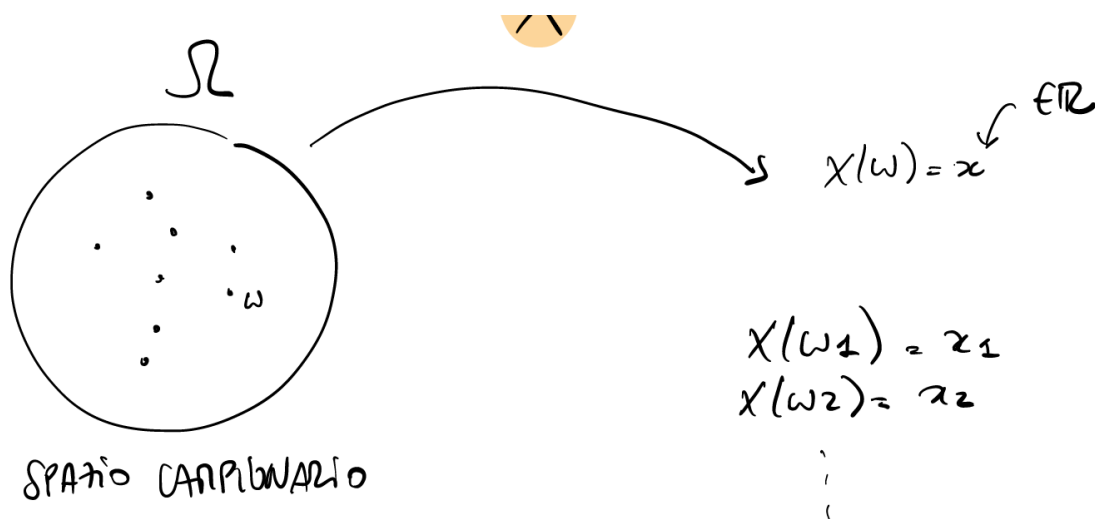


Lezione 1

1. **Statistica descrittiva** : descrivere il campione casuale aka in questo caso la variabile Luigi e' il campione casuale
2. **Statistica inferenziale** : generalizzare all'intera popolazione alla quale appartengono i soggetti

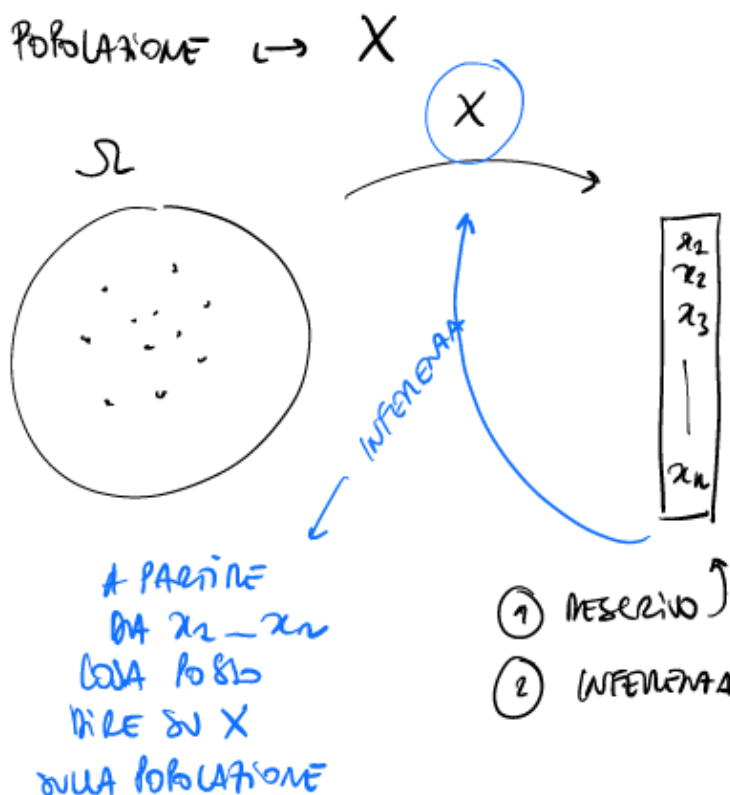
Campione casuale e cosa sono le variabili :



In uno spazio campionario (omega) con all'interno degli omega piccolo

La variabile aleatoria e' il valore che viene associato al punto omega dell'esperimento

Campionare gli elementi di omega tramite estrazione con reimbussolamento (rimetto il campione nello spazio campionario)



Le variabili del dataset x_1, x_2, x_n sono le n misurazioni di una variabile aleatoria. Le descrivo e le inferenze (passare da n osservazioni a una qualche informazione sulla popolazione)

Statistica descrittiva: tipologia delle variabili

1. **Variabili quantitative:** misurano i numeri i quali sono interessanti come quantità
 - **Discrete:** variabili di conteggio
 - **Continue:** misurano grandezze fisiche
2. **Variabili qualitative / fattori:** misurano le caratteristiche del soggetto, codificate con etichette / livelli per appartenenza di una certa classe (colore degli occhi)

Character Data: variabili identificativi e non si analizzano

Lezione 2

Statistica descrittiva univariata: lavoriamo su una variabile alla volta, divideremo le analisi secondo la tipologia di analisi

Variabili quantitative:

Calcolare degli indici → Summaries (riassuntivi)

collezione di pochi indici che siano pieni di significato (sostituiamo 130 numeri in 5)

1. Categoria degli **INDICI DI POSIZIONE:** collezione di indici che si occupano di indicarci dove troviamo i dati sull'asse reale
 - **MEDIA CAMPIONARIA** (\bar{x} sbarra sopra), e' pero influenzata da eventuali **valori estremi** (lontano dagli altri)

