

## FORMULARIO STATISTICA

---

### Funzione

Definire una funzione **`y<-function(x)(3*x+7)`**

---

### Grafico

plot(y) se voglio dire gli estremi **`plot(y, -10, 10)`**

---

### Cancellare lo script

**`rm(list=ls())`**

---

### Costruzione variabile nominale

Costruzione vettore nominale sesso di bambini (esempio)

**`sesso<-("M", "F", "M", "M", "M", "F")`**

---

### Frequenza relativa

La frequenza di un valore divisa per il numero totale di frequenze

**`str(x)`** —> per vedere la stringa "x" ovvero tutti gli elementi

---

### Tabella delle frequenze, frequenze relative o proporzionali

**`#tabella delle frequenze`**

**`freq_peso<-table(peso)`**

**`#tabella delle frequenze relative o proporzionali`** chiamata "rel peso"

**`rel_peso<- prop.table(freq_peso)`**

---

### Approssimare alla seconda decimale

**`appr_peso<- round(rel_peso, digits=2)`**

---

### Frequenze percentuali

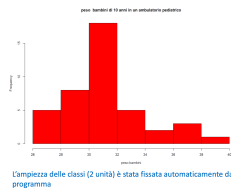
**`#frequenze percentuali`**

**`freq_percent<- round( rel_peso*100, digits=2)`**

---

## Istogrammi

***hist(peso, col="green", main="freq assolute peso", xlab= "peso bambini", ylab= "frequenze assolute")***



## Altri grafici

### A) Diagrammi a torta

***tabella<-table(allergie)*** tabella sulla console

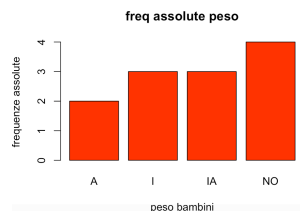
***pie(tabella)***



### B) Diagramma a barre

***#diagramma a barre***

***barplot(tabella)***



## Calcolo della media

Media campionaria =  $(x_1 + x_2 + \dots x_n)/n$

***peso\_medio <- mean(peso)***

Media ponderata =  $(x_1 * f_1 + x_2 * f_2 + \dots x_n * f_n)/n$

***wa <- sum(x\*f)/n***

## Mediana

***median(peso)***

---

## Summery

Il comando "**summery()**" fornisce:

- gli estremi del campione;
- la mediana e i quartili;
- la media campionaria.

---

## Varianza

Serve per determinare la dispersione dei dati rispetto alla media campionaria

$$s^2 = (x_k - media) * f / (n - 1)$$

**var()**

---

## Deviazione standard

a) `s2<-var()` → **sqrt (s2)** →  $\sqrt{Var(x)}$

b) **sd ()**

---

## Coefficiente di variazione

Serve per confrontare la deviazione standard e il valore assoluto della media campionaria

**coeff\_di\_variazione <- sd() / mean()**

---

## Data frame

Se voglio leggere o vedere solo una colonna della tabella (ad esempio)

**#leggere una variabile all'interno di un data set**

**inquinante\$Lago.3 → str(inquinante\$Lago.3)**

---

## Percentile

è una misura usata in statistica per indicare il minimo valore sotto al quale ricade una data percentuale degli altri elementi sotto osservazione.

**Es calcola il 10 percentuale**

**perc10<-quantile (x, 0.10)** (0.10 perchè è 10%)

Tra i percentuali assumono particolare importanza i QUARTILI sono quei valori/modalità che ripartiscono la popolazione in quattro parti di uguale.

## Ordinare il campione

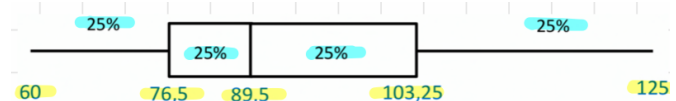
**sort()**

## Boxplot

È un **metodo grafico** per rappresentare la suddivisione in **quartili** dei dati.

Si ottiene tracciando una linea che va dal più piccolo al più grande dei dati e due suddivisioni (box) che rappresentano i quartili.

Ogni box contiene il 25% dei dati:

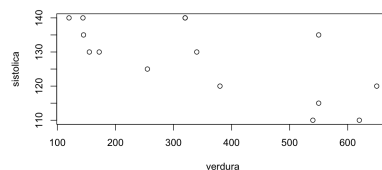


**boxplot (x, horizontal = TRUE, col="red", main="titolo")**

## Scatterplot

è un tipo di grafico in cui due variabili di un set di dati sono riportate su uno spazio cartesiano.

**plot (x,y)**



- se c'è qualche correlazione —> curva;
- se c'è correlazione lineare —> retta.

## Covarianza

Per verificare se fra due variabili statistiche c'è qualche legame lineare.

$$\text{cov}(x, y) = (x_i - x_{\text{media}})(y_i - y_{\text{media}})/(n - 1)$$

**cov (x,y)**

La covarianza può essere:

### POSITIVA

X e Y variano nella  
**STESSA DIREZIONE**

al crescere di X cresce Y  
al diminuire di X diminuisce Y  
e viceversa

### NEGATIVA

X e Y variano in  
**DIREZIONI OPPOSITE**

al crescere di X  
Y tende a diminuire

### NULLA

**nessuna**  
tendenza

$\text{cov}(X, Y) = 0$   
→ **NESSUNA**  
CORRELAZIONE

## Coefficiente di correlazione di Pearson

Serve per calcolare il grado di intensità della correlazione lineare di una coppia di campioni

**$cor(x,y)$**  → “r”

## Valutazione sull'intensità di correlazione

A) valori positivi di r

$0 < r < 0.25$	debole o nessuna
$0.25 < r < 0.75$	moderata
$0.75 < r < 1$	forte

B) valori negativi di r

$-0.25 < r < 0$	debole o nessuna
$-0.75 < r < -0.25$	moderata
$-1 < r < -0.75$	forte

Nel caso in cui vi sia una significativa correlazione lineare (FORTE) si può costruire una retta chiamata retta di regressione

