La disuguaglianza di Chebyshev può anche essere formulata nel seguente modo: Se X è una variabile aleatoria con media μ e varianza σ^2 finite, allora

$$\mathbb{P}(|X - \mu| > \eta \sigma) \leq \frac{1}{\eta^2} \quad \forall \eta > 0.$$

Ovvero: la probabilità che X disti dalla sua media μ più di una frazione η della deviazione standard σ è inferiore a $\frac{1}{\eta^2}$.

Esempio 3.2.1. Sia X_1, \dots, X_n un campione statistico di numerosità n . Supponiamo di conoscere la varianza $\sigma^2 = 4$ del campione e che la media μ sia ignota. Quanto deve essere grande n per poter affermare che

$$\mathbb{P}(|\bar{X} - \mu| > 1) \leq \frac{1}{10}?$$

Sappiamo che

$$\mathbb{P}(|\bar{X} - \mu| > 1) \leq \frac{\sigma^2}{n \cdot 1^2} = \frac{4}{n}.$$

è allora sufficiente richiedere $\frac{4}{n} \leq \frac{1}{10}$ cioè $n \geq 40$.

Dalla disuguaglianza di Chebyshev segue facilmente il seguente

Teorema 3.2.2 (Legge debole dei grandi numeri). *Sia $\{X_i\}_{i=1}^{\infty}$ una successione di v.a. indipendenti, identicamente distribuite, con media μ e varianza σ^2 finite.*

Per ogni $n \in \mathbb{N}$ sia $\bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i$. Allora

$$\lim_{n \rightarrow \infty} \mathbb{P}(|\bar{X}_n - \mu| > t) = 0 \quad \forall t > 0.$$

La legge debole dei grandi numeri ci *autorizza* a usare il valore di \bar{X}_n come sostituto della media μ della distribuzione e la disuguaglianza di Chebyshev ci dice con precisione quanto è *probabilisticamente accettabile* questa sostituzione.

Esempio 3.2.2. Ho una monetina che potrebbe essere truccata. Voglio scoprire, con un'approssimazione di ± 0.05 e con un grado di certezza del 90% quanto vale la probabilità di ottenere testa in un singolo lancio. Posso formalizzare ogni singolo lancio della monetina con una variabile aleatoria di Bernoulli di parametro p dove p è la probabilità (incognita) di

ottenere testa in un singolo lancio. Se lancio la moneta n volte ho allora un campione statistico X_1, \dots, X_n che segue la distribuzione $B(p)$. Sia \bar{X}_n la media campionaria di questo campione. Allora

$$\mathbb{E} [\bar{X}_n] = p, \quad \text{Var} [\bar{X}_n] = \frac{p(1-p)}{n}.$$

Per la disuguaglianza di Chebyshev

$$\mathbb{P} (|\bar{X}_n - p| \geq 0.05) \leq \frac{p(1-p)}{n(0.05)^2} \leq \frac{400}{4n} = \frac{100}{n}$$

Voglio

$$\mathbb{P} (|\bar{X}_n - p| \leq 0.05) \geq \frac{90}{100}$$

cioè

$$\mathbb{P} (|\bar{X}_n - p| \geq 0.05) \leq 1 - \frac{90}{100} = \frac{1}{10}$$

Basta allora avere $\frac{100}{n} \leq \frac{1}{10}$ cioè $n \geq 1000$. Dunque: tiro la moneta 1000 volte registrando il risultato ad ogni i -esimo lancio ($x_i = 1$) o croce ($x_i = 0$) vedendo questo numero come il valore assunto da una v.a. bernoulliana X_i di parametro p .

Calcolo $\bar{x} = \frac{1}{1000} \sum_{i=1}^{1000} x_i$ e lo vedo come il valore assunto dalla v.a. \bar{X} . La probabilità che il valore \bar{x} differisca da p per meno di 0.05 è maggiore-uguale del 90%.

Più in generale

Esempio 3.2.3. Sia X_1, \dots, X_n un campione statistico di numerosità n , bernoulliano di parametro (incognito) $p \in [0, 1]$. Dunque

$$\begin{aligned} \mathbb{E} [X_i] &= p & \text{Var} [X_i] &= p(1-p) \\ \mathbb{E} [\bar{X}] &= p & \text{Var} [\bar{X}] &= \frac{p(1-p)}{n} \end{aligned}$$

Allora, per la disuguaglianza di Chebyshev

$$\mathbb{P} (|\bar{X} - p| > t) \leq \frac{p(1-p)}{nt^2} \leq \frac{1}{4nt^2} \quad \forall t > 0. \quad (3.1)$$

poiché $p(1-p) \leq \frac{1}{4} \quad \forall p \in [0, 1]$.

La distribuzione gaussiana $N(\mu, \sigma^2)$ e il teorema del limite centrale

Ricordiamo che la distribuzione gaussiana di parametri $\mu \in \mathbb{R}$ e $\sigma^2 > 0$, $N(\mu, \sigma^2)$, è la distribuzione assolutamente continua associata alla densità

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), \quad x \in \mathbb{R}.$$

Se una v.a. X segue la distribuzione $N(\mu, \sigma^2)$, allora

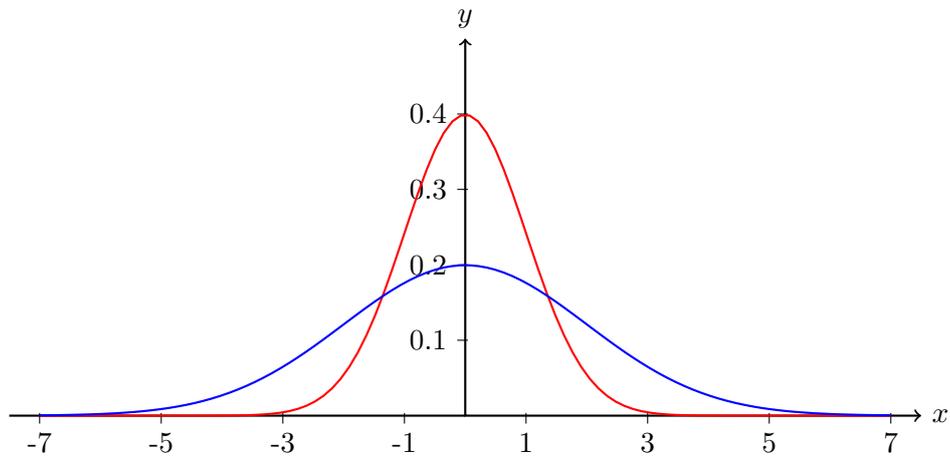


Figura 3.1: Densità associate alle distribuzioni $N(0,1)$ (in rosso) e $N(0,4)$ (in blu)

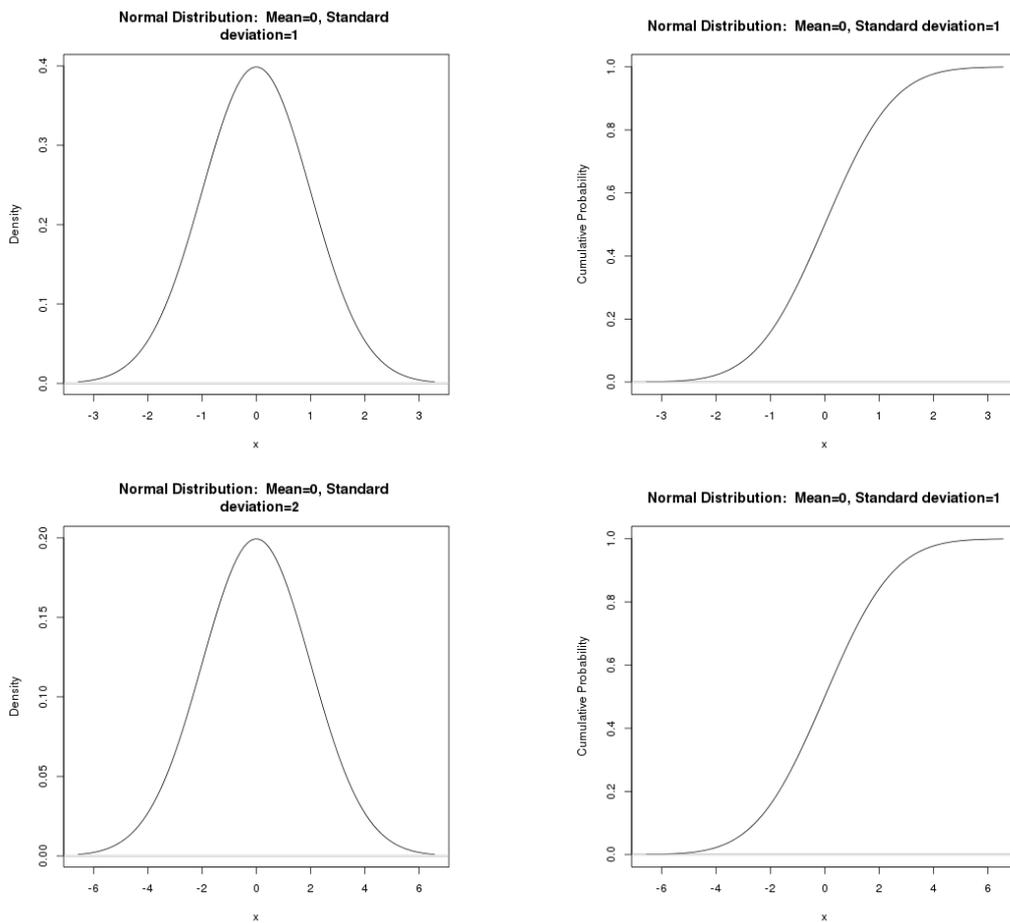


Figura 3.2: $N(0,1)$ e $N(0,4)$, densità e funzione di ripartizione

$$\mathbb{E}[X] = \mu, \quad \text{Var}[X] = \sigma^2.$$

Inoltre $f(x) > 0$ per ogni $x \in \mathbb{R}$, quindi la funzione di ripartizione $F_X(x) := \mathbb{P}(X \leq x)$ è strettamente monotona crescente. Dunque, per ogni $\alpha \in (0, 1)$ esiste uno ed un solo $x = x_\alpha \in \mathbb{R}$ tale $F_X(x_\alpha) = \alpha$. x_α si dice **quantile** di X di livello α . Inoltre, se $\mu = 0$, la densità è una funzione pari, e dunque $F_X(t) + F_X(-t) = 1$ per ogni $t \in \mathbb{R}$; in particolare $x_{1-\alpha} = -x_\alpha$.