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

(MEDIA ARITMETICA DELLE OSSERVAZIONI)

$$\bar{x} = 3.72$$
 - **MEDIANA CAMPIONARIA:** ci da l'indice di posizione del valore centrale delle osservazioni
 su 7 numeri [1,2,3,4,5,6,7] il valore centrale e' 4
 su 6 numeri [1,2,3,4,5,6] il valore centrale e' $\frac{3}{4}$
 - **QUANTILI CAMPIONARI:** **Q25=** pedice indica un numero tra 0 - 100, la percentuale di osservazioni che lasci a sinistra
 - **QUARTILI:** Q25 - Q50 - Q75 = Quartili, perche' dividono in quarti le osservazioni

2. Categoria degli **INDICI DI DISPERSIONE:**

- **Range:**
- **Varianza campionaria:** definita con s quadro, la varianza e' in un'altra scala perche' eleviamo al quadrato
 → per tornare all'unita di misura delle osservazioni $s = \sqrt{s^2}$
- **Coefficiente di variazione:**

$$CV = \frac{s}{\bar{x}} = \begin{cases} < 1 & \text{VARIAB. < EXP} \\ 1 & \text{VARIABILITA' PAR. A QUELLA DELL'ESPOSIZIONE} \\ > 1 & \text{VARIAB. > EXP} \end{cases}$$

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

MEDIA CAMPIONARIA
 STARS QUADRATIC
 OVA MEDIA CAMPIONARIA
 ... MEDIA ARITMETICA

3. categoria degli **Indici di forma:**

- ## LEZIONE 3

-

si tabula la variabile : conto quante osservazioni in ciascuna categoria, ovvero quante volte compare una certa label

Tabella di conteggio

A	B	C	D
3	7	9	1

BAR CHARTS

Fruit	Number of People
A	3
B	7
C	9
D	1

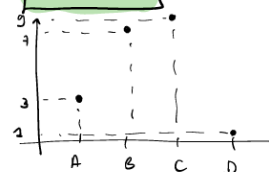
- **Bar charts:** altezza variabile di ciascuna variabile
- **Dot chart:**
- **Pie Chart:** DAI CHE LO SAI E NON LO DEVO METTERE IL DISEGNINO, che tanto sono inutili

DOT CHART

Variabile quantitativa / Variabile qualitativa

Siamo interessati alle distribuzioni delle distribuzioni nella v.quantitativa nelle diverse classi

VAR 1	VAR 2
x_1	A
x_2	B
1	A
1	1
	1
x_n	A



LEZIONE 4

V. quantitativa / V. quantitativa : che relazione hanno X e Y, hanno un qualche legame ?

- **Scatterplot**: grafico a punti delle informazioni
 - 1) **nuvola informe** : i punti si spargono in modo casuale (non c'è relazione tra X e Y)
 - 2) Osservazioni si dispongono lungo il grafico di una funzione

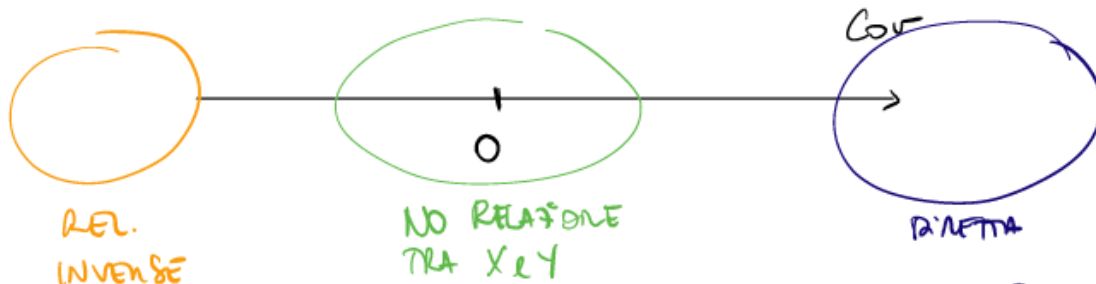
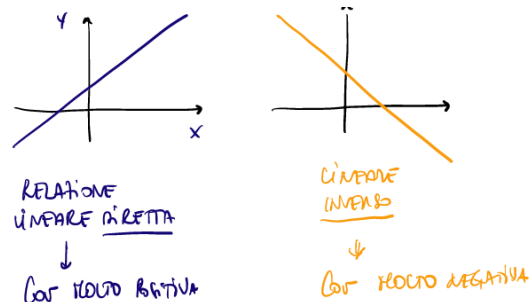
Anche degli indici che riassumano e quantifichino la relazione:

Media varianza - MOMENTI

Covarianza(x,y), media $x \cdot y$ - la media di x - la media di Y

Covarianza positiva, Covarianza Negativa

Versione Normalizzata, correlazione(X,Y) [-1, 1]

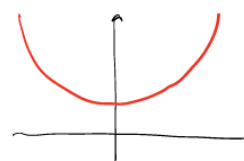


1]

LA COV. PUÒ VENIRE UGUALE A 0 ANCHE SE LE VAR. SONO LEGATE es: $Y = X^2$

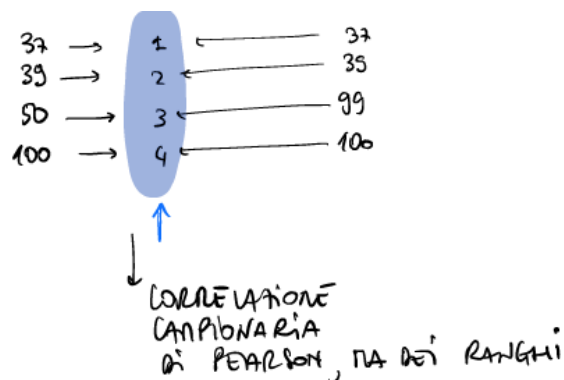
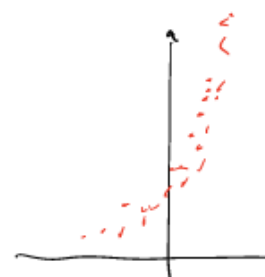
INDICE DI CORRELAZIONE CAMPIONARIA DI PEARSON:

stima la correlazione di prima appartenente all'intervallo [-1, 1], più si avvicina allo zero e meno sono correlati, più si avvicina a uno e più sono correlati, a -1 sono negativamente correlati