## Retta di regressione

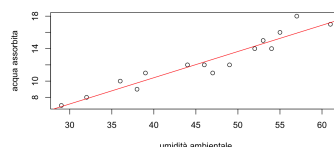
m = pendenza →  $\beta_1 = r_{xy} \frac{s_y}{s_x}$     q = intercetta →  $\beta_0 = \bar{y} - \hat{\beta}_1 \bar{x}$

Con R :  **$lm(x \sim y)$**

## Sovrapposizione tra retta di regressione e i dati reali del campione

**$plot(x \sim y)$**

**$abline(\text{retta di regressione})$**



## Probabilità:

---

### Spazio campionario

Simbolo:  $\Omega$

Significato: insieme di **tutti i possibili esiti** di un esperimento

---

### Evento

Simbolo:  $E$

Significato: **sottoinsieme** dello spazio campionario  $E \subset \Omega$

---

### Operazioni tra insiemi

A) **UNIONE**  $A \cup B = \{x \in \Omega; x \in A \text{ oppure } x \in B\}$

B) **INTERSEZIONE**  $A \cap B = \{x \in \Omega; x \in A \text{ e } x \in B\}$

C) **COMPLEMENTARE**  $A^c = \{x \in \Omega; x \notin A\}$

**Leggi di de Morgan:**  $(A \cup B)^c = A^c \cap B^c$  e  $(A \cap B)^c = A^c \cup B^c$

---

### Insieme delle parti

Diciamo “**insieme delle parti**” una famiglia  $A$  di sottoinsiemi di  $\Omega$  che soddisfi i seguenti criteri:

1.  $A$  non è vuoto  $\rightarrow A \neq \emptyset$
2. Se  $E \in A$ , allora  $\rightarrow E^c \in A$
3. Se  $E_1, E_2 \in A$  allora  $\rightarrow E_1 \cup E_2 \in A$  ( $A$  è chiuso rispetto alle operazioni di unione, intersezione)

$P(\Omega)$  = famiglia di tutti i sottoinsiemi di  $\Omega$  / insieme delle parti di  $\Omega$

---

### Definizione di probabilità

Dati  $\Omega$  spazio campionario e  $A$  sistema di eventi, diciamo **probabilità su  $\Omega$**  un'applicazione  $\rightarrow P: A \rightarrow R; E \rightarrow P(E)$

che soddisfa: 1.  $0 < P(E) < 1$

2.  $P(\Omega) = 1$

3. Se  $A, B \in A$ ,  $A \cap B = \emptyset$  allora  $P(A \cup B) = P(A) + P(B)$

## Regole del calcolo della probabilità

- A) Probabilità del complementare  $P(A^c) = 1 - P(A)$
- B) Evento impossibile  $P(\emptyset) = 0$
- C) Partizione dell'evento certo  $P(B) = P(B \cap A) + P(B \cap A^c)$
- D) Ordinamento  $A \subset C \rightarrow P(A) \leq P(B)$
- E) Unione di eventi non disgiunti  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

---

## Eventi indipendenti

A e B sono indipendenti se il verificarsi di A **non influenza** la probabilità che si verifichi B e viceversa

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

---

## Spazi di probabilità uniformi

Diciamo che  $(\Omega, P)$  è uno **"spazio di probabilità uniforme"** se

1.  $\Omega$  è finito;
2.  $\forall \omega \in \Omega, P\{\omega\} = p$  indipendentemente dal valore di  $\omega$ . (Ovvero se la probabilità è la stessa per ogni esito in  $\Omega$ )

Usiamo la notazione  $|A|$  = tutti gli elementi dell'insieme A

$$p = 1 / |\Omega|$$

Dunque  $\forall A \in P(\Omega)$  abbiamo che  $P(A) = (|A|) \cdot p = |A| / |\Omega|$

Dunque in questo caso si può dire che **Probabilità = casi favorevoli/casi possibili**

---

## Variabili aleatorie

Data una popolazione generica  $\Omega$ , diciamo **variabile aleatoria su  $\Omega$**  una variabile "X" che assume **casualmente** i suoi valori nella popolazione  $\Omega$ .

Generalmente possiamo considerare l'evento  $X = k$ :

Considerata una variabile aleatoria X, ha senso calcolare  $P(X=k)$

## Variabile aleatoria **discreta**

È una variabile aleatoria che assume valori discreti  $x=0,1,2,\dots,n$  e  $y= x_1, x_2, x_3,\dots$

Sempre in generale diciamo che  $x_k \in \text{range di } X$  se  $P(X=x_k) \neq 0$

## Densità di probabilità o funzione di distribuzione o pdf (probability distribution function)

Funzione : fornisce la probabilità di ogni evento costruito a partire da X :

$$f_x(x) = P(X=x) \text{ oppure } p_k = P(X=k)$$

Regole :

1.  $\forall k \ 0 < p_k < 1$
2.  $p_k = P(X=k) = 0$  se  $k \notin \text{range di } X$
3.  $\sum p_k = 1$

## Valore atteso

In generale il valore atteso di una variabile casuale discreta (che assuma cioè solo un numero finito o una infinità numerabile di valori) è dato dalla somma dei possibili valori di tale variabile, ciascuno moltiplicato per la probabilità di essere assunto (ossia di verificarsi), cioè è la media ponderata dei possibili risultati.

Diciamo “valore atteso” di X (media, valor medio, speranza matematica, attesa) la quantità:

$$\mu = E[X] = \sum x_i \cdot P(X=x_i)$$

Proprietà:

1.  $E[cX] = cE[X]$
2.  $E[X + Y] = E[X] + E[Y]$

## Varianza e deviazione standard di una v.a.

Rappresentano la dispersione dei possibili valori di X rispetto al valore atteso.

Data X, variabile aleatoria finita, diciamo “varianza di X” la quantità:

$$\sigma^2 = s^2 = \text{Var}(X) = \sum (x_i - \mu)^2 \cdot P(X=x_i)$$

Diciamo “deviazione standard di X”:

$$\sigma = s = \sqrt{\text{Var}(X)}$$

Proprietà:

1.  $\text{Var}(aX) = a^2 \text{Var}(X)$  ;
2.  $\text{Var}(a + X) = \text{Var}(X)$  (invarianza per traslazione);
3. Se X,Y sono indipendenti allora  $E[XY] = E[X] E[Y]$  e  $\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y)$ .



