

Brought to you by:

ASTRA

STATISTICA

2° ANNO CLEAM / CLEF

Written by
Matteo Cordaro

2022-2023 Edition

Find more at:

astrabocconi.it

This handout has no intention of substituting University material for what concerns exams preparation, as this is only additional material that does not grant in any way a preparation as exhaustive as the ones proposed by the University.

Questa dispensa non ha come scopo quello di sostituire il materiale di preparazione per gli esami fornito dall'Università, in quanto è pensato come materiale aggiuntivo che non garantisce una preparazione esaustiva tanto quanto il materiale consigliato dall'Università.

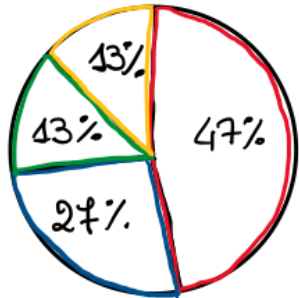
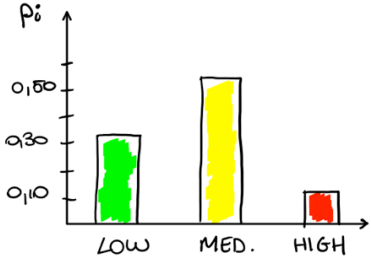
**Ciao! Sono Matteo, studente CLEAM del terzo anno.
Questa dispensa è innanzi tutto il mio metodo di studio ed in
questo modo preparo tutti i miei esami: spero che possa essere
utile anche a voi per preparare i vostri.
Vi segnalo che sono disponibile per ripetizioni nelle materie di cui
realizzo le dispense.
Se siete interessati, qui trovate il mio Instagram:
[**@_matteocordaro.**](#)**

Buono studio!



CONCETTI DI BASE DELLA STATISTICA	
Terminologia di base	<ul style="list-style-type: none"> • Popolazione (N) → numero di unità statistiche nella popolazione • Parametri → caratteristiche della popolazione • Campione (n) → sottoinsieme di unità statistiche selezionate dalla popolazione la cui <i>ampiezza</i> è inferiore rispetto a quella della popolazione ($n \ll N$) • Statistica → caratteristiche del campione • Variabilità campionaria → il valore di una statistica dipende dal campione selezionato; essa assume un valore diverso/varia da campione a campione → incertezza nelle conclusioni/decisioni • Statistica descrittiva: banca della statistica che consiste in un insieme di strumenti (tabelle, grafici, indici, misure di sintesi) volti a sintetizzare ed organizzare i dati grezzi (dataset) al fine di farne emergere caratteristiche rilevanti per il fenomeno e/o il problema decisionale affrontato • Statistica inferenziale: consiste in un insieme di logiche e procedure (stimatori, intervalli di confidenza, verifiche di ipotesi) volte, in primo luogo, ad estendere e generalizzare le conclusioni tratte dal campione all'intera popolazione e, in secondo luogo, a misurare e controllare l'incertezza nelle conclusioni, dovuta alla variabilità campionaria, ovvero l'errore statistico • Variabile statistica: aspetto, caratteristica delle unità statistiche della popolazione • Modalità: uno dei possibili valori (non necessariamente numerico) che può assumere una variabile • Tipi di variabili statistiche: <ul style="list-style-type: none"> ○ Qualitative/Categoriche → esprimono una qualità posseduta dall'unità, l'appartenenza ad una categoria; modalità: nomi, aggettivi <ul style="list-style-type: none"> ▪ Ordinali: modalità ordinabili in modo oggettivo (<i>livello di gradimento alto, medio, basso</i>) ▪ Nominali: non c'è un ordinamento logico ○ Quantitative/Numeriche → frutto della rilevazione di una quantità sull'unità; modalità: valori numerici <ul style="list-style-type: none"> ▪ Discrete → con poche modalità distinte e derivante da un processo di conteggio → <i>numero di esami superati, numero di componenti della famiglia</i> ▪ Continue → con molte modalità distinte (spesso le quantità monetarie) → qualsiasi numero reale in un intervallo → <i>tempo impiegato per arrivare all'università</i> • Scale di misurazione (da aggiungere se richiesto il tipo di variabile): <ul style="list-style-type: none"> ○ Scala nominale ○ Scala ordinale ○ Scala a livello di intervallo → importanza delle distanze tra le modalità: l'origine è arbitraria (<i>zero in senso arbitrario</i>) ○ Scala a livello di rapporto → importanza delle distanze e dei rapporti tra le modalità: l'origine è fissa (<i>zero in senso assoluto</i>)
STATISTICA DESCRITTIVA: ANALISI UNIVARIATA DI VARIABILI CATEGORICHE E NUMERICHE DISCRETE	
Tipologie di analisi	<ol style="list-style-type: none"> 1. Analisi univariata → focus su singola variabile 2. Analisi bivariata → focus su una coppia di variabili (relazione tra variabili) 3. Analisi multivariata → focus su un gruppo di variabili (relazione di dipendenza/interdipendenza in un gruppo di variabili)

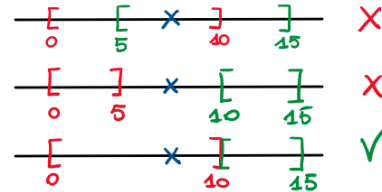


<p>Distribuzione di frequenza</p>	<p>Sia X una variabile statistica, e siano $x_1, x_2, x_i, \dots, x_k$ le k modalità distinte di X osservate negli n dati ($k < n$). Si chiama distribuzione di frequenza una tabella che permette di organizzare i dati a seconda delle modalità</p> <p>Si definiscono, $\forall i = 1, \dots, k$:</p> <ol style="list-style-type: none"> FREQUENZA ASSOLUTA della i-esima modalità, che indichiamo con f_i, il numero di volte per cui x_i si ripete, o alla stessa maniera, il numero di unità nei dati che presentano $X = x_i$ FREQUENZA RELATIVA della i-esima modalità, che indichiamo con p_i la proporzione di unità per cui x_i si ripete nei dati: $p_i = Fr\{X = x_i\} = \frac{f_i}{n}$ <p>Chiameremo FREQUENZA RELATIVA PERCENTUALE della i-esima modalità: $p_i\% = 100 * p_i$</p>																						
<p>Frequenza cumulata</p>	<p>FREQUENZA RELATIVA CUMULATA DI X_i</p> $F_i = Fr\{X \leq x_i\} = \sum_{j=1}^i p_j = p_1 + p_2 + \dots + p_i$ <p>Ed è, quindi, il peso congiunto delle prime i modalità</p>																						
<p>Proprietà delle frequenze</p>	<p>FREQUENZE ASSOLUTE</p> <ul style="list-style-type: none"> $0 \leq f_i \leq n$ $\sum_{i=1}^k f_i = f_1 + \dots + f_i + \dots + f_k = n$ <p>FREQUENZE RELATIVE</p> <ul style="list-style-type: none"> $0 \leq p_i \leq 1$ $\sum_{i=1}^k p_i = p_1 + \dots + p_i + \dots + p_k = 1$ <p>FREQUENZE RELATIVE CUMULATE (SU SCALA ORDINALE O SUPERIORE)</p> <ul style="list-style-type: none"> $0 \leq F_i \leq 1$ $F_0 = 0, F_k = 1$ $F_{i-1} \leq F_i \rightarrow$ NON DECRESCENTI $p_i = F_i - F_{i-1} \rightarrow$ dalle frequenze cumulate alle frequenze relative 																						
<p>Tabella di distribuzione delle frequenze</p>	<table border="1" style="width: 100%; border-collapse: collapse; text-align: center;"> <thead> <tr> <th style="padding: 5px;">MODALITÀ X_i</th> <th style="padding: 5px;">FREQUENZA ASSOLUTA</th> <th style="padding: 5px;">FREQUENZA RELATIVA</th> <th style="padding: 5px;">FREQUENZA CUMULATA</th> </tr> </thead> <tbody> <tr> <td style="padding: 5px;">[a, b)</td> <td style="padding: 5px;">f_1</td> <td style="padding: 5px;">p_1</td> <td style="padding: 5px;">p_1</td> </tr> <tr> <td style="padding: 5px;">[c, d)</td> <td style="padding: 5px;">f_2</td> <td style="padding: 5px;">p_2</td> <td style="padding: 5px;">$p_1 + p_2$</td> </tr> <tr> <td style="padding: 5px;">...</td> <td style="padding: 5px;">...</td> <td style="padding: 5px;">...</td> <td style="padding: 5px;">$p_1 + p_2 + \dots$</td> </tr> <tr> <td style="padding: 5px;">[m, k]</td> <td style="padding: 5px;">f_k</td> <td style="padding: 5px;">p_k</td> <td style="padding: 5px;">$p_1 + p_2 + \dots + p_k = 1$</td> </tr> </tbody> </table>			MODALITÀ X_i	FREQUENZA ASSOLUTA	FREQUENZA RELATIVA	FREQUENZA CUMULATA	[a, b)	f_1	p_1	p_1	[c, d)	f_2	p_2	$p_1 + p_2$	$p_1 + p_2 + \dots$	[m, k]	f_k	p_k	$p_1 + p_2 + \dots + p_k = 1$
MODALITÀ X_i	FREQUENZA ASSOLUTA	FREQUENZA RELATIVA	FREQUENZA CUMULATA																				
[a, b)	f_1	p_1	p_1																				
[c, d)	f_2	p_2	$p_1 + p_2$																				
...	$p_1 + p_2 + \dots$																				
[m, k]	f_k	p_k	$p_1 + p_2 + \dots + p_k = 1$																				
<p>Rappresentazioni grafiche</p>	<p>DIAGRAMMA A TORTA</p>	<p>Per variabili qualitative nominali con al più 5 modalità distinte</p> <p>Calcolo dei gradi degli angoli al centro</p> $\% : 100 = x : 360$																					
	<p>DIAGRAMMA A BARRE</p>	<p>Per variabili qualitative ordinali.</p> <p>Sull'asse verticale sono riportate le frequenze; sull'asse orizzontale le modalità.</p> <p>Le barre sono equidistante e non è orientato</p>																					

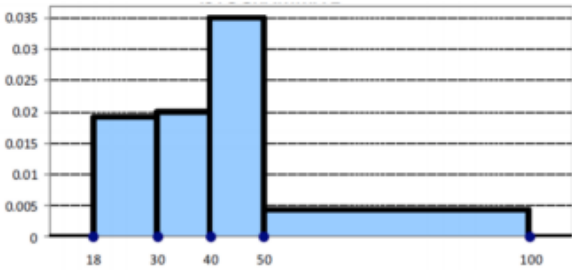
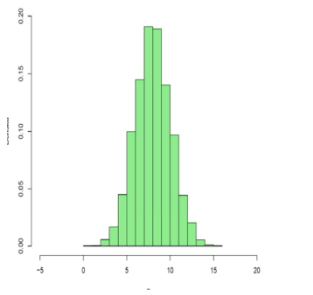
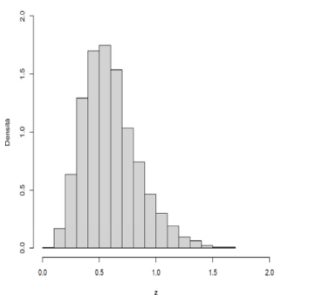
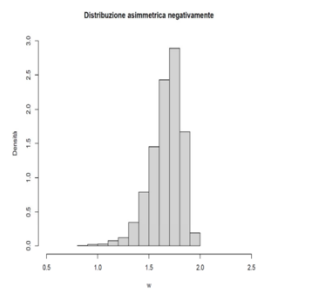
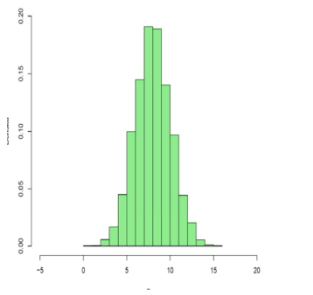
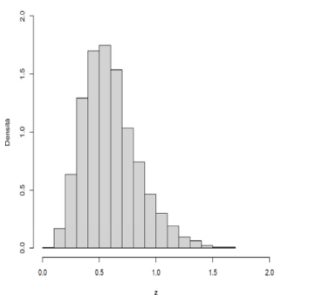
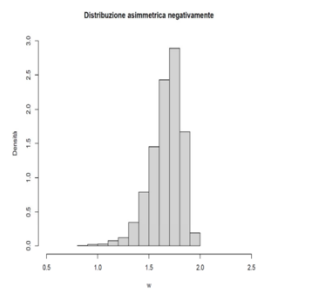
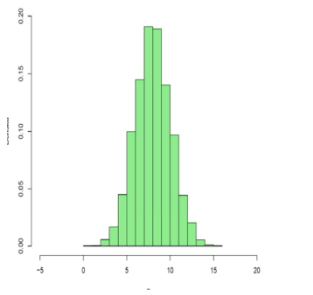
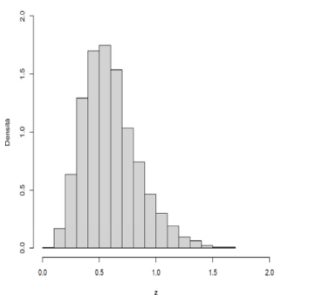
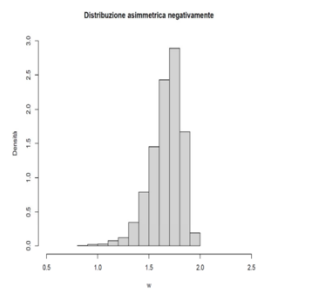


	<p>DIAGRAMMA AD ASTE</p>	<p>Per variabile quantitativa discreta. L'asta ha altezza proporzionale alla frequenza relativa, osservabile sull'asse verticale. Sull'asse orizzontale, orientato, vengono riportate i valori numerici osservati e riportati con la loro posizione corretta. Le aste sono più sottili delle barre così da permettere di mantenere la distanza tra unità</p>	
	<p>CURVA DELLE FREQUENZE CUMULATE</p>	<p>Permette di rappresentare la funzione cumulata di una variabile quantitativa discreta. È un grafico non decrescente, continuo da destra, con limite inferiore pari a 0 e limite superiore di 1 Rappresenta la relativa funzione di ripartizione La funzione è costante a tratti con punti di salto corrispondenti al valore della variabile e ampiezza di salto corrispondente alla frequenza relativa nel punto di salto</p>	
<p>Funzione della curva delle frequenze cumulate</p>	<p>È la funzione di ripartizione di equazione generica:</p> $F(x) = \begin{cases} 0 & x < x_1 \\ F_i & x_i \leq x < x_{i+1} \\ 1 & x \geq x_k \end{cases}$		
<p>Calcolo delle frequenze cumulate dalla funzione di ripartizione</p>	<ul style="list-style-type: none"> • $Fr\{X \leq a\} = F(a)$ • $Fr\{X < a\} \rightarrow$ somma le frequenze relative fino ad a escluso • $Fr\{X > a\} = 1 - Fr\{X \leq a\} = 1 - F(a)$ • $Fr\{a < x \leq b\} = F(b) - F(a) = p(a) + \dots + p(b)$ • $Fr\{a \leq x \leq b\} \rightarrow$ somme le frequenze relative da a fino a b incluso • N.B.: $Fr\{a \leq x \leq b\} \neq Fr\{a < x \leq b\}$ • In ogni caso, posso sempre sommare le frequenze relative; attento nell'uso della differenza con la funzione delle frequenze cumulate! 		
<p>STATISTICA DESCRITTIVA: ANALISI UNIVARIATA DI VARIABILI DISCRETE CONTINUE</p>			

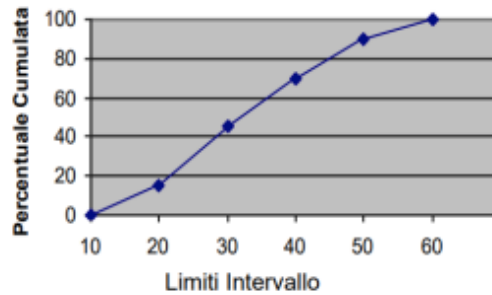


<p>Classi e Frequenza</p>	<ul style="list-style-type: none"> • Classi/Intervalli di valori → poiché le variabili rilevate sono molte, le modalità vengono raggruppate in range di valori, pur perdendo parte del dettaglio per favorire la leggibilità • Frequenza → a ciascun intervallo viene corrisposta la sua frequenza assoluta (= numero di unità per le quali il carattere in esame ha una modalità che cade nell'intervallo) e relativa (= la frequenza assoluta divisa per la dimensione del campione) 										
<p>Criteri di classificazione</p>	<p>Esistono due criteri da rispettare per definire le classi:</p> <ol style="list-style-type: none"> 1. DISGIUNTE → ogni valore deve appartenere al più ad una sola classe (no doppi conteggi) 2. ESAUSTIVE → ogni valore deve appartenere ad almeno una classe (no mancati conteggi) 										
<p>Notazione</p>	<ul style="list-style-type: none"> • k → numero di classi • $x_0, x_1, \dots, x_i, \dots, x_k$ → estremi delle classi 										
<p>Definizione degli estremi (ampiezza costante)</p>	<p>Classi di uguale ampiezza</p> <ul style="list-style-type: none"> • Scelta di k, numero di classi dell'intervallo • Ampiezza costante a tutte le classi: $w = \frac{\max - \min}{k} \rightarrow \text{arrotondata per eccesso}$ <p>Dove \min → minimo valore osservato; \max → massimo valore osservato</p> <ul style="list-style-type: none"> • Determinazione degli estremi: <ul style="list-style-type: none"> ○ $x_0 = \min$ ○ $x_1 = \min + w$ ○ $x_i = \min + i * w$ ○ $x_k = x_{k-1} + w = \min + k * w > \max$ 										
<p>Distribuzione di frequenza in classi</p>	<ul style="list-style-type: none"> • f_i → frequenza assoluta della i-esima classe = numero di unità con valore X nell'intervallo $[x_{i-1}, x_i)$ • $p_i = F\{X \in [x_{i-1}, x_i)\} = \frac{f_i}{n}$ → frequenza relativa della i-esima classe = proporzione di unità con valore X nell'intervallo $[x_{i-1}, x_i)$ • $F_i = Fr\{X < x_i\} = \sum_{j=1}^i p_j$ → frequenza cumulata = peso congiunto delle prime i classi 										
<p>Tabella</p>	<table border="1" style="width: 100%; text-align: center;"> <tr> <td>Intervallo</td> <td>f_i</td> <td>p_i</td> <td>c_i</td> <td>w_i</td> </tr> <tr> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> </tr> </table>	Intervallo	f_i	p_i	c_i	w_i
Intervallo	f_i	p_i	c_i	w_i							
...							
<p>Istogramma</p>	<p>Serie di rettangoli contigui, tali che l'i-esimo rettangolo ha:</p> <ul style="list-style-type: none"> • Come base, l'i-esima classe $[x_{i-1}, x_i) = w_i$ • Area pari alla frequenza relativa della i-esima classe p_i • Altezza pari alla densità di frequenza c_i $c_i = \frac{p_i}{w_i}$ <p>DENSITÀ DI FRQUENZA = frequenza relativa contenuta in ogni intervallo unitario nell'i-esima classe</p>										



	 <p>Come è possibile osservare, l'istogramma può avere anche intervalli di diversa ampiezza per aumentare il grado di dettaglio. Il fatto che sull'asse verticale vengano usata la densità di frequenza, al posto delle <i>frequenze assolute</i>, risiede nella necessità di non "falsare" l'analisi dal punto di vista grafico L'unico caso in cui la rappresentazione con frequenze assolute equivale a quella con densità di frequenza e il caso in cui le ampiezze delle classi siano uguali</p>									
<p>Calcolo delle frequenze a partire dall'istogramma</p>	<p>L'area descritta dall'istogramma lungo un certo intervallo descrive la frequenza relativa per quell'intervallo</p>									
<p>Forma della distribuzione</p>	<p>Gli istogrammi, a seconda della forma assunta dalla distribuzione, si classificano in: Ricorda: la simmetria si valuta osservando la densità ed osservando globalmente tutti gli intervalli</p> <table border="1" data-bbox="402 1059 1417 1435"> <thead> <tr> <th data-bbox="402 1059 740 1133">DISTRIBUZIONE SIMMETRICA</th> <th colspan="2" data-bbox="740 1059 1417 1099">DISTRIBUZIONE ASIMMETRICA</th> </tr> <tr> <td></td> <th data-bbox="740 1099 1075 1133">POSITIVAMENTE</th> <th data-bbox="1075 1099 1417 1133">NEGATIVAMENTE</th> </tr> </thead> <tbody> <tr> <td data-bbox="402 1133 740 1435">  </td> <td data-bbox="740 1133 1075 1435">  </td> <td data-bbox="1075 1133 1417 1435">  </td> </tr> </tbody> </table>	DISTRIBUZIONE SIMMETRICA	DISTRIBUZIONE ASIMMETRICA			POSITIVAMENTE	NEGATIVAMENTE			
DISTRIBUZIONE SIMMETRICA	DISTRIBUZIONE ASIMMETRICA									
	POSITIVAMENTE	NEGATIVAMENTE								
										
<p>Curva delle frequenze cumulate</p>	<p>È la funzione che, per ogni x sull'asse, ha espressione</p> $F(x) = Fr\{X \leq x\} = \text{Frequenza relativa di valori osservati inferiori ad } x$ <p>Proprietà</p> <ul style="list-style-type: none"> • $0 \leq F(x) \leq 1$ • $F(x) = 0$ se $x < x_0$ • $F(x) = 1$ se $x \geq 1$ • $F(x) \leq F(y)$ se $x < y \rightarrow$ non decrescente • È una funzione lineare a tratti • $Fr(a < X \leq b) = F(b) - F(a)$ 									
<p>Ogiva</p>	<p>Per variabili quantitative continue, è la rappresentazione della funzione cumulata È un grafico non decrescente, continuo da destra con limite inferiore pari a 0 e limite superiore pari ad 1 La pendenza di ogni tratto corrisponde alla densità della classe</p>									





→ Tipologia di esercizi: può essere chiesto il valore di x per cui la frequenza cumulata sia pari ad un valore a . Per farlo occorre:

- Individuare l'intervallo in cui cade il valore di x richiesto
- Calcolare la retta, parte dell'ogiva, passante tra i due punti dell'intervallo: prima calcolare il coefficiente angolare; poi imporre il passaggio per uno dei due estremi; esplicitare l'equazione della retta
- Sostituire il valore a al posto della y
- Risolvere l'equazione in x

STATISTICA DESCRITTIVA: ANLISI BIVARIATA DI VARIABILI CATEGORICHE

Tabella di contingenza e frequenza assoluta congiunta

X \ Y	Y_1	...	y_j	...	y_c	TOTALE
x_1	f_{11}	...	f_{1j}	...	f_{c1}	R_1
...
x_i	f_{i1}	...	f_{ij}	...	f_{ic}	R_i
...
x_r	f_{r1}	...	f_{rj}	...	f_{rc}	R_r
TOTALE	C_1	C_c	n

Con:

- r = # modalità distinte di X
- c = # modalità distinte di Y
- i = indice per le modalità di X righe
- j = indice per le modalità di Y colonne

La tabella di contingenza è ottenuta:

- Classificando le unità statistiche sulla base della coppia di modalità per variabili X e Y che esibiscono (x_i, y_j)
- **DISTRIBUZIONE CONGIUNTA.** Contando il numero di unità che esibiscono ciascuna coppia, e determinando così per ogni $i = 1, \dots, r$ e per ogni $j = 1, \dots, c$

f_{ij} = FREQUENZA ASSOLUTA CONGIUNTA della coppia (x_i, y_j)

Essa rappresenta il numero di unità per cui congiuntamente $X = x_i$ e $Y = y_j$

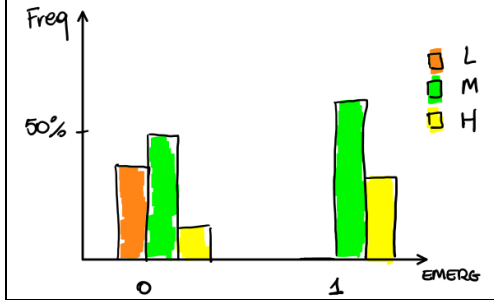
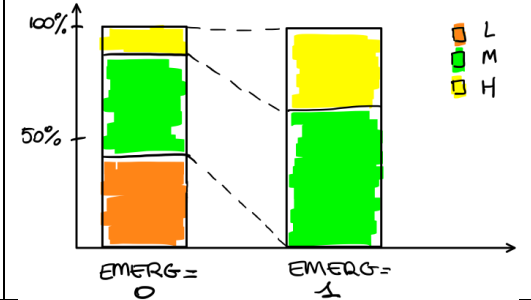
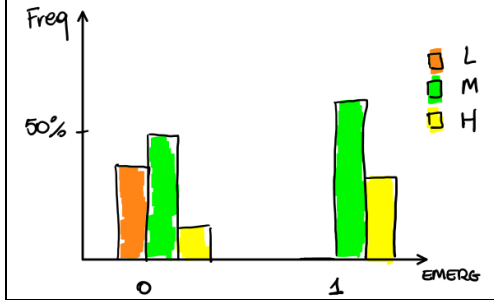
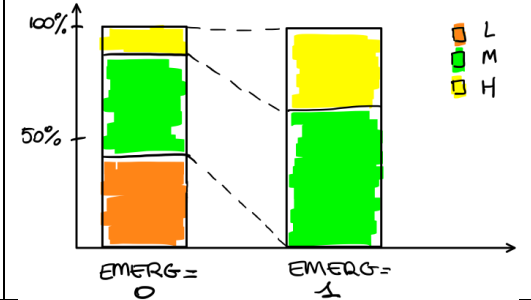
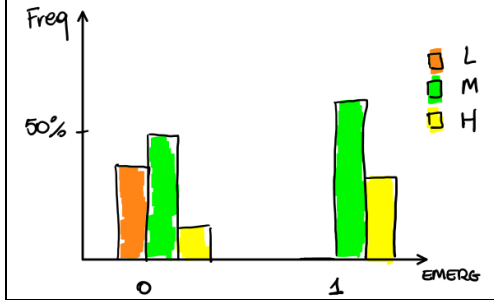
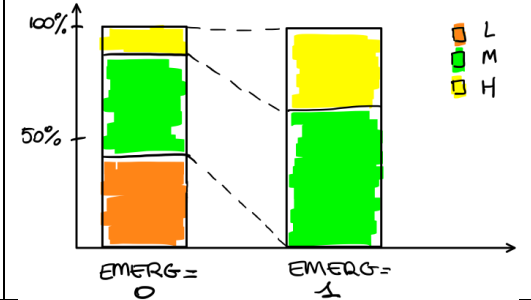
PROPRIETÀ

- $0 \leq f_{ij} \leq n$
- $\sum_{ij} f_{ij} = \sum_{i=1}^r \sum_{j=1}^c f_{ij}$


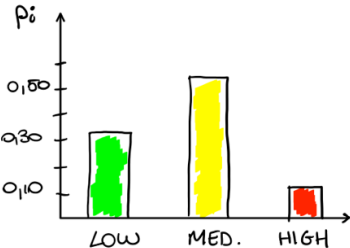
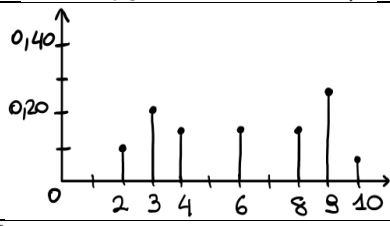
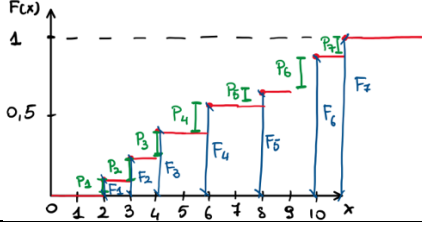
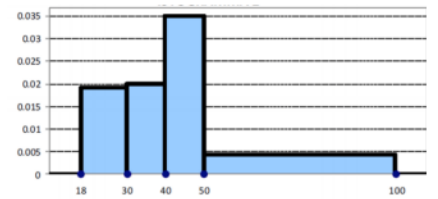