CORRELAZIONE CAMPIONARIA DI SPEARMAN: meglio di quella di Pearson, la calcolo se la relazione non è lineare ma sembra monotona

Calcolo i ranghi: posizioni di ciascuna osservazione nel vettore ordinato in maniera crescente



V.QUALI / V.QUALI :

TABELLE DI FREQUENZA

	A	B	C	D
PIPPO	32	1	0	0
PLUTO	5	5	3	21

LEZIONE 5:

Statistica inferenziale:

Spazio campionario : la popolazione

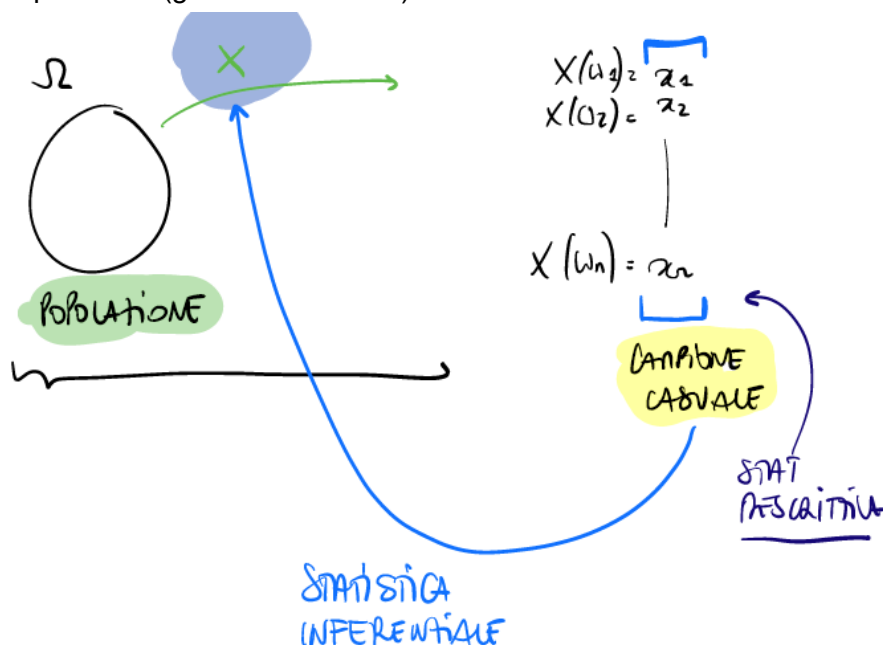
campiono gli individui : numerandoli → estrazioni con reimbussolamento

Campione casuale(dati): il nome della persona → statistica descrittiva

Parleremo di campione di misure, di popolazione delle misure

Confronto individui e misure , ci interessano le misure

A partire dal campione casuale, vorrei dare informazioni sulla popolazione, la statistica descrittiva non cambia, ti dice sempre quei numeri delle persone della popolazione, la statistica inferenziale invece va oltre, a partire dal campione casuale vorrei dare informazione sulla popolazione, ovvero dalla variabile aleatoria X dalla quale sto campionando (generalizzazione)



Statistica inferenziale parametrica :

statistica inferenziale → vogliamo informazioni su X , la popolazione

parametrica → x e' una v.aleatoria

- 1) Famiglia : v.a → x e' Bernoulli (evento con successo 1,0 insuccesso), x e' Binomiale, geometrica, Poisson, normale, esponenziale, uniforme // e' come se le conoscessi
- 2) Parametri : x e' Bernoulli(p) , binomiale(n,p), normale(μ, σ^2), esponenziale(λ)

La statistica inferenziale parametrica si occupa di dare informazioni sui parametri della popolazione X

x binomiale (n, p) = media = $n \cdot p$ e poi Varianza = $n \cdot p(1-p)$

x normale $(\mu, \text{sigma quadro})$ = media = μ e poi varianza = sigma quadro

media e varianza \rightarrow funzioni dei parametri delle famiglia

Parametri sono 3 : Media - Varianza - Proporzione

Proporzione : ci riferiamo al parametro **P** della X_n Bernoulli(p)

ci servono per codificare le variabili di tipo qualitativo

$\rightarrow x = 1 \rightarrow p$ 1=capelli rossi

$\rightarrow x = 0 \rightarrow (1 - p)$ 0=capelli non rossi

Variabilità campionaria: do informazioni sui parametri che controllano la variabilità campionaria, delle tecniche che tengano conto che il campione e' causale e che avrebbero potuto darci altri 130 numeri diversi

IDEA: Valutare quanto diversi potrebbero essere altri 130 numeri a partire dalla variabilità nei 130 che mi hanno dato

COME: 1) **calcolare una stima puntuale**(calcolato a partire dal campione casuale e che aspettiamo che sia vicino al parametro che non conosciamo) per il parametro che ci interessa ,ma lo abbiamo già fatto :

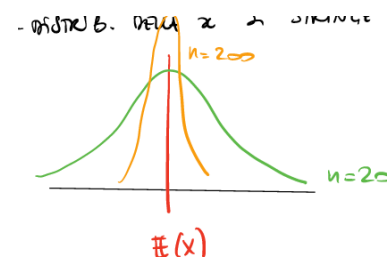
Media popolazione \rightarrow Media campionaria

Varianza popolazione \rightarrow Varianza campionaria

Proporzione \rightarrow proporzione campionaria // frequenze relative vengono chiamati stimatori, sono funzioni del campione casuale(devono essere calcolabili, senza incognite), che forniscono stime puntuali, ovvero valori ***vicini*** al valore del parametro incognito

vicini PERCHE? (non lo dice nessuno) :

- 1) la media campionaria cambia quando cambio il campione casuale, e' una v.a. ha la sua distribuzione
- 2) la distribuzione della \bar{x} cappello e' centrata la media campionaria e' uno stimatore corretto per la media della popolazione
- 3) quando n (taglia del campione) aumenta la distribuzione della media campionaria si stringe
- 4) La forma a campana c'e' sempre qualunque sia la distribuzione della popolazione x