## Processi di Bernoulli - Distribuzione Binomiale

### Prova Bernoulliana

Diciamo prova bernoulliana un esperimento che può avere solo due risultati

V successo	$P(V) = p$
F insuccesso	$P(F) = 1 - p$

N.B: p per convenzione indica sempre la probabilità di successo

### Variabile aleatoria di Bernoulli

Diciamo v.a. di Bernoulli la v.a. :

$$X \sim B(p)$$

Essa può assumere due valori:

$X = 1$	Successo ( V )	$P(X=1) = p(1) = p$
$X = 0$	Insuccesso ( F )	$P(X=0) = p(0) = 1-p$

### Valore atteso

Indica quante volte in media abbiamo successo

$$E[X] = p$$

### Varianza

$$Var(X) = p(1 - p)$$

### Processi di Bernoulli

È una sequenza (anche illimitata) di prove bernoulliane.

1. Prove indipendenti;
  2. Tutte con la stessa probabilità di successo  $p \in (0,1)$ .
- **Limitati** : sequenza di numero **fissato** n di prove bernoulliane ;
  - **Illimitati** : sequenza di numero **infinito** di prove bernoulliane .

---

## Processi di Bernoulli limitati :

- $X$  : numero di successi ottenuti in  $n$  prove;
- $X$  può assumere i valori:  $0, 1, 2, \dots, n$  ;
- $X$  è detta **variabile aleatoria binomiale** di parametri  $n$  (numero di prove) e  $p$  (probabilità di successo in ogni prova)
- Si indica  $X \sim B(n, p)$

---

## Valore atteso e varianza di un processo bernoulliano limitato

Valore atteso :  $E[X] = np$  ;

Varianza :  $Var(X) = np(1 - p)$

## FUNZIONI IN R :

---

### Fattoriale

FUNZIONE:  $n! = \text{factorial}(n)$

---

### Coefficiente binomiale

Calcoliamo coefficienti binomiali

FUNZIONE:

A) Regola matematica:  $n$  binomiale  $k \rightarrow \text{choose}(n, k)$ :  $k \leq n$

B) Funzione precostituita:  $\text{PP}_7 \leftarrow \text{dbinom}(7, \text{size}=n, \text{prob}=p)$

Per calcolare ad esempio  $P(x \leq 2) = P(x=0) + P(x=1) + P(x=2)$

usiamo la funzione  $p \leftarrow \text{sum}(\text{dbinom}(0:2, \text{size}=n, \text{prob}=p))$

---

### Setnames

Questa è una funzione comoda che imposta i nomi su un oggetto e restituisce l'oggetto. È molto utile alla fine della definizione di una funzione in cui si sta creando l'oggetto da restituire e si preferisce non memorizzarlo con un nome solo per poter assegnare i nomi.

In questo caso

Tabuliamo pdf:

**$\text{tab.pdf} \leftarrow \text{setNames}(\text{pdf}, k)$**

pdf = distribuzione e  $k$  = vettore range della variabile

## Realizzazione di una variabile

Assumiamo di aver estratto concretamente un valore  $x$  dalla popolazione:

$x$  è detto realizzazione della variabile aleatoria  $X$

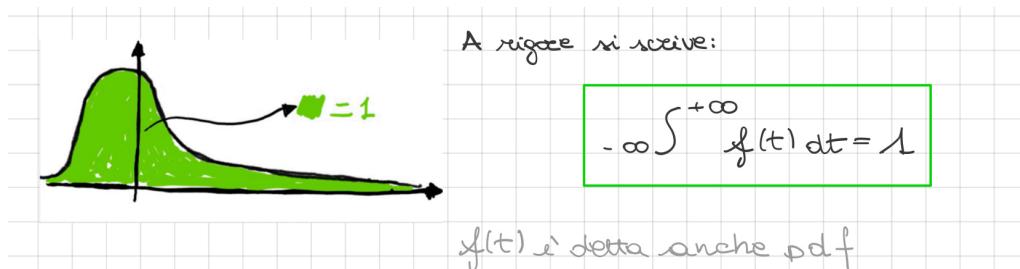
Es: scommettiamo sul lancio di due dati: scommettiamo sul 7 ma all'estrazione esce il numero 5 (realizzazione di  $X$ ). Dato di fatto: abbiamo perso qualsiasi fossero le previsioni probabilistiche di  $X$

## Variabile aleatoria discreta

Assumono valori in un intervallo di numeri reali (es.  $[0, \infty]$ ,  $[0.4, 10]$  ecc)

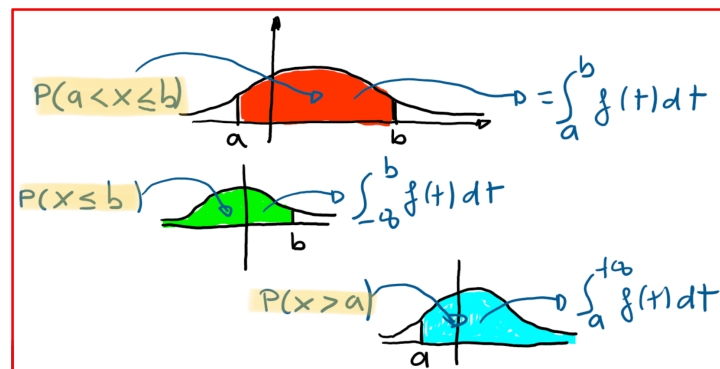
## Funzione densità PDF

È una funzione  $0 < f(t) < 1$  tale che l'area sottesa al suo grafico sia pari a 1:



## Calcolo della probabilità

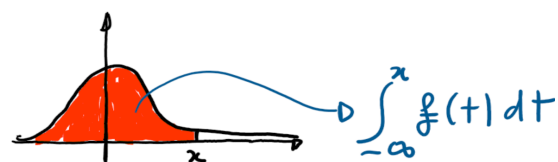
$\forall a, b \in \mathbb{R}, a < b$  la probabilità  $P(a < X < b)$  è data dall'area sottesa al grafico tra l'intervallo  $[a, b]$ :



## Funzione di ripartizione CDF (cumulative distribution function)

$$F(t) = P(X < x)$$

- 1  $F$  funzione non decrescente;
- 2  $\lim_{x \rightarrow -\infty} F(x) = 0$ ;
- 3  $\lim_{x \rightarrow +\infty} F(x) = 1$ .



