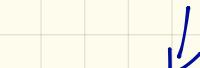


# Inferenza statistica

	Task Description	Due Date	Status
	Introduzione all'inferenza statistica: popolazione e campione, obiettivi, concetti e metodi di campionamento		
	Concetti sui principali schemi di campionamento		
	Il campionamento casuale semplice: schema con e senza ripetizione		
	Campione casuale e campione osservato. Le r.c. osservazioni campionarie		
	Statistica: definizione, principio di naturalità, sufficienza		
	Gli spazi dell'inferenza statistica: spazio parametrico, spazio campionario e spazio della statistica		
	Distribuzione campionaria di una statistica		
	Distribuzione campionaria della media campionaria		
	Distribuzione campionaria della proporzione campionaria		
	Distribuzione campionaria di una statistica: un semplice esempio		

Parte descrittiva → insieme olati POPOLAZIONE

insieme olati POPOLAZIONE



obiettivo sintesi



Parte interventuale → sottoinsieme delle POPOLAZIONE

se il campione è estratto casualmente gli strumenti della probabilità ci aiutano in questo passaggio inverso



**INFERENZA**

campione

sintesi del ]  
campiono

obiettivo  
riotto ad  
ottenere  
informazioni  
sulla popolazione

## CAMPIONAMENTO

non probabilistico

campionamento a  
scelta ragionata

campionamento oli  
comodo o oli convenienza

probabilistico

per ogni unità statistica  
è nota la probabilità di inclusione  
nel campione

campionamento casuale  
semplice

campionamento  
sistematico

campionamento  
stratificato

campionamento a  
gruppi

campionamento casuale semplice con ripetizione (più "naturale")

tutte le unità statistiche della popolazione hanno la stessa probabilità di inclusione nel campione

senza ripetizione (vantaggio: semplicità)

se  $N \rightarrow \infty$  i due schemi sono equivalenti

le unità della popolazione possono essere assimilate alle palline contenute in un'urna

da un'urna che contiene  $N$  palline  
ne estraggo casualmente  $n$  { con ripetizione  
senza ripetizione}

campionamento stratificato  $\leadsto$  la popolazione è caratterizzata da un certo numero di strati

es: genere

ricomponibili sulle differente modalità di una variabile qualitativa

es: diploma

I livelli e strati corrispondono a varie differenti

Altro il processo di campionamento cercherà di riprodurre nel campione

le proporz. degli strati presenti nella popolazione

# MEMO COMBINAZIONE LINEARE V.C.

Sia  $X_1, X_2, \dots, X_m$  una m-plo di r.c. tali che le  $X_i$  siano i.i.d. come una  $X \sim f(x_i; \theta)$ .

Sia inoltre  $\mu_x = E(X) = E(X_i)$  e  $\sigma^2 = \text{Var}(X) = \text{Var}(X_i)$ .

Sia  $W$  una combinazione lineare delle  $m$  r.c. con pesi  $a_i$ :

$$W = \sum_i a_i X_i$$

Si ha:

$$\mu_W = E(W) = \mu \sum_i a_i$$

$$\sigma_W^2 = \text{Var}(W) = \sigma^2 \sum_i a_i^2$$

Due particolari combinazioni lineari:

- somma ( $a_i = 1 \forall i$ )  $\rightarrow W = \sum_i X_i$   $\underbrace{\mu_W = m\mu}_{\sigma_W^2 = m\sigma^2}$

- media ( $a_i = \frac{1}{n} \forall i$ )  $\rightarrow W = \frac{\sum_i X_i}{m} = \bar{X}_m$   $\underbrace{\mu_W = \mu}_{\sigma_W^2 = \frac{\sigma^2}{m}}$

Cosa è possibile dire circa la distribuzione di  $W$ ?

Se  $X \sim \text{Bec}(\pi)$   $\rightarrow W = \sum X_i \sim \text{bin}(n, \pi)$

$E(W) = n\pi$   
 $\text{Var}(W) = n\pi(1-\pi)$

$W = \frac{\sum X_i}{n}$  trasformazione di una v.c. binomiale

$E(W) = \pi$   
 $\text{Var}(W) = \frac{\pi(1-\pi)}{n}$

Se  $X \sim N(\mu, \sigma^2)$   $\rightarrow W = \sum a_i X_i \sim N\left(\mu \sum a_i, \sigma_x^2 \sum a_i^2\right)$

proprietà riproduttiva v.c. normale  
(distribuzione esatta)

Se  $X \sim f(x; \theta)$   $\rightarrow W = \sum a_i X_i \xrightarrow[n \rightarrow \infty]{} N\left(\mu \sum a_i, \sigma_x^2 \sum a_i^2\right)$

