

**Parte II**

**Statistica inferenziale**



## 5. Campioni statistici

### 5.1. 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  $\mathcal{D}$  e che (molto spesso) si possono supporre indipendenti.
- La distribuzione  $\mathcal{D}$  è (parzialmente) incognita; si cercano informazioni su  $\mathcal{D}$  a partire dai dati rilevati. Le informazioni ricavate sulla distribuzione  $\mathcal{D}$  sono di natura probabilistica. Per esempio, non riusciremo ad ottenere informazioni del tipo *la media della distribuzione  $\mathcal{D}$  è 50* ma informazioni del tipo *la media della distribuzione  $\mathcal{D}$  è compresa tra 49.8 e 50.2 con probabilità del 90%*.

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

**Definizione 5.1.1** (Campione statistico). Una famiglia di variabili aleatorie

$$X_1, X_2, \dots, X_n$$

si dice un *campione statistico di numerosità  $n$*  se le v.a.  $X_1, X_2, \dots, X_n$  sono indipendenti ed identicamente distribuite.

Se  $f$  è la comune densità delle v.a.  $X_1, X_2, \dots, X_n$ , allora la v.a. vettoriale  $X := (X_1, X_2, \dots, X_n)$  ha densità congiunta

$$g_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n) = f(x_1)f(x_2)\dots f(x_n).$$

La comune distribuzione delle  $X_i$  si dice *distribuzione campionaria di  $X_1, X_2, \dots, X_n$* .

**Osservazione 5.1.1.** Poiché le v.a.  $X_1, X_2, \dots, X_n$  seguono la stessa distribuzione, esse hanno anche la stessa media e la stessa varianza (se queste quantità esistono).

**Definizione 5.1.2** (Statistica). Sia  $X_1, X_2, \dots, X_n$  un campione statistico. Una funzione (non dipendente da parametri) di  $X_1, X_2, \dots, X_n$  si dice una statistica.

**Osservazione 5.1.2.** Chiariamo cosa si intende per statistica:  $3X_1 - 2X_2$  è una statistica;  $\max\{X_1, X_2, \dots, X_n\}$  è una statistica.  $X_1 - \mu$   $\mu \in \mathbb{R}$  non è una statistica.

## 5.2. Media campionaria e varianza campionaria

**Definizione 5.2.1.** Sia  $X_1, X_2, \dots, X_n$  un campione statistico. Chiamiamo **media campionaria** di  $X_1, X_2, \dots, X_n$  la statistica

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

chiamiamo **varianza campionaria** di  $X_1, X_2, \dots, X_n$  la statistica

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

**Proposizione 5.2.1.** Sia  $X_1, X_2, \dots, X_n$  un campione statistico di numerosità  $n$  con media  $\mu$  e varianza  $\sigma^2$  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

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

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

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

### 5.2.1. 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  $\mu$  del campione.

**Teorema 5.2.1** (Disuguaglianza di Chebyshev). *Se  $X$  è una variabile aleatoria con media  $\mu$  e varianza  $\sigma^2$  finite, allora*

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

*Dimostrazione.* Consideriamo la v.a.  $Y := \begin{cases} 0 & \text{se } |X - \mu| < t, \\ t & \text{se } |X - \mu| \geq t. \end{cases}$

Sicuramente  $0 \leq Y \leq |X - \mu|$ , quindi  $Y^2 \leq (X - \mu)^2$  e dunque

$$\mathbb{E} [Y^2] \leq \mathbb{E} [(X - \mu)^2] = \text{Var} [X].$$

D'altra parte

$$\mathbb{E} [Y^2] = 0 \mathbb{P} (|X - \mu| < t) + t^2 \mathbb{P} (|X - \mu| \geq t) = t^2 \mathbb{P} (|X - \mu| \geq t),$$

da cui la tesi. □

**Osservazione 5.2.1.** La disuguaglianza di Chebyshev può anche essere formulata nel seguente modo: Se  $X$  è una variabile aleatoria con media  $\mu$  e varianza  $\sigma^2$  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  $\mu$  più di una frazione  $\eta$  della deviazione standard  $\sigma$  è inferiore a  $\frac{1}{\eta^2}$ .