2) **controllare la variabilità**: siccome non viene calcolata dalla stima puntuale,

- intervalli di confidenza

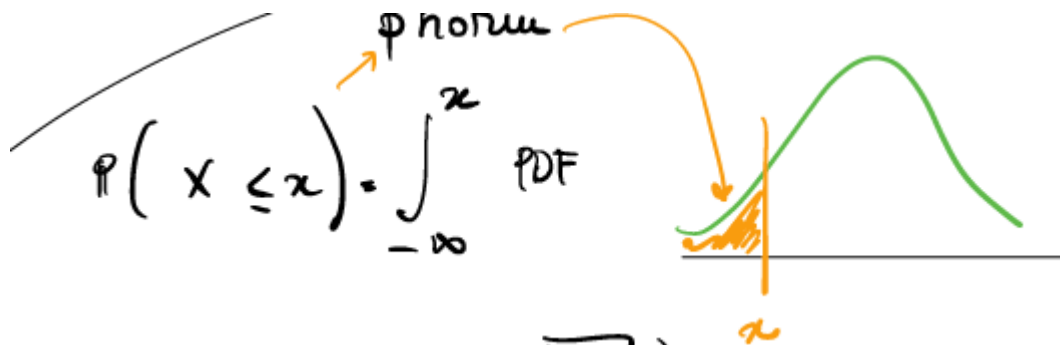
- test di ipotesi

LEZIONE 6

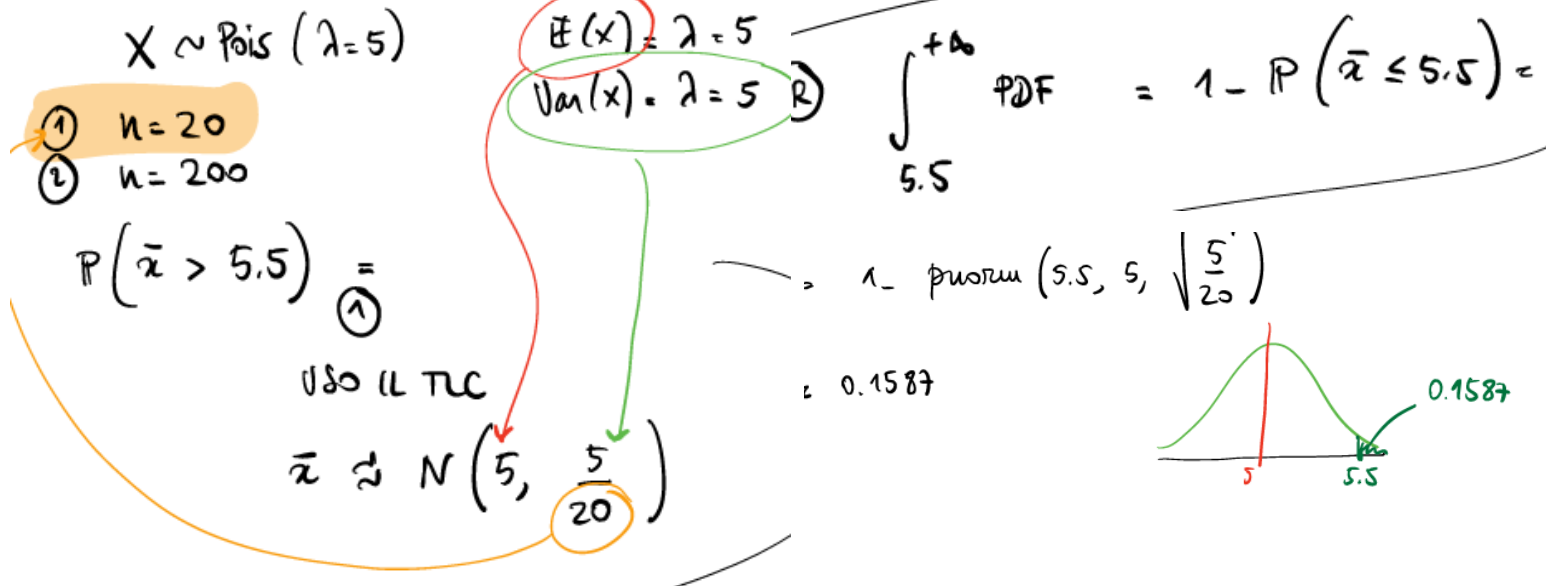
Le 4 osservazioni di prima sono riassunte in 2 teoremi:

- 1) **TEOREMA DEL LIMITE CENTRALE:** X la popolazione qualunque





pNorm calcola la parte a sx della X



$$\bar{x} \xrightarrow{n \rightarrow +\infty} E(X)$$

2) LEGGE DEI GRANDI NUMERI:

\bar{x} medio converge alla media della popolazione, sono risultati asintotici, le approssimazioni vanno bene quando n è grande

1) s quadro cambia al caviare del campione casuale, e' una v.a. anche s quadro \rightarrow ha una distribuzione, come si comporta?

2) s quadro e' corretto cioe' $E(s^2) = \text{Var}(X)$

3) c'e' una leggera asimmetria

4) se n tende a infinito (cresce) la distribuzione di s quadro si stringe cioe' var (s quadro) decresce

LGN: s quadro tende alla $\text{Var}(x)$

Per s quadro non vale il TLC cioe non posso dire che e' circa Gaussiana

Per la p cappello valgono tutti i risultati che abbiamo visto per la media campionaria

TLC per p cappello: la popolazione e' un Bernoulli

LGN per p cappello

$$\hat{p} \xrightarrow{n \rightarrow +\infty} p$$

$$\hat{p} = \bar{z}$$

METODI DI CONTROLLO DELLA VARIABILITA' → INTERVALLI DI CONFIDENZA:

due metodi che fanno la stessa cosa

$[a, b]$ → intervallo di numeri

restituiamo una confidenza indicata come $(1-\alpha) = 0.95/0.90/0.99$ → alta probabilità di essere vero

Come leggerli: la probabilità che l'intervallo $[a, b]$ contenga il valore vero del parametro (il valore vero del parametro) e' pari a $(1-\alpha)$ aka la confidenza

La confidenza non e' un risultato, ma una scelta, la fisso prima di iniziare

$(1-\alpha) = 0.95/0.90/0.99$

In pratica di tipo che il numero sta all'interno dell'intervallo con probabilità $(1-\alpha)$

$[a, b]$?