	<ul style="list-style-type: none"> • DISTRIBUZIONI MARGINALI. Nell'ultima colonna e nell'ultima riportiamo le distribuzioni marginali, così dette poiché a margine della tabella <ul style="list-style-type: none"> ○ $R_i = \text{TOTALE DELLA } i - \text{ESIMA RIGA} = \text{\#unità che presentano } X = x_i \rightarrow \text{Frequenza Assoluta della modalità } x_i$ ○ $C_j = \text{TOTALE DELLA } j - \text{ESIMA COLONNA} = \text{\#unità che presentano } Y = Y_j \rightarrow \text{Frequenza Assoluta della modalità } Y_j$ <p>PROPRIETÀ</p> <ul style="list-style-type: none"> ○ $0 \leq R_i = f_{i1} + f_{i2} + \dots + f_{ic} \leq n$ ○ $0 < C_j = c_{1j} + c_{2j} + \dots + c_{ij} \leq n$ <p>Quindi, dalla tabella di contingenza riesco a trovare le distribuzioni univariate di X e di Y</p>																																																	
<p>Frequenza relativa congiunta</p>	<p>Esiste anche la possibilità di calcolare le FREQUENZE RELATIVE CONGIUNTE: basta dividere le frequenze assolute congiunte per il totale n</p> $p_{ij} = \frac{f_{ij}}{n} = Fr\{X = x_i, Y = y_j\}$ <p>In questo caso, si avrà che $n = 1$</p> <table border="1" data-bbox="408 927 1414 1279"> <thead> <tr> <th>X \ Y</th> <th>y_1</th> <th>...</th> <th>y_j</th> <th>...</th> <th>y_c</th> <th>TOTALE</th> </tr> </thead> <tbody> <tr> <td>x_1</td> <td>p_{11}</td> <td>...</td> <td>p_{1j}</td> <td>...</td> <td>p_{1c}</td> <td>R_1/n</td> </tr> <tr> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> </tr> <tr> <td>x_i</td> <td>p_{i1}</td> <td>...</td> <td>$p_{ij} = f_{ij}/n$</td> <td>...</td> <td>p_{ic}</td> <td>R_i/n</td> </tr> <tr> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> </tr> <tr> <td>x_r</td> <td>p_{r1}</td> <td>...</td> <td>p_{rj}</td> <td>...</td> <td>p_{rc}</td> <td>R_r/n</td> </tr> <tr> <td>TOTALE</td> <td>C_1/n</td> <td>...</td> <td>C_j/n</td> <td>...</td> <td>C_c/n</td> <td>1</td> </tr> </tbody> </table>	X \ Y	y_1	...	y_j	...	y_c	TOTALE	x_1	p_{11}	...	p_{1j}	...	p_{1c}	R_1/n	x_i	p_{i1}	...	$p_{ij} = f_{ij}/n$...	p_{ic}	R_i/n	x_r	p_{r1}	...	p_{rj}	...	p_{rc}	R_r/n	TOTALE	C_1/n	...	C_j/n	...	C_c/n	1
X \ Y	y_1	...	y_j	...	y_c	TOTALE																																												
x_1	p_{11}	...	p_{1j}	...	p_{1c}	R_1/n																																												
...																																												
x_i	p_{i1}	...	$p_{ij} = f_{ij}/n$...	p_{ic}	R_i/n																																												
...																																												
x_r	p_{r1}	...	p_{rj}	...	p_{rc}	R_r/n																																												
TOTALE	C_1/n	...	C_j/n	...	C_c/n	1																																												
<p>Frequenza subordinata (o condizionata) e distribuzioni subordinate di frequenze</p>	<p>Ci permettono di calcolare le relative frequenze percentuali per ogni modalità della variabile:</p> <p>1. FREQUENZA SUBORDINATA DI Y DATA X</p> $\frac{f_{ij}}{R_i} = Fr\{Y = y_j X = x_i\}$ <p>Ovvero, la proporzione di unità che presentano $Y = y_j$ nel sottogruppo/categoria di quelle con $X = x_i$</p> <table border="1" data-bbox="435 1641 1409 1984"> <thead> <tr> <th>X \ Y</th> <th>y_1</th> <th>...</th> <th>y_j</th> <th>...</th> <th>y_c</th> <th>TOTALE</th> </tr> </thead> <tbody> <tr> <td>x_1</td> <td>f_{11}/R_1</td> <td>...</td> <td>f_{1j}/R_1</td> <td>...</td> <td>f_{1c}/R_1</td> <td>1</td> </tr> <tr> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>1</td> </tr> <tr> <td>x_i</td> <td>f_{i1}/R_i</td> <td>...</td> <td>f_{ij}/R_i</td> <td>...</td> <td>f_{ic}/R_i</td> <td>1</td> </tr> <tr> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> </tr> <tr> <td>x_r</td> <td>f_{r1}/R_r</td> <td>...</td> <td>f_{rj}/R_r</td> <td>...</td> <td>f_{rc}/R_r</td> <td>1</td> </tr> <tr> <td>TOTALE</td> <td>—</td> <td>—</td> <td>—</td> <td>—</td> <td>—</td> <td>1</td> </tr> </tbody> </table>	X \ Y	y_1	...	y_j	...	y_c	TOTALE	x_1	f_{11}/R_1	...	f_{1j}/R_1	...	f_{1c}/R_1	1	1	x_i	f_{i1}/R_i	...	f_{ij}/R_i	...	f_{ic}/R_i	1	x_r	f_{r1}/R_r	...	f_{rj}/R_r	...	f_{rc}/R_r	1	TOTALE	—	—	—	—	—	1
X \ Y	y_1	...	y_j	...	y_c	TOTALE																																												
x_1	f_{11}/R_1	...	f_{1j}/R_1	...	f_{1c}/R_1	1																																												
...	1																																												
x_i	f_{i1}/R_i	...	f_{ij}/R_i	...	f_{ic}/R_i	1																																												
...																																												
x_r	f_{r1}/R_r	...	f_{rj}/R_r	...	f_{rc}/R_r	1																																												
TOTALE	—	—	—	—	—	1																																												

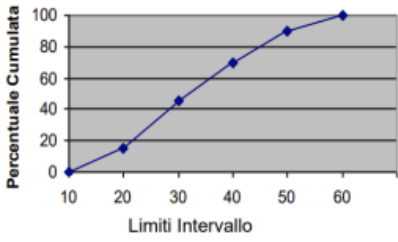
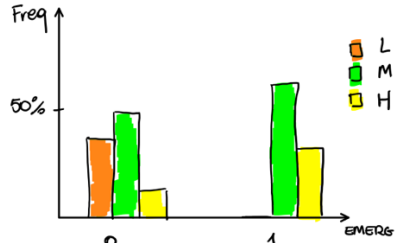
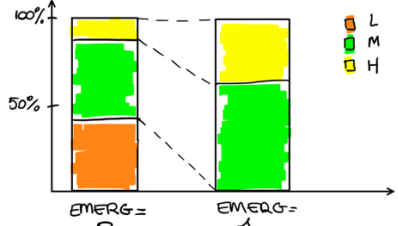


	<p>Abbiamo eseguito una normalizzazione per riga: fatto 100 la singola x_i abbiamo analizzato il contributo delle singole y_i</p> <p>2. FREQUENZA SUBORDINATA DI X DATA Y</p> $\frac{f_{ij}}{C_j} = Fr\{X = x_i Y = y_j\}$ <table border="1" data-bbox="432 481 1396 824"> <thead> <tr> <th>X \ Y</th> <th>y_1</th> <th>...</th> <th>y_j</th> <th>...</th> <th>y_c</th> <th>TOTALE</th> </tr> </thead> <tbody> <tr> <td>x_1</td> <td>f_{11}/C_1</td> <td>...</td> <td>f_{1j}/C_j</td> <td>...</td> <td>f_{1c}/C_c</td> <td>—</td> </tr> <tr> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>—</td> </tr> <tr> <td>x_i</td> <td>f_{i1}/C_1</td> <td>...</td> <td>f_{ij}/C_j</td> <td>...</td> <td>f_{ic}/C_c</td> <td>—</td> </tr> <tr> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>—</td> </tr> <tr> <td>x_r</td> <td>f_{r1}/C_1</td> <td>...</td> <td>f_{rj}/C_j</td> <td>...</td> <td>f_{rc}/C_c</td> <td>—</td> </tr> <tr> <td>TOTALE</td> <td>1</td> <td>...</td> <td>1</td> <td>...</td> <td>1</td> <td>—</td> </tr> </tbody> </table> <p>Abbiamo eseguito una normalizzazione per colonna: fatto 100 la singola y_i, abbiamo analizzato il contributo della singola x_i</p>	X \ Y	y_1	...	y_j	...	y_c	TOTALE	x_1	f_{11}/C_1	...	f_{1j}/C_j	...	f_{1c}/C_c	—	—	x_i	f_{i1}/C_1	...	f_{ij}/C_j	...	f_{ic}/C_c	—	—	x_r	f_{r1}/C_1	...	f_{rj}/C_j	...	f_{rc}/C_c	—	TOTALE	1	...	1	...	1	—
X \ Y	y_1	...	y_j	...	y_c	TOTALE																																												
x_1	f_{11}/C_1	...	f_{1j}/C_j	...	f_{1c}/C_c	—																																												
...	—																																												
x_i	f_{i1}/C_1	...	f_{ij}/C_j	...	f_{ic}/C_c	—																																												
...	—																																												
x_r	f_{r1}/C_1	...	f_{rj}/C_j	...	f_{rc}/C_c	—																																												
TOTALE	1	...	1	...	1	—																																												
<p>Dipendenza ed indipendenza statistica</p>	<ul style="list-style-type: none"> Le variabili X ed Y si dicono statisticamente indipendenti se le distribuzioni di frequenza subordinata di Y data X (o, equivalentemente di X data Y) sono tutte uguali tra loro → Ciò significa che, con riguardo alla variabile Y, le unità si comportano allo stesso modo, all'interno di ogni sottogruppo o sottopopolazione definito da una modalità di X Se ciò non accade, ovvero se almeno due frequenze subordinate sono diverse tra loro, allora X ed Y si dicono statisticamente dipendenti 																																																	
<p>Diagramma a barre accostate e diagramma a barre sovrapposte</p>	<p>Per rappresentare graficamente un'analisi descrittiva bivariata possiamo utilizzare alternativamente:</p> <table border="1" data-bbox="400 1346 1428 1861"> <thead> <tr> <th>Diagramma a barre accostate</th> <th>Diagramma a barre sovrapposte</th> </tr> </thead> <tbody> <tr> <td>Asse X → modalità della variabile presa in considerazione come condizionale Asse Y → frequenza subordinata della variabile condizionata</td> <td>Fatta cento la frequenza subordinata, viene impiegata una sola colonna ripartita nelle diverse modalità, espresse in percentuale</td> </tr> <tr> <td>  </td> <td>  </td> </tr> </tbody> </table>	Diagramma a barre accostate	Diagramma a barre sovrapposte	Asse X → modalità della variabile presa in considerazione come condizionale Asse Y → frequenza subordinata della variabile condizionata	Fatta cento la frequenza subordinata, viene impiegata una sola colonna ripartita nelle diverse modalità, espresse in percentuale																																													
Diagramma a barre accostate	Diagramma a barre sovrapposte																																																	
Asse X → modalità della variabile presa in considerazione come condizionale Asse Y → frequenza subordinata della variabile condizionata	Fatta cento la frequenza subordinata, viene impiegata una sola colonna ripartita nelle diverse modalità, espresse in percentuale																																																	
																																																		
<p>Paradosso di Simpson</p>	<p>Il paradosso di Simpson indica una situazione in cui una relazione tra due fenomeni appare modificata, o perfino invertita, a causa di variabili non presi in considerazione nell'analisi iniziale</p>																																																	



	<p>Il paradosso di Simpson permette dunque, a certe condizioni, il verificarsi di situazioni in cui il comportamento di sottogruppi è diverso dal comportamento complessivo.</p> <p>Il paradosso di Simpson avviene quando non viene inclusa nell'analisi bivariata una variabile essenziale, creando risultati fuorvianti, dettati da una errata analisi delle frequenze.</p>		
<p>Calcolo di frequenze</p>	<p>Dato un campione di n elementi, prendiamo due variabili A e B per le quali abbiamo costruito una tabella di contingenza. Prestiamo attenzione alle richieste fatte:</p> <ul style="list-style-type: none"> • Percentuale di campione che è A e B $\rightarrow Fr\{A \cap B\} = \frac{A \cap B}{n}$ • Percentuale di campione che è A tra quelli che sono B $\rightarrow Fr\{A B\}$ • Supponendo che B < A, percentuale di campione che in totale ha fatto almeno A $\rightarrow Fr\{X \leq A\} = \frac{A+B}{n}$ 		
<p>Grafici appropriati per tipologia di variabile</p>	<p>Analisi</p>	<p>Variabile</p>	<p>Grafico</p>
		<p>Qualitativa nominale</p>	<p>Grafico a torta</p> 
		<p>Qualitativa ordinale</p>	<p>Diagramma a barre</p> 
	<p>Descrittiva univariata</p>	<p>Quantitativa discreta</p>	<p>Diagramma ad aste</p> 
			<p>Curva delle frequenze cumulate</p> 
	<p>Quantitativa continua</p>	<p>Istogramma</p> 	



			Ogiva	
Descrittiva bivariata	Condizionate	Diagramma a barre accostate		
		Diagramma a barre sovrapposte		

MISURA DI TENDENZA CENTRALE

Definizione
 Sono dei **valori assunti dalle modalità che vengono calcolati con lo scopo di sintetizzare la distribuzione di una variabile mediante un'unica modalità**. Sono indicatori il più possibile rappresentativi della tendenza della variabile. Le misure di tendenza centrale (e non centrale) sono modalità, non frequenze.

Media
 La media aritmetica rappresenta il **baricentro della distribuzione**: le distanze dalla media dei valori alla sua sinistra e di quelli alla sua destra si compensano, dando luogo ad una somma di tali distanze nulla. Si utilizza **solo** per le **variabili quantitative**. È impossibile da calcolare per quelle qualitative, perché, essendo una somma, richiede necessariamente dei **valori numerici**.

PER DATI GREZZI (o DISAGGREGATI)

1) Media aritmetica semplice

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

Proprietà

- La media minimizza la funzione di errore statistico:

$$L(\bar{x}) < L(a) \rightarrow \sum_i (x_i - \bar{x})^2 < \sum_i (x_i - a)^2$$

- La somma delle deviazioni $(x_i - \bar{x})$ della media è nulla

$$\sum_+ (x_i - \bar{x}) + \sum_- (x_i - \bar{x}) = 0$$

- Preserva il totale $\rightarrow n\bar{x} = \sum_i x_i$
- È interna $\rightarrow \min < \bar{x} < \max$
- Media e trasformazioni lineari



	<p style="text-align: center;">$\forall y_i = a + bx_i \text{ con } a \text{ e } b \text{ arbitrari} \rightarrow \bar{y} = a + b\bar{x}$</p> <ul style="list-style-type: none"> • Monotonia: se x_i aumenta/diminuisce anche \bar{x} aumenta/diminuisce $\frac{d\bar{x}}{dx_i} = \frac{1}{n} > 0$ <p>2) Media aritmetica ponderata Dati x_i valori, con w_i pesi, avremo che:</p> $\bar{x}_w = \frac{\sum_{i=1}^k w_i x_i}{\sum_{i=1}^k w_i}$ <p>PER DISTRIBUZIONE DI FREQUENZE</p> <p>1) Variabile numerica discreta Con le frequenze assolute:</p> $\bar{x} = \frac{\sum_{i=1}^k f_i x_i}{\sum_{i=1}^k f_i}$ <p>Con frequenze relative:</p> $\bar{x} = \frac{\sum_{i=1}^k p_i x_i}{\sum_{i=1}^k p_i} \rightarrow \bar{x} = \sum_{i=1}^k p_i x_i$ <p>2) Variabile numerica in classi Riunendo in classi, abbiamo <i>perso</i> i veri valori del dataset. Per procedere al calcolo della media dobbiamo supporre che l'elemento rappresentativo interno di ogni classe, che chiamiamo m_i, ovvero il punto medio della classe</p> $m_i = \frac{x_{i-1} + x_i}{2}$ <p>Da cui:</p> $\bar{x} = \frac{\sum_{i=1}^k f_i m_i}{n} = \sum_{i=1}^k p_i m_i$ <p>N.B.: <i>Proprietà</i> La media preserva il totale: dati n valori, se al posto dei valori x_1, x_2, \dots, x_n, gli n individui presentassero tutti il medesimo valore \bar{x}, media della variabile X, allora il totale del carattere nel gruppo resterebbe invariato</p>
<p>Moda</p>	<p>È la modalità alla quale è associata la frequenza più elevata. La moda non tiene conto dell'aspetto numerico o meno delle modalità, ma si basa esclusivamente sulle loro frequenze</p> <p>Attenzione:</p> <ul style="list-style-type: none"> • Va bene sia per variabili qualitative che quantitative • È la modalità che esibisce la frequenza più elevata, non la frequenza più elevata • Per distribuzioni in classi si impiega la classe modale, ovvero la classe con la densità di frequenza più elevata
<p>Mediana</p>	<p>La mediana degli n (o N) valori osservati della variabile X ($x_1, x_2, \dots, x_i, \dots, x_n$) è un qualunque valore Me, che lascia alla sua sinistra e alla sua destra esattamente</p>



il 50% del totale delle osservazioni; in altri termini, Me deve soddisfare la disuguaglianza:

$$Fr\{X \leq Me\} \leq 0.5 \wedge Fr\{X \geq Me\} \geq 0.50$$

Me è quel valore che *spacca in due la distribuzione*, ovvero 50% dei dati sono a destra e 50% dei dati sono a sinistra

Per dati disaggregati

Denotando con $x_{(1)}, x_{(2)}, \dots, x_{(n)}$ la sequenza degli n valori osservati della variabile X , posti **in ordine crescente**, Me è l'osservazione di posto $0.5 \cdot (n+1)$:

$$Me = \begin{cases} x_{(0.5 \cdot (n+1))} & \text{se } n \text{ dispari} \\ \frac{x_{(0.5 \cdot n)} + x_{(0.5 \cdot (n+1))}}{2} & \text{se } n \text{ pari} \end{cases}$$

La logica sottostante è la seguente:

- Se n è dispari \rightarrow scelgo l'osservazione centrale
- Se n è pari \rightarrow non c'è un'osservazione centrale, bensì un intervallo tra due osservazioni centrali \rightarrow scelgo il valor medio di quell'intervallo

Per calcolarla procedo così:

- Calcolo la posizione della mediana nella sequenza ordinata:

$$Pos(Me) = 0.5(n + 1)$$

- Determino che tale posizione è $x_{(Pos(Me))} = x_{(0.5(n+1))}$
- Se n è dispari, il valore con posizione $Pos(Me)$ è proprio la mediana; Se n è pari, allora $Pos(Me)$ sarà un valore decimale indicante che la mediana si trova tra due valori ed, in particolare, è il punto medio tra questi
- Nel caso in cui la posizione della Me ha decimali 0,50, allora è il punto medio tra i due valori

Per distribuzione di frequenza

- Variabile numerica discreta

Me è il primo valore le cui frequenze cumulate, F_i , raggiungono o superano un valore pari a 0.5

Qualora il valore sia esattamente 0,5 allora occorre eseguire una media tra il valore che registra 0,5 e quello immediatamente successivo

- Variabile numerica continua

Me è il valore (o i valori) che soddisfano l'equazione $F(Me) = 0.50$; in particolare, se $F_{i-1} < 0.5 \leq F_i$ (le frequenze cumulate superano il valore 0.50 nella i -esima classe), Me può essere determinata risolvendo l'equazione:

$$0.5 = F_{i-1} + (Me - x_{i-1}) * c_i$$

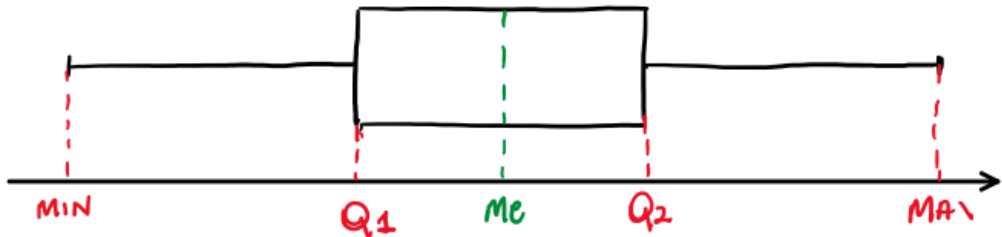
$$Me = x_{i-1} + \frac{0.5 - F_{i-1}}{c_i}$$

Quindi, prima di tutto, identifico la classe di cui fa parte la mediana (la classe che incorpora una frequenza cumulata superiore a 0.5); in seconda battuta applico la formula, ricordando che $x_{(i-1)}$ rappresenta l'estremo inferiore della classe scelta e F_{i-1} la frequenza cumulata della classe precedente



<p>Opportuni indicatori di tendenza centrale</p>	<ul style="list-style-type: none"> • Media: solo per variabili quantitative discrete o continue • Moda: sia per variabili qualitative che quantitative; tuttavia, è un indicatore eccessivamente povero di informazioni • Mediana: sia per variabili quantitative che qualitative
<p>MISURE DI TENDENZA NON CENTRALE</p>	
<p>Definizione</p>	<p>Sono valori che dividono la distribuzione in due parti non uguali e che non giacciono al centro di essa; sono estensioni dell'idea di mediana</p>
<p>Quantile di ordine p</p>	<p>Fissato p, con $0 \leq p \leq 1$, un quantile di ordine p dei valori osservati x_1, x_2, \dots, x_n è un qualunque valore x_p, tale che:</p> $Fr\{X \leq x_p\} \leq p$ $Fr\{X \geq x_p\} \geq 1 - p$ <p>Quindi il P-esimo quantile (o, analogamente, percentile) è quel valore che lascia alla sua sinistra (eventualmente includendo lo stesso valore) approssimativamente il P% di osservazioni</p> <p>Nella sequenza degli n valori osservati della variabile X, posti in ordine crescente, x_p è l'osservazione di posto $p*(n+1)$, ovvero</p> $x_p = \begin{cases} x_{(h)} & \text{con } h \text{ l'intero più vicino a } p*(n+1) \\ \frac{x_{(l)} + x_{(l+1)}}{2} & \text{se } p*(n+1) \text{ ha decimali } 0.5 \\ & l \text{ è } p(n+1) \text{ arrotondato per difetto} \end{cases}$ <p>Procedo così:</p> <ul style="list-style-type: none"> • Calcolo da posizione del quartile $\rightarrow Pos(x_p) = p*(n+1)$ • Arrotondo all'intero più vicini (per eccesso o per difetto); ma attenzione al caso in cui siamo esattamente a decimali 0.5 • Individuo x_p con la formula usata sopra e facendo attenzione ai decimali <p>Quantile di ordine p per distribuzioni di frequenza</p> <ul style="list-style-type: none"> • <u>Variabile numerica discreta</u> x_p è il primo valore le cui frequenze cumulate F_i raggiungono o superano il valore p • <u>Variabile numerica in classi</u> x_p è il valore (o i valori) che soddisfano l'equazione $F(x_p) = p$; in particolare, se $F_{i-1} < p \leq F_i$ (ovvero le frequenze cumulate superano il valore p nell'i-esima classe), x_p può essere determinato così: $p = F_{i-1} + (x_p - x_{i-1}) * c_i$ $x_p = x_{i-1} + \frac{p - F_{i-1}}{c_i}$ <p>Nel caso in cui la classe riporti esattamente l'ordine p del percentile come frequenza cumulata, allora l'estremo superiore è esattamente il percentile ricercato</p> <p>Quartili fondamentali Sono quantili particolari:</p>



	<ol style="list-style-type: none"> 1. Primo quartile o 25esimo percentile (Q_1) $\rightarrow p = 0.25$; $Pos(Q_1)=0.25(n+1)$ 2. Secondo quartile o 50esimo percentile = Mediana ($Q_2=Me$) $\rightarrow p = 0.50$; $Pos(Q_1)=0.50(n+1)$ 3. Terzo quartile o 75esimo percentile (Q_3) $\rightarrow p = 0.75$; $Pos(Q_1)=0.75(n+1)$ <p>Osservazione: Nota che tutte le relazioni possono essere espresse in funzione del primo quartile:</p> $Pos(Q_1) = 0,25(n + 1) \rightarrow Pos(Me) = 2 * 0,25(n + 1) = 2 * Pos(Q_1)$										
<p>Percentili</p>	<p>I percentili sono 99 punti che dividono la distribuzione in 100 parti con la stessa frequenza (pari a 0.01 o 1%)</p> <ul style="list-style-type: none"> • L'n-esimo percentile è il quantile di ordine $p = n\%$ • Si calcolano allo stesso modo dei quartili e i fondamentali sono esattamente quelli già citati • Ha senso calcolarli per variabili quantitative (discrete e continue) e per le qualitative ordinali; non ha senso calcolarli per variabili qualitative nominali <p>Il significato del quantile/percentile è il seguente: Ordinando in senso crescente un dataset, trovare i migliori $x\%$ risultati, vuol dire calcolare il $(1-x)$-esimo percentile, ovvero il quantile di ordine $p = (1-x)$ Quindi se viene fornita una scala di valori crescenti, e viene richiesto di trovare il 70% dei valori migliori, allora dobbiamo calcolare il quartile di ordine 0,30 (cioè, il 30esimo percentile): infatti, tale valore identifica il valore a sinistra del quale si trovano il 30% dei valori totali (quelli più bassi, i peggiori) e, riflessivamente, a destra del quale si trova il 70% dei valori (quelli più alti, i migliori) Una seconda interpretazione ci permette di dire che il valore corrispondente al 30esimo percentile è quel valore che, nel 70% dei casi, viene superato.</p>										
<p>Cinque numeri di sintesi</p>	<table border="1"> <thead> <tr> <th><i>min</i></th> <th>Q_1</th> <th><i>Me</i></th> <th>Q_3</th> <th><i>max</i></th> </tr> </thead> <tbody> <tr> <td>Valore più piccolo rilevato per la variabile campionaria</td> <td>Primo quartile $p = 0,25$</td> <td>Mediana $p = 0,50$</td> <td>Terzo quartile $p = 0,75$</td> <td>Valore più grande rilevato per la variabile campionaria</td> </tr> </tbody> </table>	<i>min</i>	Q_1	<i>Me</i>	Q_3	<i>max</i>	Valore più piccolo rilevato per la variabile campionaria	Primo quartile $p = 0,25$	Mediana $p = 0,50$	Terzo quartile $p = 0,75$	Valore più grande rilevato per la variabile campionaria
<i>min</i>	Q_1	<i>Me</i>	Q_3	<i>max</i>							
Valore più piccolo rilevato per la variabile campionaria	Primo quartile $p = 0,25$	Mediana $p = 0,50$	Terzo quartile $p = 0,75$	Valore più grande rilevato per la variabile campionaria							
<p>Box Plot</p>	<p>Il box plot o <i>diagramma a scatola e baffi</i> è la rappresentazione grafica dei cinque numeri di sintesi</p> 										



Forma della distribuzione	I cinque numeri di sintesi ed il Box Plot consentono una valutazione della simmetria o asimmetria della distribuzione dei dati, mediante i confronti di:	
	Simmetria	Numeri di sintesi
	Distribuzione approssimativamente simmetriche	$\begin{cases} Q_1 - \min \approx \max - Q_3 \\ Me - Q_1 \approx Q_3 - Me \\ Me \approx \bar{x} \end{cases}$
	Distribuzione positivamente asimmetrica (obliqua a destra)	$\begin{cases} Q_1 - \min < \max - Q_3 \\ Me - Q_1 < Q_3 - Me \\ Me < \bar{x} \end{cases}$
	Distribuzione negativamente asimmetrica (obliqua a sinistra)	$\begin{cases} Q_1 - \min > \max - Q_3 \\ Me - Q_1 > Q_3 - Me \\ Me > \bar{x} \end{cases}$
Outlier e loro identificazione	<p>OUTLIER: osservazione per cui il valore di X si discosta notevolmente dal resto dei dati, ovvero dal centro della distribuzione, influenzando significativamente l'analisi</p> <p>Cause:</p> <ul style="list-style-type: none"> • Errori di rilevazione • Errori di trascrizione • Manifestazione raramente osservata del fenomeno <p>Identificazione: Per prima cosa identifichiamo gli estremi oltre i quali cadono gli outlier</p> $\bar{L} = Q_3 + 1,5(Q_3 - Q_1) = Q_1 - 1,5IQR$ $\underline{L} = Q_1 - 1,5(Q_3 - Q_1) = Q_3 - 1,5IQR$ <p>Dove, IQR = Inter Quartile Range Quindi, tutti i valori x_i che ricadono al di fuori dell'intervallo $[\bar{L}, \underline{L}]$ sono considerati outlier</p>	
Robustezza	<p>Una misura di sintesi si dice robusta se non risente in modo significativo dei cambiamenti in una porzione limitata di valori del dataset</p> <ul style="list-style-type: none"> • La MODA e la MEDIANA sono misure ROBUSTE • La MEDIA NON è ROBUSTA 	
MISURE DI VARIABILITÀ E DISPERSIONE PER VARIABILI QUANTITATIVE UNIVARIATE		
Campo di variazione (Range)	<p>È l'ampiezza dell'intervallo che contiene il 100% dei dati:</p> $R = \max - \min$ <ul style="list-style-type: none"> • È una misura non sensibile in quanto, al variare della variabilità, il range può non cambiare 	