## Funzione di ripartizione e calcolo della probabilità

- A) Dato un intervallo  $[a, b]$   $P(a < X < b) = F(b) - F(a)$
- B) La funzione di ripartizione ci permette di calcolare  $P(a < X < b)$ , dunque  $F(x)$  descrive completamente la distribuzione della v.a.  $X$ .
- C) Calcolare gli estremi o non calcolarli non fa differenza

## Valore atteso e varianza per una variabile continua

Valore atteso:  $\mu = E[X] = \int_{-\infty}^{+\infty} t f(t) dt$

Intuitivamente può essere considerato come la media dei possibili valori di  $X$  pesati in modo continuo dalla densità  $f(t)$ .

Varianza:  $\sigma^2 = \text{Var}(X) = E[(X - \mu)^2] = \int_{-\infty}^{+\infty} (t - \mu)^2 f(t) dt$

## Distribuzione normale standard

$Z \sim N(0, 1)$

$E[Z] = 0$  e  $\text{Var}(Z) = 1$

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

In R:  $\varphi(x) = \text{dnorm}(x, \text{mean}=0, \text{sd}=1)$

## Calcolo della probabilità per distribuzione Normale Standard

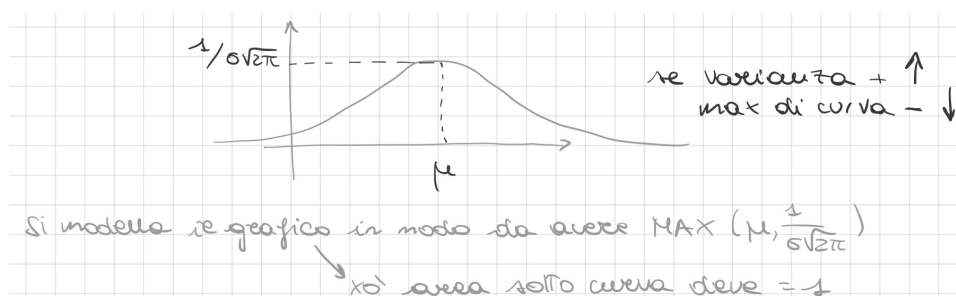
Utilizziamo R perchè non si può fare in altro modo.

## Distribuzione normale di media $\mu$ e varianza $\sigma^2$

$X \sim N(\mu, \sigma^2)$  definita da  $X = \sigma Z + \mu$

$$X \sim N(\mu, \sigma^2) \iff Z \sim N(0, 1)$$

## Curva della distribuzione normale:



## Proprietà della distr. normale

Date:

- X v.a. normale di media  $\mu_x$  e varianza  $\sigma_x^2$ ,  $X \sim N(\mu_x, \sigma_x^2)$
- Y v.a. normale di media  $\mu_y$  e varianza  $\sigma_y^2$ ,  $Y \sim N(\mu_y, \sigma_y^2)$
- X e Y indipendenti

ALLORA:  $X + Y \sim N(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$

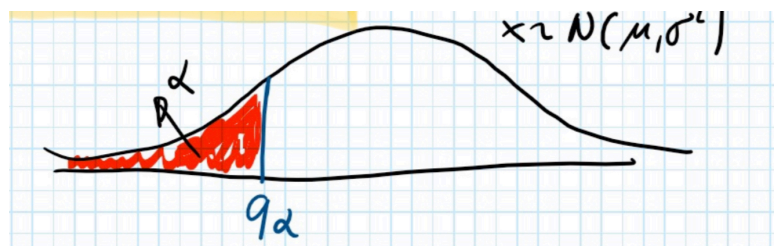
Attenzione! Le varianze si sommano, le deviazioni standard no!

## Quantili della distr. normale standard

Considerata la v.a. normale standard  $Z \sim N(0,1)$  e fissata una probabilità  $0 < \alpha < 1$

Diciamo *quantile di ordine  $\alpha$* , il valore  $z_\alpha$  tale che:  $P(Z < z_\alpha) = \alpha$

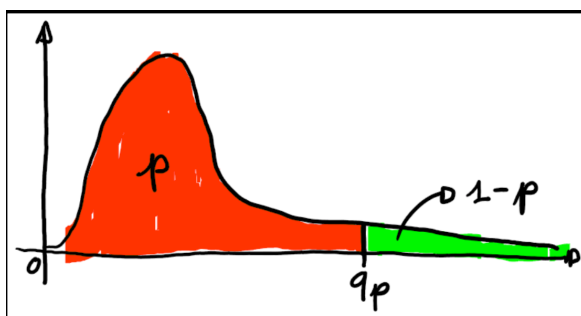
In R :  $qnorm(alpha, mean=0, sd=1) = qnorm(\alpha)$



## Quantili in generale

Data una v.a. X di densità  $f(t)$  e funzione  $F(t) = P(X < t)$ , fissata una probabilità  $0 < p < 1$  diciamo *quantile di ordine p*:

il valore  $q_p \in R$  tale che  $P(X < q_p) = p$



Le proprietà dei quantili  $z_p$  della distribuzione normale standard sono basate sulla simmetria rispetto all'asse delle ordinate della funzione gaussiana.

Dunque non sono valide in generale per i quantili  $q_p$  una distribuzione qualsiasi

## Tabella CODICI R per distribuzioni

Distribuzione	pdf (densità)	cdf (ripartizione)	quantili di ordine $\alpha$
$X \sim \text{unif}[a,b]$	<b>dunif</b> (x, min=a, max=b)	<b>punif</b> (x, min=a, max=b)	<b>qunif</b> (x, min=a, max=b)
$X \sim N(\mu, \sigma^2)$	<b>dnorm</b> (x, mean= $\mu$ , sd= $\sigma$ )	<b>pnorm</b> (x, mean= $\mu$ , sd= $\sigma$ )	<b>qnorm</b> (x, mean= $\mu$ , sd= $\sigma$ )
$X \sim B(n,p)$	<b>dbinom</b> (x, size=n, prob=p)	<b>pbinom</b> (x, size=n, prob=p)	<b>qbinom</b> (x, size=n, prob=p)

Per riassumere:

pdf (densità)	<b>d</b> "nome" = (x, parametri)
cdf (funzione di ripartizione)	<b>p</b> "nome" = (x, parametri)
Quantili di ordine $\alpha$	<b>Q</b> "nome" = (x, parametri)

## Statistica inferenziale

Vogliamo ottenere informazioni sull'**intera popolazione** a partire dai dati ottenuti attraverso il campionamento.