IC per proporzioni:

- Metodo della quantità pivotale

Passo 1) Ti danno la quantità, e' una funzione del campione casuale, che contiene il parametro incognito MA della quale conosci la distribuzione

$$\frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \sim N(0,1) \quad \text{tipo qua e' una}$$

$n \geq 30$

normale $(0,1)$

Passo 2) posso calcolare due valori q_1 e q_2 tali che:

la probabilità $(q_1 < \text{q.pivotal} < q_2) = (1-\alpha)$

$$P(q_1 < \text{Q.PIVOTALE} < q_2) = 1-\alpha$$

Praticamente sono i QUANTILI, potendo esplorare tutte le aree a pedice l'area che lasciano a sx

$(1-\alpha) = 0.95$

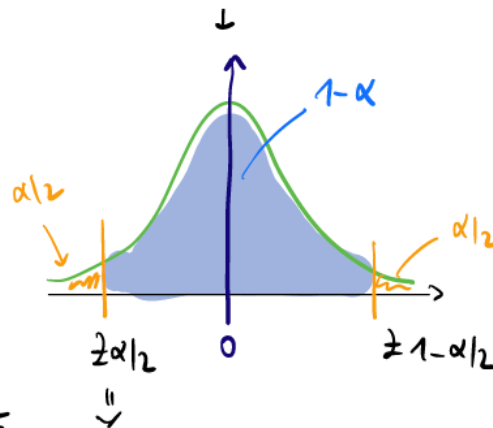
$q_2 = z_{0.95}$

Scelgo $q_1 = z_{\alpha/2}$

$q_2 = z_{1-\alpha/2}$

Cosa vuol dire? la scelta simmetrica, l'area α la dividiamo e in mezzo $1-\alpha$

Perche' e' la scelta ottimale perche si dimostra che restituisce gli IC più stretti



$$1-\alpha = 0.95$$

$$\alpha = 0.05$$

$$\alpha/2 = 0.025$$

$$1-\alpha/2 = 0.975$$

$$z_{0.025} = -1.96 = q_{\text{norm}}(0.025)$$

$$z_{0.975} = +1.96 = q_{\text{norm}}(0.975)$$

$$P(-1.96 < \frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} < 1.96) = 0.95$$

LEZIONE 7:

Metodo della quantità pivotale (continuo esercizio lezione 6), ti danno una variabile aleatoria la quale contiene un oggetto che tu sai come viene distribuito

$$P \left(\underbrace{-1.96}_{z_{1-\alpha/2}} < \frac{\hat{p} - p}{\sqrt{\hat{p}(1-\hat{p})/n}} < \underbrace{1.96}_{z_{1-\alpha/2}} \right) = \underbrace{0.95}_{1-\alpha}$$

PASSO 3:

TRASFORMARE QUELLO CHE È SCRITTO NELLA FORMA CON LA QUALE ABBIAMO VISTO L'IC

"LA P CHE IL VALORE VERO DEL PARAMETRO SIA COMPRESO TRA [a, b] È PARI A 1-α"

$$P(a < p < b) = 1 - \alpha$$

$$P \left(z_{\alpha/2} \cdot \sqrt{\hat{p}(1-\hat{p})/n} < \hat{p} - p < z_{1-\alpha/2} \cdot \sqrt{\hat{p}(1-\hat{p})/n} \right) = 1 - \alpha$$

$$P \left(z_{\alpha/2} \cdot \sqrt{\hat{p}(1-\hat{p})/n} - \hat{p} < -p < z_{1-\alpha/2} \cdot \sqrt{\hat{p}(1-\hat{p})/n} - \hat{p} \right) = 1 - \alpha$$

$$P \left(\underbrace{\hat{p} - z_{1-\alpha/2} \cdot \sqrt{\hat{p}(1-\hat{p})/n}}_a < p < \underbrace{\hat{p} - z_{\alpha/2} \cdot \sqrt{\hat{p}(1-\hat{p})/n}}_b \right) = 1 - \alpha$$

IC: $\left[\hat{p} - z_{1-\alpha/2} \cdot \sqrt{\hat{p}(1-\hat{p})/n}, \hat{p} - z_{\alpha/2} \cdot \sqrt{\hat{p}(1-\hat{p})/n} \right]$

E = errore standard, simmetrico alla p cappello.

se n aumenta Errore standard diminuisce, l'IC diventa più piccolo

se la confidenza aka (1-α) aumenta $z_{1-\alpha/2}$, allora $z_{1-\alpha/2}$ si alza

IC diventa più largo

$$E = z_{1-\alpha/2} \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = -z_{\alpha/2} \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

ERRORE STANDARD

$$[\hat{p} - E, \hat{p} + E]$$

IC per la media:

Passo 0: fisso $1-\alpha$

$$\frac{\bar{x} - \mu(x)}{\frac{s}{\sqrt{n}}}$$

Passo 1 : mi danno una quantità pivotale \rightarrow

\bar{x} medio= il valor medio

$EE(x)$ =media della popolazione

s =Deviazione Standard campionaria

se hai X una taglia maggiore di 30, $n > 30$

$$\frac{\bar{x} - \mu(x)}{\frac{s}{\sqrt{n}}} \sim t(n-1)$$

+ DI STUDENT
CON $n-1$
GRADI DI
LIB.

puoi considerarlo t di student con $n-1$

$$\frac{\bar{x} - \mu(x)}{\frac{s}{\sqrt{n}}} \sim t(n-1)$$

se X normale, puoi affermare che

Chi è una $t(n-1)$:

- 1) Una variabile aleatoria continua
- 2) PDF

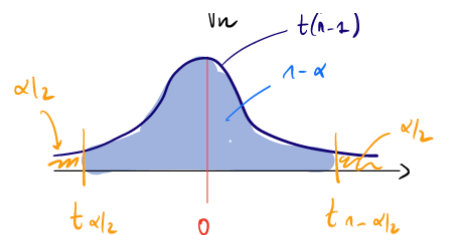
sulle code rimane più massa, ha le code alte, quando degrees of freedom (grado di libertà) aumentano, assomiglia sempre più alla $n(0,1)$

È simile alla $n(0,1)$

$$P\left(t_{\alpha/2} < \frac{\bar{x} - \mu(x)}{\frac{s}{\sqrt{n}}} < t_{1-\alpha/2}\right) = 1-\alpha$$

$- t(n-1)$

PASSO 2 :



$$1-\alpha = 0.95 \quad \alpha = 0.01 \quad \alpha/2 = 0.005$$

$$t_{\alpha/2} = t_{0.005} = \textcircled{R} = -2.68$$

$$n=50 \quad t(49)$$

$$t_{1-\alpha/2} = 2.68$$

$$IC \left[\bar{x} - t_{1-\alpha/2} \cdot \frac{s}{\sqrt{n}}, \bar{x} - t_{\alpha/2} \cdot \frac{s}{\sqrt{n}} \right]$$

Passo 3: Riscrivo, cioè

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

Errore Standard:

1) Se n cresce \rightarrow IC più piccolo

2) se $(1-\alpha)$ cresce $\rightarrow t_{1-\alpha/2} \rightarrow$ IC più grande

Prima di fare gli esercizi : Check \rightarrow 1) se $n \geq 30$

oppure

2) X normale

LEZIONE 8 :

Intervalli di confidenza per la varianza:

s quadro = varianza campionaria

var(x)

x popolazione

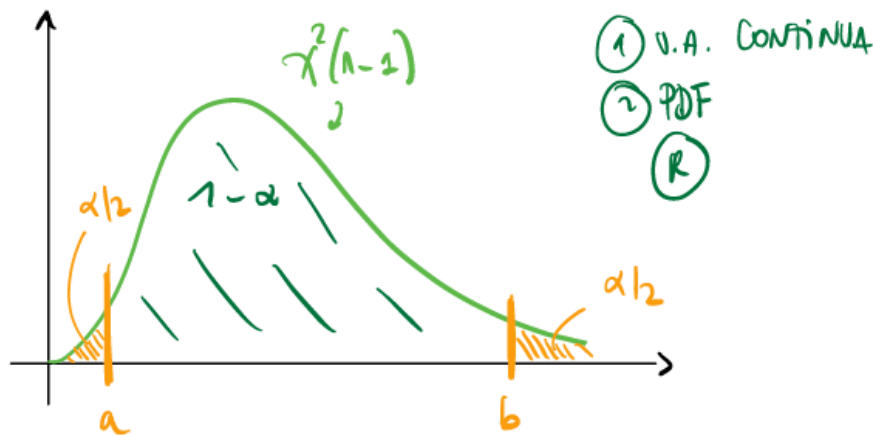
Passo 1: $(1-\alpha)$

Passo 2: ricevo la variabile Pivotal $\rightarrow \frac{(n-1)s^2}{\text{Var}(x)}$ chi quadro con $n-1$ gradi di libertà se X e'

normale $\chi^2(n-1)$

$$P(a < \text{q. PIVOTALE} < b) = 1 - \alpha$$

Passo 3 : trovare a e b tali che



la Chi quadro

sceita l'area α a metà' nelle due code :

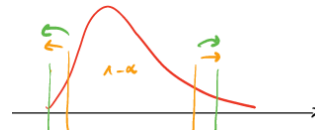
$$\begin{aligned} a &= \chi_{\alpha/2} = \chi_{0.025} = 31.5549 \\ b &= \chi_{1-\alpha/2} = \chi_{0.975} = 70.2224 \end{aligned} \quad \left| \begin{array}{l} (1-\alpha) = 0.95 \\ \alpha = 0.05 \\ \alpha/2 = 0.025 \\ n = 50 \\ \chi^2(49) \end{array} \right.$$

QUINDI:

$$P\left(\chi_{\alpha/2} < \frac{(n-1)s^2}{\text{Var}(x)} < \chi_{1-\alpha/2}\right) = 1 - \alpha$$

Passo 4 : Riscrivo →

$$\frac{EC}{PER} \rightarrow \left[\frac{(n-1)s^2}{\chi_{1-\alpha/2}}, \frac{(n-1)s^2}{\chi_{\alpha/2}} \right] \quad \text{CONFIDENZA } (1-\alpha)$$



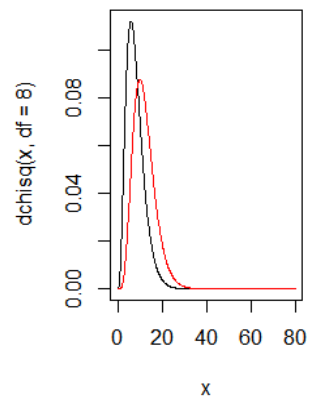
se aumento $(1 - \alpha)$: $\chi_{\alpha/2} \rightarrow$ diventa piu piccolo e $\chi_{1-\alpha/2} \rightarrow$ diventa piu grande IC si allarga

se n aumenta : $(n-1)$ aumenta ,ma IC diventano piu stretti
VarCI comando per calcolare la Varianza

Esercizio 12

Istogramma delle degenze delle neonate femmine non e' normale , X non può essere ipotizzata normale

TEST DI IPOTESI : secondo metodo di controllo della variabilita' oltre IC
Ipotesi : sono affermazioni che riguardano i parametri della popolazione



Abbiamo sempre 2 ipotesi : Ipotesi nulla $\rightarrow H_0$ & Ipotesi alternativa $\rightarrow H_1$

$$E(X) = \mu$$

$$Var(X) = \sigma^2$$

idk mo cambia gli indici ma vabb

Testiamo H_0 contro H_1 , ci chiediamo se e' opportuno abbandonare (rifiutare) H_0 a favore di H_1 ?