	<ul style="list-style-type: none"> È una misura non robusta poiché funzione unicamente del max e del min 								
Differenza interquartile (IQR)	<p>È l'ampiezza dell'intervallo che contiene il 50% centrale dei dati:</p> $IQR = Q_3 - Q_1$ <ul style="list-style-type: none"> IQR è una misura robusta poiché non risente eccessivamente della dispersione 								
Varianza	<p>Sia X una variabile quantitativa univariata e si conosca la sua media e il suo scarto quadratico medio. Si chiama varianza degli N (n) valori osservati la somma delle differenze, al quadrato, tra ciascuna osservazione e la media della popolazione (campione), divisa per la dimensione della popolazione N (ampiezza del campione n)</p> <p>Essa ci permette di misurare la variabilità dei dati ma presenta una unità di misura elevata al quadrato rispetto a quella originale.</p> <p>Per dati grezzi</p> <table border="1" style="width: 100%;"> <thead> <tr> <th>POPOLAZIONE</th> <th>CAMPIONE (VARIANZA CAMPIONARIA)</th> </tr> </thead> <tbody> <tr> <td> $\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$ <p>Formula ridotta (Devianza)</p> $\sigma^2 = \frac{\sum_{i=1}^N x_i^2}{N} - \mu^2$ </td> <td> $s^2 = \frac{\sum_{i=1}^N (x_i - \bar{x})^2}{n - 1}$ <p>Formula ridotta (Devianza)</p> $s^2 = \frac{n}{n - 1} \left[\frac{\sum_{i=1}^n x_i^2}{n} - \bar{x}^2 \right]$ </td> </tr> </tbody> </table> <p>Per distribuzioni di frequenza</p> <table border="1" style="width: 100%;"> <thead> <tr> <th>POPOLAZIONE</th> <th>CAMPIONE (VARIANZA CAMPIONARIA)</th> </tr> </thead> <tbody> <tr> <td> <p>Distribuzione discreta</p> $\sigma^2 = \frac{1}{N} \sum_{l=1}^k f_l (x_l - \mu)^2$ $= \sum_{i=1}^k p_i (x_i - \mu)^2$ <p>Formula ridotta</p> $\sigma^2 = \frac{\sum_{i=1}^k f_i x_i^2}{\sum_{i=1}^k f_i} - \mu^2 = \sum_{i=1}^k p_i x_i^2 - \mu^2$ <p>Distribuzione continua</p> $\sigma^2 = \frac{\sum_{l=1}^k f_l (m_l - \mu)^2}{N}$ $= \sum_{i=1}^k p_i (m_i - \mu)^2$ </td> <td> <p>Distribuzione discreta</p> $s^2 = \frac{1}{n - 1} \sum_{l=1}^k f_l (x_l - \bar{x})^2$ $= \frac{n}{n - 1} \sum_{i=1}^k p_i (x_i - \bar{x})^2$ <p>Formula ridotta</p> $s^2 = \frac{\sum_{i=1}^k f_i}{\sum_{i=1}^k f_i - 1} \left[\frac{\sum_{i=1}^k f_i x_i^2}{\sum_{i=1}^k f_i} - \bar{x}^2 \right]$ $= \frac{n}{n - 1} \sum_{i=1}^k p_i x_i^2 - \bar{x}^2$ <p>Distribuzione continua</p> $s^2 = \frac{\sum_{l=1}^k f_l (m_l - \bar{x})^2}{n - 1}$ </td> </tr> </tbody> </table>	POPOLAZIONE	CAMPIONE (VARIANZA CAMPIONARIA)	$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$ <p>Formula ridotta (Devianza)</p> $\sigma^2 = \frac{\sum_{i=1}^N x_i^2}{N} - \mu^2$	$s^2 = \frac{\sum_{i=1}^N (x_i - \bar{x})^2}{n - 1}$ <p>Formula ridotta (Devianza)</p> $s^2 = \frac{n}{n - 1} \left[\frac{\sum_{i=1}^n x_i^2}{n} - \bar{x}^2 \right]$	POPOLAZIONE	CAMPIONE (VARIANZA CAMPIONARIA)	<p>Distribuzione discreta</p> $\sigma^2 = \frac{1}{N} \sum_{l=1}^k f_l (x_l - \mu)^2$ $= \sum_{i=1}^k p_i (x_i - \mu)^2$ <p>Formula ridotta</p> $\sigma^2 = \frac{\sum_{i=1}^k f_i x_i^2}{\sum_{i=1}^k f_i} - \mu^2 = \sum_{i=1}^k p_i x_i^2 - \mu^2$ <p>Distribuzione continua</p> $\sigma^2 = \frac{\sum_{l=1}^k f_l (m_l - \mu)^2}{N}$ $= \sum_{i=1}^k p_i (m_i - \mu)^2$	<p>Distribuzione discreta</p> $s^2 = \frac{1}{n - 1} \sum_{l=1}^k f_l (x_l - \bar{x})^2$ $= \frac{n}{n - 1} \sum_{i=1}^k p_i (x_i - \bar{x})^2$ <p>Formula ridotta</p> $s^2 = \frac{\sum_{i=1}^k f_i}{\sum_{i=1}^k f_i - 1} \left[\frac{\sum_{i=1}^k f_i x_i^2}{\sum_{i=1}^k f_i} - \bar{x}^2 \right]$ $= \frac{n}{n - 1} \sum_{i=1}^k p_i x_i^2 - \bar{x}^2$ <p>Distribuzione continua</p> $s^2 = \frac{\sum_{l=1}^k f_l (m_l - \bar{x})^2}{n - 1}$
	POPOLAZIONE	CAMPIONE (VARIANZA CAMPIONARIA)							
	$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$ <p>Formula ridotta (Devianza)</p> $\sigma^2 = \frac{\sum_{i=1}^N x_i^2}{N} - \mu^2$	$s^2 = \frac{\sum_{i=1}^N (x_i - \bar{x})^2}{n - 1}$ <p>Formula ridotta (Devianza)</p> $s^2 = \frac{n}{n - 1} \left[\frac{\sum_{i=1}^n x_i^2}{n} - \bar{x}^2 \right]$							
POPOLAZIONE	CAMPIONE (VARIANZA CAMPIONARIA)								
<p>Distribuzione discreta</p> $\sigma^2 = \frac{1}{N} \sum_{l=1}^k f_l (x_l - \mu)^2$ $= \sum_{i=1}^k p_i (x_i - \mu)^2$ <p>Formula ridotta</p> $\sigma^2 = \frac{\sum_{i=1}^k f_i x_i^2}{\sum_{i=1}^k f_i} - \mu^2 = \sum_{i=1}^k p_i x_i^2 - \mu^2$ <p>Distribuzione continua</p> $\sigma^2 = \frac{\sum_{l=1}^k f_l (m_l - \mu)^2}{N}$ $= \sum_{i=1}^k p_i (m_i - \mu)^2$	<p>Distribuzione discreta</p> $s^2 = \frac{1}{n - 1} \sum_{l=1}^k f_l (x_l - \bar{x})^2$ $= \frac{n}{n - 1} \sum_{i=1}^k p_i (x_i - \bar{x})^2$ <p>Formula ridotta</p> $s^2 = \frac{\sum_{i=1}^k f_i}{\sum_{i=1}^k f_i - 1} \left[\frac{\sum_{i=1}^k f_i x_i^2}{\sum_{i=1}^k f_i} - \bar{x}^2 \right]$ $= \frac{n}{n - 1} \sum_{i=1}^k p_i x_i^2 - \bar{x}^2$ <p>Distribuzione continua</p> $s^2 = \frac{\sum_{l=1}^k f_l (m_l - \bar{x})^2}{n - 1}$								

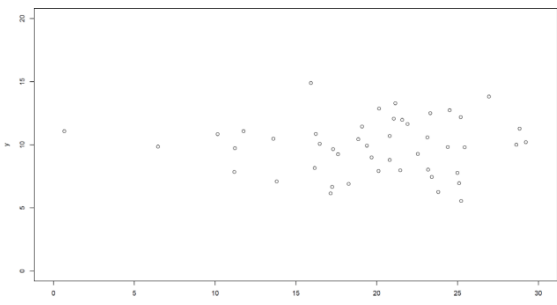
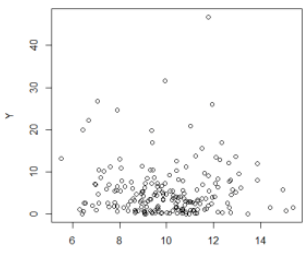
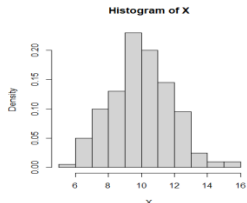
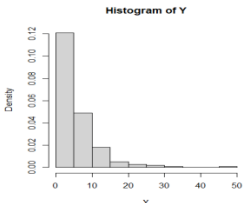


	<p><u>Formula ridotta</u></p> $\sigma^2 = \frac{\sum_{i=1}^k f_i m_i^2}{\sum_{i=1}^k f_i} - \mu^2$ $= \sum_{i=1}^k p_i m_i^2 - \mu^2$	$= \frac{n}{n-1} \sum_{i=1}^k p_i (m_i - \bar{x})^2$ <p><u>Formula ridotta</u></p> $s^2 = \frac{\sum_{i=1}^k f_i}{\sum_{i=1}^k f_i - 1} \left[\frac{\sum_{i=1}^k f_i m_i^2}{\sum_{i=1}^k f_i} - \bar{x}^2 \right]$ $= \frac{n}{n-1} \sum_{i=1}^k p_i m_i^2 - \bar{x}^2$					
<p>Scarto quadratico medio o deviazione standard</p>	<p>Sia X una variabile quantitativa univariata e si conosca la sua media e il suo scarto quadratico medio. Si chiama scarto quadratico medio o deviazione standard la radice quadrata della varianza, che permette di valutare <i>dispersione delle singole osservazioni rispetto alla media, con unità di misura uguale a quella osservata</i></p> <table border="1" data-bbox="403 987 1417 1189"> <thead> <tr> <th data-bbox="403 987 908 1030">POPOLAZIONE</th> <th data-bbox="908 987 1417 1030">CAMPIONE</th> </tr> </thead> <tbody> <tr> <td data-bbox="403 1030 908 1189"> $\sigma = \sqrt{\sigma^2} = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{N}}$ </td> <td data-bbox="908 1030 1417 1189"> $s = \sqrt{s^2} = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}$ </td> </tr> </tbody> </table>			POPOLAZIONE	CAMPIONE	$\sigma = \sqrt{\sigma^2} = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{N}}$	$s = \sqrt{s^2} = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}$
POPOLAZIONE	CAMPIONE						
$\sigma = \sqrt{\sigma^2} = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{N}}$	$s = \sqrt{s^2} = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}$						
<p>Proprietà della varianza e dello scarto quadratico medio</p>	<ul style="list-style-type: none"> • $\sigma^2 \geq 0$ e $\sigma \geq 0$ • $\sigma^2 = 0 \Leftrightarrow x_i = \mu, \forall i$ • Varianza e trasformazioni lineari: Se per ogni i, $y_i = a + bx_i$, con a e b valori reali arbitrari, allora: $\sigma_y^2 = b^2 * \sigma_x^2$ $\sigma_y = b \sigma_x$ • La varianza/scarto quadratico medio è una misura di dispersione; a trasla semplicemente la distribuzione → non modifica i due valori • Sono due misure NON ROBUSTE 						



<p>Coefficiente di Variazione</p>	<p>Sia X una variabile quantitativa univariata e si conosca la sua media e il suo scarto quadratico medio. Si chiama coefficiente di variazione l'indice che misura la variabilità di una variabile, relativamente alla media</p> <table border="1" data-bbox="403 322 1417 427"> <thead> <tr> <th>POPOLAZIONE</th> <th>CAMPIONE</th> </tr> </thead> <tbody> <tr> <td>$CV = \frac{\sigma}{ \mu }$</td> <td>$CV = \frac{s}{ \bar{x} }$</td> </tr> </tbody> </table> <p>Indica la fluttuazione (percentuale) di una variabile rispetto alla media; permette di apprezzare in termini relativi, rispetto alla media, la variazione rispetto al centro della distribuzione</p> <p>Questa misura:</p> <ul style="list-style-type: none"> • È appropriata per effettuare confronti di variabilità tra gruppi o variabili differenti • μ deve essere diverso da 0; in generale, valori di $\mu \rightarrow 0$ provocano oscillazioni molto forti • CV è un numero puro, non ha unità di misura • CV non è una misura robusta <p>N.B.: è l'indice opportuno per rispondere alla domanda: qual è la variabile con variabilità maggiore?</p>	POPOLAZIONE	CAMPIONE	$CV = \frac{\sigma}{ \mu }$	$CV = \frac{s}{ \bar{x} }$
POPOLAZIONE	CAMPIONE				
$CV = \frac{\sigma}{ \mu }$	$CV = \frac{s}{ \bar{x} }$				
<p>Disuguaglianza di Chebychev</p>	<p>Siano μ e σ la media e la deviazione standard di una popolazione qualsiasi. Fissato $k > 1$, allora:</p> $Fr\{\mu - k\sigma \leq X \leq \mu + k\sigma\} \geq 1 - \frac{1}{k^2}$ <p>La disuguaglianza ci permette di calcolare la distribuzione approssimata (ma sottostimata) dei dati in uno specifico intervallo, avendo a disposizione unicamente la media e la varianza. In particolare, permette di definire:</p> <ul style="list-style-type: none"> • La frequenza almeno osservabile in un intervallo di valori • La frequenza al più osservabile sulle code <p>Per ricavare k basta, avendo media, deviazione standard e gli estremi dell'intervallo [a,b]</p> $k = \frac{\mu - a}{\sigma}$ <p>N.B.: l'intervallo usato da Chebychev è centrato nella media ovvero, l'intervallo dato, supponiamo $[a,b]$, deve essere tale che $\frac{a+b}{2} = \bar{x}$.</p> <p>Se l'intervallo non fosse centrato nella media allora avremmo dovuto cercare l'intervallo più grande centrato nella media, contenuto nell'intervallo dato: saremo certi che la frequenza dell'intervallo dato è maggiore (o al massimo uguale) a quello centrato</p> <p>N.B.: attento a controllare tutto quello che si ha! Potrebbe capitare che l'intervallo dato corrisponda a specifici punti forniti (es.: differenza tra Max e Q_1^*): se ciò accade abbiamo una misura esatta e non approssimata (*0,75 = 75%)</p> <p>Attenzione: in generale non è possibile trovare un limite inferiore né un limite superiore per la probabilità usando il teorema dei Chebychev</p> <p>Regola empirica dei tre σ</p> <p>Se la distribuzione della variabile X è di forma campanulare e simmetrica, allora:</p>				



	<p style="text-align: center;"> $Fr\{\mu - \sigma \leq X \leq \mu + \sigma\} \approx 0.68$ $Fr\{\mu - 2\sigma \leq X \leq \mu + 2\sigma\} \approx 0.95$ $Fr\{\mu - 3\sigma \leq X \leq \mu + 3\sigma\} \approx 0.9973$ </p> <p>N.B.: prima di applicare Chebychev controlla se è possibile usare la legge empirica</p> <ol style="list-style-type: none"> 1. Verifica se la distribuzione è simmetrica, controllando la distanza tra i cinque punti di sintesi oppure la relazione tra media e mediana 2. Controlla se è campanulare, ovvero se la densità si concentra principalmente nel centro; per farlo, osserva se la distanza tra Min e Q₁ è maggiore alla distanza tra Q₁ e Me (oppure la distanza tra Q₃ e Max maggiore di Me e Q₃) 		
MISURE DI RELAZIONE LINEARE PER VARIBILE QUANTITATIVE BIVARIATE			
<p>Forma dei dati grezzi</p>	<p>Per le analisi bivariate impiegheremo due variabili numeriche X e Y raggruppate in n coppie:</p> <p style="text-align: center;">(x_n, y_n)</p>		
<p>Diagramma a dispersione</p>	<p>Il diagramma a dispersione, detto anche scatterplot, rappresenta la distribuzione delle coppie di variabili in un piano cartesiano di assi X e Y L'i-esima unità del dataset è rappresentata da un punto di coordinate (x_i, y_i) nel piano cartesiano di assi X e Y</p> <div style="text-align: center;">  </div> <p>Simmetria e correlazione con istogramma Posso ricondurre lo scatterplot a due istogrammi (per X e per Y) tracciando le proiezioni dei punti sugli assi e osservando la simmetria</p> <div style="display: flex; align-items: center; justify-content: center;"> <div style="text-align: center;">  </div> <div style="margin: 0 20px;">→</div> <div style="display: flex; gap: 20px;"> <div style="text-align: center;">  </div> <div style="text-align: center;">  </div> </div> </div>		
<p>Covarianza</p>	<p>Siano X ed Y due variabili quantitative bivariate e sia N la popolosità di una popolazione mentre n sia l'ampiezza di un campione. Si chiama covarianza l'indice di misura dell'esistenza della relazione lineare tra le due variabili X ed Y. Essa si calcola come segue:</p> <table border="1" style="width: 100%; margin-top: 10px;"> <tr> <td style="width: 50%; text-align: center;">POPOLAZIONE</td> <td style="width: 50%; text-align: center;">CAMPIONE</td> </tr> </table>	POPOLAZIONE	CAMPIONE
POPOLAZIONE	CAMPIONE		



$$\sigma_{XY} = \frac{\sum_{i=1}^N (x_i - \mu_x)(y_i - \mu_y)}{N}$$

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

$$\sigma_{XY} = Cov(X, Y) = \frac{\sum_{i=1}^N x_i y_i}{N} - \mu_x \mu_y$$

$$s_{XY} = \frac{n}{n - 1} \left[\frac{\sum_{i=1}^n x_i y_i}{n} - \bar{x} \cdot \bar{y} \right]$$

In particolare:

- **Cov > 0** → x e y tendono a muoversi nella stessa direzione (**dipendenza lineare positiva o diretta**)
- **Cov < 0** → x e y tendono a muoversi in direzioni opposte (**dipendenza lineare negativa o inversa**)
- **Cov = 0** → x e y sono non correlate, non c'è dipendenza lineare (**assenza di relazione lineare oppure relazione forte ma non monotona**)

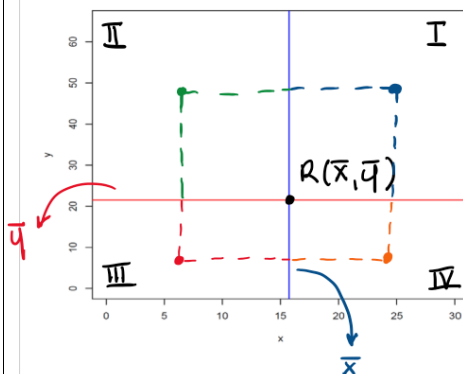
Osservazioni:

- Difficile valutare la covarianza in termini assoluti...
- Covarianza e trasformazioni lineari
Se $x_i^* = a + bx_i$ e $y_i^* = c + dy_i$ con a, b, c e d valori reali arbitrari, allora:

$$\sigma_{X^*Y^*} = bd * \sigma_{XY}$$

- Si osserva che nella trasformazione:
 - **Le variabili additive non hanno effetto**
 - **Le variabili moltiplicative hanno effetto**
- La covarianza è influenzata dalla dispersione individuale (scarto quadratico) di ciascuna variabile. Infatti, se la variabilità σ_x è alta, allora le deviazioni $x_i - \bar{x}$; conseguentemente tenderanno ad essere elevati i prodotti incrociati $(x_i - \bar{x})(y_i - \bar{y})$ e, quindi, la covarianza; quindi, basta anche solo la variabilità elevata di una variabile, inflaziona l'intero indice...
- La covarianza non è una misura robusta

Interpretazione della covarianza



QUADRANTE	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$
I	+	+	⇒ ⊕
II	-	+	⇒ ⊖
III	-	-	⇒ ⊕
IV	+	-	⇒ ⊖

Osservando i valori nei quadranti posso stimare il segno della Covarianza, in quanto posso determinare il segno di $(x_i - \bar{x})(y_i - \bar{y})$



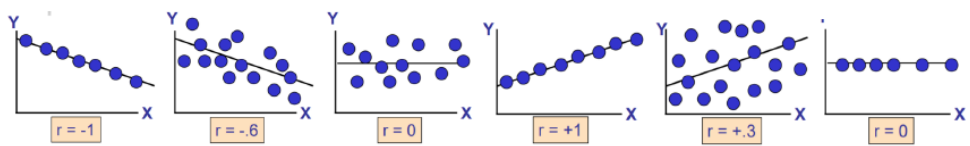
Il valore della covarianza dipende dall'unità di misura e quindi, non si tratta di un indice adeguato a valutare l'intensità della relazione lineare tra due variabili.

Siano X ed Y due variabili quantitative bivariate e si conoscano la covarianza e gli scarti quadratici medi delle due variabili. Si chiama **coefficiente di correlazione lineare di Pearson** un **indice che misura l'esistenza e l'intensità della relazione lineare** tra le due variabili X ed Y. Esso si calcola nel seguente modo:

POPOLAZIONE	CAMPIONE
$\rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$	$r_{XY} = \frac{S_{XY}}{S_X S_Y}$

Proprietà

- $-1 \leq \rho_{XY} \leq 1$
- $\rho_{XY} > 0 \rightarrow$ **correlazione lineare positiva**
- $\rho_{XY} < 0 \rightarrow$ **correlazione lineare negativa**
- $|\rho_{XY}| = 1 \rightarrow$ **correlazione lineare esatta**
 - Ovvero $y_i = a + bx_i$ con $b > 0$ se $\rho_{XY} = 1$; con $b < 0$ se $\rho_{XY} = -1$
- $\rho_{XY} \approx 0 \rightarrow$ **assenza di correlazione lineare, ma non di qualsiasi altro tipo di relazione** (ad esempio, potrebbe esserci una dipendenza quadratica)



Teorema

Indipendenza statistica \Rightarrow NON correlazione lineare

Regola pratica
Si può assumere l'esistenza di una relazione lineare se:

$$|r| = \frac{2}{\sqrt{n}}$$

ELEMENTI DI PROBABILITÀ E DISTRIBUZIONI NOTEVOLI

Funzioni aleatorie

- **Funzione di probabilità:** funzione che associa ad ogni numero reale la probabilità con cui una variabile assume quel numero reale stesso

$$p(x) = p(X = x)$$

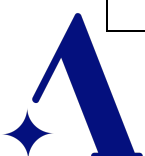
Gode delle seguenti proprietà:

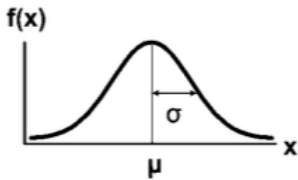
- $p(x) \geq 0$
- $\sum_{i=1}^n p(x_n) = 1$

- **Funzione di ripartizione:** funzione che associa ad ogni numero reale la probabilità con cui una variabile assume tutti i valori precedenti a tale numero, esso compreso

$$F(x) = P(X \leq x) = \sum_{u \leq x} p(u)$$


	<p>Gode delle seguenti proprietà:</p> <ul style="list-style-type: none"> ○ $0 \leq F(x) \leq 1$ ○ $F(x)$ non decrescente ○ $F(x)$ costante a tratti (discontinua)
Valore atteso	<p>Il valore atteso $E(x)$ è la media dei valori ed è il baricentro (punto centrale) della distribuzione della variabile aleatoria Si calcola come la media di X ponderata per le sue modalità:</p> $E[x] = \mu = \sum xp(x)$ <p>Interpretazione Il valore atteso $E[x]$ è approssimativamente la media aritmetica dei risultati che si ottengono replicando <i>all'infinito</i> l'esperimento</p>
Varianza e scarto quadratico medio	<p>Secondo la definizione la varianza è il valore atteso della differenza, al quadrato, tra X e la media:</p> $Var(x) = \sigma^2 = E[(x - \mu)^2] = \sum (x - \mu)^2 p(x)$ <p>Ma applichiamo la formula per il calcolo</p> $Var(x) = \sigma^2 = E[x^2] - \mu^2 = \sum [x^2 p(x)] - \mu^2$ <p>Infine, lo scarto quadratico medio</p> $\sigma = \sqrt{Var(x)}$
Covarianza	<p>La covarianza è la differenza tra il valore atteso del prodotto delle variabili aleatorie ed il prodotto dei valori attesi:</p> $Cov(X, Y) = E[XY] - E[X]E[Y]$
Trasformazioni lineari di variabili aleatorie	<p>Trasformazione lineare semplice Consideriamo una variabile aleatoria X e la sua trasformazione lineare</p> $Y = a + bX$ <p>Allora:</p> <ul style="list-style-type: none"> • $E(Y) = a + bE(x)$ • $Var(Y) = b^2 Var(x)$ • $\delta(Y) = b \delta(x)$ <p>Trasformazione di combinazione lineare Consideriamo una variabile aleatoria X e Y e la sua combinazione lineare</p> $W = aX + bY$ <p>Risulta:</p> <ul style="list-style-type: none"> ○ $E(W) = aE(X) + bE(Y)$ ○ $Var(W) = a^2 Var(X) + b^2 Var(Y) + 2abCov(X, Y)$ <p>N.B.: Se avessimo</p> $W = aX - bY$

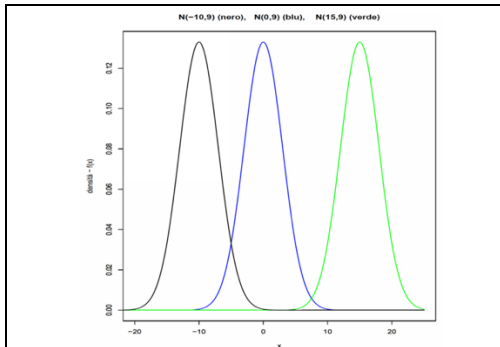


	<p>Risulterebbe</p> <ul style="list-style-type: none"> ○ $E(W) = aE(X) - bE(Y)$ ○ $Var(W) = a^2Var(X) + b^2Var(Y) - 2abCov(X,Y)$ <p>Ricorda che:</p> $\rho_{XY} = \frac{Cov(X,Y)}{\sigma_X\sigma_Y}$ <p>Se X e Y fossero indipendenti, cioè $Cov(X,Y) = 0$, risulterebbe:</p> $Var(W) = a^2Var(X) + b^2Var(Y)$
<p>Standardizzazione</p>	<p>È la trasformazione lineare di una variabile aleatoria che si ottiene sottraendo μ e dividendo per δ</p> $Z = \frac{X - \mu}{\delta}$ <p>Essa trasforma una qualsiasi variabile aleatoria X in una nuova variabile Z, detta standardizzata, tale che $E(Z) = 0$ e $Var(Z) = 1$</p>
<p>Distribuzione binomiale (variabili discrete)</p>	<p>Riporta qual è la probabilità di avere x successi su n prove</p> $X \sim Bin(n, p) = \begin{cases} \binom{n}{x} p^x (1-p)^{n-x} & x = 0, 1, 2, 3, \dots, n \\ 0 & \text{altrove} \end{cases}$ <p>Con $\binom{n}{x} = \frac{n!}{x!(n-x)!}$</p> <p>Proprietà:</p> <ul style="list-style-type: none"> • $E(x) = np$ • $Var(x) = np(1-p)$
<p>Distribuzione bernoulliana (variabili discrete)</p>	<p>Una variabile si distribuisce secondo una bernoulliana allorquando si fronteggia un esperimento dicotomico, ovvero un esperimento che può avere due soli esiti possibili (successo, fallimento). È un particolare caso di una distribuzione binomiale con n = 1</p> $X \sim Bern(p) = \begin{cases} p & x = 1 \\ 1-p & x = 0 \end{cases}$ <p>Proprietà</p> <ul style="list-style-type: none"> • $E(x) = p$ • $Var(x) = p(1-p)$
<p>Distribuzione normale o gaussiana (variabili continue)</p>	<p>Una variabile aleatori si distribuisce come una normale/gaussiana di media μ e varianza δ^2 quando:</p> $X \sim N(\mu, \sigma^2)$ <div style="text-align: right;">  </div> <ul style="list-style-type: none"> • La distribuzione ha andamento simmetrico e campanulare • Al centro della simmetria la distribuzione è $\mu = E(x)$ • Al crescere di δ^2 la densità si appiattisce e le code divengono più pesanti: questo perché, la varianza indica la dispersione della distribuzione rispetto al centro; tanto più la varianza è grande tanto più i valori saranno dispersi • $Z \sim N(0,1)$ si chiama normale standard