Nel caso in cui $\mu = 0$, $\sigma^2 = 1$, la distribuzione $N(0, 1)$ si dice *distribuzione gaussiana standard*, la funzione di ripartizione associata si indica con la lettera Φ ,

$$\Phi(x) := \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{t^2}{2}\right) dt, \quad x \in \mathbb{R}.$$

e per ogni $\alpha \in (0, 1)$ il quantile di livello α si indica z_α . Dunque

$$\Phi(x) + \Phi(-x) = 1 \quad \forall x \in \mathbb{R}, \quad z_{1-\alpha} = -z_\alpha \quad \forall \alpha \in (0, 1).$$

Ricordiamo alcune proprietà che abbiamo già visto:

Proprietà 3.2.1. 1. Se X è una v.a. gaussiana di media μ e varianza σ^2 : $\mathbb{P}_X = N(\mu, \sigma^2)$ e α, β sono due numeri reali, $\alpha \neq 0$, allora la v.a. $\alpha X + \beta$ è gaussiana di media $\alpha\mu + \beta$ e varianza $\alpha^2\sigma^2$: $\mathbb{P}_{\alpha X + \beta} = N(\alpha\mu + \beta, \alpha^2\sigma^2)$. In particolare $Y := \frac{X - \mu}{\sigma}$ è una v.a. gaussiana standard: $\mathbb{P}_Y = N(0, 1)$.

2. Siano X_1, \dots, X_n v.a. indipendenti con X_i gaussiana di media μ_i e varianza σ_i^2 : $\mathbb{P}_{X_i} = N(\mu_i, \sigma_i^2) \forall i = 1, \dots, n$. Allora la v.a. $S_n := X_1 + X_2 + \dots + X_n$ è gaussiana di media pari alla somma delle medie e varianza pari alla somma delle varianze:

$$\mathbb{P}_{S_n} = N\left(\sum_{i=1}^n \mu_i, \sum_{i=1}^n \sigma_i^2\right)$$

Teorema 3.2.3 (Teorema del limite centrale). *Sia $\{X_i\}_{i=1}^\infty$ una successione di v.a. indipendenti, identicamente distribuite, con media μ e varianza σ^2 finite. Sia $\Phi(t)$ la funzione di ripartizione associata alla distribuzione gaussiana standard $N(0, 1)$.*

Per ogni $n \in \mathbb{N}$ sia \bar{X}_n la media campionaria di X_1, \dots, X_n e sia \bar{Z}_n la sua standardizzazione:

$$\bar{Z}_n := \frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}}.$$

Allora

$$\lim_{n \rightarrow \infty} \mathbb{P}(\bar{Z}_n \leq t) = \Phi(t) \quad \forall t \in \mathbb{R}$$

ed il limite è uniforme in $t \in \mathbb{R}$.

Osservazione 3.2.2. Una formulazione equivalente della tesi del teorema del limite centrale è

$$\lim_{n \rightarrow \infty} \mathbb{P}\left(\frac{\sum_{i=1}^n X_i - n\mu}{\sigma\sqrt{n}} \leq t\right) = \Phi(t) \quad \forall t \in \mathbb{R}.$$

Esempio 3.2.4. Supponiamo di avere un campione statistico di numerosità 25 e deviazione standard 8. Qual è la probabilità che la media campionaria differisca dalla media del campione per più di 4?

Devo calcolare

$$\mathbb{P}(|\bar{X} - \mu| > 4)$$

dove $\mu = \mathbb{E}[X_i] \quad \forall i = 1, \dots, n$ e dunque è anche $\mu = \mathbb{E}[\bar{X}]$. Applicando la disuguaglianza di Chebyshev otteniamo

$$\mathbb{P}(|\bar{X} - \mu| > 4) \leq \frac{\text{Var}[\bar{X}]}{4^2} = \frac{64}{25 \cdot 16} = \frac{4}{25} = 0.16$$

Proviamo ad applicare il teorema del limite centrale. Indico con \bar{Z} la standardizzazione della media campionaria. Si ha

$$\begin{aligned} \mathbb{P}(|\bar{X} - \mu| > 4) &= \mathbb{P}\left(\frac{|\bar{X} - \mu|}{\frac{\sigma}{\sqrt{n}}} > \frac{4}{\frac{\sigma}{\sqrt{n}}}\right) = \mathbb{P}\left(|\bar{Z}| > \frac{4}{\frac{8}{\sqrt{25}}}\right) = \\ &= \mathbb{P}\left(|\bar{Z}| > \frac{5}{2}\right) = \mathbb{P}\left(\bar{Z} > \frac{5}{2}\right) + \mathbb{P}\left(\bar{Z} < -\frac{5}{2}\right) \\ &\simeq 1 - \Phi(2.5) + \Phi(-2.5) = 2(1 - \Phi(2.5)) \\ &= 2(1 - \Phi(2.5)) \simeq 2(1 - 0.9938) = 0.0124 \end{aligned}$$

Perché questa stima *sembra* tanto migliore di quella ottenuta con la disuguaglianza di Chebyshev? Perché non abbiamo un'indicazione sul significato del primo dei \simeq . In altre parole, il teorema del limite centrale è appunto un teorema di passaggio al limite e non fornisce una stima dell'errore che si compie sostituendo $\mathbb{P}(Z_n \leq t)$ con $\Phi(t)$. A tal proposito vale il seguente

Teorema 3.2.4 (Teorema di Berry–Esseen). *Sia $\{X_i\}_{i=1}^{\infty}$ una successione di v.a. indipendenti, identicamente distribuite, con media $\mu = 0$, varianza σ^2 e momento terzo $\gamma := \mathbb{E}[|X_i|^3]$ finiti. Sia $\Phi(t)$ la funzione di ripartizione associata alla distribuzione gaussiana standard $N(0, 1)$.*

Sia $C := \frac{0.8\gamma}{\sigma^3}$. Allora

$$\left| \mathbb{P}\left(\frac{\bar{X}_n}{\frac{\sigma}{\sqrt{n}}} \leq t\right) - \Phi(t) \right| \leq \frac{C}{\sqrt{n}} \quad \forall t \in \mathbb{R}.$$

Dal Teorema di Berry–Esseen, teorema 3.2.4, otteniamo dunque

$$|\mathbb{P}(\bar{Z}_n \leq t) - \Phi(t)| \leq \frac{C}{\sqrt{n}} \quad \forall t \in \mathbb{R}.$$

Alcune distribuzioni legate alla distribuzione gaussiana

Distribuzione di Pearson (o χ^2) con n gradi di libertà, χ_n^2

Si tratta della distribuzione $\Gamma(\alpha, \lambda)$ dove $\alpha = \frac{n}{2}$, $\lambda = \frac{1}{2}$. È dunque la distribuzione associata alla densità

$$f(x) := \begin{cases} \frac{1}{\Gamma(\frac{n}{2})} \left(\frac{1}{2}\right)^{\frac{n}{2}} x^{\frac{n}{2}-1} \exp\left(-\frac{x}{2}\right) & x > 0, \\ 0 & x \leq 0, \end{cases}$$

dove $\Gamma(a) := \int_0^{+\infty} x^{a-1} e^{-x} dx$, $a > 0$.

Osservazione 3.3.1. Abbiamo visto che $\forall a > 0$ si ha $\Gamma(a+1) = a\Gamma(a)$ e che $\Gamma(1) = 1$. Inoltre $\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$. Infatti (con la sostituzione $x = y^2$)

$$\Gamma\left(\frac{1}{2}\right) = \int_0^{+\infty} x^{-1/2} e^{-x/2} dx = \int_0^{+\infty} 2 e^{-y^2} dy = \int_{\mathbb{R}} e^{-y^2} dy = \sqrt{\pi}.$$

Quindi

$$\begin{aligned} \Gamma\left(\frac{3}{2}\right) &= \frac{1}{2}\Gamma\left(\frac{1}{2}\right) = \frac{1}{2}\sqrt{\pi}, & \Gamma\left(\frac{5}{2}\right) &= \frac{3}{2}\Gamma\left(\frac{3}{2}\right) = \frac{3 \cdot 1}{2 \cdot 2}\sqrt{\pi} = \frac{3!!}{2^2}\sqrt{\pi}, \\ \dots & & \Gamma\left(\frac{2k+1}{2}\right) &= \frac{(2k-1)!!}{2^k}\sqrt{\pi} \quad \text{per ogni intero non-negativo } k. \end{aligned}$$

Proprietà 3.3.1. Se X è una v.a. con distribuzione χ^2 a n gradi di libertà, $\mathbb{P}_X = \chi_n^2$, allora

$$\mathbb{E}[X] = n, \quad \text{Var}[X] = 2n.$$

Dimostrazione. Poiché una v.a. con distribuzione $\Gamma(\alpha, \lambda)$ ha valore atteso α/λ e varianza α/λ^2 , in particolare per una v.a. con distribuzione di Pearson abbiamo

$$\mathbb{E}[X] = \frac{\frac{n}{2}}{\frac{1}{2}} = n, \quad \text{Var}[X] = \frac{\frac{n}{2}}{\left(\frac{1}{2}\right)^2} = 2n.$$