$f$  potrebbe non essere normale  
oppure non si hanno le  
informazioni necessarie per  
assumere la normalità

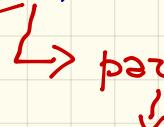
teorema del limite centrale  
(distribuzione approssimata)

POPOLAZIONE  $\rightarrow$  si è interessati a studiare un fenomeno che si assume possa essere descritto da un modello probabilistico

ATTENZIONE: 

tutta l'inferenza si basa su questa ipotesi

$$X \sim f(x; \theta)$$

 parametro di interesse

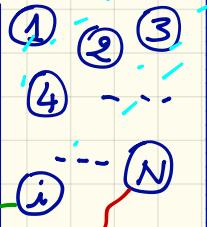
 l'inferenza parametrica mira ad ottenere informazioni su  $\theta$  a partire dai dati campionari

(H)

spazio parametrico



insieme dei possibili valori che può assumere il parametro  $\theta$



$N$ : numerosità della popolazione  
( $N \rightarrow \infty$ : popolazioni infinite)

→ [l'<sup>i</sup>-esima unità  
della popolazione]

all'<sup>i</sup>-esima posizione  
del campione potrai osservare una  
qualsiasi delle unità della popolazione

## CAMPIONE OSSERVATO

$x_1, x_2, \dots, x_i, \dots, x_m$

numerosità  
del campione

valore che la variabile  $X$   
assume sull'<sup>i</sup>-esima unità  
del campione

il valore  $x_i$  osservato è uno dei possibili  
valori che possono essere osservati all'<sup>i</sup>-esima  
posizione del campione

$x_i$  può essere considerato uno  
dei valori che una r.v.  $X$  può  
assumere

$(x_1, x_2, \dots, x_i, \dots, x_m)$ 

campione osservato (n valori)

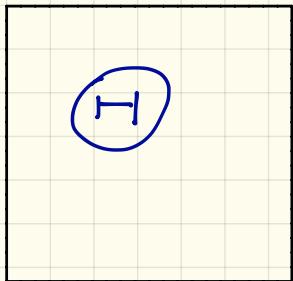
 $(X_1, X_2, \dots, X_i, \dots, X_m)$ 

campione casuale (n r.c.)

ciascuna  $X_i \sim X \sim f(x; \theta)$ inoltre se il campionamento è con ripetizione  
le  $X_i$  sono indipendenti→ le  $X_i$  sono i.i.d.  $\sim X$ parallelo con le combinazioni lineari  
di r.c. richiamato ad inizio lezionel'insieme di tutti i possibili  
campioni di ampiezza n(ovvero l'insieme di tutti i possibili  
valori che possono comporre il campione  
osservato) si chiama spazio campionario → NOTA: non confondere  
e si indica con  $\mathcal{X}_m^n$ NOTA: non confondere  
con  $\mathcal{X}$  (spazio campione)

# GLI "SPAZI" DELL'INFERENZA

SPAZIO  
PARAMETRICO

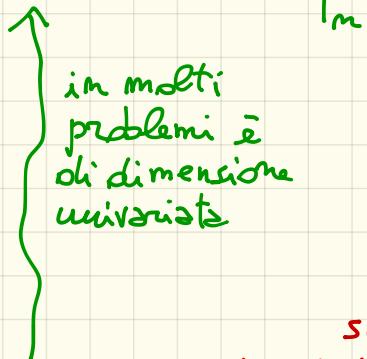


lo spazio di interesse  
rispetto cui devo "decidere"

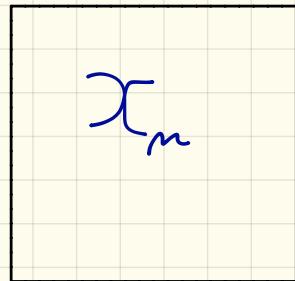
rispetto cui devo "decidere"

SPAZIO DELLA  
STATISTICA

$T_m$



SPAZIO  
CAMPIONARIO



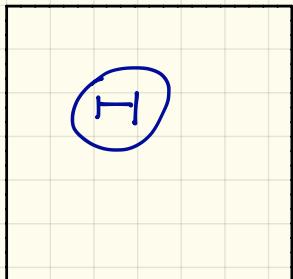
lo spazio di lavoro  
sul quale effettua la  
"decisione"