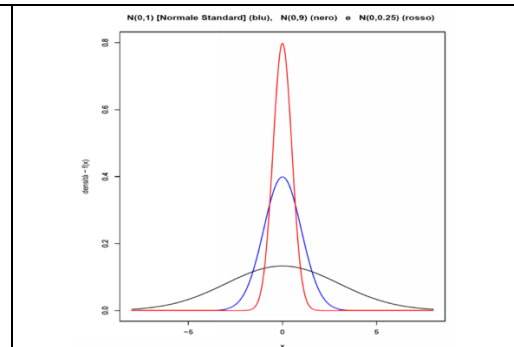


- $X \sim N(\mu, \delta^2)$ e $Y = a + bX$ allora $Y \sim (a + b\mu, b^2\delta^2)$
- **Media = Mediana = Moda**

Parametri e forma della distribuzione



Le curve hanno stessa δ^2 ma risulta:
 $\mu_{nera} < \mu_{blu} < \mu_{verde}$



Le curve hanno stessa μ ma risulta:
 $\delta^2_{rossa} < \delta^2_{blu} < \delta^2_{nera}$

Calcolo della probabilità cumulata e standardizzazione

L'area sottesa dalla curva rappresenta la probabilità cumulata; per eseguirne il calcolo si procede così:

1. Definisco la curva per il problema in oggetto
2. Standardizzo traducendo X in $Z = \frac{X-\mu}{\sigma}$
3. Riscrivo la probabilità assegnata:

$$P(X < a) = P\left(Z < \frac{a - \mu}{\sigma}\right)$$

4. Uso le tavole per ritrovare il risultato


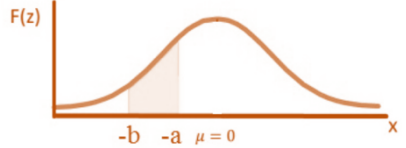
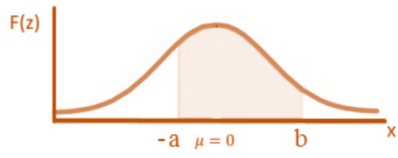
Simmetria e campanularità

- Poiché la distribuzione è simmetrica la media è uguale alla mediana e alla media interquartile

$$\mu = Me = \frac{Q_3 + Q_1}{2}$$

Casi di calcolo di probabilità con la normale	1	$P(Z < a) = F_z(a)$	
	2	$P(Z > -a) = P(Z < a) = F_z(a)$	
	3	$P(Z > a) = 1 - F_z(a)$	
	4	$P(Z < -a) = P(Z > a) = 1 - F_z(a)$	



	5	$P(a < Z < b) = F_z(b) - F_z(a)$	
	6	$P(-b < Z < -a) = F_z(b) - F_z(a)$	
	7	$P(-a < Z < b) = F_z(b) + F_z(a) - 1$	
Chebichev e normale	Se i dati forniti applicando Chebichev, su un intervallo centrato sulla media, fornita con la deviazione standard, suggerissero una forma simmetrica e campanulare , allora è possibile approssimare più precisamente la frequenza normalizzando e ricorrendo alle tavole di sintesi numerica		
Approssimazione della distribuzione binomiale con la normale	<p>La distribuzione binomiale si può approssimare alla distribuzione normale solo quando n è abbastanza grande ovvero quando $np(1 - p) > 9$</p> <p>In questo caso:</p> <ul style="list-style-type: none"> $E(x) = np$ $Var(x) = np(1 - p)$ <p>E la distribuzione binomiale $X \sim Bin(n, p)$ può essere trasformata approssimativamente in una normale:</p> $X \approx N(np, np(1 - p))$ <p>Per cui, la standardizzazione diviene:</p> $Z = \frac{X - E(x)}{Var(x)} = \frac{X - np}{\sqrt{np(1 - p)}}$		

STATISTICA INFERENZIALE	
Introduzione	La statistica inferenziale ci permette di trarre conclusioni/prendere decisioni riguardanti una popolazione sulla base dei risultati di un singolo campione
Metodi di campionamento probabilistici	<p>Campionamento casuale semplice (CCS): ogni unità della popolazione ha la stessa probabilità di essere incluso nel campione; si procede così:</p> <ol style="list-style-type: none"> Si seleziona un'unità della popolazione Si rileva il valore d'interesse Si esegue una scelta: <ol style="list-style-type: none"> Si reinserisce l'unità (campionamento con reimmissione) Non si reinserisce l'unità (campionamento senza reimmissione) Si ripetono le fasi finché non si selezionano n unità <p>Proprietà della scelta</p> <ul style="list-style-type: none"> CCS con reimmissione → le probabilità osservate da ciascuno dei due sono le medesime e, in particolare, sono uguali a quelle che caratterizzano la popolazione → V.A. INDIPENDENTI E IDENTICAMENTE DISTRIBUITE: $X_1, X_2, \dots \sim iid f_x \rightarrow Cov(X, Y) = 0$



	<ul style="list-style-type: none"> • CCS senza reimmissione con n piccolo → le probabilità osservata da ciascuno dei due sono diverse ma si può dimostrare che sono uguali a quelle che caratterizzano la popolazione → V.A. DIPENDENTI E IDENTICAMENTE DISTRIBUITE - $X_1, X_2, \dots \sim^{id} f_x \rightarrow Cov(X, Y) \neq 0$ • CCS senza reimmissione con n grande → le probabilità osservata da ciascuno dei due sono <i>minimamente diverse</i> e dipendono <i>debolmente</i> da quanto precedentemente osservato → V.A. INDIPENDENTI E IDENTICAMENTE DISTRIBUITE: $X_1, X_2, \dots \approx^{id} f_x \rightarrow Cov(X, Y) = 0$
<p>Media campionaria</p>	<p>Si consideri una campione causale di ampiezza n da una popolazione con media μ e varianza δ^2; sia tale campione IID. Allora si chiama media campionaria la media di n variabili aleatorie:</p>
	$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ <p>Proprietà:</p> <ol style="list-style-type: none"> 1. Il valore atteso della media campionaria è uguale alla media della popolazione (la media campionaria è uno stimatore non distorto per la media della popolazione, vedi oltre) $E[\bar{X}] = \mu$ <p>Dimostrazione:</p> $E[\bar{X}] = E\left[\frac{1}{n} \sum x_i\right] = \frac{1}{n} E\left[\sum x_i\right] = \frac{1}{n} \sum E[x_i] = \frac{1}{n} * n\mu = \mu$ <ol style="list-style-type: none"> 2. La varianza della media campionaria è la varianza della popolazione divisa per la popolosità del campione: $Var[\bar{X}] = \frac{\sigma^2}{n}$ <p>Dimostrazione:</p> $Var[\bar{X}] = Var\left[\frac{1}{n} \sum x_i\right] = \left(\frac{1}{n}\right)^2 Var\left[\sum x_i\right] = \frac{1}{n^2} n\sigma^2 = \frac{\sigma^2}{n}$ <ol style="list-style-type: none"> 3. La deviazione standard risulta: $\sigma = \frac{\sigma}{\sqrt{n}}$ <p>N.B.: se il campionamento è SENZA REIMMISSIONE e risulta $n > 0.05N$ allora:</p> $Var[\bar{X}] = \frac{\sigma^2}{n} * \frac{N-n}{N-1}$ <p>Dove $\frac{N-n}{N-1}$ è il fattore di correzione per popolazione finita</p> <p>Distribuzione della media campionaria</p> <ol style="list-style-type: none"> 1. POPOLAZIONE NORMALE <p>Teorema:</p>



Sia X un V.A. IID che si distribuisce come una normale $X \sim N(\mu, \sigma^2)$. Sotto queste ipotesi, **se la popolazione ha distribuzione normale allora il campione ha distribuzione ancora normale**

a. **CAMPIONE CON REIMMISSIONE**

- $E(\bar{X}) = \mu$
- $Var(\bar{X}) = \frac{\sigma^2}{n}$
- $\delta(\bar{X}) = \frac{\sigma}{\sqrt{n}}$ (**standard error**)
- Distribuzione:

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

- Standardizzazione $\rightarrow Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$

b. **CAMPIONE SENZA REIMMISSIONE** (con $n > 0.05N \rightarrow$ **fattore di correzione per popolazioni finite**)

- $E(\bar{X}) = \mu$
- $Var(\bar{X}) = \frac{\sigma^2}{n} * \frac{N-n}{N-1}$
- $\delta(\bar{X}) = \frac{\sigma}{\sqrt{n}} * \sqrt{\frac{N-n}{N-1}}$ (**standard error**)
- Distribuzione:

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n} * \frac{N-n}{N-1}\right)$$

- Standardizzazione $\rightarrow Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}} * \sqrt{\frac{N-n}{N-1}}}$

N.B.: è possibile ridurre l'errore standard (deviazione) aumentando n .

2. **POPOLAZIONE QUALSIASI**

Teorema centrale del limite

Sia X un campione casuale di ampiezza n da una popolazione **arbitraria** con $E(x) = \mu$ e $Var(X) = \sigma^2$; sia n sufficientemente elevato (**approssimativamente $n > 25$**).

Allora **la media** (o **la somma**) di variabili aleatorie IID si distribuisce **approssimativamente come una normale** con parametri:

$$\bar{X} \approx N\left(\mu, \frac{\sigma^2}{n}\right) \Leftrightarrow \sum_i X_i \approx N(n\mu, n\sigma^2)$$

Standardizzando:

$$\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \approx N(0,1) \Leftrightarrow \frac{\sum_i X_i - n\mu}{\sigma\sqrt{n}} \approx N(0,1)$$

3. **DISTRIBUZIONE BERNOULLIANA**



	<p>Si consideri un campione di una popolazione che si distribuisce secondo una Bernoulliana di parametro p, quindi con media μ e varianza $p(1-p)$. Se risulta $np(1-p) > 9$ allora è possibile applicare il teorema centrale del limite e P si distribuisce approssimativamente:</p> $\hat{P} \approx N\left(p, \frac{p(1-p)}{n}\right)$
Distribuzione della proporzione campionaria	<p>Data una distribuzione Bernoulliana. Sia p la proporzione della popolazione che possiede la caratteristica oggetto di studio. Si dice proporzione campionaria \hat{P} la stima della proporzione del successo di carattere osservato p, sul totale del campione in esame:</p> $\hat{P} = \frac{\text{\#popolazione successo}}{\text{dimensione campione}} = \frac{\sum_{i=1}^n x_i}{n}$ <p>Con $0 \leq \hat{P} \leq 1$</p> <p>La \hat{P} ha lo stesso significato di \bar{X}, ma si riferisce al caso particolare della bernoulliana; in particolare, per n sufficientemente grande [$np(1-p) > 9$] risulta che:</p> $\hat{P} \approx N\left(p, \frac{p(1-p)}{n}\right)$ <p>Standardizzazione:</p> $Z = \frac{\hat{P} - p}{\sqrt{\frac{p(1-p)}{n}}}$ <p>E deviazione standard (standard error)</p> $\sigma = \sqrt{\frac{p(1-p)}{n}}$ <p>Un modo alternativo, ma del tutto equivalente, per affrontare questo tipo di problemi, è quello di valutare la corrispondente distribuzione binomiale. Data una distribuzione bernoulliana $X \sim \text{Bern}(p)$, la somma di n distribuzioni bernoulliane si distribuiscono come una binomiale $X \sim \text{Bin}(n, p)$. La binomiale rappresenta la ripetizione di n esperimenti bernoulliani. Se $np(1-p) > 9$, può essere approssimata ad una normale con media e varianza pari a:</p> $X \approx N(np, np(1-p))$ <p>X rappresenta il totale di elementi considerati non sotto forma di proporzione. Da cui segue la standardizzazione:</p> $Z = \frac{X - np}{\sqrt{np(1-p)}}$



<p>Esercizi tipici sulla media e proporzione campionaria</p>	<ol style="list-style-type: none"> 1. <u>Calcolo dell'errore standard della media o della proporzione campionaria</u> Sono le rispettive deviazioni standard 2. <u>Calcolo della probabilità con media campionaria</u> Possono riguardare diverse tipologie, ma occorre partire sempre: <ul style="list-style-type: none"> ○ Verificare l'ampiezza del campione e controllare se si necessita del fattore di correzione (campionamento senza reimmissione con $n > 0.05N$; se N non è fornito supporremo la popolazione infinita e non si necessiterà del fattore di correzione ○ Scrivere la distribuzione della media campionaria secondo la normale: il suo valore atteso è uguale a quello della popolazione; la sua varianza dipende dalla presenza o meno del fattore di correzione ○ Non <i>spaventarsi</i> se non vengono forniti n o $Var(X)$ o $E(X)$: in questi casi l'esercizio chiede, probabilmente, di determinare proprio uno di quei valori; scriviamo, i parametri della normale mantenendo incogniti i valori non dati e ricerchiamoli 3. <u>Calcolo dei percentili</u> Viene fornito il percentile (p) o il suo complementare ($1-p$); in questi casi occorre prima standardizzare e poi cercare sulle tavole della normale <i>all'incontrario</i>: individuo il valor di probabilità che si avvicina di più a quello dato (se intermedio, eseguo una media) e risalgo al valore corrispondente 4. <u>Considerazioni circa la probabilità al variare di n</u> Un risultato fondamentale afferma che all'aumentare di n, la varianza diminuisce e viceversa. La varianza è una misura della dispersione dei dati del campione; se stiamo considerando un intervallo che comprende la media, allora una riduzione della varianza e la conseguente concentrazione maggiore dei dati attorno alla media, provocherà un aumento della probabilità; se, viceversa, stiamo considerando un intervallo che include le <i>code</i> della normale, una riduzione della varianza, provocherà un minore appiattimento della curva e, quindi, code più sottili: la probabilità diminuisce 5. <u>Calcolo della probabilità con la proporzione campionaria</u> La proporzione campionaria indica un rapporto (e quindi una percentuale) dei successi sul totale della popolazione. Uso la proporzione campionaria quando viene richiesto la probabilità circa una determinata percentuale; uso la distribuzione della somma dei successi, nel caso venga richiesto una probabilità di numero. Stesse considerazioni circa le variazioni di n, valgono anche per la proporzione campionaria 6. <u>Superamento e variazioni della media della popolazione da parte della media campionaria</u> In questo tipo di esercizi viene fornita unicamente la deviazione standard (o la varianza) ma non la media (inutile); distinguiamo due casi modello (possono esserci piccole variazioni ma il ragionamento è il medesimo): <ul style="list-style-type: none"> ○ Qual è la probabilità che la media campionaria superi la media della popolazione per più di a $P(\bar{X} > \mu + a) = P\left(Z > \frac{\mu + a - \mu}{\sigma}\right) = P\left(Z > \frac{a}{\sigma}\right)$ <p>Ora calcolo normalmente con le tavole</p> <p>Analogamente alla richiesta di quale sia la probabilità che la media campionaria sia inferiore alla media della popolazione per più di a si risponde come segue:</p> $P = P(\bar{X} < \mu - a) = P\left(Z < \frac{\mu - a - \mu}{\sigma}\right) = P\left(Z < -\frac{a}{\sigma}\right)$
--	---



	<ul style="list-style-type: none"> Qual è la probabilità che la media campionaria differisca dalla media della popolazione per più di a $P(\bar{X} - \mu > a) = P\left(\left \frac{\bar{X} - \mu}{\sigma}\right > \frac{a}{\sigma}\right) = P(Z > \frac{a}{\sigma})$ <p>Scompongo il valore assoluto</p> $P\left(Z > \frac{a}{\sigma}\right) + P\left(Z < -\frac{a}{\sigma}\right) = 2\left[1 - F\left(\frac{a}{\sigma}\right)\right]$ <p>Qualora fosse stata richiesta la probabilità che la media campionaria differisse per meno di a allora:</p> $P(\bar{X} - \mu < a) = P(\mu - a < \bar{X} < \mu + a)$ <p>N.B.: in un esercizio può capitare che sia richiesto di calcolare la probabilità di un parametro per <i>la n-esima unità specifica estratta dal campione</i>. In questo caso NON dobbiamo valutare la media campionaria ma la distribuzione nel suo complesso poiché la <i>specifica unità</i> è parte della popolazione intera. Infatti, <i>la singola unità del campione ha la stessa distribuzione del campione stesso</i></p>
<p>Esempio riepilogativo</p>	<p><i>Si consideri un campione casuale di 50 incassi settimanali (tra loro indipendenti) per un ristorante. Sapendo che l'incasso medio settimanale sia pari a 25 e che la varianza è pari a 4, quale distribuzione è una ragionevole descrizione per l'incasso complessivo delle 50 settimane?</i></p> <p>Sia X l'incasso settimanale della i-esima settimana; l'incasso complessivo delle 50 settimane considerate è dato dalla somma dei 50 incassi settimanali, (indipendenti e identicamente distribuiti) ovvero da $\sum_{i=1}^{50} X_i$</p> <p>Poiché il campione è sufficientemente ampio, possiamo applicare <u>teorema centrale del limite</u>: la somma di variabili aleatorie IID si distribuisce approssimativamente come una normale con i seguenti parametri:</p> $\sum_i X_i \approx N(n\mu, n\sigma^2) \rightarrow \sum_i X_i \approx N(25 * 50, 50 * 4) \rightarrow \sum_i X_i \approx N(1250, 200)$ <p><i>Sia 30 il primo quartile. Supponete che in una qualunque settimana il ristorante considerato raggiunga il pareggio solo se l'incasso è superiore a 30. In quante settimane su 50 ci si può attendere che il ristorante sia in pareggio?</i></p> $P(\text{Incasso settimanale} \geq 30) = 0.25$ <p>Ogni settimana rappresenta una prova bernoulliana Y_i, con probabilità di raggiungere il pareggio pari a $p = 0.25$</p> <p>Il numero di settimane su 50 in cui il pareggio sarà raggiunto corrisponde alla ripetizione per 50 volte dell'esperimento bernoulliano e quindi ad una binomiale con i seguenti parametri:</p> $Y_i \sim \text{Bern}(p = 0.25) \rightarrow \sum_i Y_i \sim \text{Bin}(n = 50; p = 0.25)$ <p>Da cui segue che:</p> $E\left[\sum_i Y_i\right] = 50 * p = 50 * 0.25 = 12.5$ <p>Il pareggio è raggiunto in 12.5 settimane su 50</p>



Supponiamo che la probabilità che il ristorante raggiunga il pareggio è pari a $p=0.40$, con quale probabilità il ristorante raggiungerà il pareggio almeno in 15 settimane su 50?

Ripartiamo dalla distribuzione binomiale: poiché $np(1-p)>9$, il teorema centrale del limite ci consente di approssimare la binomiale ad una normale; quindi, il numero di settimane in cui verrà raggiunto il pareggio su 50 ha distribuzione:

$$\sum_i Y_i \approx N(np; np(1-p)) = N(20,12)$$

Possiamo calcolare la probabilità che il ristorante raggiunga il pareggio in almeno 15 settimane su 50:

$$P\left(\sum_i Y_i \geq 15\right) = P\left(Z \geq \frac{15-20}{\sqrt{12}}\right) = P(Z > -1.44) = P(Z < 1.44) = 0.9251$$

GLI STIMATORI PUNTUALI													
Definizioni	<p>Uno stimatore di un parametro θ è:</p> <ul style="list-style-type: none"> • Una variabile aleatoria che dipende dall'informazione contenuta nella popolazione • Il cui valore fornisce un'approssimazione del valore sconosciuto del parametro (a posteriori, dopo l'estrazione del campione) • Uno specifico valore dallo stimatore viene chiamato stima puntuale <p>Un intervallo di confidenza fornisce ulteriori informazioni circa la variabilità della stima ed indica un <i>intervallo di valori plausibili</i> per il parametro delle popolazioni</p> <div style="text-align: center;"> </div>												
Stime puntuali	<p>Come possiamo trovare uno stimatore per un parametro d'interesse?</p> <ul style="list-style-type: none"> • STIME PER ANALOGIE <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th>Stima</th> <th>Parametro</th> </tr> </thead> <tbody> <tr> <td>\bar{X}</td> <td>μ</td> </tr> <tr> <td>S^2</td> <td>σ^2</td> </tr> <tr> <td>S</td> <td>σ</td> </tr> <tr> <td>S_{XY}</td> <td>σ_{XY}</td> </tr> <tr> <td>r_{XY}</td> <td>ρ_{XY}</td> </tr> </tbody> </table> <ul style="list-style-type: none"> • STIMATORI DEI MINIMI QUADRATI • STIMATORI B.L.U.E 	Stima	Parametro	\bar{X}	μ	S^2	σ^2	S	σ	S_{XY}	σ_{XY}	r_{XY}	ρ_{XY}
Stima	Parametro												
\bar{X}	μ												
S^2	σ^2												
S	σ												
S_{XY}	σ_{XY}												
r_{XY}	ρ_{XY}												
Criteri di valutazione degli stimatori	<p>1. NON DISTORSIONE Uno stimatore puntuale T per un parametro θ si dirà non distorto se:</p> $E[T] = \theta$												

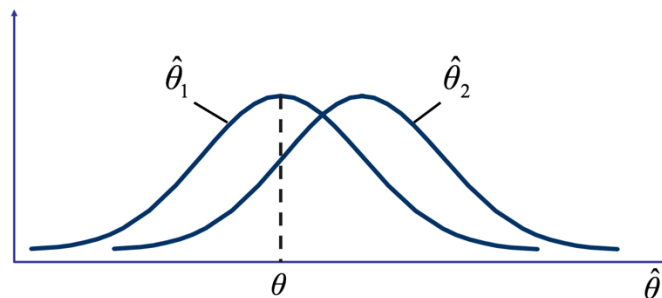


Se T fosse distorto, allora è possibile misurare la **distorsione dello stimatore $D(T)$** :

$$D(T) = E[T] - \theta$$

Infine, uno stimatore è non distorto se e solo se ha distorsione nulla:

$$T \text{ non distorto} \Leftrightarrow D(T) = 0$$



Nella figura, lo stimatore $\hat{\theta}_1$ è non distorto; $\hat{\theta}_2$ è distorto

Da qui segue che:

- La media campionaria è uno stimatore non distorto per la media
- La varianza campionaria è uno stimatore non distorto per la varianza
- La proporzione campionaria è uno stimatore non distorto per la proporzione

Sovrastima, sottostima e non distorsione

- $E[T] > \theta \text{ o } D(T) > 0 \Rightarrow$ Sovrastima
- $E[T] < \theta \text{ o } D(T) < 0 \Rightarrow$ Sottostima
- $E[T] = \theta \text{ o } D(T) = 0 \Rightarrow$ Non distorsione (correttezza)

Più in generale, dato uno stimatore nella forma:

$$T_w = w_1 X_1 + w_2 X_2 + \dots + w_n X_n$$

Con $w_n \in \mathbb{R}$, risulta che:

- $\sum w_i > 1 \Rightarrow$ Sovrastima
- $\sum w_i < 1 \Rightarrow$ Sottostima
- $\sum w_i = 1 \Rightarrow$ Non distorsione

Osservazioni utili alla **verifica** della **distorsione** di uno stimatore

- $T_1 = \frac{\sum x_i}{n} \rightarrow w_i = \frac{1}{n}$
- $E\left[\frac{\sum x_i}{n}\right] = \frac{1}{n} E[\sum x_i]$
- $E[X] = \mu$
- Per verificare la distorsione di uno stimatore T del parametro θ occorre valutare il suo $E[T]$ e verificare che questo sia pari esattamente a θ

2. NON DISTORSIONE ASINTOTICA

Dato un campione, uno stimatore T si dirà **asintoticamente non distorto** per un parametro θ se $E[T] \rightarrow \theta$ quando $n \rightarrow +\infty$ o, equivalentemente, se $D(T) \rightarrow 0$



3. EFFICIENZA

Quale stimatore scegliere se entrambi sono non distorti? Lo stimatore **più efficiente** è lo stimatore **con la varianza più piccola**:

$$\text{Var}[T_1] < \text{Var}[T_2] \rightarrow T_1 \text{ più efficiente di } T_2$$

È possibile anche definire l'**efficienza relativa di T_1 rispetto a T_2** come:

$$ER[T_1|T_2] = \frac{\text{Var}[T_2]}{\text{Var}[T_1]}$$

E T_1 si dirà più efficiente di T_2 quando $ER[T_1|T_2] > 1$

Stimatore più efficiente (stimatore BLUE)

Uno stimatore T^* non distorto per un parametro θ si dirà il più efficiente se:

$$\text{Var}[T^*] < \text{Var}[T]$$

Per qualunque altro stimatore non distorto T

4. STANDARD ERROR

Se T è uno stimatore non distorto per $\theta \rightarrow$ il suo standard error è dato:

$$\text{Std. Error} = \sqrt{\text{Var}[T]} = \sqrt{E[T - \theta]^2}$$

Esso rappresenta la deviazione standard della distribuzione campionario dello stimatore: ovvero, l'errore di stima

Esempi notevoli:

- **Standard error per media campionaria** come stimatore non distorto della media ($n < 0,05N$)

$$\text{Std. Error}(\bar{X}) = \begin{cases} \frac{\sigma}{\sqrt{n}} & \text{se } \sigma \text{ è nota} \\ \frac{s}{\sqrt{n}} & \text{se } \sigma \text{ è incognita} \end{cases}$$

Nel secondo caso ho **stimato la varianza con la varianza campionaria**

- Standard error per proporzione campionaria

$$\text{Std. Error}(\bar{P}) = \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

5. ERRORE QUADRATICO MEDIO

L'errore quadratico medio di uno stimatore T (potenzialmente distorto) è dato da:

$$EQM[T] = E[(T - \theta)^2] = \text{Var}[T] + D(T)^2$$

Quindi, se T è non distorto e $D(T)=0$, risulta:

$$EQM[T] = \text{Var}[T]$$



	Dati due stimatori, lo stimatore con minor EQM risulta essere quello maggiormente preciso ed è, quindi, da preferire
Distribuzione della media campionaria	<p>Proprietà La media campionaria è uno stimatore non distorto per la media</p> <p style="text-align: center;">$E[\bar{X}] = \mu$</p> <p>Dimostrazione:</p> $E[\bar{X}] = E\left[\frac{1}{n} \sum x_i\right] = \frac{1}{n} E\left[\sum x_i\right] = \frac{1}{n} \sum E[x_i] = \frac{1}{n} * n\mu = \mu$ <p>Distribuzione: <i>Vedi sopra (parte relativa al primo parziale)</i></p>
Distribuzione della proporzione campionaria	<i>Vedi sopra (parte relativa al primo parziale)</i>
Distribuzione della varianza campionaria	<p>Si consideri un campione casuale di ampiezza n da una popolazione con media μ e varianza σ^2. La varianza campionaria:</p> $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ <p>Proprietà:</p> <ul style="list-style-type: none"> • Per una qualsiasi popolazione arbitraria si dimostra che la varianza campionaria è uno stimatore non distorto per la varianza $\rightarrow E[S^2] = \sigma^2$ • Per una popolazione normale, in particolare, si può dimostrare che la varianza della varianza campionaria risulta: $Var[S^2] = \frac{2\sigma^4}{n-1}$ <p>È possibile, inoltre, ricondurre la varianza campionaria ad una distribuzione chi-quadrato con (n-1) gradi di libertà nel modo seguente</p> $\frac{(n-1)S^2}{\sigma^2} \sim \chi^2_{(n-1)}$ <p>E quindi, la varianza campionaria si distribuisce:</p> $S^2 \sim \frac{\sigma^2}{n-1} \cdot \chi^2_{(n-1)}$
Distribuzione chi-quadrato	<p>Se V ha una distribuzione chi-quadrato con k gradi di libertà, ovvero $V \sim \chi^2_{(k)}$ allora:</p> <ul style="list-style-type: none"> ▪ $E[V] = k$ ▪ $Var[V] = 2k$