□

Teorema 3.3.1. Se X e Y sono due variabili di Pearson indipendenti, $\mathbb{P}_X = \chi_n^2$, $\mathbb{P}_Y = \chi_k^2$, allora la v.a. $X + Y$ segue la distribuzione di Pearson a $n + k$ gradi di libertà:

$$\mathbb{P}_{X+Y} = \chi_{n+k}^2.$$

Dimostrazione. Sappiamo che la distribuzione di $X + Y$ è a.c. con densità $h(x)$ data dal prodotto di convoluzione delle densità associate a χ_n^2 e χ_k^2 . Dunque $h(x) = 0$ per $x \leq 0$. Per $x > 0$ abbiamo invece

$$\begin{aligned} h(x) &= \int_0^x \frac{\lambda^\alpha}{\Gamma(\alpha)} y^{\alpha-1} e^{-\lambda y} \frac{\lambda^\beta}{\Gamma(\beta)} (x-y)^{\beta-1} e^{-\lambda(x-y)} dy \\ &= e^{-\lambda x} \frac{\lambda^{\alpha+\beta}}{\Gamma(\alpha)\Gamma(\beta)} \int_0^x y^{\alpha-1} (x-y)^{\beta-1} dy = && \text{(sostituisco } y = xt) \\ &= e^{-\lambda x} \frac{\lambda^{\alpha+\beta} x^{\alpha+\beta-1}}{\Gamma(\alpha)\Gamma(\beta)} \int_0^1 t^{\alpha-1} (1-t)^{\beta-1} dt = C x^{\alpha+\beta-1} e^{-\lambda x} \end{aligned}$$

dove $C = \frac{\lambda^{\alpha+\beta}}{\Gamma(\alpha)\Gamma(\beta)} \int_0^1 t^{\alpha-1} (1-t)^{\beta-1} dt$. Poiché h deve essere una densità di probabilità può solo essere $C = \frac{\lambda^{\alpha+\beta}}{\Gamma(\alpha+\beta)}$. Scegliendo $\alpha = \frac{n}{2}$, $\beta = \frac{k}{2}$, $\lambda = \frac{1}{2}$, si ottiene la tesi. □

Il seguente teorema dà un legame tra la distribuzione gaussiana e le distribuzioni χ^2 :

Teorema 3.3.2. *Se X è una v.a. gaussiana standard, $\mathbb{P}_X = N(0, 1)$, allora X^2 segue la distribuzione di Pearson ad un grado di libertà, $\mathbb{P}_{X^2} = \chi_1^2$.*

Dimostrazione. Sappiamo che $\mathbb{P}_X = N(0, 1) = f(x)dx$ con $f(x) = \frac{1}{\sqrt{2\pi}}e^{-x^2/2}$. Dunque $\mathbb{P}_{X^2} = g(x)dx$ con

$$g(x) = \begin{cases} 0 & x \leq 0, \\ \frac{1}{\sqrt{2\pi}}x^{-1/2}e^{-x/2} & x > 0, \end{cases}$$

cioè $\mathbb{P}_{X^2} = \chi_1^2$. □

Teorema 3.3.3. *Se X_1, \dots, X_n sono v.a. indipendenti e gaussiane, con X_i di media μ_i e varianza σ_i^2 , $\forall i = 1, \dots, n$, allora la v.a. $\sum_{i=1}^n \left(\frac{X_i - \mu_i}{\sigma_i}\right)^2$ segue la distribuzione di Pearson a n gradi di libertà, χ_n^2 .*

Dimostrazione. Poiché la v.a. $\frac{X_i - \mu_i}{\sigma_i}$ ha distribuzione gaussiana standard, applicando i teoremi 3.3.2 e 3.3.1 ed il principio di induzione si ottiene la tesi. □

Corollario 3.3.4. *Se X_1, \dots, X_n è un campione statistico gaussiano, con media μ e varianza σ^2 , allora la v.a. $\sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma}\right)^2$ segue una distribuzione χ^2 con n gradi di libertà.*

Esempio 3.3.1. Si vuole localizzare un oggetto puntiforme, misurandone le tre coordinate cartesiane rispetto ad un prefissato sistema di riferimento. L'errore sperimentale, misurato in millimetri per ciascuna delle tre coordinate è una v.a. gaussiana di media 0 e deviazione standard 2.

Supponendo che i tre errori siano v.a. indipendenti, calcolare la probabilità che la distanza tra la posizione misurata e la posizione reale sia inferiore a 1.2 mm.

Soluzione. Indico con X_1, X_2, X_3 , gli errori commessi nella misurazione delle tre coordinate. Per il Teorema di Pitagora la distanza tra le due posizioni è

$$D = \sqrt{X_1^2 + X_2^2 + X_3^2}$$

Vogliamo calcolare $\mathbb{P}(D < 1.2) = \mathbb{P}(D^2 < 1.44) = \mathbb{P}(X_1^2 + X_2^2 + X_3^2 < 1.44)$.

Pongo $Z_i := \frac{X_i}{\sigma} = \frac{X_i}{2}$, $i = 1, 2, 3$, da cui $X_i^2 = 4Z_i^2$ e dunque

$$\begin{aligned} \mathbb{P}(D < 1.2) &= \mathbb{P}(X_1^2 + X_2^2 + X_3^2 < 1.44) = \mathbb{P}(4(Z_1^2 + Z_2^2 + Z_3^2) < 1.44) \\ &= \mathbb{P}(Z_1^2 + Z_2^2 + Z_3^2 < .36). \end{aligned}$$

Basterà dunque controllare (vedi ultima riga del listato a seguire) il valore della funzione di ripartizione delle v.a. di distribuzione χ_3^2 nel punto 0.36 che è (circa) 0.052.

```
> setwd("/home/laura/Documents/didattica/2017-18_analisi_reale/alcuni_appunti")
> .x <- seq(0.015, 18.015, length.out=100)
```

```
> plot(.x, dchisq(.x, df=3), xlab="x", ylab="Density",
+ main=paste("ChiSquared Distribution: Degrees of freedom=3"), type="l")
> plot(.x, pchisq(.x, df=3), xlab="x", ylab="Density",
+ main=paste("ChiSquared Distribution: Degrees of freedom=3"), type="l")
> abline(h=0.36, col="red")
> pchisq(c(0.36), df=3, lower.tail=TRUE)
[1] 0.05162424
```

Il seguente teorema raccoglie alcune importanti proprietà dei campioni statistici gaussiani e delle loro media e varianza campionarie.

Teorema 3.3.5. *Sia X_1, \dots, X_n un campione statistico gaussiano di numerosità n , valore atteso μ e varianza σ^2 .*

Allora, la media campionaria \bar{X} e la varianza campionaria S^2 sono v.a. indipendenti.

Sia Z_1, Z_2, \dots, Z_n la standardizzazione del campione statistico X_1, \dots, X_n i.e.

$$Z_i := \frac{X_i - \mu}{\sigma} \quad \forall i = 1, \dots, n$$

e sia \bar{Z} la media campionaria del campione normalizzato Z_1, \dots, Z_n .

Allora $\bar{Z} = \frac{\bar{X} - \mu}{\sigma}$ e la v.a. $\sum_{i=1}^n (Z_i - \bar{Z})^2$ sono indipendenti e quest'ultima segue una distribuzione χ^2 con $n - 1$ gradi di libertà.

Dimostrazione. 1. n = 2. Sappiamo che $\mathbb{P}_{X_1+X_2} = N(2\mu, 2\sigma^2)$ e $\mathbb{P}_{\bar{X}} = N(\mu, \sigma^2/2)$. Inoltre

$$S^2 = (X_1 - \bar{X})^2 + (X_2 - \bar{X})^2 = \frac{1}{2}(X_1 - X_2)^2.$$

Dunque \bar{X} e S^2 sono indipendenti se e solo se $X_1 + X_2$ e $X_1 - X_2$ sono indipendenti. Poiché $\mathbb{P}_{-X_2} = N(-\mu, \sigma^2)$ abbiamo che $\mathbb{P}_{X_1-X_2} = N(0, 2\sigma^2)$.

Per provare che $U := X_1 + X_2$ e $V := X_1 - X_2$ sono indipendenti ne calcoliamo la densità congiunta e mostriamo che è uguale al prodotto delle densità marginali. Abbiamo già visto che $\mathbb{P}_{X_1+X_2} = N(2\mu, 2\sigma^2)$. Inoltre, poiché $\mathbb{P}_{-X_2} = N(-\mu, \sigma^2)$ abbiamo che $\mathbb{P}_{X_1-X_2} = N(0, 2\sigma^2)$. Posto

$$\varphi: (x, y) \in \mathbb{R}^2 \mapsto (x + y, x - y) \in \mathbb{R}^2$$

abbiamo

$$(U, V) = \varphi \circ (X_1, X_2)$$

dunque, per ogni funzione boreliana non-negativa $\psi: \mathbb{R}^2 \rightarrow \mathbb{R}$ abbiamo

$$\begin{aligned} \int_{\mathbb{R}^2} \psi(u, v) \mathbb{P}_{U, V}(dudv) &= \int_{\mathbb{R}^2} \psi(x + y, x - y) \mathbb{P}_{X_1, X_2}(dxdy) \\ &= \int_{\mathbb{R}^2} \psi(x + y, x - y) \frac{1}{2\pi\sigma^2} \exp\left(\frac{-(x - \mu)^2}{2\sigma^2}\right) \exp\left(\frac{-(y - \mu)^2}{2\sigma^2}\right) dx dy \end{aligned}$$

con il cambiamento di variabile $u = x + y, v = x - y$

$$= \int_{\mathbb{R}^2} \psi(u, v) \frac{1}{2\pi(\sqrt{2}\sigma)^2} \exp\left(\frac{-(u - 2\mu)^2}{2(\sqrt{2}\sigma)^2}\right) \exp\left(\frac{-v^2}{2(\sqrt{2}\sigma)^2}\right) dudv$$