**Esempio 5.2.1.** Sia  $X_1, X_2, \dots, X_n$  un campione statistico di numerosità  $n$ . Supponiamo di conoscere la varianza  $\sigma^2 = 4$  del campione e che la media  $\mu$  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 5.2.2** (Legge debole dei grandi numeri). *Sia  $\{X_i\}_{i=1}^{\infty}$  una successione di v.a. indipendenti, identicamente distribuite, con media  $\mu$  e varianza  $\sigma^2$  finite.*

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

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

*Dimostrazione.* Poiché  $\mathbb{E}[\bar{X}_n] = \mu$  e  $\text{Var}[\bar{X}_n] = \frac{\sigma^2}{n}$ , per la disuguaglianza di Chebyshev si ha

$$\mathbb{P}(|\bar{X}_n - \mu| > t) \leq \frac{\sigma^2}{nt^2} \quad \forall n \in \mathbb{N}.$$

La tesi segue passando a limite. □

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

**Esempio 5.2.2.** Ho una monetina che potrebbe essere truccata. Voglio scoprire, con un'approssimazione di  $\pm 0.05$  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 monetina  $n$  volte ho allora un campione statistico  $X_1, X_2, \dots, X_n$  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 monetina 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 5.2.3.** Sia  $X_1, X_2, \dots, X_n$  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

$$\mathbb{P}(|\bar{X} - p| > t) \leq \frac{p(1-p)}{n t^2} \leq \frac{1}{4n t^2} \quad \forall t > 0. \quad (5.1)$$

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

### 5.2.2. La distribuzione gaussiana $\mathcal{N}(\mu, \sigma^2)$ e il teorema del limite centrale

Ricordiamo che la distribuzione gaussiana di parametri  $\mu \in \mathbb{R}$  e  $\sigma^2 > 0$ ,  $\mathcal{N}(\mu, \sigma^2)$ , 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  $\mathcal{N}(\mu, \sigma^2)$ , allora

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

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  $F_X(x_\alpha) = \alpha$ .  $x_\alpha$  si dice **quantile** di  $X$  di livello  $\alpha$ . 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  $\mathcal{N}((0), 1)$  si dice *distribuzione gaussiana standard*, la funzione di ripartizione associata si indica con la lettera  $\Phi$ ,

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

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

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

Si possono inoltre dimostrare le seguenti proprietà

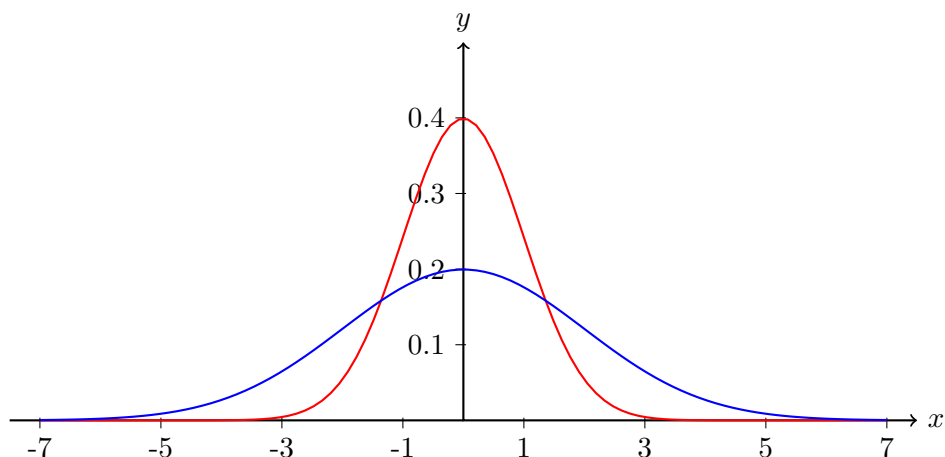


Figura 5.1: Densità associate alle distribuzioni  $\mathcal{N}((0), 1)$  (in rosso) e  $\mathcal{N}((0), 4)$  (in blu)

**Proprietà 5.2.1.** 1. Se  $X$  è una v.a. gaussiana di media  $\mu$  e varianza  $\sigma^2$ :  $X \sim \mathcal{N}(\mu, \sigma^2)$  e  $\alpha, \beta$  sono due numeri reali,  $\alpha \neq 0$ , allora la v.a.  $\alpha X + \beta$  è gaussiana di media  $\alpha\mu + \beta$  e varianza  $\alpha^2\sigma^2$ :  $\alpha X + \beta \sim \mathcal{N}(\alpha\mu + \beta, \alpha^2\sigma^2)$ . In particolare  $Y := \frac{X - \mu}{\sigma}$  è una v.a. gaussiana standard:  $Y \sim \mathcal{N}((0), 1)$ .