---

### Procedimento

- 1) Fissiamo un campione aleatorio
- 2) Dallo studio teorico otteniamo informazioni sulla distribuzione di probabilità
- 3) Realizzazione:
  - estraiamo campione di dati;
  - Con strumenti teorici (B) effettuiamo stime su  $\mu$  e  $\sigma$

---

### Campione aleatorio **indipendente identicamente distribuito**

Chiamiamo *campione aleatorio i.i.d.* una famiglia di variabili aleatorie che soddisfino:

1. Indipendenti
2. Stessa distribuzione ( = pdf )
3. Stessa varianza ( =  $\sigma^2$  )
4. Stesso valore atteso ( =  $E[X]$  )

---

### Valore atteso $E[X_{media}]$

$$E[X_{media}] = \mu$$

Non dipende da  $n$  (ampiezza)

---

### Varianza $Var(X_{media})$

$$\sigma^2 / n$$

Dipende da  $n$  (  $\uparrow n$  ;  $\downarrow Var$  )

Se la popolazione è normale di media  $\mu$  e varianza  $\sigma^2$

$$X_{media} \sim N(\mu, \sigma^2 / n)$$

## Teorema del limite centrale

Se:

- campione i.i.d.
- media  $\mu$
- varianza  $\sigma^2$
- **$n \geq 30$**

Allora: Assumiamo che la popolazione abbia **distribuzione normale**:

$$\underline{X_{media} \sim N(\mu, \sigma^2 / n)}$$

---

## Realizzazione di un campione aleatorio

Un campione casuale semplice estratto concretamente dalla popolazione

Su R comando **r+nome variabile(...)**

$N(\mu, \sigma^2)$	<code>rnorm(n, mean=mu, sd=sigma)</code>
$unif([a,b])$	<code>runif(n, max=b, min=a)</code>
$B(n,p)$	<code>rbinom(n, size=N, prob=p)</code>

## Intervallo di confidenza per media nota con varianza non nota

- a) Estraiamo un campione casuale di ampiezza  $n$
- b) Fissiamo un livello di fiducia/confidenza  $CL=1 - \alpha$

**`cl<-0.99`**

**`alpha<-1-cl`**

**`alpha`**

- c) Calcoliamo la media campionaria

**`xbar=mean(x)`**

- d) Intervallo di confidenza per la media, al livello di fiducia  $1-\alpha$  è **`[xbar-E, xbar+E]`**

**`zstar<-qnorm(1-alpha/2)`**

**`E<-zstar * sigma/sqrt(n)`**

**`IC<- xbar+c(-1,+1)*E`**

Risposta: - l'intervallo di c. è [estremo sx, estremo dx]

- il livello medio di rumorosità è  $\mu = xbar \pm E$



Se la popolazione ha **media e varianza entrambe NON note?**

$\sigma^2$  non nota  $\rightarrow$  non utilizzabile per costruire intervallo di confidenza

Introduciamo quindi:

- **Varianza campionaria**  $S^2 = 1/(n-1) \sum (x_j - \bar{x})^2$
- **Deviazione standard campionaria**  $S = \sqrt{S^2}$

Esse sono chiamate “*variabili aleatorie statistiche*”  $\rightarrow$  servono per stimare  $\mu$  e  $\sigma$ .

Si verifica che  $(\bar{x} - \mu) / (S / \sqrt{n}) = t_{n-1}$  = distribuzione t di student con n-1 gradi di libertà

## Distribuzione t-di student $X \sim t_n$

Descritta da densità  $f(t)$  con grafico molto simile alla normale standard (campana simmetrico rispetto ad origine) ma è più precisa della normale

$$E[X] = 0$$

$$\text{Var}(X) = n / (n-2) \rightarrow \text{sempre } \text{Var}(X) > 1$$

$$\text{In R: } f(t) = dt(t, df=n)$$

---

Quantili della t di student

$$t^* = qt(alpha, df=n)$$

## Intervallo di confidenza per media nota con varianza non nota

- consideriamo un campione casuale di ampiezza n
- Calcoliamo media e varianza campionarie  $\bar{x} \leftarrow \text{mean}(x)$   $s^2 \leftarrow \text{var}(x)$
- Intervallo di confidenza al livello di confidenza  $CL = 1 - \alpha$

$$[\bar{x} - E, \bar{x} + E]$$

$$E \leftarrow t^* S / \sqrt{n}$$

$$t^* = qt(alpha, df=n)$$

Intervallo di confidenza è dunque:  $\mu = \bar{x} \pm E$

$$\mu \in [\text{estremo sx; estremo dx}]$$

## Proporzione di successi in una popolazione bernoulliana

- Consideriamo la popolazione bernoulliana (successo/insuccesso) con probabilità di successo  $p$  in ogni prova ( $p$  ignota)
- Consideriamo un campione aleatorio i.i.d.
- Variabile aleatoria che conta il numero di successi del campione:

$$\hat{p} = (\text{numero di successi di campione di ampiezza } n) / n$$

$$\text{SE } n \cdot p \geq 5 \text{ e } n \cdot (1-p) \geq 5 \rightarrow \text{ALLORA } [(\hat{p} - p) / \sqrt{p \cdot (1-p)/n}] \sim N$$

(Ha distribuzione normale standard)

## Intervalli di confidenza per proporzione di successi in pop. Bernoulliana

- popolazione Bernoulliana
- Fissiamo livello di confidenza  $CL = 1 - \alpha$
- Estraiamo campione casuale
- Calcoliamo  $\hat{p}$
- Verifichiamo che  $n \cdot p \geq 5$  e  $n \cdot (1-p) \geq 5$
- Calcoliamo errore statistico  $E = z_{(1-\alpha/2)} * [\hat{p}(1-p) / n]$

$$z_{(1-\alpha/2)} = qnorm(1-\alpha/2)$$

Intervallo di confidenza è dunque:  $[\hat{p} - E ; \hat{p} + E]$

$$p = \hat{p} \pm E$$

Come si stima la varianza di una popolazione?