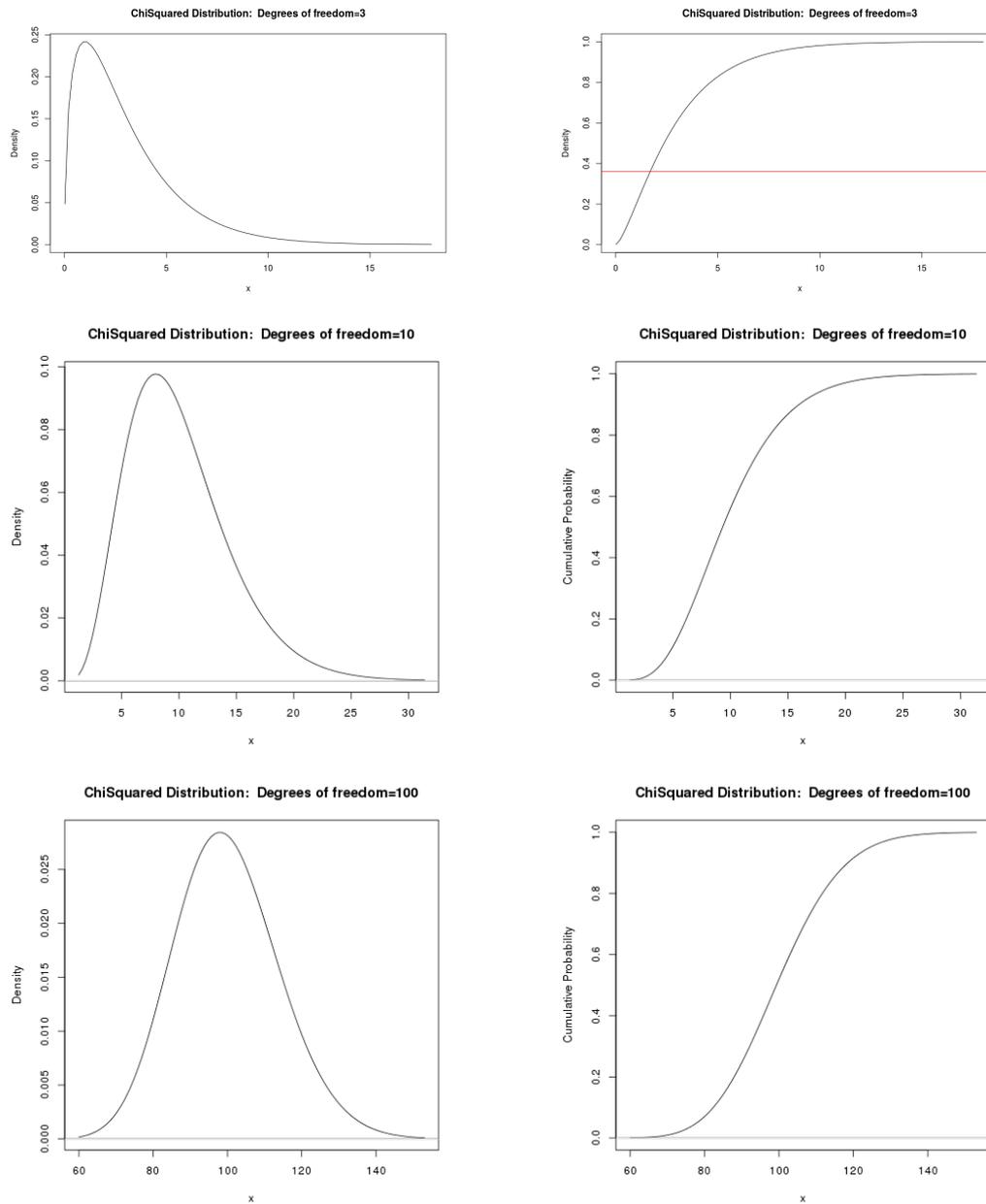


Figura 3.3: χ^2_3 , χ^2_{10} e χ^2_{100} , densità e funzione di ripartizione

ovvero la densità congiunta è il prodotto delle densità marginali

$$f_{X_1+X_2}(u) = \frac{1}{\sqrt{2\pi(\sqrt{2}\sigma)^2}} \exp\left(\frac{-(u-2\mu)^2}{2(\sqrt{2}\sigma)^2}\right), \quad f_{X_1-X_2}(v) = \frac{1}{\sqrt{2\pi(\sqrt{2}\sigma)^2}} \exp\left(\frac{-v^2}{2(\sqrt{2}\sigma)^2}\right).$$

Inoltre, se Z_1 e Z_2 sono gaussiane standard indipendenti abbiamo:

$$(Z_1 - \bar{Z})^2 + (Z_2 - \bar{Z})^2 = \frac{1}{2}(Z_1 - Z_2)^2 = \left(\frac{Z_1 - Z_2}{\sqrt{2}}\right)^2.$$

La v.a. $Z_1 - Z_2$ ha distribuzione $N(0, 2)$, dunque $\frac{Z_1 - Z_2}{\sqrt{2}}$ ha distribuzione $N(0, 1)$. Applicando il Teorema 3.3.2 otteniamo la tesi.

2. $n \geq 3$. Procediamo per induzione, supponendo che \bar{X}_{n-1} e S_{n-1}^2 siano indipendenti. Osserviamo che

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{n} ((n-1)\bar{X}_{n-1} + X_n) = \frac{n-1}{n} \bar{X}_{n-1} + \frac{1}{n} X_n \quad (3.2)$$

e dunque

$$\bar{X}_n - \bar{X}_{n-1} = \frac{1}{n} (X_n - \bar{X}_{n-1}).$$

Abbiamo dunque

$$\begin{aligned} S_n^2 &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_{n-1} + \bar{X}_{n-1} - \bar{X}_n)^2 \\ &= \frac{1}{n-1} \left(\sum_{i=1}^n (X_i - \bar{X}_{n-1})^2 + 2 \sum_{i=1}^n (\bar{X}_{n-1} - \bar{X}_n) (X_i - \bar{X}_{n-1}) + \sum_{i=1}^n (\bar{X}_{n-1} - \bar{X}_n)^2 \right) \\ &= \frac{1}{n-1} \left((n-2)S_{n-1}^2 + (X_n - \bar{X}_{n-1})^2 + 2(\bar{X}_{n-1} - \bar{X}_n) n(\bar{X}_n - \bar{X}_{n-1}) + n(\bar{X}_{n-1} - \bar{X}_n)^2 \right) \\ &= \frac{1}{n-1} \left((n-2)S_{n-1}^2 + (X_n - \bar{X}_{n-1})^2 - \frac{2}{n} (X_n - \bar{X}_{n-1})(X_n - \bar{X}_{n-1}) + \frac{1}{n} (X_n - \bar{X}_{n-1})^2 \right) \\ &= \frac{1}{n-1} \left((n-2)S_{n-1}^2 + \frac{n-1}{n} (X_n - \bar{X}_{n-1})^2 \right) \quad (3.3) \end{aligned}$$

Per la (3.2) e l'ipotesi di induzione \bar{X}_n è indipendente da S_{n-1}^2 . Avremo dunque che S_n^2 e \bar{X}_n sono indipendenti se e solo se \bar{X}_n e $X_n - \bar{X}_{n-1}$ sono indipendenti.

Sappiamo che $\mathbb{P}_{X_n} = N\left(\mu, \frac{\sigma^2}{n}\right)$, dunque

$$\mathbb{P}_{\bar{X}_n} = N\left(\mu, \frac{\sigma^2}{n}\right), \quad \mathbb{P}_{\bar{X}_{n-1}} = N\left(\mu, \frac{\sigma^2}{n-1}\right), \quad \mathbb{P}_{X_n - \bar{X}_{n-1}} = N\left(0, \sigma^2 \frac{n}{n-1}\right),$$

Devo provare che $U := \frac{n-1}{n} \bar{X}_{n-1} + \frac{1}{n} X_n$ e $V = X_n - \bar{X}_{n-1}$ sono indipendenti. Osserviamo che

$$(U, V) = \varphi \circ (\bar{X}_{n-1}, X_n), \quad \varphi(x, y) = \left(\frac{n-1}{n}x + \frac{1}{n}y, y - x\right).$$

Sia dunque $\psi: \mathbb{R}^2 \rightarrow \mathbb{R}$ una funzione di Borel non negativa. Abbiamo

$$\begin{aligned} \int_{\mathbb{R}^2} \psi(u, v) \mathbb{P}_{U, V}(dudv) &= \int_{\mathbb{R}^2} \psi\left(\frac{n-1}{n}x + \frac{1}{n}y, y-x\right) \mathbb{P}_{\bar{X}_{n-1}, X_n} dx dy \\ &= \int_{\mathbb{R}^2} \psi\left(\frac{n-1}{n}x + \frac{1}{n}y, y-x\right) \frac{\sqrt{n-1}}{2\pi\sigma^2} \exp\left(\frac{-(n-1)(x-\mu)^2 - (y-\mu)^2}{2\sigma^2}\right) dx dy \end{aligned}$$

con il cambiamento di variabile $u = \frac{n-1}{n}x + \frac{1}{n}y$, $v = y-x$

$$\begin{aligned} &= \int_{\mathbb{R}^2} \psi(u, v) \frac{\sqrt{n-1}}{2\pi\sigma^2} \exp\left(\frac{-(u-\mu)^2 (\sqrt{n})^2}{2\sigma^2}\right) \exp\left(\frac{-v^2 \left(\sqrt{\frac{n-1}{n}}\right)^2}{2\sigma^2}\right) dudv \\ &= \int_{\mathbb{R}^2} \psi(u, v) \frac{1}{\sqrt{2\pi\frac{\sigma^2}{n}}} \exp\left(\frac{-(u-\mu)^2}{2\left(\frac{\sigma}{\sqrt{n}}\right)^2}\right) \frac{1}{\sqrt{2\pi\sigma^2\frac{n}{n-1}}} \exp\left(\frac{-v^2}{2\left(\sigma\sqrt{\frac{n-1}{n}}\right)^2}\right) dudv \end{aligned}$$

ovvero la densità congiunta è il prodotto delle densità marginali. Questo prova l'indipendenza di U e V e dunque la prima parte della tesi.