2. Siano  $X_1, X_2, \dots, X_n$  v.a. indipendenti. Supponiamo che ciascuna v.a.  $X_i$  sia gaussiana di media  $\mu_i$  e varianza  $\sigma_i^2$ :  $X_i \sim \mathcal{N}(\mu_i, \sigma_i^2) \forall i = 1, \dots, n$ . Allora la v.a.  $X_1 + X_2 + \dots + X_n$  è gaussiana di media pari alla somma delle medie e varianza pari alla somma delle varianze:

$$\sum_{i=1}^n X_i \sim \mathcal{N}\left(\sum_{i=1}^n \mu_i, \sum_{i=1}^n \sigma_i^2\right)$$

**Teorema 5.2.3** (Teorema del limite centrale). Sia  $\{X_i\}_{i=1}^\infty$  una successione di v.a. indipendenti, identicamente distribuite, con media  $\mu$  e varianza  $\sigma^2$  finite. Sia  $\Phi(t)$  la funzione di ripartizione associata alla distribuzione gaussiana standard  $\mathcal{N}(0, 1)$ .

Per ogni  $n \in \mathbb{N}$  sia  $\bar{X}_n$  la media campionaria di  $X_1, X_2, \dots, X_n$  e sia  $\bar{Z}_n$  la sua standardizzazione:

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

Allora

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

**Osservazione 5.2.2.** Una formulazione equivalente della tesi del teorema del limite centrale è

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



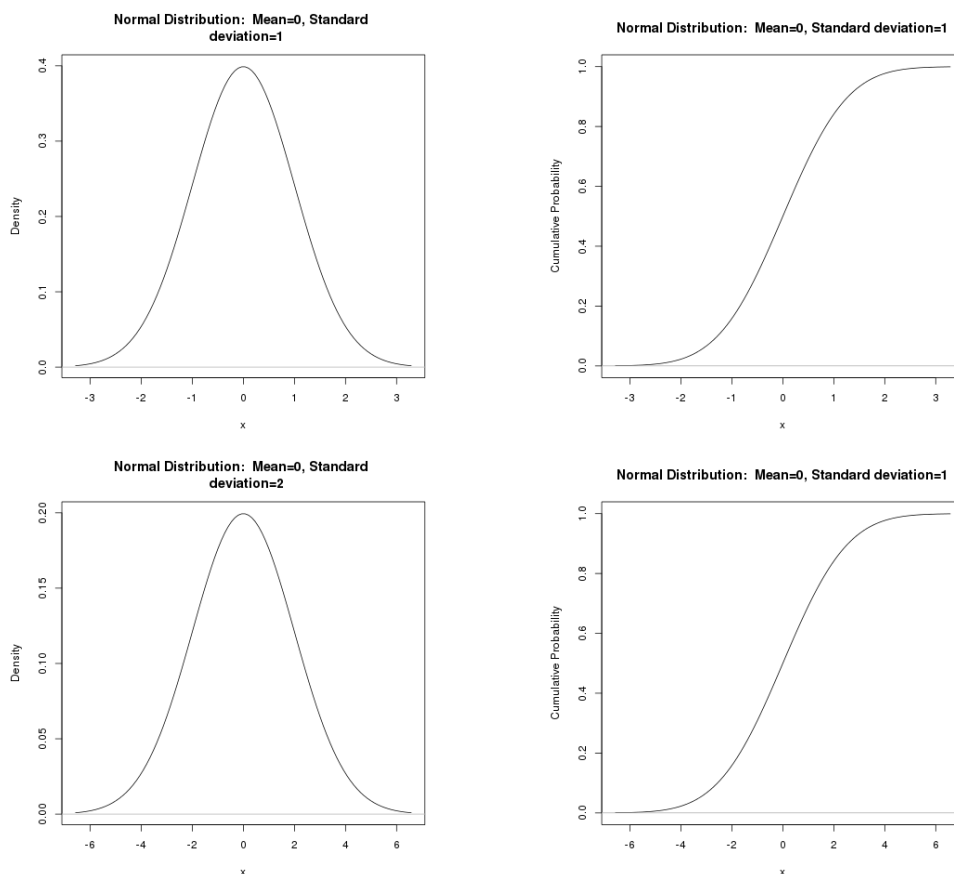


Figura 5.2:  $\mathcal{N}((0), 1)$  e  $\mathcal{N}((0), 4)$ , densità e funzione di ripartizione

**Esempio 5.2.4.** 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

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

Perché questa stima *sembra* tanto migliore di quella ottenuta con la disuguaglianza di Chebyshev? Perché non abbiamo un'indicazione sul significato di quei  $\simeq$ . 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 5.2.4** (Teorema di Berry–Esseen). *Sia  $\{X_i\}_{i=1}^{\infty}$  una successione di v.a. indipendenti, identicamente distribuite, con media  $\mu = 0$ , varianza  $\sigma^2$  e momento terzo  $\gamma := \mathbb{E}[|X_i|^3]$  finiti. Sia  $\Phi(t)$  la funzione di ripartizione associata alla distribuzione gaussiana standard  $\mathcal{N}((0), 1)$ .*

*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 5.2.4, otteniamo dunque

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

### 5.3. Alcune distribuzioni legate alla distribuzione gaussiana

#### 5.3.1. Distribuzione di Pearson (o $\chi^2$ ) con $n$ gradi di libertà, $\chi_n^2$

Si chiama così la distribuzione associata alla densità

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

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

**Osservazione 5.3.1.** Si può dimostrare che  $\forall a > 0$  si ha  $\Gamma(a+1) = a\Gamma(a)$  e che  $\Gamma(1) = 1$ ,  $\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$ . Quindi

$\Gamma(2) = 1 \cdot 1$ ,  $\Gamma(3) = 2 \Gamma(2) = 2 \cdot 1 = 2!$  ...  $\Gamma(n) = (n-1)!$  per ogni intero positivo  $n$

mentre

$$\Gamma\left(\frac{3}{2}\right) = \frac{1}{2}\Gamma\left(\frac{1}{2}\right) = \frac{1}{2}\sqrt{\pi}, \quad \Gamma\left(\frac{5}{2}\right) = \frac{3}{2}\Gamma\left(\frac{3}{2}\right) = \frac{3 \cdot 1}{2 \cdot 2}\sqrt{\pi} = \frac{3!!}{2^2}\sqrt{\pi},$$

$$\dots \quad \Gamma\left(\frac{2k+1}{2}\right) = \frac{(2k-1)!!}{2^k}\sqrt{\pi} \quad \text{per ogni intero non-negativo } k.$$

Si possono calcolare media e varianza di una v.a. che segua una distribuzione di Pearson:

**Proprietà 5.3.1.** Se  $X$  è una v.a. con distribuzione  $\chi^2$  a  $n$  gradi di libertà,  $X \sim \chi_n^2$ , allora

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

**Proprietà 5.3.2.** Se  $X$  e  $Y$  sono due variabili di Pearson indipendenti,  $X \sim \chi_n^2$ ,  $Y \sim \chi_k^2$ , allora la v.a.  $X + Y$  segue la distribuzione di Pearson a  $n + k$  gradi di libertà:

$$X + Y \sim \chi_{n+k}^2.$$

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

**Teorema 5.3.1.** Se  $X_1, X_2, \dots, X_n$  sono v.a. indipendenti e gaussiane, con  $X_i$  di media  $\mu_i$  e varianza  $\sigma_i^2$ ,  $\forall i = 1, \dots, n$ , 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,  $\chi_n^2$ .

**Corollario 5.3.2.** Se  $X_1, X_2, \dots, X_n$  è un campione statistico gaussiano, con media  $\mu$  e varianza  $\sigma^2$ , allora la v.a.  $\chi^2 := \sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma}\right)^2$  segue una distribuzione  $\chi^2$  con  $n$  gradi di libertà.