Credo in H_0 , raccolgo un campione casuale e mi domando se cio che e' contenuto nei numeri mi costringe ad abbandonare H_0 a favore di H_1

$H_0=5 - H_1 >5$ O $H_1 <5$ Test unilaterali a una coda // cerco valori solo da una parte

$H_0=5 H_1 \neq 5$ Test bilaterali a due code

Test di ipotesi sulle proporzioni:

passo 1 : $H_0 : p=0.5$ $H_1 : p>0.5$ scrivi direttamente H_0 e H_1

passo 2 : Fisso la significativa alfa piccolo : 0.05 - 0.01 - 0.1

Errore di 1 specie: rifiuti H_0 quando H_0 e' vera, alfa e' la soglia
Max per l'errore di prima specie

Errore di 2 specie: Non rifiuti H_0 quando e' falsa

Matrice confusione	H_0 vera	H_0 falsa
Rifiuto H_0	✗	✓
Non rifiuto H_0	✓	✗

Passo 3 : calcolo la statistica del test, sono quantita' aleatorie che dipendono dal campione e dal parametro p che non e' incognito per chi fa il test perche CREDO in H_0 (la mia ipotesi)

$$\frac{\hat{p} - p}{\sqrt{p(1-p)}} \approx N(0,1) \quad (n \geq 30)$$

$$\frac{\frac{83}{125} - 0.5}{\sqrt{\frac{0.5 \cdot 0.5}{125}}} = 3.6672$$

0.5 lo abbiamo deciso noi ,83 e' il numero di successi ed $n=125$

Passo 4 : e' un valore ragionevole con ciò che mi aspettavo ?

H1,sarebbe stato meglio sotto alternativa H1?

$$P \left(\begin{array}{c} \text{statistica} \\ \text{del} \\ \text{test} \end{array} > 3.6672 \right) = \underline{\text{PVALUE}}$$

$N(0,1)$

Calcolo la probabilita' aka statistica del test :

> e' il nostro H1

$1 - p(\text{stat test} < 3.6672) = \text{pnorm}(3.6672) = 0.0001$

PValue < α ? mi fissa la soglia dell'errore di prima specie

se e' piu' piccolo rifiuto H0 a favore di H1

LEZIONE 9 → inizio lezione con riassunto lezione 8

Non possono essere vere entrambe

binom.test.

$\alpha = 0.01$ succ=83 p.value= 0.0002 < $\alpha = 0.01$?

Si rifiuto H0 a favore di H1 ovvero $p > 0.5$ posso affermare con significativa $\alpha = 0.01$ che la popolazione e' d'accordo con me

OSS.supponiamo

$H_0 = p=0.5$ $H_1 = p > 0.5$ $n=125$ succ=20 $p \text{ cappello} = 2/125 = 0.16$

pvalue = 1 < $\alpha = 0.01$

No non posso rifiutare H0 in favore di H1 – Pvalue = 1

ESERCIZIO 9.10 : The variable sat.m in th data set stud.recs (Using R) contains math SAT scores for a group of students sampled from a larger population. Test the null hypothesis that the population mean score is 500 against a two-sided alternative.

Would you accept or reject at a 0.05 signification level? Pvalue

$H_0 = \mu=500$

$H_1 = \mu \neq 500$

$\alpha = 0.05$

`t.test(stud.recs$sat.m, $\mu = 500$, alternative = "two.sided")`

otteniamo che `t = -2.5731, df = 159 (n casi meno 1), p-value = 0.01099`
`IC = 475.1437 sx 496.7313 dx`

P Value = 0.011 < 0.05 ?

SI) rifiuto H0 a favore di H1 $\mu \neq 500$

SE $\alpha = 0.05$ - P Value = 0.011 < 0.01 ? NO) Non posso rifiutare H0 a favore di H1

ESERCIZIO 9.11: In the babies (Using R) data set, the variable dht records the father's height for the sampled cases. Do a significance test of the null hypothesis that the population mean height is 68 inches against an alternative that it is taller. Remove the values of 99 from the data, as these indicate missing data

$$H_0 = \mu = 68$$

$$H_1 = \mu > 68$$

$$\alpha = 0.01$$

t.test(altezze, $\mu =$ soglia, alternative = "greater")

otteniamo che $t = 20.796$, $df = 743$, $p\text{-value} < 2.2e-16$

IC = 177.8755 Inf

P value $< 2.2e-16 < 0.01$

SI) rifiuto H_0 a favore di H_1 , affermo che la popolazione e i padri ha altezza media maggiore di 68 inches

Test su differenza di medie: Proporzioni medie e varianze

2 casi

1) Campioni appaiati: Le diete, soggetti con pesi iniziali, e soggetti con pesi dopo la dieta

vorrei dire che la media dei pesi dopo, e' inferiore alla media dei pesi prima

Ipotesi $\mu_{prima} > \mu_{dopo}$

$$H_0: \mu_p = \mu_d$$

$$H_1: \mu_p > \mu_d$$

$$H_0: \mu_p - \mu_d = 0$$

$$H_1: \mu_p - \mu_d > 0$$

sulla stessa riga ho lo stesso soggetto

APPAIATI?

PRIMA	DOPO
.	x

2) Campioni indipendenti: sono interessato al peso e alla differenza di peso tra maschio e femmina, mi interessa sapere se il medio dei maschi e' maggiore del peso medio delle femmine

$$H_0: \mu_m - \mu_f = 0$$

$$H_1: \mu_m > \mu_f$$

$$H_0: \mu_m - \mu_f = 0$$

$$H_1: \mu_m - \mu_f > 0 \rightarrow \text{differenza di medie}$$

sulla stessa riga stesso soggetto, avro 2 var una quantitativa Sex

INDIPENDENTI?