Per dimostrare la seconda parte della tesi, osserviamo che essa è sicuramente vera per $n-1$, grazie al Teorema 3.3.2. Procediamo per induzione e riconsideriamo ora la formula (3.3) e supponiamo che essa non sia relativa al campione X_1, \dots, X_n ma alla sua versione standardizzata Z_1, \dots, Z_n :

$$\sum_{i=1}^n (Z_i - \bar{Z}_n)^2 = (n-1)S_n^2 = (n-2)S_{n-1}^2 + \left(\sqrt{\frac{n-1}{n}}(Z_n - \bar{Z}_{n-1})\right)^2.$$

Poiché il campione Z_1, \dots, Z_n è gaussiano standard, $\mathbb{P}_{Z_n - \bar{Z}_{n-1}} = N\left(0, \frac{n}{n-1}\right)$ dunque la

v.a. $\sqrt{\frac{n-1}{n}}(Z_n - \bar{Z}_{n-1})$ è gaussiana standard e quindi il suo quadrato segue una distribuzione di Pearson con un grado di libertà. D'altra parte, per induzione, $\sum_{i=1}^{n-1} (Z_i - \bar{Z}_{n-1})^2 = (n-2)S_{n-1}^2(Z)$ segue una distribuzione di Pearson a $n-2$ gradi di libertà. Per il Teorema 3.3.1 otteniamo la tesi. \square

Corollario 3.3.6. *Sia X_1, \dots, X_n un campione statistico gaussiano di numerosità n , media μ e varianza σ^2 e sia S^2 la sua varianza campionaria. Allora la v.a. $V := (n-1)\frac{S^2}{\sigma^2}$ segue una distribuzione χ^2 con $n-1$ gradi di libertà.*

Dimostrazione. Si ha infatti

$$V = (n-1)\frac{S^2}{\sigma^2} = \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{\sigma^2} \sum_{i=1}^n ((\mu + \sigma Z_i) - (\mu + \sigma \bar{Z}))^2 = \sum_{i=1}^n (Z_i - \bar{Z})^2$$

\square

Distribuzione t di Student con n gradi di libertà, $t(n)$

Si chiama così la distribuzione associata alla densità

$$\tau_n(x) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\Gamma\left(\frac{n}{2}\right)\sqrt{n\pi}} \left(1 + \frac{x^2}{n}\right)^{-\frac{(n+1)}{2}} \quad x \in \mathbb{R}.$$

Proprietà 3.3.2. Se X è una v.a. con distribuzione t di Student a n gradi di libertà, allora

$$\mathbb{E}[X] = 0, \quad \text{Var}[X] = \begin{cases} \frac{n}{n-2} & \text{se } n \geq 3, \\ +\infty & \text{se } n = 1, 2. \end{cases}$$

Osservazione 3.3.2. Il quantile di livello $\alpha \in (0, 1)$ associato alla distribuzione $t(n)$ si indica $t_{n,\alpha}$. Poiché la densità τ_n è una funzione pari, se $X \sim t(n)$, allora $F_X(x) + F_X(-x) = 1$. Dunque per i quantili della distribuzione $t(n)$ si ha $t_{n,\alpha} = -t_{n,1-\alpha}$ per ogni $\alpha \in (0, 1)$.

Teorema 3.3.7. Se Z è una v.a. gaussiana standard, $\mathbb{P}_Z = N(0, 1)$, se Y segue la distribuzione χ^2 con n gradi di libertà, $\mathbb{P}_Y = \chi_n^2$ e se Z e Y sono indipendenti, allora la v.a. $T := \frac{Z\sqrt{n}}{\sqrt{Y}}$ segue la distribuzione t di Student a n gradi di libertà: $\mathbb{P}_T = t(n)$.

Dimostrazione. Possiamo scrivere $T = \varphi \circ (Y, Z)$ dove $\varphi: (y, z) \in \mathbb{R}^2 \mapsto \begin{cases} \frac{z\sqrt{n}}{y} & y > 0 \\ 0 & y \leq 0 \end{cases} \in \mathbb{R}$.

Sia dunque $\psi: \mathbb{R} \rightarrow \mathbb{R}$ una funzione di Borel non negativa.

$$\begin{aligned} \int_{\mathbb{R}} \psi(t) \mathbb{P}_T(dt) &= \int_{y>0, z \in \mathbb{R}} \psi\left(\frac{z\sqrt{n}}{\sqrt{y}}\right) \mathbb{P}_{Y,Z}(dydz) \\ &= \int_{y>0, z \in \mathbb{R}} \psi\left(\frac{z\sqrt{n}}{\sqrt{y}}\right) \frac{1}{\sqrt{2\pi}} \frac{1}{\Gamma\left(\frac{n}{2}\right)} \left(\frac{1}{2}\right)^{\frac{n}{2}} y^{\frac{n}{2}-1} \exp\left(\frac{-y}{2}\right) \exp\left(\frac{-z^2}{2}\right) dydz \end{aligned}$$

con il cambio di variabile $t = \frac{z\sqrt{n}}{\sqrt{y}}$, $z = \frac{t\sqrt{y}}{\sqrt{n}}$, $dz = \frac{\sqrt{y}}{\sqrt{n}} dt$,

$$= \int_{\mathbb{R}} \psi(t) \frac{1}{\sqrt{2n\pi}} \frac{1}{\Gamma\left(\frac{n}{2}\right)} \left(\frac{1}{2}\right)^{\frac{n}{2}} \left(\int_0^{+\infty} y^{\frac{1}{2}} y^{\frac{n}{2}-1} \exp\left(\frac{-y}{2}\right) \exp\left(\frac{-yt^2}{2n}\right) dy\right) dt$$

con il cambio di variabile $u = \frac{y}{2} \left(1 + \frac{t^2}{n}\right)$, $y = 2u \left(1 + \frac{t^2}{n}\right)^{-1}$, $dy = 2 \left(1 + \frac{t^2}{n}\right)^{-1} du$,

$$= \int_{\mathbb{R}} \psi(t) \frac{1}{\sqrt{2n\pi}} \frac{1}{\Gamma\left(\frac{n}{2}\right)} \left(\frac{1}{2}\right)^{\frac{n}{2}} \left(\int_0^{+\infty} (2u)^{\frac{n+1}{2}-1} \exp(-u) \left(1 + \frac{t^2}{n}\right)^{-\frac{(n+1)}{2}} du\right) dt$$

$$= \int_{\mathbb{R}} \psi(t) \frac{1}{\sqrt{2n\pi}} \frac{1}{\Gamma\left(\frac{n}{2}\right)} \left(\frac{1}{2}\right)^{\frac{n}{2}} \left(1 + \frac{t^2}{n}\right)^{-\frac{(n+1)}{2}} \Gamma\left(\frac{n+1}{2}\right) dt$$

da cui la tesi. □

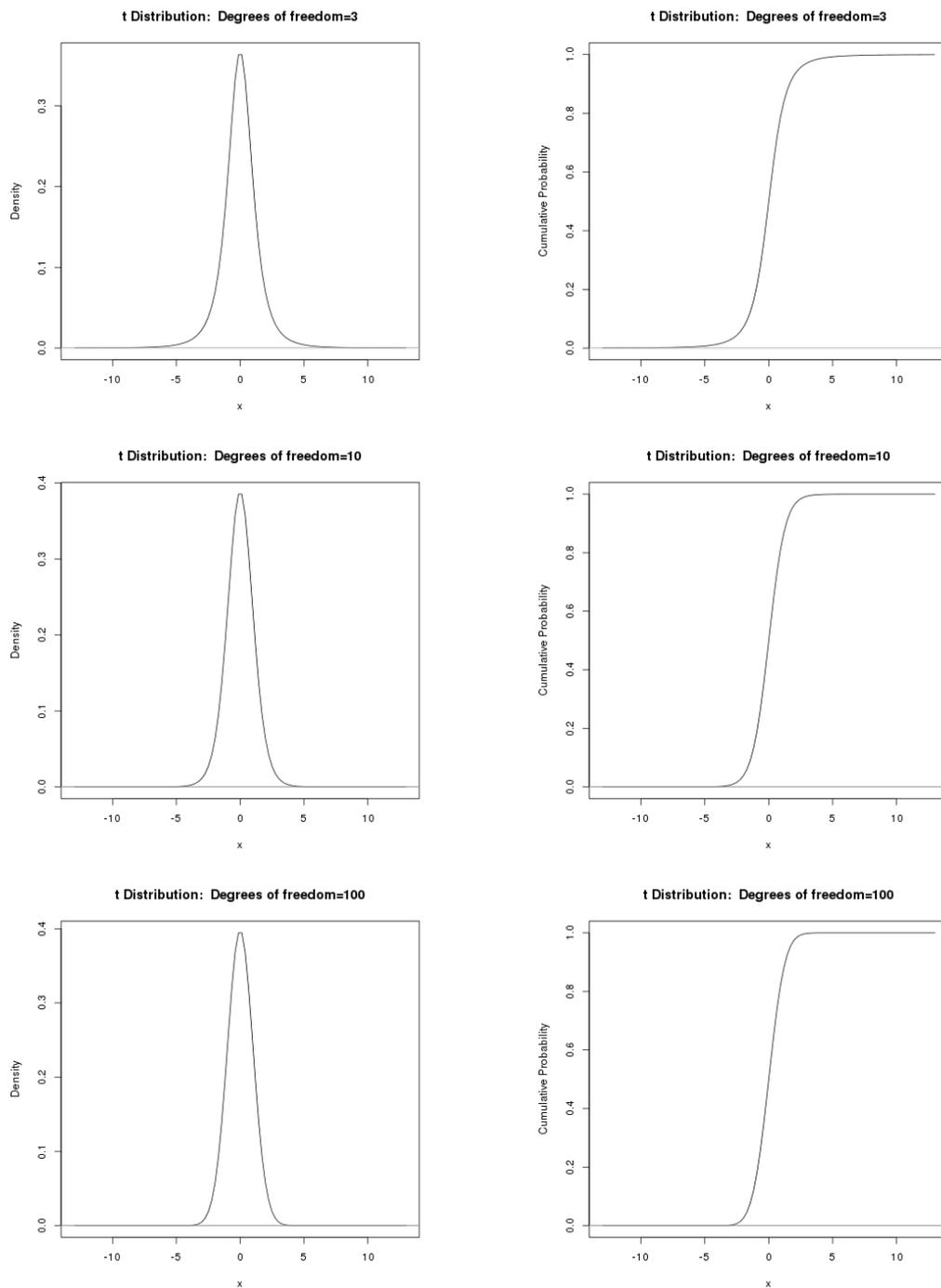


Figura 3.4: $t(3)$, $t(10)$, $t(100)$, densità e funzione di ripartizione

Corollario 3.3.8. *Se X_1, \dots, X_n è un campione statistico gaussiano di numerosità n , valore atteso μ e varianza σ^2 , allora*

$$T := \frac{(\bar{X} - \mu) \sqrt{n}}{S}$$

segue la distribuzione *t* di Student con $n - 1$ gradi di libertà: $\mathbb{P}_T = t(n - 1)$.