**Esempio 5.3.1.** 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  $\chi_3^2$  nel punto 0.36 che è (circa) 0.052.

```
> setwd("/home/laura/Documents/didattica/2012-13_elaborazioni_B194")

> .x <- seq(0.015, 17.73, length.out=100)

> plot(.x, dchisq(.x, df=3), xlab="x", ylab="Density",
      main=paste("ChiSquared Distribution: Degrees of freedom=3"), type="l")

> abline(h=0, col="gray")

> remove(.x)

> dev.copy(png, 'densitachiquadro3.png'); dev.off()
png
  3
X11cairo
  2

> .x <- seq(0.015, 17.73, length.out=100)

> plot(.x, pchisq(.x, df=3), xlab="x", ylab="Cumulative Probability",
      main=paste("ChiSquared Distribution: Degrees of freedom=3"), type="l")

> abline(h=0, col="gray")

> remove(.x)

> pchisq(c(0.36), df=3, lower.tail=TRUE)
[1] 0.05162424

> .x <- seq(0.015, 17.73, length.out=100)

> plot(.x, dchisq(.x, df=3), xlab="x", ylab="Density",
      main=paste("ChiSquared Distribution: Degrees of freedom=3"), type="l")

> abline(h=0, col="gray")

> remove(.x)

> dev.copy(png, 'densitachiquadro3.png'); dev.off()
png
  3
```

```

X11cairo
  2

> .x <- seq(0.015, 17.73, length.out=100)

> plot(.x, pchisq(.x, df=3), xlab="x", ylab="Cumulative Probability",
  main=paste("ChiSquared Distribution: Degrees of freedom=3"),type="l")

> abline(h=0, col="gray")

> remove(.x)

> dev.copy(png,'ripartizionechiquadro3.png');dev.off()
png
  3
X11cairo
  2

> pchisq(c(0.36), df=3, lower.tail=TRUE)
[1] 0.05162424

```

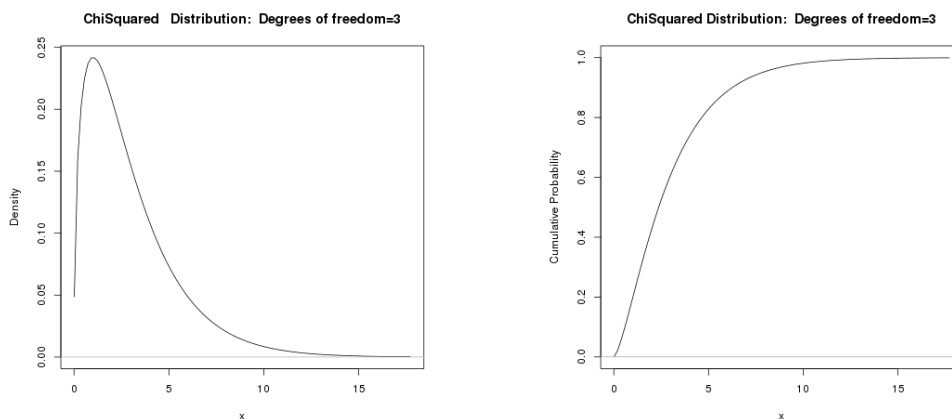
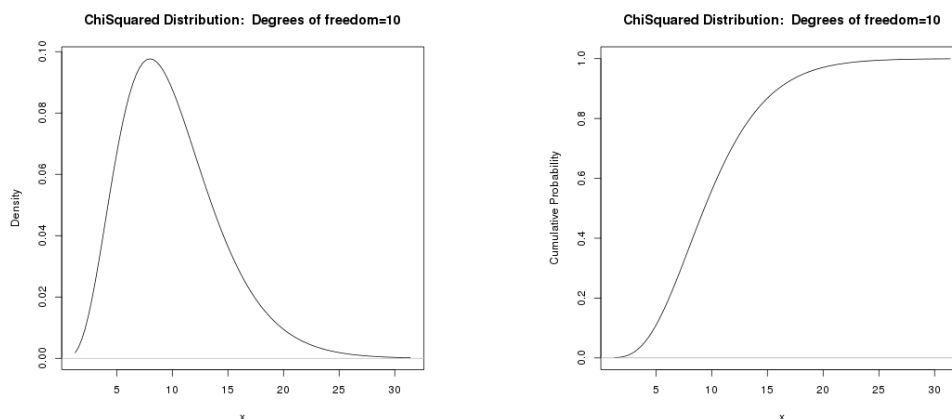
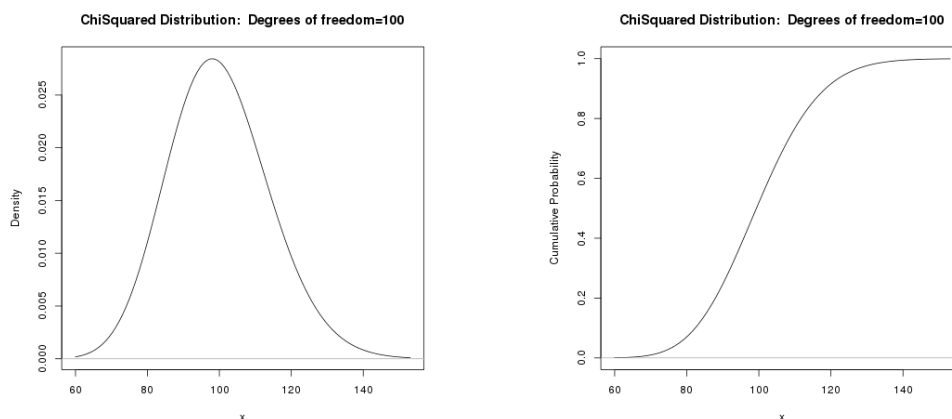