PESSO	SEX

Passo 1 : il mio problema ha campioni appaiati o indipendenti?

LEZIONE 10 :

PESO PRIMA	PESO DOPO	DIFFER.
x	o	$x - o$

$$\begin{aligned} \text{DIFFERENZA} &= \\ &= \text{PESO PRIMA} - \text{PESO DOPO} \end{aligned}$$

1) **campioni appaiati:**

$$E(X - Y) = E(X) - E(Y)$$

media della differenza e' uguale alla differenza delle medie

$$H_0: \mu_p - \mu_d = 0$$

$$H_1: \mu_p - \mu_d > 0$$

$$H_0: \mu_{diff} = 0$$

$H_1: \mu_{diff} > 0 \rightarrow$ quando i campioni sono appaiati, ha senso costruire la variabile aleatoria e testare la differenza di medie equivale a testare la media della differenza
→ TEST SULLA MEDIA

2) Campioni indipendenti: invece non ha senso calcolare le differenze del peso di una donna e di un uomo, magari ho preso 2 campioni femmine,

Passo 1: Scrivo bene H0 e H1

H0: $\mu_M - \mu_F = 0$ H1: $\mu_M - \mu_F > 0$
($\mu_M > \mu_F$)

$$\frac{\bar{x}_M - \bar{x}_F - (\mu_M - \mu_F)}{\sqrt{\frac{\sigma_M^2}{n_M} + \frac{\sigma_F^2}{n_F}}}$$

Passo 2: fisso $\alpha = 0.01 / 0.05 / 0.1$

Passo 3: calcolo la statistica del test
se $n_M, n_F > 30$ oppure x_M & x_F Normali
Calcolo il valore campionario e sotto H0

Passo 4: probabilità che la statistica del test assuma valori pari al valore campionario o più estremi nel verso dell'alternativa

$P(\text{stat del test} > \text{valore campionario}) = P\text{-value}$

Passo 5: $P\text{-value} < \alpha$?

SE SI rifiuto H0 in favore di H1

SE NO non posso rifiutare H0

Esercizio 9.31: For the babies (UsingR) data set, the variable age contains the recorded mom's age and dage contains the dad's age for several different cases in the sample. Do a significance test of the null hypothesis of equal ages against a one-sided alternative that the dads are older in the sampled population.

H0 : $\mu_d - \mu_m = 0$ H1: $\mu_d - \mu_m > 0$ ($\mu_d - \mu_m$) $\alpha = 0.05$

`t = 28.092, df = 1226, p-value < 2.2e-16 IC = 2.943902 Inf`
`mean difference = 3.127139 // MEDIA CAMPIONARIA DELLE DIFFERENZE, i padri`
`sono più anziani delle madri di 3 anni`

$P\text{-value} = 2.2e-16 < \alpha = 0.01$ SI) posso rifiutare H0 a favore di H1 ovvero che i padri siano più vecchi delle madri

Controllo ipotesi

n>30 : $df = 1226$ $n = df + 1 = 1227$ 

oppure

`Xnormale = differenze <- babies$dage - babies$age`

Esercizio 9.32 : The data set normtemp(UsingR) contains body measurements for 130 healthy ,randomly selected individuals from some parent population. The variable temperature contains normal body temperature data and the variable gender contains gender information,with male coded as 1 and female as 2. Is the sample difference across the two groups statistically significant?

H0 : $\mu_f = \mu_m$

H1: $\mu_f \neq \mu_m$

$\mu_f - \mu_m = 0$

$\mu_f - \mu_m \neq 0 \rightarrow$ modo formale per scrivere il test

$\alpha = 0.05$

Campioni indipendenti

```

normtemp$temp <- (normtemp$temperature - 32)/1.8      # trasformazione gradi
tempf <- normtemp$temp[normtemp$gender == 2]          # temperature femmine
tempm <- normtemp$temp[normtemp$gender == 1]          # temperature maschi

hist(tempf)      # istogrammi rispettivi
hist(tempm)

boxplot(tempf, tempm)      # boxplot affiancati per vedere differenze

t.test(tempf, tempm, mu = 0, alternative = "two.sided", paired = FALSE) # t di student per
# femmine meno maschi , x - y
# mu = 0 ipotesi nulla
# alternativa con una two sided
# gli dobbiamo dire se sono appaiati
t = 2.2854, df = 127.51, p-value = 0.02394
95 percent confidence interval: 0.02156277 0.29980476
mean of x mean of y : 36.88547 36.72479

appaiati P value : 0.02394 < 0.05 ? SI) rifiuto H0 a favore di H1
OSS:  $\mu_f > 30$ ?  $\mu_m > 30$ ? SI oppure  $X_m$  normale ?  $X_f$  Normale ?

```

Esempio 9.36 : The Galton (HisData) data set contains data used by Francis Galton in 1885. Each data point contains a child's height and an average of his or her parents heights. Assuming the data is a random sample for a population of interest, perform a t-test to see if there is a difference in the population mean height. assume the paired t-test is appropriate

```

c= child          p=parent

H0 :  $\mu_c - \mu_p = 0$           H1:  $\mu_c - \mu_p > 0$            $\alpha = 0.01$ 

```

Campioni appaiati (testo)

```

boxplot(Galton$parent, Galton$child)

t.test(Galton$child, Galton$parent, mu = 0, paired = TRUE,
       alternative = "two.sided")
# child - parent ,x-y
# mu = 0
# paired= True (detto dal testo)
# alternative = 2 sided
t = -2.8789, df = 927, p-value = 0.004082
95 percent confidence interval: -0.36949983 -0.06993982
mean difference: -0.2197198

```

Pvalue= 0.004082 < $\alpha = 0.01$? SI) posso rifiutare H0 a favore di H1
 Controllo ipotesi : $n > 30$? SI) OK , ma così non fosse dovrei controllare la Normale guardando le altezze nell'istogramma, come qui sotto


```
difference <- Galton$child - Galton$parent  
hist(difference) # se a forma di campana allora e' corretto
```