INTERVALLI DI CONFIDENZA E STIMATORI PER INTERVALLO	
Stimatore per intervallo	<p>Uno stimatore per intervallo con livello di confidenza $1 - \alpha$ per un parametro θ della popolazione è una coppia di variabili aleatorie, A e B, funzioni del campione casuale, tali che:</p> <ul style="list-style-type: none"> ▪ $A < B$ ▪ $P(A < \theta < B) = 1 - \alpha \in (0, 1) \forall \theta$ <p>$1 - \alpha$ è detto livello di confidenza Un intervallo di confidenza è un intervallo numerico (a,b) calcolato in corrispondenza di una specifica realizzazione campionaria</p> $IC_{1-\alpha}(\theta) = (a, b)$ <p>Interpretazione frequentista $1 - \alpha$ è la frequenza di lungo periodo con cui parametro viene osservato nell'intervallo (a,b) in un gran numero di prove dell'esperimento</p> <p>Interpretazione dell'intervallo di confidenza Dire che, ad esempio per la media, $IC_{1-\alpha}(\mu) = (a, b)$ vuol dire che siamo confidenti al $1 - \alpha * 100$ (%) che il parametro θ cada in un range di valori compreso tra a e b (ATTENZIONE: è <i>scorretto</i> parlare di probabilità) Inoltre, NON siamo certi che il parametro cada all'interno dell'intervallo (ma abbiamo solo un certo grado di confidenza che ciò accada); siamo, invece, certi che la media campionaria cada nell'intervallo poiché è su essa centrato</p> <p>Esiste una relazione diretta tra livello di confidenza e ampiezza dell'intervallo: tanto maggiore è il livello di confidenza, tanto maggiore è l'ampiezza dell'intervallo</p>
Metodo delle quantità pivotali	<p>Il procedimento utilizzato per la costruzione di un intervallo di confidenza si basa su una quantità detta quantità pivotale o pivot Dato un campione iid da una distribuzione f_x con parametro θ, una quantità pivotale V è una variabile aleatoria:</p> <ul style="list-style-type: none"> ▪ Funzione $V = g(X_1, X_2, \dots, X_n; \theta)$ del campione e del parametro θ ▪ La distribuzione di V non dipende dal parametro incognito θ ▪ $V = g(X_1, X_2, \dots, X_n; \theta)$ invertibile
INTERVALLI DI CONFIDENZA PER LA MEDIA DELLA POPOLAZIONE	
Intervalli di confidenza per la media della popolazione	<p>Analizzeremo ora gli intervalli di confidenza per la media a partire da:</p> <ul style="list-style-type: none"> ▪ POPOLAZIONE NORMALE E... <ul style="list-style-type: none"> ○ VARIANZA NOTA ○ VARIANZA INCOGNITA (sempre, a meno che diversamente segnalato) ▪ POPOLAZIONE ARBITRARIA ▪ POPOLAZIONE BERNOULLIANA
Intervalli di confidenza per la media di una popolazione normale e varianza nota	<p>Consideriamo un campione casuale semplice da una popolazione normale e varianza nota σ^2</p> $X \sim N(\mu, \sigma^2)$ <p>Avremo che:</p> <ul style="list-style-type: none"> ▪ $\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$ ▪ Pivot: $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0,1)$ <div style="border: 1px solid black; background-color: yellow; padding: 5px; width: fit-content; margin: 10px auto;"> $IC_{1-\alpha}(\mu) = \left(\bar{X} \pm z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \right) = (LCL, UCL)$ </div>



	<ul style="list-style-type: none"> ▪ Ampiezza: $w = UCL - LCL = 2z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$ ▪ Margine di errore: $ME = \frac{w}{2} = z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$ <p>N.B.: ricorda di calcolare correttamente α; $z_{\frac{\alpha}{2}}$ è il quantile della $N(0,1)$ di ordine $1 - \frac{\alpha}{2}$</p> <p>Proprietà dell'intervallo:</p> <ul style="list-style-type: none"> ▪ Affidabilità $\rightarrow 1 - \alpha$ rappresenta il livello di confidenza, ovvero quanto siamo confidenti di trovare μ in quell'intervallo. Attenzione: siamo certi di trovare la media campionaria (è il centro dell'intervallo) ma non la media della popolazione (per la quale siamo solo confidenti al livello $1 - \alpha$ di trovarla) ▪ Precisione \rightarrow quanto sono precise le informazioni fornite dall'intervallo circa la media? Ci viene detto dal margine d'errore e dall'ampiezza dell'intervallo; <ul style="list-style-type: none"> ○ Nel caso di varianza nota, la precisione aumenta all'aumentare di n del campione e, quindi, l'intervallo diviene più piccolo ○ Nel caso di varianza incognita, questo non è più vero poiché, al variare di n, varia anche la varianza campionaria (e non sappiamo in che direzione): non vi è certezza ○ N.B.: anche per la bernoulliana non è vero dal momento che lo <i>standard error</i> dipende da P cappuccio <p>Le due condizioni rappresentano un trade off: all'aumentare dell'uno si riduce necessariamente l'altro</p> <p>Fattori che influenzano il margine d'errore:</p> $ME = z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$ <ul style="list-style-type: none"> ▪ Maggiore è σ (cioè la variabilità di \bar{X}) maggiore sarà ME ▪ Maggiore è n, minore è ME ▪ Maggiore è $1 - \alpha$, maggiore è ME (minore precisione) <p>Determinazione dell'ampiezza campionaria per ottenere ME desiderato Per un fissato livello di confidenza $1 - \alpha$ quale valore di n garantisce che si raggiunga un ME^* desiderato</p> $ME = z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \leq ME^*$ <p>Basta risolvere l'equazione per n e si trova facilmente:</p> $n \geq \frac{\left(z_{\frac{\alpha}{2}}\right)^2 \sigma^2}{ME^{*2}}$
Intervalli di confidenza per la media di una	<p>Consideriamo un campione casuale semplice da una popolazione normale varianza non nota σ^2</p> $X \sim N(\mu, \sigma^2)$

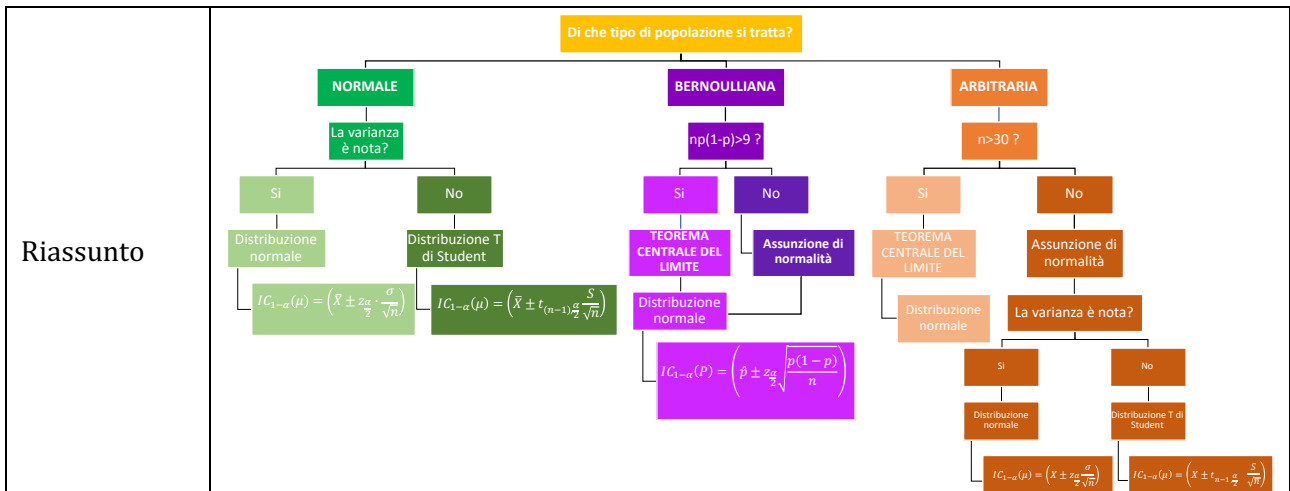


<p>popolazione normale e varianza non nota</p>	<p>Avremo che:</p> <ul style="list-style-type: none"> ▪ $\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$ ▪ Pivot: $Z = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim T_{(n-1)}$ <div style="border: 1px solid black; background-color: yellow; padding: 5px; margin: 10px 0;"> $IC_{1-\alpha}(\mu) = \left(\bar{X} \pm t_{(n-1), \frac{\alpha}{2}} \cdot \frac{S}{\sqrt{n}} \right) = (LCL, UCL)$ </div> <ul style="list-style-type: none"> ▪ Ampiezza: $w = UCL - LCL = 2t_{(n-1), \frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$ ▪ Margine di errore: $ME = \frac{w}{2} = t_{(n-1), \frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$ <p>$T_{(n-1)}$ è la distribuzione T di Student con n - 1 gradi di libertà</p> <p>Proprietà: se T ha una distribuzione T di Student con k gradi di libertà, e scriveremo $T \sim T_{(k)}$, allora:</p> $E[T] = 0 \text{ se } k > 1 \quad e \quad Var[T] = \frac{k}{k-2} \text{ se } k > 2$ <p>N.B.: nelle tavole della T di Student:</p> <ul style="list-style-type: none"> • α indica i quantili della distribuzione verso destra (indica la probabilità delle code) • v indicano i gradi di libertà <p style="background-color: yellow;">È L'UNICO CASO IN CUI NASCE LA T DI STUDENT</p>
<p>Intervalli di confidenza per la media di una popolazione arbitraria</p>	<p>Consideriamo un campione casuale semplice da una popolazione arbitraria e varianza non nota. Sia $n \geq 30$. Avremo:</p> <ul style="list-style-type: none"> ▪ $\bar{X} \approx N\left(\mu, \frac{\sigma^2}{n}\right)$ in base al teorema centrale del limite ▪ Pivot: $Z = \frac{\bar{X} - \mu}{S/\sqrt{n}} \approx N(0,1)$ ▪ Intervallo di confidenza di livello (approssimativamente) $1 - \alpha$ <div style="border: 1px solid black; background-color: yellow; padding: 5px; margin: 10px 0;"> $IC_{1-\alpha}(\mu) = \left(\bar{X} \pm z_{\frac{\alpha}{2}} \cdot \frac{S}{\sqrt{n}} \right) = (LCL, UCL)$ </div> <ul style="list-style-type: none"> ▪ Ampiezza: $w = UCL - LCL = 2z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$ ▪ Margine di errore: $ME = \frac{w}{2} = z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$
<p>Intervalli di confidenza per la proporzione campionaria</p>	<p>Consideriamo un campione casuale semplice da una popolazione bernoulliana di parametro $p \in [0, 1]$. Sia $n \geq 30$ (sufficientemente elevato). Avremo:</p> <ul style="list-style-type: none"> ▪ $\hat{P} \approx N\left(p, \frac{p(1-p)}{n}\right)$ in base al teorema centrale del limite ▪ Pivot: $Z = \frac{\hat{P} - p}{\sqrt{\frac{\hat{P}(1-\hat{P})}{n}}} \approx N(0,1)$ ▪ Intervallo di confidenza di livello (approssimativamente) $1 - \alpha$ <div style="border: 1px solid black; background-color: yellow; padding: 5px; margin: 10px 0;"> $IC_{1-\alpha}(p) = \left(\hat{P} \pm z_{\frac{\alpha}{2}} \cdot \sqrt{\frac{\hat{P}(1-\hat{P})}{n}} \right) = (LCL, UCL)$ </div>



	<ul style="list-style-type: none"> Ampiezza: $w = UCL - LCL = 2z_{\alpha/2} \sqrt{\frac{\hat{P}(1-\hat{P})}{n}}$ Margine di errore: $ME = \frac{w}{2} = z_{\alpha/2} \sqrt{\frac{\hat{P}(1-\hat{P})}{n}}$ <p>Determinazione dell'ampiezza campionaria per ottenere ME desiderato Per un fissato livello di confidenza $1 - \alpha$ quale valore di n garantisce che si raggiunga un ME^* desiderato</p> $ME = z_{\alpha/2} \sqrt{\frac{\hat{P}(1 - \hat{P})}{n}} \leq ME^*$ <p>Si dimostra</p> $n \geq \frac{\left(\frac{z_{\alpha}}{2}\right)^2 \cdot 0.25}{ME^{*2}}$ <p>(da arrotondare sempre per eccesso)</p> <p>0,25 rappresenta valore massimo che $\hat{P}(1 - \hat{P})$ può raggiungere: nè il punto di massimo della parabola $\hat{P}(1 - \hat{P}) = 0$ nel piano $[\hat{P}, \hat{P}(1 - \hat{P})]$</p>				
<p>Intervalli di confidenza con assunzione di normalità</p>	<p>Consideriamo un campione casuale semplice da una popolazione arbitraria. Sia $n < 30$: n non è abbastanza grande per usare il teorema centrale del limite Per poter procedere occorre assumere la normalità del campione; così facendo ci riconduciamo al caso di popolazione normale e varianza nota/non nota a seconda del caso richiesto:</p> <table border="1" data-bbox="375 1243 1404 1467"> <tr> <td data-bbox="375 1243 885 1355"> <p>Varianza nota</p> </td> <td data-bbox="885 1243 1404 1355"> $IC_{1-\alpha}(\mu) = \left(\bar{X} \pm z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \right)$ </td> </tr> <tr> <td data-bbox="375 1355 885 1467"> <p>Varianza non nota</p> </td> <td data-bbox="885 1355 1404 1467"> $IC_{1-\alpha}(\mu) = \left(\bar{X} \pm t_{(n-1), \alpha/2} \cdot \frac{S}{\sqrt{n}} \right)$ </td> </tr> </table>	<p>Varianza nota</p>	$IC_{1-\alpha}(\mu) = \left(\bar{X} \pm z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \right)$	<p>Varianza non nota</p>	$IC_{1-\alpha}(\mu) = \left(\bar{X} \pm t_{(n-1), \alpha/2} \cdot \frac{S}{\sqrt{n}} \right)$
<p>Varianza nota</p>	$IC_{1-\alpha}(\mu) = \left(\bar{X} \pm z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \right)$				
<p>Varianza non nota</p>	$IC_{1-\alpha}(\mu) = \left(\bar{X} \pm t_{(n-1), \alpha/2} \cdot \frac{S}{\sqrt{n}} \right)$				
<p>Fattore di correzione per popolazioni finte</p>	<p>Occorre ricordare che qualora la popolazione presentasse un numero di individui finito N occorre applicare allo standard errore il fattore di correzione</p> $n > 0.05N \Rightarrow SD * \frac{N - n}{N - 1}$ <p>Questo risultato vale sempre ed è bene ricordarlo!</p>				





Riassunto

INTERVALLI DI CONFIDENZA PER LA DIFFERENZA TRA MEDIE DI DUE POPOLAZIONI

Intervalli di confidenza per la differenza tra le medie di due popolazioni

Consideriamo due popolazioni; vogliamo confrontarle in termini di media di una variabile di interesse

- POPOLAZIONE 1

$$X_1, \dots, X_{n_x} \sim iid f_x \quad E[X] = \mu_x \quad Var[X] = \sigma_x^2$$

- POPOLAZIONE 2

$$Y_1, \dots, Y_{n_y} \sim iid f_y \quad E[Y] = \mu_y \quad Var[Y] = \sigma_y^2$$

Diamo due importanti definizioni:

- **CAMPIONI DIPENDENTI:** le stesse unità vengono analizzate prima e dopo il trattamento $\rightarrow Cov(X_i, Y_i) \neq 0$
- **CAMPIONI INDIPENDENTI:** unità (potenzialmente) differenti vengono analizzate prima e dopo il trattamento $\rightarrow Cov(X_i, Y_i) = 0$

Analizzeremo i seguenti casi:

- **CAMPIONI DIPENDENTI** (varianza campionaria S_D^2)
 - POPOLAZIONE NORMALE E VARIANZE NON NOTE E DIVERSE
 - POPOLAZIONE ARBITRARIA ($n > 30$) E VARIANZE NON NOTE MA ASSUNTE UGUALI
 - POPOLAZIONE ARBITRARIA ($n > 30$) E VARIANZE NON NOTE E DIVERSE (SOFTWARE R)
- **CAMPIONI INDIPENDENTI** (varianza pooled S_p^2)
 - POPOLAZIONE NORMALE
 - VARIANZE INCOGNITE MA ASSUNTE UGUALI
 - VARIANZE INCOGNITE E DIVERSE (SOFTWARE R)

Intervalli di confidenza per la differenza tra le medie di due popolazioni normali dipendenti (varianze non note)

Consideriamo due campioni casuali semplici **dipendenti** ($Cov(X_i, Y_i) \neq 0$) da **popolazioni normali** (alternativamente, assumiamo la normalità) **varianza non nota** σ^2

Definito $D_i = X_i - Y_i$ avremo che:

$$D_1, \dots, D_n \sim N(\mu_D, \sigma_D^2) \text{ con}$$

$$\mu_D = \mu_x - \mu_y$$

$$S_D^2 = S_x^2 + S_y^2 + (-2r_{xy})$$

- Stimatore di $\mu_D \rightarrow \bar{D} = \frac{1}{n} \sum D_i \sim N\left(\mu_D, \frac{\sigma_D^2}{n}\right)$



	<ul style="list-style-type: none"> • Pivot $\rightarrow T = \frac{\bar{D} - \mu_D}{S_D/\sqrt{n}} \sim T_{(n-1)}$ • Intervallo di confidenza di livello $1 - \alpha$ $IC_{1-\alpha}(\mu_D) = IC_{1-\alpha}(\mu_X - \mu_Y) = \left(\bar{D} \pm t_{(n-1), \frac{\alpha}{2}} \frac{S_D}{\sqrt{n}} \right)$ <p>Proprietà di D_i</p> <ul style="list-style-type: none"> • $E[D_i] = E[X_i - Y_i] = E[X_i] - E[Y_i] = \mu_X - \mu_Y = \mu_D$ • $Var[D_i] = Var[X_i - Y_i] = Var[X_i] + Var[Y_i] - 2Cov[X_i, Y_i]$ <p>Proprietà dello stimatore \bar{D}</p> <ul style="list-style-type: none"> • $E[\bar{D}] = \mu_D$ • $Var[\bar{D}] = \frac{\sigma_D^2}{n} = \frac{\sigma_X^2 + \sigma_Y^2 - 2\sigma_{XY}}{n}$
<p>Intervallo di confidenza per la differenza tra le medie di due popolazioni normali indipendenti (varianze note)</p>	<p>Consideriamo due campioni casuali semplici indipendenti ($Cov(X_i, Y_i) = 0$) da popolazioni normali e varianza nota σ^2</p> <ul style="list-style-type: none"> • Stimatore di $\Delta\mu = \mu_X - \mu_Y \rightarrow \bar{X} - \bar{Y} \sim N\left(\mu_X - \mu_Y, \frac{\sigma_X^2}{n_X} + \frac{\sigma_Y^2}{n_Y}\right)$ • Pivot $\rightarrow Z = \frac{\bar{X} - \bar{Y} - (\mu_X - \mu_Y)}{\sqrt{\frac{\sigma_X^2}{n_X} + \frac{\sigma_Y^2}{n_Y}}} \sim N(0,1)$ • Intervallo di confidenza di livello $1 - \alpha$ $IC_{1-\alpha}(\mu_X - \mu_Y) = \left((\bar{X} - \bar{Y}) \pm z_{\frac{\alpha}{2}} \sqrt{\frac{\sigma_X^2}{n_X} + \frac{\sigma_Y^2}{n_Y}} \right)$ <p>Proprietà di $\bar{X} - \bar{Y}$</p> <ul style="list-style-type: none"> • $E[\bar{X} - \bar{Y}] = E[\bar{X}] - E[\bar{Y}]$ • $Var[\bar{X} - \bar{Y}] = Var[\bar{X}] + Var[\bar{Y}]$ <p>N.B.: se una popolazione fosse arbitraria e n non abbastanza grande per usare il teorema centrale del limite, dobbiamo assumere la normalità</p>
<p>Intervallo di confidenza per la differenza tra le medie di due popolazioni normali indipendenti (varianze non note ma assunte uguali)</p>	<p>Consideriamo due campioni casuali semplici indipendenti ($Cov(X_i, Y_i) = 0$) da popolazioni normali con media μ non nota e varianze non note ma assunte uguali</p> <ul style="list-style-type: none"> • Considerando che $\sigma_X^2 = \sigma_Y^2$, le due varianze campionarie sono stimatori naturali e non distorti; tuttavia, si dimostra che la varianza pooled è uno stimatore non distorto e più efficiente delle varianze campionarie <i>semplici</i> $S_P^2 = \frac{(n_X - 1)S_X^2 + (n_Y - 1)S_Y^2}{n_X + n_Y - 2}$ <ul style="list-style-type: none"> • Possiamo, quindi, stimare la varianza di $\bar{X} - \bar{Y}$: $Var[\bar{X} - \bar{Y}]_{stima} = S_P^2 \left(\frac{1}{n_X} + \frac{1}{n_Y} \right)$ <ul style="list-style-type: none"> • Pivot $\rightarrow T = \frac{\bar{X} - \bar{Y} - (\mu_X - \mu_Y)}{S_P \sqrt{\frac{1}{n_X} + \frac{1}{n_Y}}} \sim T_{(n_X + n_Y - 2)}$ • Intervallo di confidenza di livello $1 - \alpha$



	$IC_{1-\alpha}(\mu_X - \mu_Y) = \left((\bar{X} - \bar{Y}) \pm t_{n_X+n_Y-2, \frac{\alpha}{2}} S_P \sqrt{\frac{1}{n_X} + \frac{1}{n_Y}} \right)$
<p>Intervalli di confidenza per la differenza tra le medie di due popolazioni normali indipendenti (varianze non note e differenti)</p>	<p>Consideriamo due campioni casuali semplici indipendenti ($Cov(X_i, Y_i) = 0$) da popolazioni normali con media μ non nota e varianze non note e differenti</p> <ul style="list-style-type: none"> Possiamo, quindi, stimare la varianza di $\bar{X} - \bar{Y}$: $Var[\bar{X} - \bar{Y}]_{stima} = \frac{S_X^2}{n_X} + \frac{S_Y^2}{n_Y}$ Pivot $\rightarrow Z = \frac{\bar{X} - \bar{Y} - (\mu_X - \mu_Y)}{\sqrt{\frac{S_X^2}{n_X} + \frac{S_Y^2}{n_Y}}} \approx T_{(v)}$ Intervallo di confidenza di livello approssimativamente pari a $1 - \alpha$ $IC_{1-\alpha}(\mu_X - \mu_Y) = \left((\bar{X} - \bar{Y}) \pm t_{v, \frac{\alpha}{2}} \sqrt{\frac{S_X^2}{n_X} + \frac{S_Y^2}{n_Y}} \right)$
<p>Intervalli di confidenza di popolazioni arbitrarie dipendenti (grandi campioni) con varianze non note ma uguali</p>	<p>Consideriamo due campioni sufficientemente grandi ($n \geq 30$) casuali semplici dipendenti ($Cov(X_i, Y_i) \neq 0$) da popolazioni arbitrarie e varianze non note ma uguali</p> <ul style="list-style-type: none"> $\bar{D} = \frac{1}{n} \sum D_i \sim N\left(\mu_D, \frac{\sigma_D^2}{n}\right)$ $Z = \frac{\bar{D} - \mu_D}{S_D / \sqrt{n}} \approx N(0,1)$ Intervallo di confidenza di livello $1 - \alpha$ $IC_{1-\alpha}(\mu_D) = IC_{1-\alpha}(\mu_X - \mu_Y) = \left(\bar{D} \pm z_{\frac{\alpha}{2}} \frac{S_D}{\sqrt{n}} \right)$
<p>Intervalli di confidenza di popolazioni arbitrarie indipendenti (grandi campioni) con varianze non note e differenti</p>	<p>Consideriamo due campioni sufficientemente grandi ($n \geq 30$) casuali semplici indipendenti ($Cov(X_i, Y_i) = 0$) da popolazioni arbitrarie e varianze non note e differenti</p> <ul style="list-style-type: none"> Pivot $\rightarrow Z = \frac{\bar{X} - \bar{Y} - (\mu_X - \mu_Y)}{\sqrt{\frac{S_X^2}{n_X} + \frac{S_Y^2}{n_Y}}} \approx N(0,1)$ Intervallo di confidenza di livello approssimativamente pari a $1 - \alpha$ $IC_{1-\alpha}(\mu_X - \mu_Y) = \left((\bar{X} - \bar{Y}) \pm z_{\frac{\alpha}{2}} \sqrt{\frac{S_X^2}{n_X} + \frac{S_Y^2}{n_Y}} \right)$
<p>Intervalli di confidenza di popolazioni arbitrarie (campioni non sufficientemente grandi)</p>	<p>Assunzioni:</p> <ul style="list-style-type: none"> Campioni indipendenti Varianze (non note) ma assunte uguali Normalità della popolazione <p>Mi riconduco al caso di campioni indipendenti con varianze non note ma assunte uguali \rightarrow calcolo della varianza pooled</p>



	$IC_{1-\alpha}(\mu_X - \mu_Y) = \left((\bar{X} - \bar{Y}) \pm t_{n_X+n_Y-2, \frac{\alpha}{2}} S_P \sqrt{\frac{1}{n_X} + \frac{1}{n_Y}} \right)$
<p>Riassunto</p>	<pre> graph TD A[Come sono i campioni] --> B[DIPENDENTI le stesse unità vengono analizzate prima e dopo un trattamento] A --> C[INDIPENDENTI unità differenti vengono analizzate prima e dopo un trattamento] B --> B1["Varianza campionaria S_D^2 = S_X^2 + S_Y^2 + 2Cov(X,Y)"] C --> C1["Varianza pooled S_P^2 = ((n_X - 1)S_X^2 + (n_Y - 1)S_Y^2) / (n_X + n_Y - 2)"] D[La varianza della popolazione è nota?] --> E[VARIANZA NOTA] D --> F[VARIANZA INCOGNITA] E --> E1[Distribuzione normale (Z_α/2)] F --> F1[Distribuzione T di Student (t_{n_X+n_Y-2, α/2})] G[Com'è distribuita la popolazione?] --> H[Popolazione normale] G --> I[Popolazione arbitraria] I --> J["np(1-p) > 9?"] J --> K[SI -> Teorema Centrale del Limite] J --> L[NO -> Assunzioni] L --> L1[Campioni indipendenti Varianze incognite ma assunte uguali] </pre>

TEST D'IPOTESI	
<p>Concetti di base dell'ipotesi statistica</p>	<p>Si definisce ipotesi statistica un'affermazione circa il valore di un parametro incognito della popolazione</p> <p>Tipi di ipotesi:</p> <ul style="list-style-type: none"> • Semplice: specifica un singolo valore per il parametro della popolazione considerato • Composta: specifica uno o più intervalli di valori per il parametro della popolazione considerato <ul style="list-style-type: none"> ○ Unilaterale: considera tutti i possibili valori del parametro della popolazione a destra o a sinistra rispetto all'ipotesi fatta ○ Bilaterale: considera tutti i possibili valori diversi dall'ipotesi nulla, <p>Ipotesi nulla e alternativa</p> <ul style="list-style-type: none"> • IPOTESI NULLA H₀ → ipotesi da considerarsi vera a meno di ottenere prove evidenti della validità del suo contrario (rappresenta lo status quo) • IPOTESI ALTERNATIVA H₁ → ipotesi contro la quale viene verificata l'ipotesi nulla e che viene considerata vera se l'ipotesi nulla è considerata falsa (affermazione da porre a verifica e per la quale si ricerca l'evidenza) <p>Relazione tra test d'ipotesi bilaterali e intervalli di confidenza</p> <p>Un intervallo di confidenza identifica una coppia di valori a e b che con una confidenza del 100(1 - α)% contiene il valore del parametro θ.</p> <p>Se l'intervallo include il 100(1 - α)% di valori più plausibili per θ, esso esclude il 100α% di campioni meno plausibili.</p> <p>Se il valore posto in ipotesi nulla θ₀ è fuori dall'intervallo perché inferiore ad a o superiore a b sarà considerato un valore non verosimile rispetto alla realizzazione campionaria</p>
<p>Decisione, statistica test ed errori</p>	<p>Il processo decisionale usa una statistica test costruita a partire dagli stimatori puntuali del parametro oggetto di studio; essa assume una distribuzione campionaria nota: questo ci permette di determinare i valori della statistica test che avrebbero una bassa probabilità di verificarsi se l'ipotesi nulla fosse vera.</p> <p>Se la statistica test assume uno tra questi valori, rifiutiamo l'ipotesi nulla; altrimenti non rifiutiamo l'ipotesi nulla</p>