Figura 5.3:  $\chi_3^2$ , densità e funzione di ripartizione

A titolo di confronto, visualizziamo anche densità e funzione di ripartizione delle distribuzioni  $\chi_{10}^2$  e  $\chi_{100}^2$ . Il seguente teorema raccoglie alcune importanti proprietà dei campioni statistici gaussiani e delle loro media e varianza campionarie.

**Teorema 5.3.3.** *Sia  $X_1, X_2, \dots, X_n$  un campione statistico gaussiano di numerosità  $n$ , media  $\mu$  e varianza  $\sigma^2$ .*


 Figura 5.4:  $\chi_{10}^2$ , densità e funzione di ripartizione

 Figura 5.5:  $\chi_{100}^2$ , densità e funzione di ripartizione

Allora, la media campionaria  $\bar{X} := \frac{1}{n} \sum_{i=1}^n X_i$  e la varianza campionaria

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

Sia  $Z_1, Z_2, \dots, Z_n$  la standardizzazione del campione statistico  $X_1, X_2, \dots, X_n$  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  $Z_1, Z_2, \dots, Z_n$ .

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

**Corollario 5.3.4.** Sia  $X_1, X_2, \dots, X_n$  un campione statistico gaussiano di numerosità  $n$ , media  $\mu$  e varianza  $\sigma^2$  e sia  $S^2$  la sua varianza campionaria. Allora la v.a.  $V := (n-1) \frac{S^2}{\sigma^2}$  segue una distribuzione  $\chi^2$  con  $n-1$  gradi di libertà.

*Dimostrazione.* Si ha infatti

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

□

### 5.3.2. 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 5.3.3.** 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 5.3.2.** Il quantile di livello  $\alpha \in (0, 1)$  associato alla distribuzione  $t(n)$  si indica  $t_{n,\alpha}$ . Poiché la densità  $\tau_n$  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 5.3.5.** *w* Se  $Z$  è una v.a. gaussiana standard,  $Z \sim \mathcal{N}(0, 1)$ , se  $Y$  segue la distribuzione  $\chi^2$  con  $n$  gradi di libertà,  $Y \sim \chi_n^2$  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:

$$T := \frac{Z\sqrt{n}}{\sqrt{Y}} \sim t(n).$$

**Corollario 5.3.6.** Se  $X_1, X_2, \dots, X_n$  è un campione statistico gaussiano di numerosità  $n$ , media  $\mu$  e varianza  $\sigma^2$ , allora

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

segue la distribuzione  $t$  di Student con  $n-1$  gradi di libertà:

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

*Dimostrazione.* Basta applicare il teorema 5.3.5 con  $Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$  e  $Y = (n-1) \frac{S^2}{\sigma^2}$ . □

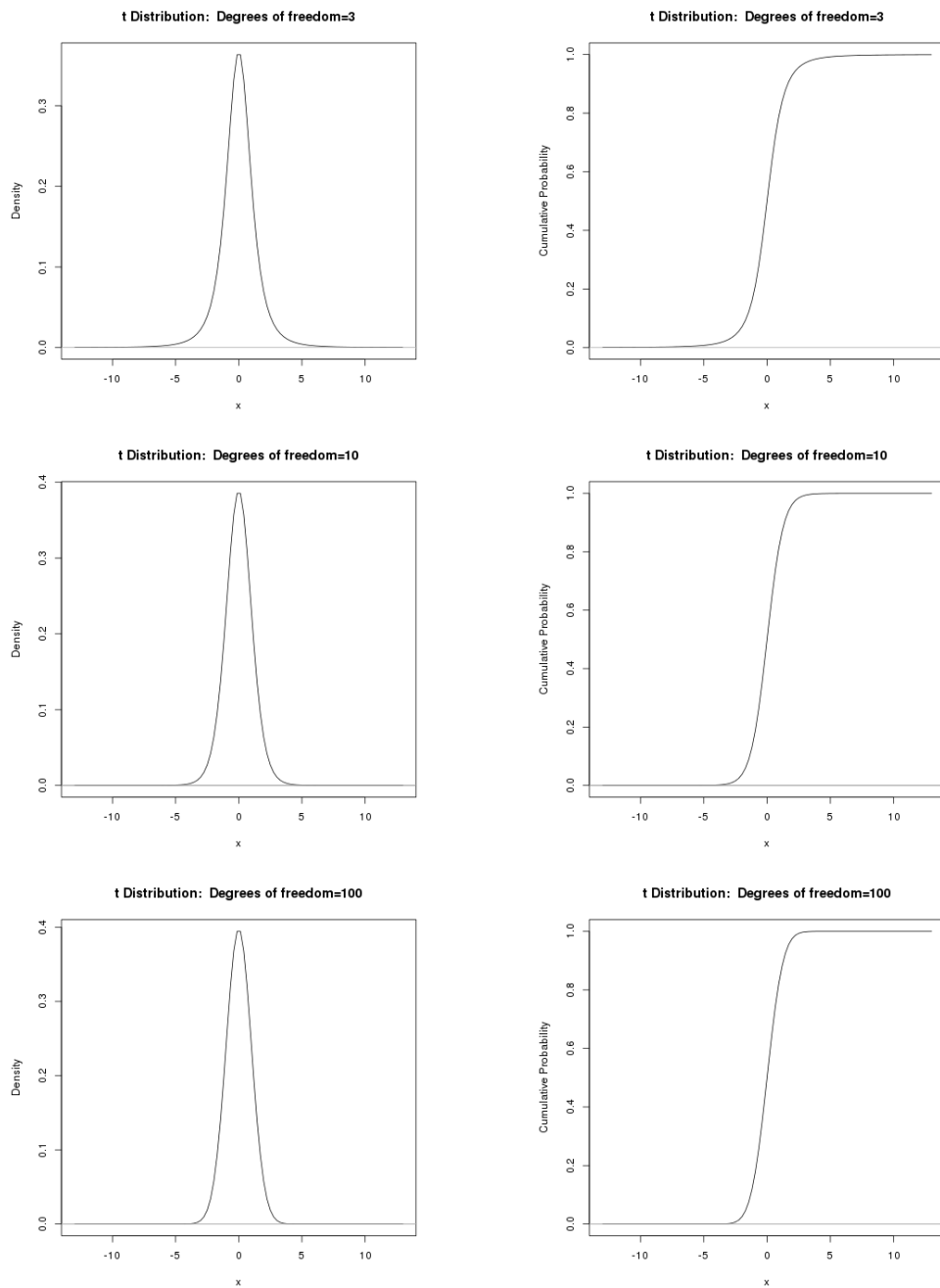


Figura 5.6:  $t(3)$ ,  $t(10)$ ,  $t(100)$ , densità e funzione di ripartizione