Dimostrazione. Basta applicare il teorema 3.3.7 con $Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$ e $Y = V = (n - 1) \frac{S^2}{\sigma^2}$. \square

Stimatori di massima versosimiglianza

Sia X_1, \dots, X_n un campione statistico e sia $Y = \varphi(X_1, \dots, X_n)$ una sua statistica. Se Y ha lo scopo di stimare un parametro θ della distribuzione del campione, diciamo che Y è uno *stimatore del parametro* θ .

Supponiamo di conoscere la distribuzione del campione a meno di un parametro θ e supponiamo che tale distribuzione sia discreta o assolutamente continua e dunque dotata di densità (discreta o meno). Tale densità dipenderà dal parametro θ e la indico col simbolo $g(x|\theta)$. La distribuzione congiunta si indica col simbolo $f(x_1, \dots, x_n|\theta)$ e sappiamo che

$$f(x_1, \dots, x_n|\theta) = g(x_1|\theta) \cdot \dots \cdot g(x_n|\theta) = \prod_{i=1}^n g(x_i|\theta).$$

Interpreto $f(x_1, \dots, x_n|\theta)$ come la *plausibilità* che la n -upla x_1, \dots, x_n si realizzi nel campione empirico quando il parametro incognito prende il valore θ . Sappiamo infatti che, se f è continua nel punto $(x_1, \dots, x_n, \theta)$, allora

$$\begin{aligned} & \mathbb{P} \left(\|X_1 - x_1\| < \frac{\delta}{2}, \dots, \|X_n - x_n\| < \frac{\delta}{2} \right) \\ &= \mathbb{P} \left((X_1, \dots, X_n) \in \prod_{i=1}^n \left(x_i - \frac{\delta}{2}, x_i + \frac{\delta}{2} \right) \right) \simeq f(x_1, \dots, x_n|\theta) \delta^n \end{aligned}$$

Dunque: dato il campione empirico x_1, \dots, x_n , cerco $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$ che massimizza la funzione $f(x_1, \dots, x_n|\theta)$. La statistica $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$ si dirà *stimatore di massima verosimiglianza del parametro* θ .

Osservazione 4.0.1. Poiché la funzione $\ln: (0, +\infty) \rightarrow \mathbb{R}$ è strettamente monotona crescente, massimizzare $f(x, n_1, \dots, x, n|\theta) = \prod_{i=1}^n g(x_i|\theta)$ equivale a massimizzare la funzione $\ln f(x, n_1, \dots, x, n|\theta) = \sum_{i=1}^n \ln g(x_i|\theta)$ e si ha

$$\frac{\partial}{\partial \theta} \ln f(x, n_1, \dots, x, n|\theta) = \frac{\partial}{\partial \theta} \sum_{i=1}^n \ln g(x_i|\theta) = \sum_{i=1}^n \frac{\partial}{\partial \theta} \ln g(x_i|\theta) = \sum_{i=1}^n \frac{1}{g(x_i|\theta)} \frac{\partial g(x_i|\theta)}{\partial \theta}$$

Distribuzione di Bernoulli

Sappiamo che la distribuzione di Bernoulli dipende dal solo parametro $p = \mathbb{P}X = 1$. Sia dunque X_1, \dots, X_n un campione statistico di Bernoulli di parametro incognito $p \in [0, 1]$.

Realizzo n prove di Bernoulli e ottengo il campione empirico $x_1, \dots, x_n, x_i \in \{0, 1\}$.

$$f(x_1, \dots, x_n|p) = \mathbb{P}(X_1 = x_1, \dots, X_n = x_n) = p^k(1-p)^{n-k},$$

$$k = k(x_1, \dots, x_n) := \sum_{i=1}^n x_i.$$

Abbiamo

$$\begin{aligned} \frac{\partial f}{\partial p} &= kp^{k-1}(1-p)^{n-k} - (n-k)p^k(1-p)^{n-k-1} \\ &= p^{k-1}(1-p)^{n-k-1}(k-np) \geq 0 \iff k-np \geq 0 \iff p \leq \frac{k}{n}. \end{aligned}$$

Poiché $k = \sum_{i=1}^n x_i$, lo stimatore di massima verosimiglianza per il parametro p è $\frac{\sum_{i=1}^n X_i}{n}$ cioè la media campionaria \bar{X} .

Distribuzione di Poisson

La distribuzione di Poisson è concentrata sugli interi nonnegativi e dipende da un solo parametro:

$$g(x|\lambda) = e^{-\lambda} \frac{\lambda^x}{x!}, \quad x = 0, 1, \dots$$

e dunque

$$f(x_1, \dots, x_n|\lambda) = \prod_{i=1}^n \left(e^{-\lambda} \frac{\lambda^{x_i}}{x_i!} \right)$$

$$\begin{aligned} \ln f(x_1, \dots, x_n|\lambda) &= \sum_{i=1}^n \ln \left(e^{-\lambda} \frac{\lambda^{x_i}}{x_i!} \right) \\ &= \sum_{i=1}^n (-\lambda + x_i \ln(\lambda) - \ln(x_i!)) = -n\lambda + n\bar{x} \ln(\lambda) - \sum_{i=1}^n \ln(x_i!) \end{aligned}$$

Da cui

$$\frac{\partial}{\partial \lambda} \ln f(x_1, \dots, x_n|\lambda) = n \left(-\lambda + \frac{\bar{x}}{\lambda} \right) \geq 0 \iff \lambda \leq \bar{x}.$$

Quindi anche in questo caso lo stimatore di massima verosimiglianza per il parametro λ è la media campionaria \bar{X} .

Distribuzione gaussiana

In questo caso la densità dipende da due parametri, $\mu \in \mathbb{R}$ e $\sigma > 0$:

$$\begin{aligned} f(x_1, \dots, x_n|\mu, \sigma) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(\frac{-(x_i - \mu)^2}{2\sigma^2} \right) \\ &= (2\pi)^{-\frac{n}{2}} (\sigma)^{-n} \exp \left(\frac{-\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2} \right) \end{aligned}$$

cosicché

$$\begin{aligned}\ln f(x_1, \dots, x_n | \mu, \sigma) &= -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \\ &= -\frac{n}{2} \ln(2\pi) - n \ln(\sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2.\end{aligned}$$

Si ha quindi

$$\begin{aligned}\frac{\partial}{\partial \mu} \ln f(x_1, \dots, x_n | \mu, \sigma) &= \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = n(\bar{x} - \mu), \\ \frac{\partial}{\partial \sigma} \ln f(x_1, \dots, x_n | \mu, \sigma) &= \frac{-n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \mu)^2 = \frac{1}{\sigma^3} \left(-n\sigma^2 + \sum_{i=1}^n (x_i - \mu)^2 \right).\end{aligned}$$

Dunque le due derivate parziali si annullano contemporaneamente se e solo se

$$\mu = \bar{x}, \quad \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{n-1}{n} \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2.$$

Dunque la media campionaria \bar{X} è uno stimatore di massima verosimiglianza per il valore atteso μ mentre $\frac{n-1}{n} S^2$ è uno stimatore di massima verosimiglianza per la varianza σ^2 .

Distribuzione uniforme su un intervallo

Se (a, b) è l'intervallo, allora la densità del campione è

$$g(x|a, b) = \begin{cases} \frac{1}{b-a} & x \in [a, b], \\ 0 & \text{altrimenti} \end{cases}$$

da cui

$$f(x_1, \dots, x_n | a, b) = \begin{cases} \frac{1}{(b-a)^n} & x_i \in [a, b] \quad \forall i = 1, \dots, n, \\ 0 & \text{altrimenti.} \end{cases}$$

Devo massimizzare $\frac{1}{(b-a)^n}$ con il vincolo $a \leq x_i \leq b$ per ogni $i = 1, \dots, n$. Devo dunque minimizzare la lunghezza dell'intervallo $b - a$ con il vincolo $a \leq x_i \leq b$ per ogni $i = 1, \dots, n$. È dunque

$$a = \min \{x_1, \dots, x_n\}, \quad b = \max \{x_1, \dots, x_n\}.$$

Dunque

$$\min \{X_1, \dots, X_n\}, \quad \max \{X_1, \dots, X_n\}$$

sono stimatori di massima verosimiglianza rispettivamente per l'estremo inferiore e per l'estremo superiore dell'intervallo.

Intervalli di confidenza

La media campionaria e la varianza campionaria ci offrono una stima dei parametri valore atteso e varianza del campione statistico in esame. Abbiamo però bisogno di sapere *quanto ci si possa fidare di questa stima* ovvero quale sia la probabilità che il *vero* valore del parametro incognito non sia *troppo distante* dalla stima trovata.

Diamo perciò la seguente definizione:

Definizione 5.0.1 (Intervallo di confidenza). Sia X_1, \dots, X_n un campione statistico e sia θ un parametro (ignoto) che caratterizza la distribuzione del campione.

Siano $L_i = l_i(X_1, \dots, X_n)$ e $L_s = l_s(X_1, \dots, X_n)$ due statistiche del campione e sia $\alpha \in (0, 1)$. Dico che l'intervallo (L_i, L_s) è un *intervallo di confidenza* (o di fiducia) di livello $1 - \alpha$ se $\mathbb{P}(\theta \in (L_i, L_s)) \geq 1 - \alpha$, ovvero che (L_i, L_s) è un intervallo di confidenza (o di fiducia) di errore α se $\mathbb{P}(\theta \notin (L_i, L_s)) \leq \alpha$.

Dico che la semiretta $(L_i, +\infty)$ è un *intervallo di confidenza unilaterale superiore* di livello $1 - \alpha$ se $\mathbb{P}(\theta > L_i) \geq 1 - \alpha$

Dico che la semiretta $(-\infty, L_s)$ è un *intervallo di confidenza unilaterale inferiore* di livello $1 - \alpha$ se $\mathbb{P}(\theta < L_s) \geq 1 - \alpha$

Osservazione 5.0.1. 1. La scelta dei nomi delle due statistiche non è casuale: L_i sta per limitazione inferiore mentre L_s sta per limitazione superiore.