## Stima della varianza

- Consideriamo una popolazione con varianza NON nota

SE consideriamo la varianza campionaria  $S^2 \rightarrow$  ALLORA si verifica che  $E[S^2] = \sigma^2$

$S^2$  è un buon stimatore di  $\sigma^2$

Si verifica inoltre che la variabile aleatoria  $(n-1) * S^2 / \sigma^2$  ha distribuzione “**chi quadro con n-1 gradi di libertà**”

## Distribuzione chi-quadro con n-1 gradi di libertà ( $\chi$ )

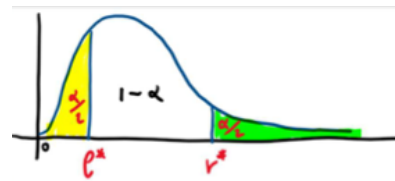
In R: **dchisq(x, df=n)**

---

Quantili della chi-quadro

Sinistro : **lstar <- qchisq(alpha/2, df=n-1)**

Destro : **rstar <- qchisq(1-alpha/2, df=n-1)**



## Intervalli di confidenza per la varianza

- Popolazione normale con varianza NON nota
- Fissiamo il livello di confidenza CL = 1-alpha
- Estraiamo campione casuale
- Media e varianza campionarie  $\bar{x}$  e  $S^2$
- Calcoliamo i quantili lstar e rstar

Intervallo di confidenza è dunque:  $\left( \frac{(n-1)s^2}{r^*} ; \frac{(n-1)s^2}{l^*} \right)$

## Test di ipotesi

Procedura per assumere con un certo livello di significatività statistica se l'affermazione  $H_0$  debba essere rifiutata o non rifiutata.

### A) Ipotesi nulla $H_0$

Afferma che un parametro della popolazione è uguale ad un valore teorico fissato.

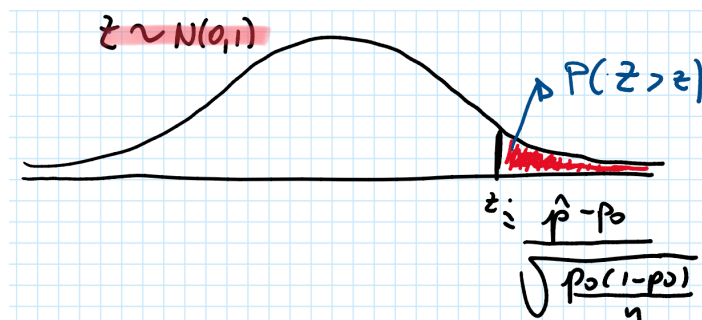
### B) Ipotesi alternativa $H_A$

È l'intervallo di valori che il parametro assume quando l'ipotesi nulla è FALSA.

Come stabilire se rifiutare o no  $H_0$ ?

### Statistica di test z

$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$$



**$P(Z > z) = \text{p-value}$**  (  $Z$ =normale standard;  $z$ =statistica di test)

**$\alpha \in (0,1) = \text{livello di significatività}$**

( In generale  $0.01 < \alpha < 0.1 \rightarrow$  tra l'1% e il 10% )

Quindi:

SE	ALLORA
<b><math>\text{p-value} \leq \alpha</math></b>	Rifiutiamo $H_0$
<b><math>\text{p-value} &gt; \alpha</math></b>	Non rifiutiamo $H_0$

## Test di ipotesi sulla proporzione p di successi in una popolazione Bernoulliana

- Popolazione bernoulliana con  $p \in (0,1)$
- Fissiamo un valore di confronto  $p_0$  e un livello di significatività  $\alpha$
- Estraiamo un campione in modo che  $n \cdot p_0 > 5$  e  $n \cdot (1-p_0) > 5$
- Calcoliamo la proporzione campionaria  $\hat{p} = \text{num di successi} / n$
- Calcoliamo statistica di test  $z = (\hat{p} - p_0) / \sqrt{p_0 (1-p_0) / n}$

Test ad una coda a destra	Test a due code	Test ad una coda a sinistra
$H_0: p=p_0$ ; $H_A: p>p_0$	$H_0: p=p_0$ ; $H_A: p \neq p_0$	$H_0: p=p_0$ ; $H_A: p<p_0$
p-value = $P(Z > z)$	p-value = $2P(Z >  z )$	p-value = $P(Z < z)$

***prop.test(x,n,p=p0,alternative="greater");***

***prop.test(x,n,p=p0,alternative="two.sided");***

***prop.test(x,n,p=p0,alternative="less").***

## Statistica di test per la costruzione di processi decisionali

Tipologia di test	Statistica di test	Distribuzioni delle statistica di test
Test sulla popolazione $H_0: p=p_0$	$z = (\hat{p} - p_0) / \sqrt{p_0 (1-p_0) / n}$	N (0,1)
Test sulla media $H_0: \mu=\mu_0$	a) sigma nota $z = (\bar{x} - \mu_0) / (\sigma / \sqrt{n})$ b) sigma non nota $T = (\bar{x} - \mu_0) / (s / \sqrt{n})$	a) sigma nota N (0,1) b) sigma non nota $t_{(n-1)}$
Test sulla deviazione standard $H_0: \sigma=\sigma_0$	$\text{chiquadro} = [(n-1) \cdot S^2] / \sigma_0^2$	chiquadro^2 con n-1 gradi di libertà

# Test di ipotesi per media e varianza

## Test per la media con varianza NOTA

- Popolazione con media  $\mu$  incognita e  $\sigma^2$  nota
- Fissiamo  $H_0: \mu = \mu_0$
- Fissiamo  $\alpha (0,1)$
- Estraiamo campione con ampiezza  $n$
- Calcoliamo media campionaria  $\bar{x} = \text{mean}(x)$
- Costruiamo  $z = (\bar{x} - \mu_0) / (\sigma / \sqrt{n})$