	<p>Per scegliere tra ipotesi nulla e ipotesi alternativa dobbiamo sviluppare una <i>regola di decisione</i> chiamata test per verificare un'ipotesi contro l'altra Tuttavia, possiamo anche commettere degli errori:</p> <table border="1" data-bbox="383 324 1412 537"> <thead> <tr> <th rowspan="2"></th> <th colspan="2">IPOTESI VERA</th> </tr> <tr> <th>H₀</th> <th>H₁</th> </tr> </thead> <tbody> <tr> <th>RIFIUTO H₀</th> <td style="color: red;">ERRORE DI PRIMA SPECIE</td> <td style="color: green;">NESSUN ERRORE</td> </tr> <tr> <th>NON RIFIUTO H₀</th> <td style="color: green;">NESSUN ERRORE</td> <td style="color: red;">ERRORE DI SECONDA SPECIE</td> </tr> </tbody> </table> <p>Quindi diremo di aver commesso:</p> <ul style="list-style-type: none"> • Errore di I specie: rifiuto H₀ quando è vera H₀ <ul style="list-style-type: none"> ○ La nostra regola di decisione sarà definita in modo da ottenere una probabilità di commettere l'errore di primo tipo pari ad α e detta livello di significatività $\alpha = P(\text{rifiutare } H_0 H_0 \text{ vera} = h_0)$ <ul style="list-style-type: none"> ○ I livelli di significatività usuali sono $\alpha = (0.05, 0.01, 0.10)$ • Errore di II specie: rifiuto H₁ quando è vera H₁ <ul style="list-style-type: none"> ○ La probabilità che si verifichi l'errore di secondo tipo è indicata con β ○ $1 - \beta$ viene detta potenza del test ed è la probabilità di non rifiutare H₁ quando H₁ è vera, ovvero di commettere l'errore di secondo tipo $\beta = P(\text{non rifiutare } H_0 H_1 \text{ vera} = h_1)$ <p>Esiste un trade off tra queste due probabilità: ogni riduzione della probabilità dell'errore di primo tipo α implica un aumento della probabilità dell'errore di secondo tipo β, e viceversa</p>		IPOTESI VERA		H ₀	H ₁	RIFIUTO H ₀	ERRORE DI PRIMA SPECIE	NESSUN ERRORE	NON RIFIUTO H ₀	NESSUN ERRORE	ERRORE DI SECONDA SPECIE
	IPOTESI VERA											
	H ₀	H ₁										
RIFIUTO H ₀	ERRORE DI PRIMA SPECIE	NESSUN ERRORE										
NON RIFIUTO H ₀	NESSUN ERRORE	ERRORE DI SECONDA SPECIE										
<p>Problemi tipo nella verifica di ipotesi per la media di una popolazione</p>	<p>Sia μ_0 il valore <i>noto</i> della media che vogliamo testare. Distinguiamo tre casi:</p> <table border="1" data-bbox="383 1377 1412 1489"> <thead> <tr> <th>CASO I</th> <th>CASO II</th> <th>CASO III</th> </tr> </thead> <tbody> <tr> <td>$H_0: \mu = \mu_0$ $H_1: \mu \neq \mu_0$</td> <td>$H_0: \mu \leq \mu_0$ $H_1: \mu > \mu_0$</td> <td>$H_0: \mu \geq \mu_0$ $H_1: \mu < \mu_0$</td> </tr> </tbody> </table>	CASO I	CASO II	CASO III	$H_0: \mu = \mu_0$ $H_1: \mu \neq \mu_0$	$H_0: \mu \leq \mu_0$ $H_1: \mu > \mu_0$	$H_0: \mu \geq \mu_0$ $H_1: \mu < \mu_0$					
CASO I	CASO II	CASO III										
$H_0: \mu = \mu_0$ $H_1: \mu \neq \mu_0$	$H_0: \mu \leq \mu_0$ $H_1: \mu > \mu_0$	$H_0: \mu \geq \mu_0$ $H_1: \mu < \mu_0$										
<p><i>p-value</i></p>	<p>Definiamo livello di significatività osservato o p-value la probabilità di ottenere un valore della statistica test più estremo del valore osservato, ovvero quel livello di significatività al quale l'ipotesi nulla può essere rifiutata, dato il valore osservato della statistica campionaria L'importanza del p-value è legata al fatto che esso fornisce l'esatta probabilità di rifiuto dell'ipotesi nulla (livello di significatività al quale rifiutiamo l'ipotesi nulla) derivante dall'osservazione dei dati campionari. Il p-value ci dice <i>quanto lontani o vicini siamo dal rifiutare l'ipotesi nulla</i></p>											
<p>Regione di rifiuto e regione di accettazione</p>	<p>La regione di rifiuto, complementare a quella di accettazione, è una regione indicabile con:</p> $R = \{x_1, x_2, \dots, x_n: \text{regola di rifiuto di } H_0\}$ <p>E, analogamente, la regione di accettazione si indica con:</p>											



	$R = \{x_1, x_2, \dots, x_n: \text{regola di non rifiuto di } H_0\}$
--	--

VERIFICA DI IPOTESI PER LA MEDA DI UNA POPOLAZIONE															
Verifica di ipotesi sulla media di una popolazione normale con varianza nota → distribuzione normale	<table border="1" style="width: 100%;"> <tr> <td style="width: 30%;">Statistica test</td> <td style="text-align: center;">$Z = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} \sim \mu = \mu_0 N(0,1)$</td> </tr> <tr> <td colspan="2" style="text-align: center;"> CASO I $H_0: \mu = \mu_0$ $H_1: \mu \neq \mu_0$ </td> </tr> <tr> <td>Regola di rifiuto</td> <td> <ul style="list-style-type: none"> • $\bar{X} < \mu_0 - z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$ oppure $\bar{X} > \mu_0 + z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$ • $Z = \left \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} \right > z_{\frac{\alpha}{2}}$ • $p - \text{value} = 2P(Z > z_{\text{oss}}) < \alpha$ </td> </tr> <tr> <td colspan="2" style="text-align: center;"> CASO II $H_0: \mu \leq \mu_0$ $H_1: \mu > \mu_0$ </td> </tr> <tr> <td>Regola di rifiuto</td> <td> <ul style="list-style-type: none"> • $\bar{X} > \mu_0 + z_{\alpha} \frac{\sigma}{\sqrt{n}}$ • $Z = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} > z_{\alpha}$ • $p - \text{value} = P(Z > z_{\text{oss}}) < \alpha$ </td> </tr> <tr> <td colspan="2" style="text-align: center;"> CASO III $H_0: \mu \geq \mu_0$ $H_1: \mu < \mu_0$ </td> </tr> <tr> <td>Regola di rifiuto</td> <td> <ul style="list-style-type: none"> • $\bar{X} < \mu_0 - z_{\alpha} \frac{\sigma}{\sqrt{n}}$ • $Z = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} < -z_{\alpha}$ • $p - \text{value} = P(Z < z_{\text{oss}}) < \alpha$ </td> </tr> </table>	Statistica test	$Z = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} \sim \mu = \mu_0 N(0,1)$	CASO I $H_0: \mu = \mu_0$ $H_1: \mu \neq \mu_0$		Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{X} < \mu_0 - z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$ oppure $\bar{X} > \mu_0 + z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$ • $Z = \left \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} \right > z_{\frac{\alpha}{2}}$ • $p - \text{value} = 2P(Z > z_{\text{oss}}) < \alpha$ 	CASO II $H_0: \mu \leq \mu_0$ $H_1: \mu > \mu_0$		Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{X} > \mu_0 + z_{\alpha} \frac{\sigma}{\sqrt{n}}$ • $Z = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} > z_{\alpha}$ • $p - \text{value} = P(Z > z_{\text{oss}}) < \alpha$ 	CASO III $H_0: \mu \geq \mu_0$ $H_1: \mu < \mu_0$		Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{X} < \mu_0 - z_{\alpha} \frac{\sigma}{\sqrt{n}}$ • $Z = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} < -z_{\alpha}$ • $p - \text{value} = P(Z < z_{\text{oss}}) < \alpha$
	Statistica test	$Z = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} \sim \mu = \mu_0 N(0,1)$													
	CASO I $H_0: \mu = \mu_0$ $H_1: \mu \neq \mu_0$														
	Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{X} < \mu_0 - z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$ oppure $\bar{X} > \mu_0 + z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$ • $Z = \left \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} \right > z_{\frac{\alpha}{2}}$ • $p - \text{value} = 2P(Z > z_{\text{oss}}) < \alpha$ 													
	CASO II $H_0: \mu \leq \mu_0$ $H_1: \mu > \mu_0$														
	Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{X} > \mu_0 + z_{\alpha} \frac{\sigma}{\sqrt{n}}$ • $Z = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} > z_{\alpha}$ • $p - \text{value} = P(Z > z_{\text{oss}}) < \alpha$ 													
	CASO III $H_0: \mu \geq \mu_0$ $H_1: \mu < \mu_0$														
Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{X} < \mu_0 - z_{\alpha} \frac{\sigma}{\sqrt{n}}$ • $Z = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} < -z_{\alpha}$ • $p - \text{value} = P(Z < z_{\text{oss}}) < \alpha$ 														
Verifica di ipotesi sulla media di una popolazione normale con varianza non nota → distribuzione T di Student	<table border="1" style="width: 100%;"> <tr> <td style="width: 30%;">Statistica test</td> <td style="text-align: center;">$T = \frac{\bar{X} - \mu_0}{\frac{S}{\sqrt{n}}} \sim \mu = \mu_0 T_{(n-1)}$</td> </tr> <tr> <td colspan="2" style="text-align: center;"> CASO I $H_0: \mu = \mu_0$ $H_1: \mu \neq \mu_0$ </td> </tr> <tr> <td>Regola di rifiuto</td> <td> <ul style="list-style-type: none"> • $\bar{X} < \mu_0 - t_{n-1, \frac{\alpha}{2}} \frac{S}{\sqrt{n}}$ oppure $\bar{X} > \mu_0 + t_{n-1, \frac{\alpha}{2}} \frac{S}{\sqrt{n}}$ • $T = \left \frac{\bar{X} - \mu_0}{\frac{S}{\sqrt{n}}} \right > t_{n-1, \frac{\alpha}{2}}$ • $p - \text{value} = 2P(T > t_{\text{oss}}) < \alpha$ </td> </tr> <tr> <td colspan="2" style="text-align: center;"> CASO II $H_0: \mu \leq \mu_0$ $H_1: \mu > \mu_0$ </td> </tr> <tr> <td>Regola di rifiuto</td> <td> <ul style="list-style-type: none"> • $\bar{X} > \mu_0 + t_{n-1, \alpha} \frac{S}{\sqrt{n}}$ • $Z = \frac{\bar{X} - \mu_0}{\frac{S}{\sqrt{n}}} > t_{n-1, \alpha}$ • $p - \text{value} = P(T > t_{\text{oss}}) < \alpha$ </td> </tr> <tr> <td colspan="2" style="text-align: center;"> CASO III $H_0: \mu \geq \mu_0$ $H_1: \mu < \mu_0$ </td> </tr> </table>	Statistica test	$T = \frac{\bar{X} - \mu_0}{\frac{S}{\sqrt{n}}} \sim \mu = \mu_0 T_{(n-1)}$	CASO I $H_0: \mu = \mu_0$ $H_1: \mu \neq \mu_0$		Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{X} < \mu_0 - t_{n-1, \frac{\alpha}{2}} \frac{S}{\sqrt{n}}$ oppure $\bar{X} > \mu_0 + t_{n-1, \frac{\alpha}{2}} \frac{S}{\sqrt{n}}$ • $T = \left \frac{\bar{X} - \mu_0}{\frac{S}{\sqrt{n}}} \right > t_{n-1, \frac{\alpha}{2}}$ • $p - \text{value} = 2P(T > t_{\text{oss}}) < \alpha$ 	CASO II $H_0: \mu \leq \mu_0$ $H_1: \mu > \mu_0$		Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{X} > \mu_0 + t_{n-1, \alpha} \frac{S}{\sqrt{n}}$ • $Z = \frac{\bar{X} - \mu_0}{\frac{S}{\sqrt{n}}} > t_{n-1, \alpha}$ • $p - \text{value} = P(T > t_{\text{oss}}) < \alpha$ 	CASO III $H_0: \mu \geq \mu_0$ $H_1: \mu < \mu_0$			
	Statistica test	$T = \frac{\bar{X} - \mu_0}{\frac{S}{\sqrt{n}}} \sim \mu = \mu_0 T_{(n-1)}$													
	CASO I $H_0: \mu = \mu_0$ $H_1: \mu \neq \mu_0$														
	Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{X} < \mu_0 - t_{n-1, \frac{\alpha}{2}} \frac{S}{\sqrt{n}}$ oppure $\bar{X} > \mu_0 + t_{n-1, \frac{\alpha}{2}} \frac{S}{\sqrt{n}}$ • $T = \left \frac{\bar{X} - \mu_0}{\frac{S}{\sqrt{n}}} \right > t_{n-1, \frac{\alpha}{2}}$ • $p - \text{value} = 2P(T > t_{\text{oss}}) < \alpha$ 													
	CASO II $H_0: \mu \leq \mu_0$ $H_1: \mu > \mu_0$														
Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{X} > \mu_0 + t_{n-1, \alpha} \frac{S}{\sqrt{n}}$ • $Z = \frac{\bar{X} - \mu_0}{\frac{S}{\sqrt{n}}} > t_{n-1, \alpha}$ • $p - \text{value} = P(T > t_{\text{oss}}) < \alpha$ 														
CASO III $H_0: \mu \geq \mu_0$ $H_1: \mu < \mu_0$															



		Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{X} < \mu_0 - t_{n-1,\alpha} \frac{S}{\sqrt{n}}$ • $Z = \frac{\bar{X} - \mu_0}{\frac{S}{\sqrt{n}}} < -t_{n-1,\alpha}$ • $p - value = P(T < t_{oss}) < \alpha$ 	
<p>Verifica di ipotesi sulla media di una popolazione arbitraria con varianza non nota → grandi campioni (n>30)</p>	<p>Per $n \geq 30$ applichiamo il Teorema centrale del limite e impieghiamo una distribuzione normale</p>			
	Statistica test	$Z = \frac{\bar{X} - \mu_0}{\frac{S}{\sqrt{n}}} \approx^{\mu=\mu_0} N(0,1)$		
	<p>CASO I $H_0: \mu = \mu_0$ $H_1: \mu \neq \mu_0$</p>			
	Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{X} < \mu_0 - z_{\frac{\alpha}{2}} \frac{S}{\sqrt{n}}$ oppure $\bar{X} > \mu_0 + z_{\frac{\alpha}{2}} \frac{S}{\sqrt{n}}$ • $Z = \left \frac{\bar{X} - \mu_0}{\frac{S}{\sqrt{n}}} \right > z_{\frac{\alpha}{2}}$ • $p - value = 2P(Z > z_{oss}) < \alpha$ 		
	<p>CASO II $H_0: \mu \leq \mu_0$ $H_1: \mu > \mu_0$</p>			
	Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{X} > \mu_0 + z_{\alpha} \frac{S}{\sqrt{n}}$ • $Z = \frac{\bar{X} - \mu_0}{\frac{S}{\sqrt{n}}} > z_{\alpha}$ • $p - value = P(Z > z_{oss}) < \alpha$ 		
	<p>CASO III $H_0: \mu \geq \mu_0$ $H_1: \mu < \mu_0$</p>			
Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{X} < \mu_0 - z_{\alpha} \frac{S}{\sqrt{n}}$ • $Z = \frac{\bar{X} - \mu_0}{\frac{S}{\sqrt{n}}} < -z_{\alpha}$ • $p - value = P(Z < z_{oss}) < \alpha$ 			



<p>Verifica di ipotesi sulla media di una proporzione con varianza non nota → grandi campioni ($n > 30$)</p>	<p>Se n è estratto da una popolazione grande a sufficienza da poter affermare che $np(1-p) > 9$ allora vale il teorema centrale del limite e applichiamo una distribuzione normale</p>		
	<table border="1"> <tr> <td data-bbox="370 324 662 443">Statistica test</td> <td data-bbox="662 324 1444 443"> $Z = \frac{\bar{P} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \approx^{p=p_0} N(0,1)$ </td> </tr> </table>	Statistica test	$Z = \frac{\bar{P} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \approx^{p=p_0} N(0,1)$
	Statistica test	$Z = \frac{\bar{P} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \approx^{p=p_0} N(0,1)$	
	<p style="text-align: center;">CASO I $H_0: \mu = \mu_0$ $H_1: \mu \neq \mu_0$</p>		
	Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{P} < p_0 - z_{\frac{\alpha}{2}} \sqrt{\frac{p_0(1-p_0)}{n}}$ oppure $\bar{X} > p_0 + z_{\frac{\alpha}{2}} \sqrt{\frac{p_0(1-p_0)}{n}}$ • $Z = \left \frac{\bar{P} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \right > z_{\frac{\alpha}{2}}$ • $p - value = 2P(Z > z_{OSS}) < \alpha$ 	
	<p style="text-align: center;">CASO II $H_0: \mu \leq \mu_0$ $H_1: \mu > \mu_0$</p>		
Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{P} > p_0 + z_{\alpha} \sqrt{\frac{p_0(1-p_0)}{n}}$ • $Z = \frac{\bar{P} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} > z_{\alpha}$ • $p - value = P(Z > z_{OSS}) < \alpha$ 		
<p style="text-align: center;">CASO III $H_0: \mu \geq \mu_0$ $H_1: \mu < \mu_0$</p>			
Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{P} > p_0 - z_{\alpha} \sqrt{\frac{p_0(1-p_0)}{n}}$ • $Z = \frac{\bar{P} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} < -z_{\alpha}$ • $p - value = P(Z < z_{OSS}) < \alpha$ 		
<p>Potenza di un test</p>	<p>Quando rifiutiamo l'ipotesi nulla deduciamo che vi sia una forte evidenza empirica a favore della nostra conclusione</p> <p>Quando non rifiutiamo l'ipotesi nulla è perché:</p> <ul style="list-style-type: none"> • L'ipotesi nulla è vera • L'ipotesi nulla è falsa ma noi non la rifiutiamo (errore di secondo tipo) <p>Qual è la probabilità di commettere un errore di secondo tipo?</p> <p>La regola di rifiuto permette di determinare i valori della media campionaria che permettono di non rifiutare l'ipotesi nulla, ovvero la regione di accettazione: essa è infatti la complementare della regione di rifiuto</p> <p>La probabilità dell'errore di secondo tipo β è pari alla probabilità che la media campionaria appartenga alla regione di accettazione</p> <p>Il suo complementare $1 - \beta$ si chiama potenza</p> <p>Fissato un livello di significatività pari ad α e fornito il vero valore del parametro in test dobbiamo:</p> <ul style="list-style-type: none"> • Definire la regola per cui $\beta = P(\text{non rifiutare } H_0 H_1 \text{ vera})$ • Applicare la regola di non rifiuto: cerco il valore critico x_c che separa la regione di rifiuto da quella di accettazione. • Calcolare la probabilità normalizzato: si ottiene β • Alla fine, calcola la potenza come $1 - \beta$ 		



Per **aumentare la potenza del test** occorre:

- Aumentare l'**ampiezza campionaria**
- Aumentare α , ossia la probabilità di rifiutare H_0

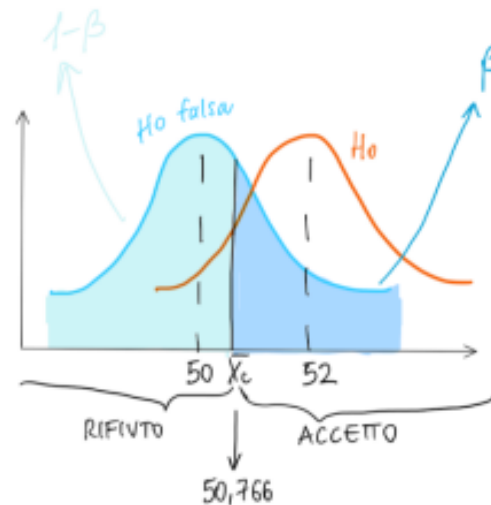
Esempio numerico con grafico

Consideriamo il seguente test unilaterale sulla media della popolazione:

$$H_0: \mu \geq 52$$

$$H_1: \mu < 52$$

Assumiamo che la vera media sia $\mu^* = 50$ (ovvero, H_1 vera); assumiamo che la varianza sia nota, e pari a 6, e che la popolosità del campione sia $n = 64$. Consideriamo $\alpha = 0.05$



Come primo step dobbiamo trovare β , e per farlo occorre calcolr il valore critico di rifiuto x_c che separa la regione di rifiuto da quella di accettazione. Essendo la vera media (50) inferiore a quella fornita dall'ipotesi nulla (52), allora la regione di rifiuto di H_0 risulta essere quella minore di x_c

Perciò si ha:

$$\beta = P(X \geq x_c | \mu^* = 50)$$

Per trovare x_c , nel seguente caso, applico:

$$x_c = \mu_0 - z_\alpha * \frac{\sigma}{\sqrt{n}} = 50,766$$

Da cui:

$$\beta = P(X \geq 50,766) = P\left(X \geq \frac{50,766 - 50}{6/\sqrt{64}}\right) = P(Z \geq 1,02) = 0,1539$$

Segue

$$1 - \beta = 0,8461 \rightarrow \text{Potenza}$$



Rappresentiamo graficamente e confrontiamo l'ipotesi nulla e l'ipotesi alternativa. A titolo esemplificativo, consideriamo il seguente test con la relativa regione di rifiuto:

$$H_0: \mu < \mu_0$$

$$H_1: \mu > \mu_0 = \mu_1$$

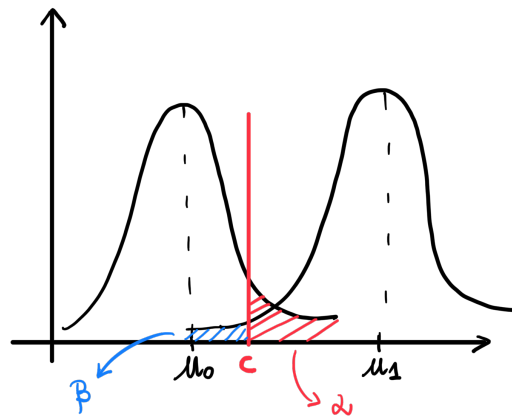
$$R = \{\bar{x} > c\}$$

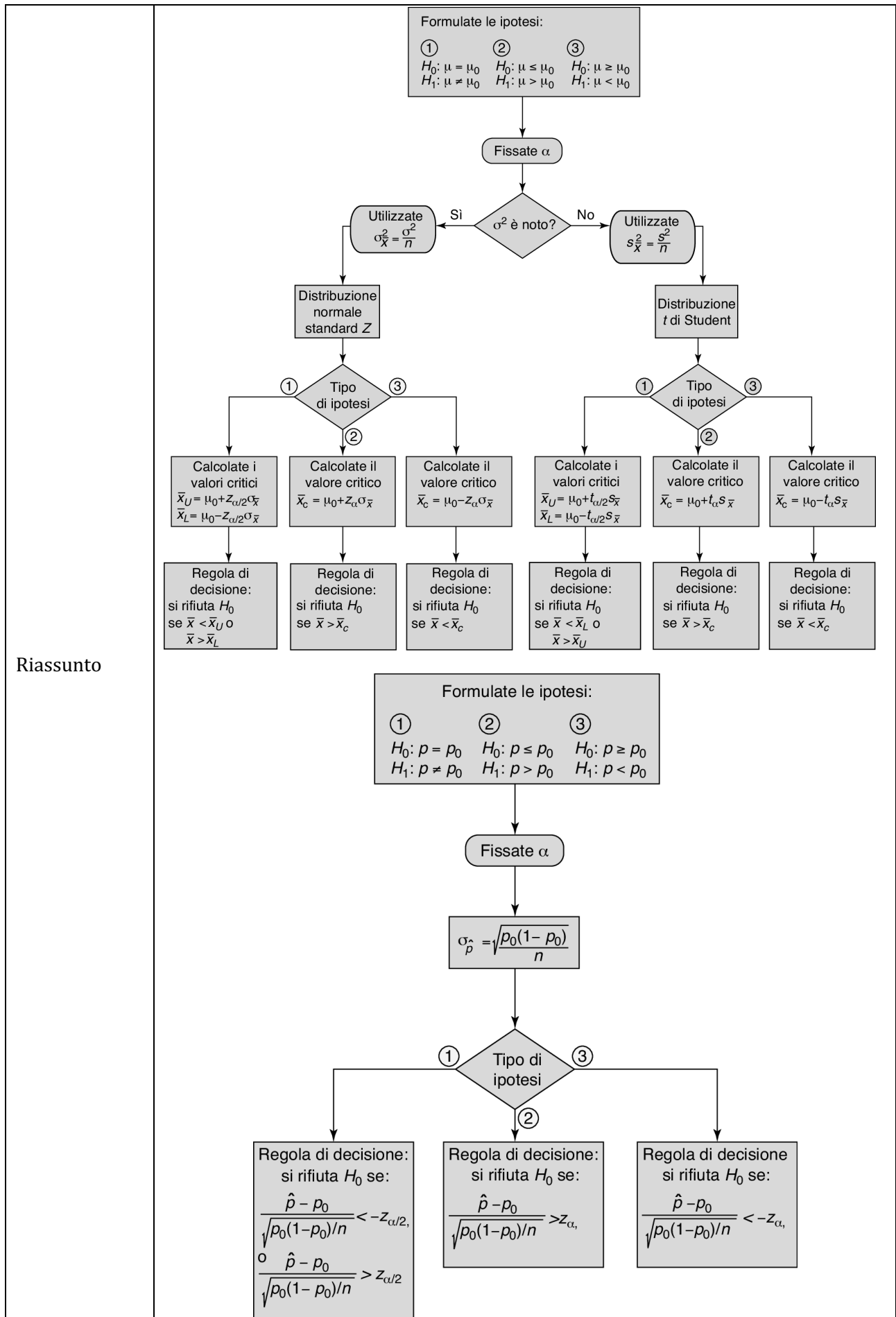
Consideriamo:

- α = probabilità di rifiutare H_0 quando H_0 è vera (area a destra di c sotto la curva di H_0)
- β = probabilità di non rifiutare H_0 quando H_1 è vera (area a sinistra di c sotto la curva di H_1)

Relazione grafica tra α e β

Può essere chiesto di descrivere la relazione tra le aree per determinati cambiamenti di c , μ_0 , μ_1 , n (N.B.: aumentando n , la varianza diminuisce e la curva si appiattisce, concentrandosi sul centro)





VERIFICA DI IPOTESI SULLA MEDIA DI DUE POPOLAZIONI			
Problemi tipo	CASO I	CASO II	CASO III
	$H_0: \mu_x - \mu_y = d_0$ $H_1: \mu_x - \mu_y \neq d_0$	$H_0: \mu_x - \mu_y \leq d_0$ $H_1: \mu_x - \mu_y > d_0$	$H_0: \mu_x - \mu_y \geq d_0$ $H_1: \mu_x - \mu_y < d_0$
	Occorre distinguere due casi fondamentali: <ul style="list-style-type: none"> • Campioni dipendenti → campioni uguali, provenienti dalla medesima popolazione, vengono analizzati prima e dopo un trattamento • Campioni indipendenti → campioni diversi vengono analizzati prima e dopo un trattamento 		
Popolazioni normali - campioni dipendenti popolazione normale	Sia $D_i = X_i - Y_i$ con $D_1, \dots, D_n \sim^{iid} N(\mu_D, \sigma_D^2)$ ed in particolare: <ul style="list-style-type: none"> • $\bar{D} = \bar{x} - \bar{y}$ • $S_D^2 = S_X^2 + S_Y^2 - 2r_{XY}$ 		
	Statistica test	$T = \frac{\bar{D} - d_0}{S_D/\sqrt{n}} \approx^{\mu_D=d_0} T_{(n-1)}$	
	CASO I $H_0: \mu_x - \mu_y = d_0$ $H_1: \mu_x - \mu_y \neq d_0$		
	Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{D} < d_0 - t_{n-1, \frac{\alpha}{2}} \frac{S_D}{\sqrt{n}}$ oppure $\bar{D} > d_0 + t_{n-1, \frac{\alpha}{2}} \frac{S_D}{\sqrt{n}}$ • $T = \left \frac{\bar{D} - d_0}{\frac{S_D}{\sqrt{n}}} \right > t_{n-1, \frac{\alpha}{2}}$ • $p - value = 2P(T > t_{oss}) < \alpha$ 	
	CASO II $H_0: \mu_x - \mu_y \leq d_0$ $H_1: \mu_x - \mu_y > d_0$		
	Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{D} > d + t_{n-1, \alpha} \frac{S}{\sqrt{n}}$ • $T = \frac{\bar{D} - d_0}{\frac{S_D}{\sqrt{n}}} > t_{n-1, \alpha}$ • $p - value = P(Z > t_{oss}) < \alpha$ 	
	CASO III $H_0: \mu_x - \mu_y \geq d_0$ $H_1: \mu_x - \mu_y < d_0$		
Regola di rifiuto	<ul style="list-style-type: none"> • $\bar{D} < d + t_{n-1, \alpha} \frac{S}{\sqrt{n}}$ • $T = \frac{\bar{D} - d_0}{\frac{S_D}{\sqrt{n}}} < -t_{n-1, \alpha}$ • $p - value = P(Z < t_{oss}) < \alpha$ 		
Varianza pooled	La principale differenza del caso in cui le varianze siano incognite ma assunte uguali con campioni indipendenti è l'utilizzo della varianza pooled S_P $S_P^2 = \frac{(n_X - 1)S_X^2 + (n_Y - 1)S_Y^2}{n_X + n_Y - 2}$		
Popolazioni normali - campioni indipendenti con varianze incognite ma	Statistica test	$T = \frac{(\bar{X} - \bar{Y}) - d_0}{S_P \sqrt{\frac{1}{n_X} + \frac{1}{n_Y}}} \sim^{\mu_x - \mu_y = d_0} T_{(n_X + n_Y - 2)}$	
	CASO I $H_0: \mu_x - \mu_y = d_0$ $H_1: \mu_x - \mu_y \neq d_0$		