2. Di solito si è interessati a *piccoli* valori di α , più precisamente a $\alpha \in (10^{-2}, 10^{-1})$.

3. La disuguaglianza di Chebyshev ci ha fornito un intervallo di confidenza per il valore atteso μ del campione nel caso in cui la varianza σ^2 sia nota

$$\mathbb{P}(|\bar{X} - \mu| \geq t) \leq \frac{\sigma^2}{t^2} \quad \forall t > 0$$

ovvero

$$\mathbb{P}(|\bar{X} - \mu| < t) \geq 1 - \frac{\sigma^2}{t^2} \quad \forall t > 0$$

cioè

$$\mathbb{P}(\bar{X} - t < \mu < \bar{X} + t) \geq 1 - \frac{\sigma^2}{t^2} \quad \forall t > 0.$$

Fissato $\alpha \in (0, 1)$ scelgo $t = \frac{\sigma}{\sqrt{\alpha}}$. La disuguaglianza di Chebyshev si legge allora

$$\mathbb{P}\left(\bar{X} - \frac{\sigma}{\sqrt{\alpha}} < \mu < \bar{X} + \frac{\sigma}{\sqrt{\alpha}}\right) \geq 1 - \alpha \quad \forall \alpha \in (0, 1).$$

Dunque l'intervallo $\left(\bar{X} - \frac{\sigma}{\sqrt{\alpha}}, \bar{X} + \frac{\sigma}{\sqrt{\alpha}}\right)$ è un intervallo di confidenza di livello $1 - \alpha$ per il valore atteso μ del campione.

Stima per intervalli del valore atteso di campioni gaussiani

Campione gaussiano di cui è nota la varianza

Intervallo bilaterale

Sia X_1, \dots, X_n un campione gaussiano di valore atteso μ incognita e varianza σ^2 nota.

Sia Z una v.a. gaussiana standard e sia $\alpha \in (0, 1)$. Calcolo $\mathbb{P}\left(|Z| \leq z_{1-\frac{\alpha}{2}}\right)$:

$$\begin{aligned} \mathbb{P}\left(|Z| \leq z_{1-\frac{\alpha}{2}}\right) &= \mathbb{P}\left(-z_{1-\frac{\alpha}{2}} \leq Z \leq z_{1-\frac{\alpha}{2}}\right) \\ &= \mathbb{P}\left(Z \leq z_{1-\frac{\alpha}{2}}\right) - \mathbb{P}\left(Z \leq -z_{1-\frac{\alpha}{2}}\right) = \mathbb{P}\left(Z \leq z_{1-\frac{\alpha}{2}}\right) - \mathbb{P}\left(Z \leq z_{\frac{\alpha}{2}}\right) \\ &= \Phi\left(z_{1-\frac{\alpha}{2}}\right) - \Phi\left(z_{\frac{\alpha}{2}}\right) = 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha. \end{aligned} \quad (5.1)$$

Sappiamo che $\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$ e che dunque $\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$. Applichiamo quindi la disuguaglianza (5.1) a $\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$. Si ha:

$$\begin{aligned} 1 - \alpha &= \mathbb{P}\left(\frac{|\bar{X} - \mu|}{\frac{\sigma}{\sqrt{n}}} \leq z_{1-\frac{\alpha}{2}}\right) = \mathbb{P}\left(-z_{1-\frac{\alpha}{2}} \leq \frac{\mu - \bar{X}}{\frac{\sigma}{\sqrt{n}}} \leq z_{1-\frac{\alpha}{2}}\right) \\ &= \mathbb{P}\left(\frac{-\sigma z_{1-\frac{\alpha}{2}}}{\sqrt{n}} \leq \mu - \bar{X} \leq \frac{\sigma z_{1-\frac{\alpha}{2}}}{\sqrt{n}}\right) \\ &= \mathbb{P}\left(\bar{X} - \frac{\sigma z_{1-\frac{\alpha}{2}}}{\sqrt{n}} \leq \mu \leq \bar{X} + \frac{\sigma z_{1-\frac{\alpha}{2}}}{\sqrt{n}}\right) \end{aligned}$$

L'intervallo

$$\left(\bar{X} - \frac{\sigma z_{1-\frac{\alpha}{2}}}{\sqrt{n}}, \bar{X} + \frac{\sigma z_{1-\frac{\alpha}{2}}}{\sqrt{n}}\right)$$

è dunque un intervallo di confidenza di livello $1 - \alpha$ per il valore atteso μ del campione.

Osservazione 5.1.1 (Dimensionamento del campione). Fissato il livello di confidenza $1 - \alpha$, supponiamo di voler controllare l'ampiezza dell'intervallo di confidenza $L_s - L_i$. Nel caso in esame l'ampiezza dell'intervallo di confidenza è $\frac{2\sigma z_{1-\frac{\alpha}{2}}}{\sqrt{n}}$. Se fissiamo una limitazione superiore 2δ per l'ampiezza di tale intervallo, deve dunque essere

$$\frac{2\sigma z_{1-\frac{\alpha}{2}}}{\sqrt{n}} \leq 2\delta$$

ovvero

$$n \geq \left(\frac{\sigma z_{1-\frac{\alpha}{2}}}{\delta}\right)^2.$$

Intervallo unilaterale superiore

Sia $Z \sim N(0, 1)$. Sappiamo che

$$\mathbb{P}(Z \leq t) = 1 - \alpha \quad \text{se e solo se} \quad t = z_{1-\alpha}.$$

Abbiamo dunque

$$1 - \alpha = \mathbb{P} \left(\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \leq z_{1-\alpha} \right) = \mathbb{P} \left(\bar{X} - \mu \leq \frac{\sigma z_{1-\alpha}}{\sqrt{n}} \right) = \mathbb{P} \left(\mu \geq \bar{X} - \frac{\sigma z_{1-\alpha}}{\sqrt{n}} \right).$$

Quindi la semiretta

$$\left(\bar{X} - \frac{\sigma z_{1-\alpha}}{\sqrt{n}}, +\infty \right)$$

è un intervallo di confidenza unilaterale superiore di livello $1 - \alpha$.

Intervallo unilaterale inferiore

Sia $Z \sim N(0, 1)$. Sappiamo che

$$\mathbb{P}(Z \geq t) = 1 - \alpha \quad \text{se e solo se} \quad \mathbb{P}(Z \leq t) = \alpha \quad \text{se e solo se} \quad t = z_\alpha.$$

Abbiamo dunque

$$1 - \alpha = \mathbb{P} \left(\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \geq z_\alpha \right) = \mathbb{P} \left(\bar{X} - \mu \geq \frac{\sigma z_\alpha}{\sqrt{n}} \right) = \mathbb{P} \left(\mu \leq \bar{X} - \frac{\sigma z_\alpha}{\sqrt{n}} \right).$$

Quindi la semiretta

$$\left(-\infty, \bar{X} - \frac{\sigma z_\alpha}{\sqrt{n}} \right) = \left(-\infty, \bar{X} + \frac{\sigma z_{1-\alpha}}{\sqrt{n}} \right)$$

è un intervallo di confidenza unilaterale inferiore di livello $1 - \alpha$.

Campione gaussiano di cui non è nota la varianza

Intervallo bilaterale

Sia X_1, \dots, X_n un campione gaussiano di valore atteso μ varianza σ^2 , entrambe incognite.

Sappiamo che la v.a. $T := \frac{(\bar{X} - \mu)\sqrt{n}}{S}$ segue la distribuzione t di Student con $n - 1$ gradi di libertà:

$$T \sim t(n - 1).$$

Sia $t_{n-1, 1-\frac{\alpha}{2}}$ il relativo quantile di livello $1 - \frac{\alpha}{2}$:

$$\mathbb{P} \left(T \leq t_{n-1, 1-\frac{\alpha}{2}} \right) = 1 - \frac{\alpha}{2}.$$

Calcolo $\mathbb{P} \left(|T| \leq t_{n-1, 1-\frac{\alpha}{2}} \right)$:

$$\begin{aligned} \mathbb{P} \left(|T| \leq t_{n-1, 1-\frac{\alpha}{2}} \right) &= \mathbb{P} \left(-t_{n-1, 1-\frac{\alpha}{2}} \leq T \leq t_{n-1, 1-\frac{\alpha}{2}} \right) \\ &= \mathbb{P} \left(T \leq t_{n-1, 1-\frac{\alpha}{2}} \right) - \mathbb{P} \left(T \leq -t_{n-1, 1-\frac{\alpha}{2}} \right) \\ &= \mathbb{P} \left(T \leq t_{n-1, 1-\frac{\alpha}{2}} \right) - \mathbb{P} \left(T \leq t_{n-1, \frac{\alpha}{2}} \right) = 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha. \end{aligned}$$

Abbiamo dunque

$$\begin{aligned} 1 - \alpha &= \mathbb{P} \left(|T| \leq t_{n-1, 1-\frac{\alpha}{2}} \right) = \mathbb{P} \left(\frac{|\bar{X} - \mu| \sqrt{n}}{S} \leq t_{n-1, 1-\frac{\alpha}{2}} \right) \\ &= \mathbb{P} \left(|\bar{X} - \mu| \leq \frac{S}{\sqrt{n}} t_{n-1, 1-\frac{\alpha}{2}} \right) \\ &= \mathbb{P} \left(\frac{-S}{\sqrt{n}} t_{n-1, 1-\frac{\alpha}{2}} \leq \mu - \bar{X} \leq \frac{S}{\sqrt{n}} t_{n-1, 1-\frac{\alpha}{2}} \right) \\ &= \mathbb{P} \left(\bar{X} - \frac{S}{\sqrt{n}} t_{n-1, 1-\frac{\alpha}{2}} \leq \mu \leq \bar{X} + \frac{S}{\sqrt{n}} t_{n-1, 1-\frac{\alpha}{2}} \right) \end{aligned}$$

L'intervallo

$$\left(\bar{X} - \frac{S t_{n-1, 1-\frac{\alpha}{2}}}{\sqrt{n}}, \bar{X} + \frac{S t_{n-1, 1-\frac{\alpha}{2}}}{\sqrt{n}} \right)$$

è dunque un intervallo di confidenza di livello $1 - \alpha$ per il valore atteso μ del campione.