Test ad una coda a destra	Test a due code	Test ad una coda a sinistra
$H_0: \mu = \mu_0 ; H_A: \mu > \mu_0$	$H_0: \mu = \mu_0 ; H_A: \mu \neq \mu_0$	$H_0: \mu = \mu_0 ; H_A: \mu < \mu_0$
p-value = $P( Z > z )$	p-value = $2P( Z >  z  )$	p-value = $P( Z < z )$

**$p\text{-value} < -pnorm(z) ; \quad p\text{-value} < 2*pnorm(z) ; \quad p\text{-value} < 1-pnorm(z) .$**

**$prop.test(x,n,p=p0,alternative="greater");$**

**$prop.test(x,n,p=p0,alternative="two.sided");$**

**$prop.test(x,n,p=p0,alternative="less").$**

## Test per la media con varianza NON NOTA

(primi 4 punti uguali a prima)

- Calcoliamo media campionaria  $\bar{x}$  e varianza campionaria  $S^2$
- Costruiamo  $t = (\bar{x} - \mu_0) / (s / \sqrt{n})$

Test ad una coda a destra	Test a due code	Test ad una coda a sinistra
$H_0: \mu = \mu_0 ; H_A: \mu > \mu_0$	$H_0: \mu = \mu_0 ; H_A: \mu \neq \mu_0$	$H_0: \mu = \mu_0 ; H_A: \mu < \mu_0$
p-value = $P( T > t )$	p-value = $2P( T >  t  )$	p-value = $P( T < t )$

**$p\text{-value} < -pt(t, df=n-1) ; \quad p\text{-value} < 2*pt( abs(t), df=n-1, lower.tail=FALSE) ;$**

**$p\text{-value} < 1-pt(t, df=n-1) .$**

oppure

**$t.test(x,\mu=\mu_0, alternative="greater");$**

**$t.test(x,\mu=\mu_0,alternative="two.sided");$**

**$t.test(x,\mu=\mu_0,alternative="less")$**

## Test per la varianza

- Popolazione normale con  $\sigma^2$  incognita
- Fissiamo  $H_0$ :  $\sigma^2 = \sigma_0^2$  e fissiamo  $\alpha$  (0,1)
- Estraiamo campione con ampiezza  $n$
- Calcoliamo varianza campionaria  $S^2$
- Costruiamo chiquadro =  $[(n-1) \cdot S^2] / \sigma_0^2$

Test ad una coda a destra	Test a due code	Test ad una coda a sinistra
$H_A$ : $\sigma > \sigma_0$	$H_A$ : $\sigma \neq \sigma_0$	$H_A$ : $\sigma < \sigma_0$
p-value = $P(X > \chi^2)$	p-value = $2 \min[ P(X < \chi^2); P(X > \chi^2) ]$	p-value = $P(X < \chi^2)$

***p-value <- pchisq(chi, df=n-1, lower.tail=FALSE) ;***

***p-value <- 2\*min[pchisq(chi,df=n-1,lower.tail=TRUE),  
pchisq(chi,df=n-1,lower.tail=FALSE)]***

***p-value <- pchisq(chi,df=n-1,lower.tail=TRUE)***

## [ Test non parametrici: test di wilcoxon per la mediana

- Consideriamo  $H_0$ : mediana=m
- Fissiamo  $\alpha$  (0,1)
- Estraiamo campione casuale

Test ad una coda a destra	Test a due code	Test ad una coda a sinistra
$H_A$ : mediana > m	$H_A$ : mediana $\neq$ m	$H_A$ : mediana < m
<b><i>wilcox.test(x, mu=m, alternative="greater")</i></b>	<b><i>wilcox.test(x, mu=m, alternative="two.sided")</i></b>	<b><i>wilcox.test(x, mu=m, alternative="less")</i></b>

Il test non accetta TIES, ovvero ripetizioni nel vettore di dati.

Si aggiunge "exact=FALSE", ovvero R non calcola il p-value preciso, ma lo stima in modo che possa funzionare lo stesso. ]

## Confronto tra parametri (media, dev.standard etc) tra due popolazioni diverse

### Test sul confronto tra due popolazioni bernoulliane

- Abbiamo popolazione 1 con proporzione  $p_1$  di successo e  
popolazione 2 con proporzione  $p_2$  di successo
  - Ampiezza dei campioni  $n_1$  e  $n_2$ , non per forza uguali
  - Estraiamo campioni X e Y
  - Calcoliamo  **$\hat{p}_1 = \text{numero successi in } X / n_1$**  e  
 **$\hat{p}_2 = \text{numero successi in } Y / n_2$**
  - Costruiamo la statistica di test così:
    - $\hat{p} = \text{proporzione totale di successi} = (\hat{p}_1 * n_1 + \hat{p}_2 * n_2) / (n_1 + n_2)$
    - Statistica di test =  **$z = (\hat{p}_1 - \hat{p}_2) / \sqrt{\hat{p} * (1 - \hat{p}) * (1/n_1 + 1/n_2)}$**
- Si distribuisce come una normale standard  $N(0,1)$

Codifica R:

- Vettore numero di successi  **$vp = c(n_1 * \hat{p}_1, n_2 * \hat{p}_2)$**
- Vettore ampiezze dei campioni  **$vn = (n_1, n_2)$**
- Effettuiamo il prop.test  **$prop.test(vp, vn, alternative = "greater, two.sided, less")$**



# Test di confronto tra medie, mediane e varianze

## Test di confronto tra medie di due popolazioni

- Abbiamo popolazione X e popolazione Y
- Ampiezza dei campioni  $n_x$  e  $n_y$ , non per forza uguali
- Estraiamo campioni X e Y
- Consideriamo medie e varianze REALI  $\mu_x, \mu_y, \sigma^2_x, \sigma^2_y$
- $H_0: \mu_x = \mu_y$
- Calcoliamo  $\bar{x}$ ,  $\bar{y}$  (medie campionarie)
- Calcoliamo  $S_x^2$  e  $S_y^2$  varianze campionarie

Distinguiamo i test sulla base delle informazioni su  $\sigma^2_x$  e  $\sigma^2_y$