<p>assunte uguali</p>		<p>Regola di rifiuto</p>	<ul style="list-style-type: none"> $\bar{X} - \bar{Y} < d_0 - t_{n_x+n_y-2, \frac{\alpha}{2}} S_P \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}$ oppure $\bar{X} - \bar{Y} > d_0 + t_{n_x+n_y-2, \frac{\alpha}{2}} S_P \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}$ $T = \left \frac{\bar{X} - \bar{Y} - d_0}{S_P \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}} \right > t_{n_x+n_y-2, \frac{\alpha}{2}}$ $p - value = 2P(T > t_{OSS}) < \alpha$ 		
		<p>CASO II $H_0: \mu_x - \mu_y \leq d_0$ $H_1: \mu_x - \mu_y > d_0$</p>			
		<p>Regola di rifiuto</p>	<ul style="list-style-type: none"> $\bar{X} - \bar{Y} > d_0 + t_{n_x+n_y-2, \frac{\alpha}{2}} S_P \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}$ $T = \frac{\bar{X} - \bar{Y} - d_0}{S_P \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}} > t_{n_x+n_y-2, \alpha}$ $p - value = P(Z > t_{OSS}) < \alpha$ 		
		<p>CASO III $H_0: \mu_x - \mu_y \geq d_0$ $H_1: \mu_x - \mu_y < d_0$</p>			
		<p>Regola di rifiuto</p>	<ul style="list-style-type: none"> $\bar{X} - \bar{Y} < d_0 - t_{n_x+n_y-2, \frac{\alpha}{2}} S_P \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}$ $T = \frac{\bar{X} - \bar{Y} - d_0}{S_P \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}} < -t_{n_x+n_y-2, \alpha}$ $p - value = P(Z < t_{OSS}) < \alpha$ 		
<p>Popolazioni arbitrarie con grandi campioni ($n \geq 30$)</p>	<p>Applicazione del teorema:</p> <ul style="list-style-type: none"> Teorema centrale del limite → distribuzione normale <p>a) Campioni dipendenti con varianze non note e uguali</p> <p>Statistica test: $Z = \frac{\bar{D} - d_0}{\frac{S_D}{\sqrt{n}}} \approx \mu_d = d_0 \ N(0,1)$</p> <p>Le regole di rifiuto sono uguali a quelle viste per la singola popolazione normale varianza incognita, cambiando però i quantili con quelli della normale standard ($t_{n-1, \alpha} \rightarrow z_\alpha$)</p> <p>b) Campioni indipendenti con varianze non note e differenti</p> <p>Statistica test: $Z = \frac{(\bar{X} - \bar{Y}) - d_0}{\sqrt{\frac{s_P^2}{n_x} + \frac{s_P^2}{n_y}}} \approx \mu_{X-Y} = d_0 \ N(0,1)$</p> <p>Le regole di rifiuto sono uguali a quelle viste per la singola popolazione normale varianza incognita, cambiando però i quantili con quelli della normale standard ($t_{v, \alpha} \rightarrow z_\alpha$)</p>				
<p>Popolazioni arbitrarie</p>	<p>Assunzioni:</p> <ul style="list-style-type: none"> Indipendenza delle popolazioni Uguaglianza delle varianze (non note) delle popolazioni Assunzione di normalità (ampiezza campionaria non sufficiente per applicare il teorema centrale del limite) <p>Ci possiamo ricondurre al caso di popolazioni indipendenti con varianze assunte uguali → calcolo la <u>varianza pooled</u> e distribuzione <u>t di student</u></p>				



<p>Test sulla bontà di adattamento</p>	<p>Consideriamo un campione casuale di n osservazioni che possano essere classificate in K categorie o classi di misura. Indichiamo con O_1, O_2, \dots, O_K il numero di osservazioni effettive.</p> <p>Formuliamo un'ipotesi H_0 che ci permetta di descrivere i dati: in particolare supponiamo che l'ipotesi nulla specifichi p_i con $i = 1, 2, \dots, K$ t. c. $\sum_{i=1}^k p_i = 1$ come la probabilità che un'osservazione appartenga all'i-esima categoria. Sia infine:</p> $E_i = np_i$ <p>La frequenza attesa della i-esima classe sotto l'ipotesi H_0. Verificare che i dati effettivamente riscontrati si adattano ai valori attesi, sotto l'ipotesi H_0, vuol dire condurre un test sulla bontà di adattamento</p> $H_0: p_i = p_i \text{ per ogni } i$ $H_1: p_i \neq p_i \text{ per almeno un } i$ <p>Abbiamo la tabella:</p> <table border="1" data-bbox="383 772 1412 1052"> <tr> <th>Categoria</th> <th>1</th> <th>2</th> <th>...</th> <th>K</th> <th>Totale</th> </tr> <tr> <td>Frequenze osservate</td> <td>O_1</td> <td>O_2</td> <td>...</td> <td>O_K</td> <td>n</td> </tr> <tr> <td>Probabilità (sotto H_0)</td> <td>p_1</td> <td>p_2</td> <td>...</td> <td>p_k</td> <td>1</td> </tr> <tr> <td>Frequenze attese (sotto H_0)</td> <td>$E_1 = np_1$</td> <td>$E_2 = np_2$</td> <td>...</td> <td>$E_K = np_K$</td> <td>n</td> </tr> </table> <p>Per n elevato (grandi campioni) e $E_i > 5 \forall i$ avremo la seguente statistica test</p> $\chi_{oss}^2 = \sum_{i=1}^K \frac{(O_i - E_i)^2}{E_i} \approx_{H_0} \chi_{(k-1)}^2$ <p>E presenta la seguente regola di rifiuto:</p> <ul style="list-style-type: none"> $\chi_{oss}^2 > \chi_{k-1, \alpha}^2$ $p\text{-value} = P(\chi^2 > \chi_{oss}^2) < \alpha$ 	Categoria	1	2	...	K	Totale	Frequenze osservate	O_1	O_2	...	O_K	n	Probabilità (sotto H_0)	p_1	p_2	...	p_k	1	Frequenze attese (sotto H_0)	$E_1 = np_1$	$E_2 = np_2$...	$E_K = np_K$	n																									
Categoria	1	2	...	K	Totale																																													
Frequenze osservate	O_1	O_2	...	O_K	n																																													
Probabilità (sotto H_0)	p_1	p_2	...	p_k	1																																													
Frequenze attese (sotto H_0)	$E_1 = np_1$	$E_2 = np_2$...	$E_K = np_K$	n																																													
<p>Test chi-quadrato sulla indipendenza (o sulla associazione di due variabili)</p>	<p>Si considerino le variabili categoriche X ed Y. Vogliamo testare le seguenti ipotesi:</p> $H_0: X \text{ e } Y \text{ indipendenti}$ $H_1: X \text{ e } Y \text{ dipendenti}$ <p>Sulla base dei dati campionari raccolti su n unità e sintetizzati in una tabella di contingenza</p> <table border="1" data-bbox="542 1612 1268 1892"> <tr> <th>$X \setminus Y$</th> <th>y_1</th> <th>...</th> <th>y_j</th> <th>...</th> <th>y_c</th> <th>Tot.</th> </tr> <tr> <td>x_1</td> <td>O_{11}</td> <td>...</td> <td>O_{1j}</td> <td>...</td> <td>O_{1c}</td> <td>R_1</td> </tr> <tr> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> </tr> <tr> <td>x_i</td> <td>O_{i1}</td> <td>...</td> <td>O_{ij}</td> <td>...</td> <td>O_{ic}</td> <td>R_i</td> </tr> <tr> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> </tr> <tr> <td>x_r</td> <td>O_{r1}</td> <td>...</td> <td>O_{rj}</td> <td>...</td> <td>O_{rc}</td> <td>R_r</td> </tr> <tr> <td>Tot.</td> <td>C_1</td> <td>...</td> <td>C_j</td> <td>...</td> <td>C_c</td> <td>n</td> </tr> </table> <p>Dove O_{ij} è la frequenza congiunta assoluta osservata di unità nel campione. Consideriamo la frequenza attesa della cella definibile come:</p>	$X \setminus Y$	y_1	...	y_j	...	y_c	Tot.	x_1	O_{11}	...	O_{1j}	...	O_{1c}	R_1	x_i	O_{i1}	...	O_{ij}	...	O_{ic}	R_i	x_r	O_{r1}	...	O_{rj}	...	O_{rc}	R_r	Tot.	C_1	...	C_j	...	C_c	n
$X \setminus Y$	y_1	...	y_j	...	y_c	Tot.																																												
x_1	O_{11}	...	O_{1j}	...	O_{1c}	R_1																																												
...																																												
x_i	O_{i1}	...	O_{ij}	...	O_{ic}	R_i																																												
...																																												
x_r	O_{r1}	...	O_{rj}	...	O_{rc}	R_r																																												
Tot.	C_1	...	C_j	...	C_c	n																																												



	$E_{ij} = \frac{R_i C_j}{n}$ <p>Sotto l'ipotesi nulla di assenza di dipendenza, ci aspettiamo, infatti, di osservare una distribuzione <i>proporzionale</i> al totale delle osservazioni totali per riga e per colonna. Se n elevato e $E_{ij} < 5$ al massimo nel 20% delle celle, possiamo costruire la seguente statistica test:</p> $\chi^2 = \sum_{i=1}^n \sum_{j=1}^k \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \approx_{H_0} \chi^2_{(r-1)(c-1)}$ <p>Avremo la seguente regola di rifiuto</p> <ul style="list-style-type: none"> • $\chi^2 > \chi^2_{(r-1)(c-1), \alpha}$ • $p - value = P(\chi^2 > \chi^2_{oss}) < \alpha$
--	---

<p>Test sulla correlazione lineare (r)</p>	<p>Si considerino due variabili numeriche X ed Y estratte da un campione casuale semplice di ampiezza n la cui distribuzione sia normale bidimensionale; si vuole indagare circa l'eventuale esistenza e forza, nella popolazione, di una relazione lineare tra le variabili considerate: l'obiettivo dell'analisi è il coefficiente di correlazione lineare nella popolazione e valutare se esista una correlazione lineare tra le due variabili</p> $\rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y} = \frac{Cov(X, Y)}{\sqrt{Var[X] * Var[Y]}}$ <p>Poiché disponiamo unicamente di un campione dobbiamo riferirci al coefficiente di correlazione lineare campionario</p> $r_{XY} = \frac{S_{XY}}{S_X S_Y} = \frac{\sum_i (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_i (X_i - \bar{X})^2 * \sum_i (Y_i - \bar{Y})^2}}$ <p>Si può dimostrare che è possibile ottenere la seguente statistica test</p> $T = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} \sim_{\rho=0} T_{(n-2)}$ <p>E avremo le seguenti regole di rifiuto:</p> <table border="1" style="width: 100%; text-align: center;"> <thead> <tr> <th>CASO</th> <th>REGOLA DI RIFIUTO</th> </tr> </thead> <tbody> <tr> <td>$H_0: \rho = 0$ $H_1: \rho \neq 0$</td> <td>$t_{oss} > t_{n-2, \frac{\alpha}{2}}$ $p - value = 2P(T > t_{oss}) < \alpha$</td> </tr> <tr> <td>$H_0: \rho = 0$ $H_1: \rho > 0$</td> <td>$t_{oss} > t_{n-2, \alpha}$ $p - value = 2P(T > t_{oss}) < \alpha$</td> </tr> <tr> <td>$H_0: \rho = 0$ $H_1: \rho < 0$</td> <td>$t_{oss} < -t_{n-2, \alpha}$ $p - value = 2P(T < t_{oss}) < \alpha$</td> </tr> </tbody> </table>	CASO	REGOLA DI RIFIUTO	$H_0: \rho = 0$ $H_1: \rho \neq 0$	$ t_{oss} > t_{n-2, \frac{\alpha}{2}}$ $p - value = 2P(T > t_{oss}) < \alpha$	$H_0: \rho = 0$ $H_1: \rho > 0$	$t_{oss} > t_{n-2, \alpha}$ $p - value = 2P(T > t_{oss}) < \alpha$	$H_0: \rho = 0$ $H_1: \rho < 0$	$t_{oss} < -t_{n-2, \alpha}$ $p - value = 2P(T < t_{oss}) < \alpha$
CASO	REGOLA DI RIFIUTO								
$H_0: \rho = 0$ $H_1: \rho \neq 0$	$ t_{oss} > t_{n-2, \frac{\alpha}{2}}$ $p - value = 2P(T > t_{oss}) < \alpha$								
$H_0: \rho = 0$ $H_1: \rho > 0$	$t_{oss} > t_{n-2, \alpha}$ $p - value = 2P(T > t_{oss}) < \alpha$								
$H_0: \rho = 0$ $H_1: \rho < 0$	$t_{oss} < -t_{n-2, \alpha}$ $p - value = 2P(T < t_{oss}) < \alpha$								

MODELLO DI REGRESSIONE LINEARE SEMPLICE	
Terminologie e modello statistico	<p>Il modello di regressione lineare fornisce il valore atteso della variabile aleatoria Y quando X assume un particolare valore, aggiungendovi un margine di errore ϵ:</p> $y = \beta_0 + \beta_1 x + \epsilon_i$



	<p>Essa necessita alcune assunzioni di base:</p> <ol style="list-style-type: none"> 1. La x è costante o è una realizzazione di una variabile aleatoria X, indipendente dalle componenti aleatorie di errore ε_i 2. Il valore atteso della variabile aleatoria Y è una funzione lineare della variabile aleatoria esplicativa X 3. I termini di errore sono variabili aleatorie con media 0 e varianza costante σ^2. La seconda condizione è detta omoschedasticità $E[\varepsilon_i] = 0 \quad Var[\varepsilon_i] = \sigma^2$ <ol style="list-style-type: none"> 4. Gli errori aleatori non sono correlati tra loro e quindi $E[\varepsilon_i \varepsilon_j] = 0$ <p>Da queste assunzioni, seguono alcune importanti osservazioni:</p> <ul style="list-style-type: none"> • L'ipotesi di media nulla implica $Cov(X, \varepsilon) = \mathbf{0} \rightarrow X$ e ε non sono correlati • L'ipotesi di varianza costante significa che la variabilità di ε non dipende dal valore di X <p>Determinazione dei parametri del modello</p> <ul style="list-style-type: none"> • $E[Y X = x] = E[b_0 + b_1x + \varepsilon X = x] = E[b_0 + b_1x X = x] + E[\varepsilon X = x] = b_0 + b_1x$ • $b_0 = E[Y X = 0] =$ medio di Y nella sottopopolazione con $X = 0 \rightarrow$ è l'intercetta all'origine • $b_1 =$ variazione nella media di Y se X aumenta di una unità; rappresenta il coefficiente angolare $\rightarrow E[Y X = x + 1] - E[Y X = x] = b_0 + b_1(x + 1) - (b_0 + b_1x) = b_1$ <p>Ulteriore conseguenza:</p> <ul style="list-style-type: none"> • $Var[Y X = x] = \sigma^2 \rightarrow$ la varianza di Y è costante in ogni sottopopolazione
<p>Metodo dei minimi quadrati</p>	<p>Per determinare, da n coppie di osservazioni, gli stimatori dei coefficienti incogniti del modello, β_0 e β_1, possiamo ricorrere al modello dei minimi quadrati. Dato un campione di ampiezza n:</p> $(x_1, y_1) \dots (x_n, y_n)$ <p>Stimiamo β_0 e β_1 in modo tale che minimizzino:</p> $S(\beta_0, \beta_1) = \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)]^2$ <p>Tali stime hanno espressione esplicita:</p> $b_1 = \frac{s_{XY}}{s_X^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = r_{XY} \frac{s_Y}{s_X}$ $b_0 = \bar{y} - b_1 \bar{x}$ <p>Interpretazione dei coefficienti della regressione (se esplicativi)</p> <ul style="list-style-type: none"> • b_1 indica quanto in media aumenta y all'aumentare di una unità di x <ul style="list-style-type: none"> ○ Il segno positivo indica una relazione diretta (retta positivamente inclinata) ○ Il segno negativo indica una relazione inversa (retta negativamente inclinata) • b_0 indica quanto vale in media y, per $x = 0$; occorre valutare il suo senso osservando i dati del problema



Vogliamo sviluppare delle misure per indicare **quanto efficacemente la variabile X spiega il comportamento di Y** nel modello lineare.

Partiamo dalla retta di regressione stimata:

$$\hat{y} = b_0 + b_1 x$$

Questa consente di determinare:

- **I valori previsti**

$$\hat{y}_i = b_0 + b_1 x_i$$

- **I valori residui**

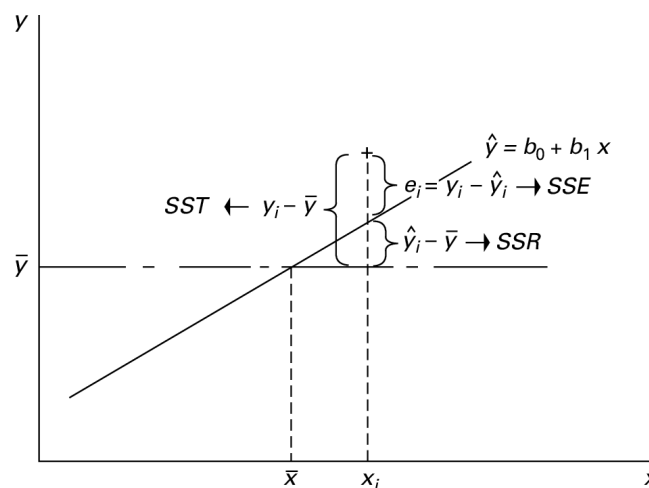
$$e_i = y_i - \hat{y}_i = y_i - (b_0 + b_1 x_i)$$

Come è facile osservare i valori osservati y_i si discostano dai valori previsti \hat{y}_i per un residuo e_i :

$$y_i = \hat{y}_i + e_i$$

Vi è quindi una parte che il modello non riesce a spiegare, ovvero e_i

Capacità
esplicativa
della retta di
regressione



Quello che stiamo facendo si chiama **scomposizione della devianza**. Questa può essere condotta anche considerando:

- $y_i - \bar{y}$ scarto della singola osservazione dalla media
- $\hat{y}_i - \bar{y}$ scarto del singolo valore previsto dalla media
- $y_i - \hat{y}_i$ scarto della singola osservazione dal corrispondente valore previsto

La loro analisi ci permette di scomporre lo scarto di una singola osservazione di Y dalla media nella somma algebrica degli altri due:

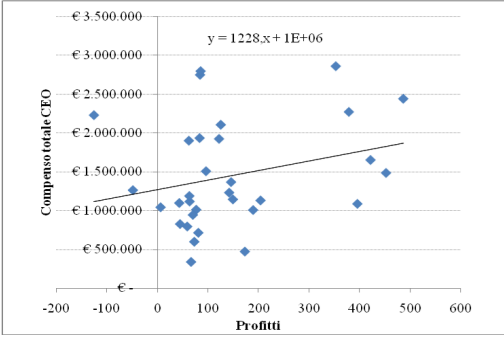
$$y_i - \bar{y} = (\hat{y}_i - \bar{y}) + (y_i - \hat{y}_i)$$

Elevano al quadrato ciascun membro dell'uguaglianza e sommando gli n addendi si ottiene:



	$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2$ $SST = SSR + SSE$ <p>Note le varianze della variabile risposta Y e della variabile esplicativa X possiamo giungere alle seguenti formule:</p> $SSR = b_1^2(n - 1)s_X^2 = b_1^2 \sum (x_i - \bar{x})^2$ $SST = (n - 1)s_Y^2 = \sum (y_i - \bar{y})^2$ $SSE = SST - SSR = (n - 1)s_Y^2 - b_1^2(n - 1)s_X^2$
<p>Coefficiente di determinazione</p>	<p>Si chiama coefficiente di determinazione:</p> $R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$ <p>Rappresenta la percentuale di variabilità di Y spiegata dal modello di regressione. R^2 aumenta in modo direttamente proporzionale alla dispersione della variabile indipendente X.</p> <p>Varia tra 0 ed 1: valori più elevati indicano una migliore bontà del modello, in quanto meglio capace di spiegare la variabilità di Y</p> <p>N.B.: la bontà del modello mediante R^2 è sempre relativa, nel senso che è corretto valutare due R^2, ed il maggiore indicherà il modello migliore; non è corretto valutare R^2 in senso assoluto</p> <p>Proprietà:</p> <ul style="list-style-type: none"> • $0 \leq R^2 \leq 1$ • Maggiore è il suo valore, migliore è l'adattamento del modello ai dati <p>Esiste una relazione tra R^2 ed il coefficiente di correlazione lineare</p> $R^2 = r^2 = \left(\frac{Cov(x, y)}{\sqrt{Var(x)} * \sqrt{Var(y)}} \right)^2$
<p>Stima della varianza del modello</p>	<p>Chiamiamo stimatore dei residui o varianza dell'errore o varianza dei residui:</p> $s_e^2 = \frac{\sum_i e_i^2}{n - 2} = \frac{\sum_i (y_i - \hat{y}_i)^2}{n - 2} = \frac{SSE}{n - 2}$ <p>Si chiama errore standard del modello:</p> $s_e = \sqrt{s_e^2}$



<p>Interpretazione grafica del modello (esempio)</p>	<p>Considerando la distribuzione dei valori osservati per la variabile dipendente e quella esplicativa riportata nel diagramma a dispersione, è possibile notare come i punti non sembrano suggerire l'esistenza di una relazione lineare tra le due variabili. Essendo, inoltre, la distanza tra i punti e la retta piuttosto ampia (accentuata dalla presenza di alcuni valori anomali) ci si aspetta di commettere un errore di stima piuttosto elevato e di ottenere, di conseguenza, un coefficiente di determinazione R quadrato non molto elevato.</p>	
<p>Proprietà degli stimatori dei minimi quadrati</p>	<p>Teorema di GAUSS-MARKOV Sotto le assunzioni di base del modello di regressione lineare, gli stimatori dei minimi quadrati sono i migliori stimatori lineari e non distorti per il modello di regressione lineare (stimatori BLUE) Valgono perciò le seguenti proprietà (σ^2 è non noto)</p> <p><u>Valore atteso</u></p> $E[b_0] = \beta_0$ $E[b_1] = \beta_1$ <p><u>Varianza</u></p> $Var[b_0] = S_{b_0}^2 = s_e^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_i (x_i - \bar{x})^2} \right)$ $Var[b_1] = S_{b_1}^2 = \frac{s_e^2}{\sum_i (x_i - \bar{x})^2} = \frac{s_e^2}{(n-1)s_{\bar{x}}^2}$ <p>Inoltre, fissa l'assunzione di Normalità degli errori, ovvero $\epsilon_1, \dots, \epsilon_n \sim iid N(0, \sigma^2)$, gli stimatori dei minimi quadrati hanno distribuzione normale:</p> $b_0 \sim N(\beta_0, \sigma_{b_0}^2) \quad b_1 \sim N(\beta_1, \sigma_{b_1}^2)$ <p><u>Efficienza degli stimatori</u> Uno stimatore è tanto più efficiente tanto più la sua varianza è piccola. La varianza dello stimatore del coefficiente angolare è pari a $Var[b_1] = \frac{s_e^2}{\sum_i (x_i - \bar{x})^2}$, per cui, affinché $Var[b_1]$ sia piccolo è necessario che $\sum_i (x_i - \bar{x})^2$ sia elevato e cioè che la variabile indipendente sia dispersa intorno alla media.</p>	
<p>Inferenza sui parametri</p>	<p>Intervalli di confidenza (assunzione di normalità degli errori) Dato che σ^2 è incognito, stimiamo la varianza con s^2, la varianza campionaria; per b_1 e b_0 si ha che:</p> $S_{B_1} = \sqrt{\frac{s_e^2}{\sum_i (x_i - \bar{x})^2}} = \sqrt{\frac{s_e^2}{(n-1)Var(x)}}$ $S_{B_0} = \sqrt{s^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_i (x_i - \bar{x})^2} \right)} = \sqrt{s^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{(n-1)Var(x)} \right)}$	



Ricorda che $\sum_i(x_i - \bar{x})^2 = (n - 1)Var(x)$

Da cui determiniamo i pivot:

$$\frac{B_1 - \beta_1}{S_{B_1}} \sim T_{(n-2)} \quad \frac{B_0 - \beta_0}{S_{B_0}} \sim T_{(n-2)}$$

Da cui otteniamo gli intervalli di confidenza, a livello $1 - \alpha$

$$IC_{1-\alpha}(\beta_1) = \left(b_1 \pm t_{n-2, \frac{\alpha}{2}} S_{B_1} \right)$$

$$IC_{1-\alpha}(\beta_0) = \left(b_0 \pm t_{n-2, \frac{\alpha}{2}} S_{B_0} \right)$$

Test d'ipotesi (assunzione di normalità degli errori)

	REGOLA DI RIFIUTO
<p>CASO I $H_0: \beta_1 = b_1^*$ $H_1: \beta_1 \neq b_1^*$</p>	$b_1 < b_1^* - t_{n-2, \frac{\alpha}{2}} S_{B_1}$ oppure $b_1 > b_1^* + t_{n-2, \frac{\alpha}{2}} S_{B_1}$ $ T = \left \frac{b_1 - b_1^*}{S_{B_1}} \right > t_{n-2, \frac{\alpha}{2}}$ $p - value = 2P(T > t_{oss}) < \alpha$
<p>CASO II $H_0: \beta_1 \leq b_1^*$ $H_1: \beta_1 > b_1^*$</p>	$b_1 > b_1^* + t_{n-2, \alpha} S_{B_1}$ $T = \frac{b_1 - b_1^*}{S_{B_1}} > t_{n-2, \alpha}$ $p - value = P(T > t_{oss}) < \alpha$
<p>CASO III $H_0: \beta_1 \geq b_1^*$ $H_1: \beta_1 < b_1^*$</p>	$b_1 < b_1^* - t_{n-2, \alpha} S_{B_1}$ $T = \frac{b_1 - b_1^*}{S_{B_1}} < -t_{n-2, \alpha}$ $p - value = P(T < t_{oss}) < \alpha$