Intervallo unilaterale superiore

Sappiamo che

$$\mathbb{P}(T \leq t) = 1 - \alpha \quad \text{se e solo se} \quad t = t_{n-1, 1-\alpha}.$$

Abbiamo dunque

$$\begin{aligned} 1 - \alpha &= \mathbb{P} \left(\frac{(\bar{X} - \mu) \sqrt{n}}{S} \leq t_{n-1, 1-\alpha} \right) = \mathbb{P} \left(\bar{X} - \mu \leq \frac{S t_{n-1, 1-\alpha}}{\sqrt{n}} \right) \\ &= \mathbb{P} \left(\mu \geq \bar{X} - \frac{S t_{n-1, 1-\alpha}}{\sqrt{n}} \right). \end{aligned}$$

Quindi la semiretta

$$\left(\bar{X} - \frac{S t_{n-1, 1-\alpha}}{\sqrt{n}}, +\infty \right)$$

è un intervallo di confidenza unilaterale superiore di livello $1 - \alpha$.

Intervallo unilaterale inferiore

Sappiamo che

$$\mathbb{P}(T \geq t) = 1 - \alpha \quad \text{se e solo se} \quad \mathbb{P}(T \leq t) = \alpha \quad \text{se e solo se} \quad t = t_{n-1, \alpha}.$$

Abbiamo dunque

$$1 - \alpha = \mathbb{P} \left(\frac{(\bar{X} - \mu) \sqrt{n}}{S} \geq t_{n-1, \alpha} \right) = \mathbb{P} \left(\bar{X} - \mu \geq \frac{S t_{n-1, \alpha}}{\sqrt{n}} \right) = \mathbb{P} \left(\mu \leq \bar{X} - \frac{S t_{n-1, \alpha}}{\sqrt{n}} \right).$$

Quindi la semiretta

$$\left(-\infty, \bar{X} - \frac{S t_{n-1, \alpha}}{\sqrt{n}} \right) = \left(-\infty, \bar{X} + \frac{S t_{n-1, 1-\alpha}}{\sqrt{n}} \right)$$

è un intervallo di confidenza unilaterale inferiore di livello $1 - \alpha$.

Stima per intervalli della varianza di campioni gaussiani

Intervallo bilaterale

Sia X_1, \dots, X_n un campione gaussiano di valore atteso μ (incognita o nota) e varianza σ^2 incognita.

Sappiamo che la v.a. $V := (n-1)\frac{S^2}{\sigma^2}$ segue la distribuzione χ^2 a $n-1$ gradi di libertà. Per ogni $\alpha \in (0, 1)$ indico con $\chi_{n-1, \alpha}^2$ il quantile di livello α della v.a. V :

$$F_V(\chi_{n-1, \alpha}^2) = \alpha \quad \forall \alpha \in (0, 1).$$

Osservazione 5.2.1. $\chi_{n-1, \alpha}^2 > 0$ per ogni $\alpha \in (0, 1)$.

Calcolo $\mathbb{P}\left(\chi_{n-1, \frac{\alpha}{2}}^2 < V < \chi_{n-1, 1-\frac{\alpha}{2}}^2\right)$:

$$\begin{aligned} \mathbb{P}\left(\chi_{n-1, \frac{\alpha}{2}}^2 < V < \chi_{n-1, 1-\frac{\alpha}{2}}^2\right) &= \mathbb{P}\left(V < \chi_{n-1, 1-\frac{\alpha}{2}}^2\right) - \\ &\quad - \mathbb{P}\left(V < \chi_{n-1, \frac{\alpha}{2}}^2\right) = 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha. \end{aligned}$$

Dunque

$$\begin{aligned} 1 - \alpha &= \mathbb{P}\left(\chi_{n-1, \frac{\alpha}{2}}^2 < (n-1)\frac{S^2}{\sigma^2} < \chi_{n-1, 1-\frac{\alpha}{2}}^2\right) \\ &= \mathbb{P}\left(\frac{1}{\chi_{n-1, 1-\frac{\alpha}{2}}^2} < \frac{\sigma^2}{(n-1)S^2} < \frac{1}{\chi_{n-1, \frac{\alpha}{2}}^2}\right) = \mathbb{P}\left(\frac{(n-1)S^2}{\chi_{n-1, 1-\frac{\alpha}{2}}^2} < \sigma^2 < \frac{(n-1)S^2}{\chi_{n-1, \frac{\alpha}{2}}^2}\right) \end{aligned}$$

Quindi l'intervallo

$$\left(\frac{(n-1)S^2}{\chi_{n-1, 1-\frac{\alpha}{2}}^2}, \frac{(n-1)S^2}{\chi_{n-1, \frac{\alpha}{2}}^2}\right)$$

è un intervallo di confidenza di livello $1 - \alpha$ per la varianza σ^2 del campione.

Intervallo unilaterale superiore

Sappiamo che

$$\mathbb{P}(V \leq t) = 1 - \alpha \quad \text{se e solo se} \quad t = \chi_{n-1, 1-\alpha}^2.$$

Dunque

$$1 - \alpha = \mathbb{P}\left((n-1)\frac{S^2}{\sigma^2} < \chi_{n-1, 1-\alpha}^2\right) = \mathbb{P}\left(\sigma^2 > (n-1)\frac{S^2}{\chi_{n-1, 1-\alpha}^2}\right).$$

Quindi la semiretta

$$\left(\frac{(n-1)S^2}{\chi_{n-1, 1-\alpha}^2}, +\infty\right)$$

è un intervallo di confidenza di livello $1 - \alpha$ per la varianza σ^2 del campione.

Intervallo unilaterale inferiore

Sappiamo che

$$\mathbb{P}(V \geq t) = 1 - \alpha \quad \text{se e solo se} \quad \mathbb{P}(V \leq t) = \alpha \quad \text{se e solo se} \quad t = \chi_{n-1, \alpha}^2.$$

Dunque

$$1 - \alpha = \mathbb{P}\left((n-1)\frac{S^2}{\sigma^2} > \chi_{n-1, \alpha}^2\right) = \mathbb{P}\left(\sigma^2 \leq (n-1)\frac{S^2}{\chi_{n-1, \alpha}^2}\right).$$

Quindi l'intervallo

$$\left(0, \frac{(n-1)S^2}{\chi_{n-1, \alpha}^2}\right)$$

è un intervallo di confidenza di livello $1 - \alpha$ per la varianza σ^2 del campione.

Esempio 5.2.1. Calcoliamo gli intervalli di confidenza per il carattere Totpor dei dati tratti da [2], nell'ipotesi che si tratti della realizzazione di v.a. normali.

```
> setwd("~/Documents/didattica/2017-18_analisi_reale/alcuni_appunti/esempio_statistica")
>
> library(readr)
>
> table2 <- read_delim("~/Documents/didattica/2017-18_analisi_reale/alcuni_appunti/
table2.csv", "\t", escape_double = FALSE, trim_ws = TRUE)
Parsed with column specification:
cols(
  Code = col_character(),
  Totpor = col_double(),
  PRA = col_double(),
  PV = col_double(),
  Densi = col_double(),
  TenStr = col_double(),
  CO2SBW = col_double(),
  FirTemp = col_integer()
)
>
> ## definisco la funzione che calcola l'intervallo bilaterale con varianza nota
>
> bilat.norm = function(x, sigma, conf) { n = length(x); xbar=mean(x);
+ alpha = 1 - conf;
+ zstar = qnorm(1-alpha/2);
+ SE = sigma/sqrt(n);
+ xbar + c(-zstar*SE, zstar*SE)}
>
> # definisco la funzione che calcola l'intervallo bilaterale con varianza ignota
>
> bilat.stud = function(x, conf) { n = length(x);
+ m = n-1;
+ xbar=mean(x);
+ alpha = 1 - conf;
+ zstar = qt(1-alpha/2, m, lower.tail=TRUE);
```

```

+ SE = sd(x)/sqrt(n);
+ xbar + c(-zstar*SE,zstar*SE)
+ }
>
> # definisco la funzione che calcola l'intervallo bilaterale per la varianza
>
> bilat.chi = function(x,conf) {
+   n = length(x);
+   m = n-1;
+   alpha = 1 - conf;
+   zsup = qchisq(alpha/2, m, lower.tail=TRUE);
+   zinf = qchisq(1 - alpha/2, m, lower.tail=TRUE);
+   SE = sd(x)*sd(x)*m;
+   c(SE/zinf,SE/zsup)
+ }
>
>
> numSummary(table2[,c("Totpor", "PRA", "PV", "Densi", "TenStr", "CO2SBW", "FirTemp")],
+ statistics=c("mean", "sd", "quantiles"), quantiles=c(0,.25,.5,.75,1))
      mean      sd      0%      25%      50%      75%      100%  n NA
Totpor  40.1193548  7.0371760  26.850  36.0550  40.900  44.4200  54.640 31  0
PRA      0.6732581  0.4760389   0.158   0.4220   0.622   0.7305   2.657 31  0
PV       55.3290323 28.5498417  10.200  30.4500  59.400  80.7000  88.600 31  0
Densi    1.6929032  0.1701214   1.340   1.5600   1.680   1.8150   2.020 31  0
TenStr   0.6092258  0.3143682   0.143   0.4065   0.527   0.7165   1.405 31  0
CO2SBW   0.5816667  0.5259152   0.050   0.2900   0.390   0.4950   1.960 30  1
FirTemp 764.8387097 52.9698636 730.000 740.0000 740.000 750.0000 960.000 31  0
>
> bilat.norm(table2$Totpor, 7.04, .9)
[1] 38.03957 42.19914
> bilat.norm(table2$Totpor, 7.04, .95)
[1] 37.64113 42.59758
>
> bilat.stud(table2$Totpor, .9)
[1] 37.97416 42.26455
> bilat.stud(table2$Totpor, .95)
[1] 37.53810 42.70061
>
> bilat.chi(table2$Totpor, .9)
[1] 33.94002 80.33757
> bilat.chi(table2$Totpor, .95)
[1] 31.62366 88.48047
>

```