### 1) Varianze note

Statistica di test:  $z = (\bar{x} - \bar{y}) / \sqrt{[\sigma^2_x/n_x + \sigma^2_y/n_y]}$

Coda a dx: **`p-value<-pnorm(z,lower.tail=FALSE);`**

Due code: **`p-value<-2*pnorm(abs(z), lower.tail=FALSE );`**

Coda a sx: **`p-value<-1-pnorm(z, lower.tail=TRUE).`**

### 2) Varianze NON note ma uguali $\sigma^2_x = \sigma^2_y$

Statistica di test:  $t = (\bar{x} - \bar{y}) / \sqrt{[1/n_x + 1/n_y]}$

Coda a dx: **`t.test<-(x, y, alternative="greater", var.equal=TRUE)`**

Due code: **`t.test(x, y, alternative="two.sided", var.equal=TRUE)`**

Coda a sx: **`t.test(x, y, alternative="less",var.equal=TRUE)`**

### 3) Varianze NON note e ignote $\sigma^2_x \neq \sigma^2_y$

Statistica di test:  $t = (\bar{x} - \bar{y}) / \sqrt{[S_x^2/n_x + S_y^2/n_y]}$

Coda a dx: **`t.test<-(x, y, alternative="greater", var.equal=FALSE)`**

Due code: **`t.test(x, y, alternative="two.sided", var.equal=FALSE)`**

Coda a sx: **`t.test(x, y, alternative="less",var.equal=FALSE)`**

## Test di confronto tra mediane di due popolazioni

- $H_0$ : medianaA = medianaB

Utilizziamo il test di Wilcoxon:

Coda a dx: **wilcox.test<-(x, y, alternative="greater")**

Due code: **wilcox.test(x, y, alternative="two.sided")**

Coda a sx: **wilcox.test(x, y, alternative="less")**

## Test di confronto tra varianze

- $H_0$ :  $\sigma^2_A = \sigma^2_B$
- Calcoliamo varianze campionarie  $S^2_A$  e  $S^2_B$

Statistica di test:  $f = S^2_A / S^2_B$  distribuzione di Fisher

### [ Distribuzione di Fisher

Codifica R:

- Quantili: qf ( $\alpha$ , n, m)
- $P(F < f) = \text{pf}(\alpha, n, m, \text{lower.tail}=\text{TRUE})$
- $P(F > f) = \text{pf}(\alpha, n, m, \text{lower.tail}=\text{FALSE})$  ]

Coda a dx: **var.test<-(x, y, alternative="greater")**

Due code: **var.test(x, y, alternative="two.sided")**

Coda a sx: **var.test(x, y, alternative="less")**

## Test di indipendenza

Serve per verificare se due variabili sono indipendenti o dipendono l'una dall'altra.

- $H_0$ : variabili sono indipendenti
- $H_A$ : variabili non indipendenti

Codifica R:

- Calcolare la matrice  $x=(25,13; 12,15)$  **`x<-rbind( c(25,13), c(12,15) )`**
- Eseguiamo il test chiquadro **`chisq.test(x)`**

Attenzione: se le frequenze in una casella sono meno di 5 appare un Warning

allora **`chisq.test(x,simulate.p.value=TRUE)`**

## Test di adattamento ad un modello

Serve per verificare se un modello probabilistico ipotizzato per la popolazione sia o meno compatibile con i dati campionari che abbiamo.

### Test chi-quadro di adattamento

- n categorie A1, A2 etc
- Vettore delle proporzioni teoriche di elementi per categoria è  
 $pt = p_1, p_2, \dots, p_k$  ( $p_1 + p_2 + \dots + p_k = 1$ )
- Estraiamo campione casuale
- Frequenze per categoria  $z=z_1, z_2, \dots, z_k$

Eseguiamo il test: **`chisq.test(z,p=pt)`**

## Test chi-quadro per le distribuzioni continue non note

Consideriamo popolazione con distribuzione X non nota

Serve per verificare se la popolazione si adatta ad un modello di distribuzione  $X_0$

$$H_0: X=X_0 \text{ e } H_A: X \neq X_0$$

Si usa il **Test di Kolmogorov-Smirnov**: **`ks.test(x, "pchisq"/"pnorm"/"punif", df=n)`**

## Confronto tra le distribuzioni di due popolazioni

- Abbiamo popolazione X e popolazione Y

Serve per stabilire se le distribuzioni delle due popolazioni sono uguali oppure diverse.

$$H_0: X=Y \text{ e } H_A: X \neq Y$$

- Estraiamo i campioni x e y

Eseguiamo il test: ***ks.test(x,y)***

## Verifica della normalità di una popolazione

- Abbiamo popolazione X

Serve per stabilire se la distribuzione di probabilità della popolazione X è normale

$$H_0: X \text{ ha distribuzione normale; } H_A: X \text{ non ha distribuzione normale}$$

- Estraiamo un campione casuale x

Eseguiamo il **Test di Shapiro-Wilk**: ***shapiro.test(x)***

## Test ANOVA (Analysis of variance)

È un test di ipotesi sulle medie che utilizza la varianza delle distribuzioni per arrivare al calcolo del p-value.

- Consideriamo n popolazioni indipendenti:  $X_1, X_2, \dots, X_n$
- Assumiamo che in tutte le popolazioni:
  - Distribuzione normale
  - Varianze tutte uguali (tutte la stessa varianza)

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k ; H_A: \text{almeno una } \neq \text{ dalle altre}$$

- Estraiamo i campioni:  $x_1, x_2, x_n$
- Costruiamo una lista con tutti i campioni: **`l<-list(maggio<-may, ...)`**
- Costruiamo il dataframe: **`d<-stack(l)`**
- Guardiamo struttura del dataframe: **`str(d)`**

Eseguiamo il test: **`oneway.test (values~ind, data=d, var.equal=TRUE)`**

Se NON possiamo assumere la normalità della distribuzione?

## Test di Kruskal-Wallis

- Consideriamo n popolazioni indipendenti, tutte con la stessa distribuzione (non necessariamente nota)

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k ; H_A: \text{almeno una } \neq \text{ dalle altre}$$

- Eseguiamo le altre operazioni descritte precedentemente

Eseguiamo il test: **`kruskal.test (values~ind, data=d)`**