Test di significatività

Un particolare test d'ipotesi da condurre è il **test di significatività**, atto a verificare se **X sia una variabile significativa o meno**; se rifiuto l'ipotesi nulla, X è da considerarsi una variabile significativa. Tale test prevede le seguenti ipotesi:

$$H_0: \beta_1 = 0 \text{ (non significativa)}$$

$$H_1: \beta_1 \neq 0 \text{ (significativa)}$$

Per cui vale la seguente statistica test:

$$T = t - value = \frac{b_1}{S_{b_1}} \sim_{b_1=0} T_{(n-2)}$$

Valgono le seguenti regole di rifiuto:

$$|T| = \left| \frac{b_1 - 0}{S_{b_1}} \right| = |t - value| > t_{n-2, \frac{\alpha}{2}}$$

$$p - value = 2P(T_{n-2} > |t_{oss}|) < \alpha$$



	<p>In caso di rifiuto si dirà che:</p> <ul style="list-style-type: none"> • b_1 è significativamente diverso da 0 • X è una variabile esplicativa statisticamente significativa <p><u>Relazione tra test di significatività e test di correlazione lineare (ρ)</u></p> $\rho = 0 \Leftrightarrow \beta_1 = 0$
Previsione	<p>Sono possibili due tipi di previsione:</p> <ul style="list-style-type: none"> • Previsione puntuale del valore y_{n+1} risultate da una singola osservazione, x_{n+1} • Previsione del valore atteso condizionato $E[Y_{n+1} X = x_{n+1}]$, cioè della media della variabile risposta quando la variabile esplicativa assume il valore prefissato x_{n+1} <p>a) Stima puntuale La stima è uguale per entrambi e, in particolare, avremo che:</p> $y_{n+1} = b_0 + b_1 x_{n+1}$ <p>b) Stima per intervallo <u>Intervallo di confidenza per il valore atteso della stima</u></p> $IC_{1-\alpha}(E[Y_{n+1} X = x_{n+1}]) = \left(y_{n+1} \pm t_{n-2, \frac{\alpha}{2}} s_e \sqrt{\frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum_i (x_i - \bar{x})^2}} \right)$ <ul style="list-style-type: none"> • Viene fornito x_{n+1} • Calcolo il centro dell'intervallo $y_{n+1} = b_0 + b_1 x_{n+1}$ • Calcolo $s_e = \sqrt{\frac{SSE}{n-2}}$ • Considero il termine sotto radice: <ul style="list-style-type: none"> ○ Il denominatore della seconda frazione può essere riscritto come $\sum_i (x_i - \bar{x})^2 = (n-1)Var(X)$ <p><u>Intervallo di previsione per il valor singolo stimato</u></p> $IP_{1-\alpha}(Y_{n+1}) = \left(y_{n+1} \pm t_{n-2, \frac{\alpha}{2}} s_e \sqrt{1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum_i (x_i - \bar{x})^2}} \right)$ <p>Si osserva che l'intervallo di confidenza è più piccolo dell'intervallo di previsione ed entrambi sono centrati sul valore previsto y: con l'intervallo di confidenza prevediamo un valore medio; con l'intervallo di previsione una variabile aleatoria, che non ci permette di avere un intervallo più stretto a causa della sua aleatorietà</p> <p>Si possono ottenere ulteriori informazioni studiando la forma generali degli intervalli e tenendo presente il fatto che tanto maggiore è l'ampiezza intervallare, tanto minore è la previsione puntuale:</p> <ol style="list-style-type: none"> 1. Più l'ampiezza n del campione è elevata, tanto minore è l'ampiezza dell'intervallo 2. Tanto maggiore è s_e^2, tanto più ampio è l'intervallo di confidenza: infatti tale misura è una stima della varianza degli errori



3. Tanto maggiore è $\sum_i(x_i - \bar{x})^2$ e tanto minore risulta l'ampiezza dell'intervallo: infatti, si avranno più informazioni circa i valori più lontani dalla media
4. Se il valore per il quale vogliamo fare la previsione x_{n+1} si allontana dalla media, allora otterremo intervalli più ampi

Call:
lm(formula = Y..Retail.Sales ~ X..Income, data = RetailSales)

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	559.4600	1450.6975	0.39	0.7
X..Income	0.3815	0.0253	15.08	2.2e-12

Residual standard error: 148 on 20 degrees of freedom
Multiple R-squared: 0.919, Adjusted R-squared: 0.915
F-statistic: 228 on 1 and 20 DF, p-value: 2.17e-12

Analysis of Variance Table

Response: Y..Retail.Sales	Df	Sum Sq	Mean Sq	F value	Pr(>F)
X..Income	1	4961434	4961434	228	2.2e-12
Residuals	20	436127	21806		

SSE, Somma dei quadrati degli errori

(a)

OUTPUT RIEPILOGO

Statistica della regressione	
R multiplo	0.95875
R al quadrato	0.91920
R al quadrato corretto	0.91516
Errore standard	147.66972
Osservazioni	22

ANALISI VARIANZA					
	gdl	SQ	MQ	F	Significatività F
Regressione	1	4961434.406	4961434.406	227.523	2.17134E-12
Residuo	20	436126.913	21806.346		
Totale	21	5397561.318			

	Coefficienti	Errore Standard	Stat t	Valore di significatività	Inferiore 95%	Superiore 95%
Intercetta	559.460014	1450.69753	0.385648974	0.703828029	-2466.64201	3585.562035
X Income	0.38151672	0.02529306	15.08384918	2.17134E-12	0.328756319	0.434277121

(b)

IL MODELLO DI REGRESSIONE MULTIPLA

Il modello teorico della regressione multipla definisce la relazione tra una **variabile risposta Y** e un **insieme di variabili esplicative X_k** con $k = 1, \dots, n$; n rappresenta il numero totale delle osservazioni, mentre k rappresenta il numero totale delle variabili esplicative

Avremo infatti che:

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \varepsilon_i$$

Ci aspettiamo che la varianza di ε del modello di regressione semplice sia maggiore dell' ε del modello di regressione multipla

Assunzioni di base valide per ogni i



	<p>1. Le x_{ij} sono costanti o realizzazioni di una variabile aleatoria X, indipendenti dalle componenti aleatorie di errore ε_i</p> <p>2. Il valore atteso della variabile aleatorie Y è una funzione lineare delle variabili indipendenti $X_i \rightarrow E[Y X = x_k, k = 1,2, \dots]$</p> <p>3. I termini di errore sono variabili aleatorie con media 0 e varianza costante σ^2. La seconda condizione è detta omoschedasticità</p> $E[\varepsilon_i] = 0 \quad Var[\varepsilon_i] = \sigma^2$ <p>4. Gli errori aleatori non sono correlati tra loro e quindi $E[\varepsilon_i \varepsilon_j] = 0$</p> <p>5. <u>Non è possibile trovare un insieme di coefficienti</u>, tutti non nulli, c_1, c_2, \dots, c_n tali che:</p> $c_0 + c_1 x_1 + c_2 x_2 + \dots + c_n x_n = 0$ <p>Ovvero non esiste una relazione lineare tra le X_i</p>
<p>Stima dei minimi quadrati</p>	<p>Anche nel caso della regressione multipla, le stime dei parametri vengono determinati mediante gli stimatori dei minimi quadrati; tali stimatori sono BLUE, ovvero sono i migliori stimatori lineari non distorti</p>
<p>Interpretazione della stima dei coefficienti</p>	<p>Il coefficiente di regressione stimato di una variabile, <u>se significativo</u>, rappresenta <u>l'effetto sulla media della variabile Y dovuto ad una variazione unitaria di quella variabile</u>, a <u>parità</u> dei valori <u>delle altre variabili</u> esplicative. È importante specificare che la variabile è “significativa” e l'effetto è presente “in media” e “a parità delle altre variabili”, altrimenti la conclusione non è corretta</p>
<p>Scomposizione della varianza e indice R^2</p>	<p>Vale la scomposizione</p> $SST = SSR + SSE$ <p>Ovvero:</p> $\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2$ <p>E l'indice R^2 (o coefficiente di determinazione) è definito come:</p> $R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SSR} \text{ con } 0 \leq R^2 \leq 1$ <p>R^2 misura la variabilità spiegata dal modello rispetto alla variabilità totale e fornisce un'informazione circa la bontà del modello stesso: tanto maggiore è R^2 tanto migliore sarà la spiegazione fornita dal modello</p> <p>Vale la relazione per cui:</p> $R^2 = \text{corr}(y, \hat{y})^2$ <p>Ovvero il coefficiente di determinazione è uguale al quadrato della correlazione campionaria, al quadrato, tra i valori previsti ed osservati per y</p>
<p>Stima della varianza del modello</p>	<p>Una stima non distorta della varianza degli errori aleatori ε_i è chiamata varianza dell'errore o varianza dei residui è:</p> $s_e^2 = \frac{SSE}{n - k - 1} = \frac{\sum_{i=1}^n e_i^2}{n - k - 1}$ <p>Dove K rappresenta il numero di variabili indipendenti del modello.</p>



	<p>Si chiama errore standard della stima</p> $s_e = \sqrt{s_e^2}$
<p>Coefficiente di determinazione corretto</p>	<p>R^2 non può diminuire se nuove variabili esplicative vengono incluse nel modello: il confronto tra modelli basati sul valore di R^2 porterà sempre a scegliere il modello che include tutte le possibili variabili esplicative disponibili:</p> $R^2_{multipla} \geq R^2_{semplice}$ <p>Per confrontare due modelli di regressione multipla (o un modello di regressione multipla con un modello di regressione semplice) con un diverso numero di variabili esplicative ($K \neq K^*$) è possibile usare il coefficiente di determinazione corretto, che tiene conto della significatività statistica delle variabili inserite nel modello</p> $\overline{R^2} = 1 - \frac{SSE/(n - k - 1)}{SST/(n - 1)}$ <p>Con una formula alternativa risulta che:</p> $\overline{R^2} = 1 - (1 - R^2) * \frac{n - 1}{n - k - 1}$ <p>Si utilizza questo indice per bilanciare la riduzione, anche se minima, della somma dei quadrati degli errori che si determina con l'aggiunta di variabili esplicative non rilevanti. L'indice R^2 fornisce quindi un migliore strumento per confrontare modelli di regressione multipla con un numero diverso di variabili indipendenti.</p> <p>Vale sempre che $\overline{R^2} < R^2$</p>
<p>Coefficiente di correlazione multipla</p>	$R = r(\hat{Y}, Y) = \sqrt{\overline{R^2}}$ <p>Ci permette di misurare la relazione lineare tra la variabile dipendente e le variabili indipendenti</p>
<p>Interpretazione dei parametri e inferenza sui parametri singolarmente considerati</p>	<p>Sia dato il modello teorico di regressione multipla e siano b_0, b_1, \dots, b_k le stime dei minimi quadrati dei parametri del modello teorico e siano $s_{b_0}, s_{b_1} \dots$ le stime degli errori standard degli stimatori dei minimi quadrati. Allora possiamo costruire la seguente statistica test:</p> $T_{b_i} = \frac{b_i - \beta_i^*}{s_{b_i}} \sim T_{(n-k-1)}$ <p>Avremo che l'intervallo di confidenza:</p> $IC_{1-\alpha}(\beta_i) = \left(b_i \pm t_{n-k-1, \frac{\alpha}{2}} s_{b_i} \right)$ <p>Sarà anche possibile condurre il test d'ipotesi, supponendo la normalità degli errori, sui parametri singolarmente considerati (T-test)</p> $H_0: \beta_i = \beta_i^* \quad H_1: \beta_i \neq \beta_i^*$ <p>Valgono le seguenti regole di rifiuto:</p>



	$ T = \left \frac{b_1}{S_{B_1}} \right = t - value > t_{n-2, \frac{\alpha}{2}}$ $p - value = 2P(T^{(n-k-1)} > t_{oss}) < \alpha$ <p>I casi II e III sono analogamente risolvibili N.B.: t-value viene spesso indicato con Stat t</p> <p>Un test di significatività del parametro singolarmente considerato si conduce semplicemente imponendo:</p> $H_0: \beta_i = 0 \text{ (non significativo)}$ $H_1: \beta_i \neq 0 \text{ (significativo)}$ <p>Le regole di rifiuto sono le medesime del modello generale</p>
<p>Verifica della significatività di tutti i parametri congiuntamente e considerati (cenni)</p>	<p>Conduciamo un test di significatività generale del modello (F-test):</p> $H_0: \beta_1, \beta_2, \dots = 0 \text{ (il modello non è significativo)}$ $H_1 = \text{almeno un } \beta_i \neq 0 \text{ (almeno una variabile esplicativa è rilevante} \rightarrow \text{il modello è, nel complesso, significativo)}$ <p>La statistica test (statistica F) è:</p> $F = \frac{SSR/k}{SSE/(n-k-1)} = \frac{MSR}{s_e^2} = \frac{MSR}{MSE} \sim^{H_0} F_{k, n-k-1}$ <p>E la regola di rifiuto è:</p> $F > F_{k, n-k-1, \alpha}$ $p - value < \alpha$
<p>Situazioni problematiche nell'analisi di regressione: l'approccio dell'analisi grafica dei residui</p>	<p>Dopo aver costruito il modello di regressione lineare è buona norma condurre l'analisi grafica dei residui che consente di determinare la bontà del modello e per verificare se sono soddisfatte le ipotesi standard della regressione. Condurre tale analisi significa costruire:</p> <ul style="list-style-type: none"> • Un istogramma che permetta di verificare se la distribuzione dei residui sia simmetrica e uniforme • Costruire un diagramma a dispersione che permette di osservare eventuali osservazioni ad alto leverage o outliers <p>Osservazioni ad alto leverage</p>



	<p>osservazione ad alto leverage che tuttavia non è influente</p> <p>osservazione ad alto leverage che è anche influente</p>
	<p>Un'osservazione ad alto leverage è quindi differente/distante dalle altre in termini di valori delle variabili esplicative: presentano un'ordinata che si discosta notevolmente dalle altre.</p> <p>Per verificare che sia anche influente è possibile stimare i modelli con o senza l'osservazione</p> <p>Outlier nella regressione È una osservazione con valore elevato del corrispondente residuo, ovvero non ben descritta dal modello; in particolare, essi sono punti la cui ascissa si discosta notevolmente</p> <p>La sua influenza può essere valutata come nel caso delle osservazioni ad alto leverage</p>
<p>Multicollinearità</p>	<p>È l'esistenza di forti (ma non perfette) relazioni lineari tra le variabili esplicative nel dataset osservato. Come ricordiamo (è la condizione aggiuntiva) non esistono dei coefficienti tali che:</p> $c_0 + c_1x_1 + \dots + c_kx_k = 0$



	<p>Escludendo l'esistenza di relazioni lineari esatte tra le variabili esplicative. Tuttavia, relazioni di forte intensità, seppur non esatte, generano problemi:</p> $c_0 + c_1x_1 + \dots + c_kx_k \approx 0$ <p>Alcuni indizi di multicollinearità sono:</p> <ul style="list-style-type: none"> • I coefficienti di regressione sono molto diversi, anche dal punto di vista del segno, da quelli che ci si potrebbe attendere secondo le teorie economiche o l'esperienza (logica) • I coefficienti delle variabili che si ritengono rilevanti hanno una statistica T molto bassa, trattandosi quindi di valori nulli • Le statistiche T di tutti i coefficienti sono basse (i coefficienti singolarmente considerati sono ritenuti non significativi), ma il valore della statistica F del modello indica la complessiva significatività statistica del modello stesso • Vi sono forti correlazioni tra le coppie di variabili indipendenti: in un modello di regressione semplice una variabile era ritenuta significativa; inserendolo nel modello di regressione multipla la stessa variabile perde la sua significatività <p>Effetti Gli stimatori dei minimi quadrati sono ancora non distorti ma la varianza e lo standard error sono inflazionati, ovvero di valore manifestamente elevato. L'aggiunta di variabili collineari porta a creare un modello più complesso ma non maggiormente esplicativo (ridondante)</p> <p>Come eliminarla Un <i>suggerimento generale</i> può essere quello di eliminare semplicemente una delle variabili ritenute tra loro collineari</p>
--	--

R/RADIANT	
Attribuzione di nome ad un oggetto	Oggetto semplice <code>x <- log(3)</code> Oggetto complesso <code>z <- c(2, 4, 3, 1)</code>
Funzioni	Media <code>mean(z)</code> Matrice <code>matrix(z, nrow=2, ncol=2, byrow=TRUE)</code>
Caricare un dataset	Selezionare la cartella di lavoro File → Cambia Directory Caricare il file <code>load(file="nome_file.RData")</code>
Funzioni da file	Funzione str: permette di indagare la natura dell'oggetto <code>str(nome_file)</code> Funzione media <code>mean(nome_file\$nome_variabibile)</code> Funzione summary: ci permette di visualizzare gli indicatori di sintesi che il software ritiene siano appropriati per la specifica variabile) <code>summary(nome_file\$nome_variabibile)</code>



Rappresentazioni grafiche	<p>Funzione table: costruisce una distribuzione di frequenza assoluta <code>table(nome_file\$nome_variabile)</code></p> <p>Grafico a torta: funzione pie(table) <code>pie(table(nome_file\$nome_variabile))</code> <code>pie(table(nome_file\$nome_variabile)/length(nome_file\$nome_variabile))</code></p> <p>Boxplot per una singola variabile: funzione boxplot <code>boxplot(nome_file\$nome_variabile)</code></p>
Radiant: operazioni di base	<p>Radiant è caratterizzato da un'interfaccia di menu a tendina. Mediante il pulsante <i>load</i> è possibile caricare i dataset che intendiamo analizzare</p> <p>Possiamo usare:</p> <ul style="list-style-type: none"> • View per analizzare globalmente il dataset • Visualize per creare le rappresentazioni grafiche; dal menu che viene visualizzato possiamo settare tutte le caratteristiche della rappresentazione • Pivot per generare le distribuzioni di frequenza assolute; cliccando si <i>normalize by</i>: <ul style="list-style-type: none"> ○ <i>Total</i> possiamo generare le frequenze congiunte relative ○ <i>Row</i> possiamo generare le frequenze condizionate per riga ○ <i>Column</i> possiamo generare le frequenze condizionate per colonna • Explore per generare gli indici di sintesi delle diverse variabili, dal sottomenu <i>apply function</i>
Analisi univariata e istogrammi	<p>Radiant permette di costruire istogrammi con <u>classi di uguale ampiezza</u>: Data → Visualize:</p> <ol style="list-style-type: none"> 1. Plot type → distribution 2. X-variabile → selezionare la variabile d'interesse 3. Create plot <p>Lo slider Number of bins permette di scegliere il numero di classi da usare</p> <p>Radiant permette di costruire anche istogrammi per <u>distribuzioni condizionate</u> utilizzando il comando <i>Facet column</i> e selezionando la variabile condizionante</p> <p>Con R/RStudio è possibile costruire istogrammi con <u>classi di diversa ampiezza</u>: occorre usare la funzione hist() che si compone dei seguenti argomenti:</p> <ul style="list-style-type: none"> • nome_file\$variabile • main = "titolo istogramma" • xlab = "nome asse x" • ylab = "nome asse y" • breaks = c(0, estremo1, estremo2, estremo3, ..., estremoN) <p>R ci permette di creare rapidamente anche i boxplot relativi ad <u>una sola variabile</u> con la funzione boxplot <code>boxplot(nome_file\$nome_variabile)</code></p>
Analisi bivariata	<p><u>Variabili categoriche</u></p> <p>Radiant ci permette di creare <u>tabelle a doppia entrata</u> mediante il tab <i>Pivot</i>, selezionando le due variabili da analizzare nel box <i>Categorical variables</i> e cliccando sul pulsante <i>Create pivot table</i>: abbiamo ottenuto le distribuzioni assolute</p> <ul style="list-style-type: none"> • Con il comando Normalize by: total possiamo ottenere la distribuzione relativa; aggiungendo una spunta accanto a <i>Percentage</i> possiamo visualizzare la distribuzione percentuale • Con il comando Normalize by: column/row otteniamo la distribuzione condizionata secondo la variabile disposta in colonna/riga



	<ul style="list-style-type: none"> • Cliccando su Show plot è possibile visualizzare il diagramma a barre accostate; con <i>plot type</i> → <i>fill</i> otteniamo il grafico a barre sovrapposte <p><u>Variabili numeriche</u> Dal menu Basics → Correlation possiamo selezionare le variabili da analizzare e cliccare sul checkbox <i>Show covariance matrix</i>. Nella tab <i>Plot</i> sono visualizzabili gli scatterplot relativi</p>
Probabilità	<p>Le probabilità possono essere valutate mediante il menu <i>Basics</i></p> <p>In Probability calculator possiamo calcolare le probabilità da varie distribuzioni selezionabili in <i>Distribution</i> Possiamo lavorare con <i>Input type</i>:</p> <ul style="list-style-type: none"> • Values che definisce i valori estremi dell'intervallo; da qui possiamo calcolare l'area sottesa alla curva tra gli estremi dell'intervallo, ovvero la probabilità. In pratica, ci permette di trovare la probabilità p per cui: $P(\alpha \leq X \leq \beta) = p$ <ul style="list-style-type: none"> • Probability ci permette di definire i percentili a cui è associata l'area sottesa sulla coda sinistra. In pratica, ci permette di trovare la x per cui: $P(X \leq x) = 0,95$
Intervalli di confidenza	<p>È possibile trovare gli intervalli di confidenza utilizzando, dal menu <i>Basics</i>, la sezione relativa al parametro d'interesse (Single mean, Compare means, proportions...) Occorre eseguire un test Two Sided ed impostare il livello di confidenza con lo slider. Il risultato si può osservare sotto le percentuali nel tabulato, che indicano gli estremi dell'intervallo</p>
Inferenza sulla media di una popolazione normale	<p>R e Radiant mettono a disposizione comandi solo nel caso di varianza non nota; per farlo andiamo nel menu Basics → Single mean A questo punto:</p> <ul style="list-style-type: none"> • Selezioniamo la variabile da analizzare nel box <i>Variable</i> • Scegliere il tipo di ipotesi alternativa da testare nel box <i>Alternative hypothesis</i> con valori possibili <i>Two sided, Less than, Greater than</i> • Scegliamo il livello di confidenza in <i>Confidence level</i> • Indichiamo il valore da testare sotto ipotesi nulla in <i>Comparison value</i> <p>Ci viene restituito sia l'intervallo di confidenza sia il p-value; valori molto piccoli di p-value vengono indicati con la dicitura < 0.001</p>
Inferenza sulla proporzione di successi da una popolazione bernoulliana	<p>Basics → Single proportion</p> <ul style="list-style-type: none"> • Selezioniamo la variabile da analizzare nel box <i>Variable</i> • Scegliere il tipo di ipotesi alternativa da testare nel box <i>Alternative hypothesis</i> con valori possibili <i>Two sided, Less than, Greater than</i> • Scegliamo il livello di confidenza in <i>Confidence level</i> • Indichiamo il valore da testare sotto ipotesi nulla in <i>Comparison value</i> <p>Occorre sempre selezionare il test type: Z-test</p>
Inferenza sul confronto tra medie di popolazioni normali	<p>Ci concentreremo sul caso di varianze non note e diverse</p> <p><u>Campioni indipendenti</u> Basics → Compare means</p> <ul style="list-style-type: none"> • <i>Select a factor o numeric variable</i> inseriamo la variabile che identifica i gruppi di cui vogliamo confrontare le medie (tipo factor)



	<ul style="list-style-type: none"> • <i>Numeric variable</i> selezioniamo la variabile numerica di cui vogliamo confrontare le medie • <i>Alternative hypothesis</i> indichiamo che tipo di ipotesi vogliamo testare (Two sided, less/greater than) • Inseriamo il <i>Confidence level</i> • Spuntiamo il box <i>Show additional statistics</i> • <i>Sample type</i> → <i>independent</i> • <i>Test type</i> → <i>t-test</i> <p><u>Campioni dipendenti</u> La procedura è uguale ma dobbiamo selezionare il box <i>Sample type</i> → <i>paired</i></p>
Indice di correlazione lineare	<p>L'indice di correlazione lineare costituisce lo strumento principale per valutare l'intensità dell'associazione lineare tra due variabili numeriche.</p> <p>Basics → Correlation</p> <p>Ci vengono riportati l'indice di correlazione lineare e il p-value: bassi valori del p-value implicano un'associazione altamente significativa</p>
Bontà di adattamento	<p>Consente di valutare se una distribuzione di probabilità osservata per una variabile aleatoria discreta (con supporto finito) possa essere considerata "compatibile" con una distribuzione di probabilità teorica definita a priori</p> <p>Basics → Goodness of fit</p> <ul style="list-style-type: none"> • La variabile nel data set che contiene i dati con cui effettuare l'analisi (box <i>Select a categorical variable</i>) • Le probabilità che definiscono la distribuzione teorica con cui confrontare quella osservata (box <i>Probabilities</i>); queste probabilità devono essere inserite come numeri compresi tra 0 e 1, la cui somma deve essere pari a 1 e il cui numero deve essere pari alle categorie osservate per la variabile indicata nel box precedente • Il tipo di output da mostrare (box <i>Observed</i> per le frequenze osservate, <i>Expected</i> per quelle attese sotto l'ipotesi nulla, <i>Chi-squared</i> per il contributo al calcolo dell'indice χ^2 di ogni categoria)
Test d'indipendenza	<p>Basics → Cross-tabs</p> <ul style="list-style-type: none"> • Il nome delle variabili che contengono i dati su cui vogliamo effettuare il test di indipendenza (due box successivi, entrambi denominati <i>Select a categorical variable</i>; consigliamo di selezionare la variabile di riga nel primo box e quella di colonna nel secondo) • Il tipo di output da mostrare, ovvero le frequenze osservate (<i>Observed</i>), quelle attese sotto l'ipotesi nulla di indipendenza (<i>Expected</i>), i contributi al calcolo dell'indice χ^2 di ogni cella della tabella (<i>Chi-squared</i>), le radici quadrate dei medesimi contributi (<i>Deviation std.</i>), le frequenze condizionate date le righe, quelle condizionate date le colonne o quelle congiunte relative (rispettivamente <i>Row percentages</i>, <i>Column percentages</i> e <i>Table percentages</i>).
Modello di regressione lineare semplice	<p>La funzione principale per stimare il modello di regressione lineare in R è lm() che presenta la seguente sintassi:</p> $reg1 <- lm(y \sim x, data = nome_dataset)$ <p>Possiamo assegnare un nome a questa funzione e successivamente impiegare la funzione</p>



	<p style="text-align: center;"><i>summary(reg1)</i></p> <p>Questa funzione ci restituisce:</p> <ul style="list-style-type: none"> • Nella sezione Residuals una sintesi dei residui del modello • Nella sezione Coefficients fornisce le stime dei coefficienti del modello insieme ai relativi p-value • Nelle ultime tre righe vengono forniti: <ul style="list-style-type: none"> ○ La stima di σ (Residual standard error) ○ L'indice R^2 (Multiple R-squared) ○ Il test F (ultima riga) <p>Per ottenere una visualizzazione grafica possiamo usare il seguente codice:</p> <p style="text-align: center;"><i>plot(y~x, data = nome_dataset)</i> <i>abline(reg1, lwd = 2, col = "blue")</i></p> <p>R ci fornisce solamente i p-value per i test sui due coefficienti ma non gli intervalli di confidenza per entrambi i coefficienti; questi possiamo ottenerli:</p> <p style="text-align: center;"><i>confint(reg1, level = α)</i></p> <p>Per ottenere la scomposizione della varianza usiamo</p> <p style="text-align: center;"><i>anova(reg1predict)</i></p>
Modello di regressione lineare multipla	<p>Ancora una volta usiamo la funzione <i>lm()</i> separando le variabili con un +</p> <p style="text-align: center;"><i>reg2 <- lm(y~x₁ + x₂ + ... + x_n, data = comp)</i> <i>summary(reg2)</i></p> <p>Come nel modello semplice possiamo usare le funzioni <i>plot()</i> e <i>anova()</i></p>
Inferenza sui coefficienti del modello	<p>Usiamo la funzione <i>confint()</i> per ottenere l'intervallo di confidenza (livello standard = 0.95) per i parametri β_0 e β_1</p> <pre>lm.fiori <- lm(height ~ distance, data=fiori) confint(lm.fiori, level = 0.9)</pre>
Previsioni per un modello di regressione lineare	<p>Uno dei motivi della popolarità del modello di regressione lineare consiste nella possibilità di calcolare previsioni in relazione a scenari che non sono stati necessariamente osservati nel campione.</p> <p>La funzione da usare in R è <i>predict()</i> con i seguenti argomenti:</p> <ul style="list-style-type: none"> • object, ovvero l'oggetto restituito dalla funzione <i>lm()</i> contenente i risultati della stima • newdata, che indica un nuovo data frame contenente i valori delle variabili indipendenti da utilizzare per il calcolo delle previsioni • interval, che può assumere solo tre possibili valori, ovvero confidence se si desidera calcolare gli intervalli di confidenza per la previsione del valore medio, prediction se si desidera calcolare gli intervalli di previsione per i valori individuali, none se non si vuole calcolare nessun intervallo • level, il quale consente di specificare il livello di confidenza per il calcolo degli intervalli di confidenza o previsione <p><u>Previsione di Y dato X</u></p>



```
predict(object, newdata = data.frame(variabile = numero))
```

N.B.: per valutare se una previsione è affidabile dobbiamo valutare la variabile esplicativa richiesta è all'interno del range di valori della variabile: per farlo, controlliamo gli estremi con la funzione *summary(variabile)*

Previsione intervallare per singolo valore

```
predict(object, newdata = data.frame(variabile = numero), interval  
= "prediction", level = 0.90 )
```

Confidenza intervallare per media dei valori

```
predict(object, newdata = data.frame(variabile = numero), interval  
= "confidence", level = 0.90 )
```