di solito non si lavora direttamente su  $X_m$  ma su una sua opportuna sintesi

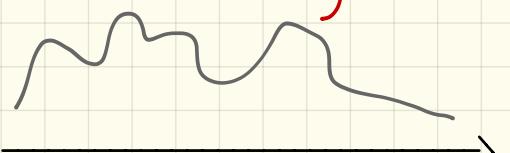
STATISTICA

# GLI "SPAZI" DELL'INFERENZA

SPAZIO  
PARAMETRICO

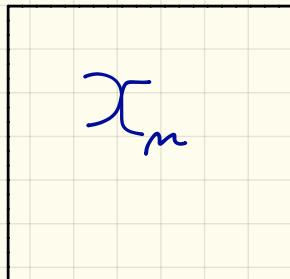


SPAZIO DELLA  
STATISTICA



la distribuzione di probabilità di  $T_m$  si chiama  
distribuzione campionaria

SPAZIO  
CAMPIONARIO



e che quindi può essere calcolata

qualunque funzione del campione reale che non dipende da parametri incogniti

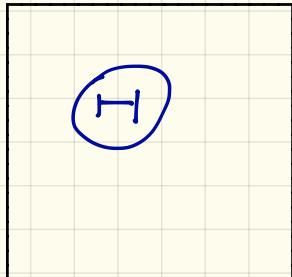
$$T_m = T(X_1, X_2, \dots, X_i, \dots, X_n)$$

è una r.c. (funzione di m r.c.)

$$t_m = T(y_1, y_2, \dots, y_i, \dots, y_n)$$

valore assunto dalla statistica in corrispondenza del campione osservato

## SPAZIO PARAMETRICO



## SPAZIO DELLA STATISTICA



## SPAZIO CAMPIONARIO



in molti problemi di interesse è possibile fare riferimento a statistiche caratterizzate da distribuzioni campionarie normali

$$X \sim \text{Bin}(\pi)$$

$$\bar{x}_m = \hat{\pi} = \frac{\sum x_i}{m} \rightarrow N\left(\pi, \frac{\pi(1-\pi)}{m}\right)$$

distribuzione approssimata  
(TLC)

$$X \sim N(\mu, \sigma^2)$$

$$\bar{X}_m = \frac{\sum x_i}{m} \sim N\left(\mu, \frac{\sigma^2}{m}\right)$$

distrib. eratta (proprietà riproduttiva normale)

$$X \sim f(x; \theta)$$

$$\bar{X}_m \rightarrow N\left(\mu, \frac{\sigma^2}{m}\right)$$

distribuzione approssimata (TLC)

# SULLA SCELTA DELLA "MIGLIORE" STATISTICA $T_m$

↳

esistono una serie di proprietà formali che consentono di scegliere fra statistiche alternative

( tali proprietà dipendono anche dal particolare utilizzo che si fa di  $T_m$  )

- ↳ problemi di stima  $\rightarrow T_m$  si chiama stimatore
- ↳ verifica di ipotesi  $\rightarrow T_m$  si chiama statistica test

due requisiti essenziali:

- NATURALITÀ applico la stessa statistica che aspetterei sulla popolazione all'unico campione disponibile
- SUFFICIENZA  $\rightarrow$  il requisito essenziale è che  $T_m$  sia in grado di catturare tutte le informazioni utili su  $\theta$  contenute nel campione

una volta calcolata la statistica  $T_m$  non è necessario conservare il campione

## Inferenza su $\mu$ [ DISTRIBUZIONE MÉDIA CAMPIONARIA ]



$$T_m = \bar{X}_m = \frac{\sum X_i}{n} \quad \text{si usa la statistica media campionaria}$$

- se  $X \sim N(\mu, \sigma^2)$   $\Rightarrow \bar{X}_m \sim N\left(\mu, \frac{\sigma^2}{m}\right)$  proprietà riproducibile r.c. normale

- se  $X \sim f(x; \theta)$   $\Rightarrow \bar{X}_m \rightarrow N\left(\mu, \frac{\sigma^2}{m}\right)$  teorema del limite centrale

$\downarrow$

$X$  non è normale  
oppure non è possibile  
assumere la normalità

$m \rightarrow \infty$

$\downarrow$

REGOLA EMPIRICA :

se la dimensione del campione ( $n$ )  
 $n \geq 30$  posso sfruttare il TLC

## Inferenza su $\pi$ [ DISTRIBUZIONE PROPORTIONE CAMPIONARIA ]



$$T_m = \bar{J}_m = \frac{\sum X_i}{m} \quad \text{si usa la statistica proporzione campionaria}$$

$\left( \begin{array}{l} \frac{1}{m} \bar{X}_m \quad \text{si tratta di una media di osservazioni} \\ \{0, 1\} \end{array} \right)$

$\xrightarrow{\quad}$  possesso delle caratteristiche  
di interesse (successo)

In questo caso:  $X \sim \text{Ber}(\pi)$



$\xrightarrow{\quad}$  proporzione nella popolazione

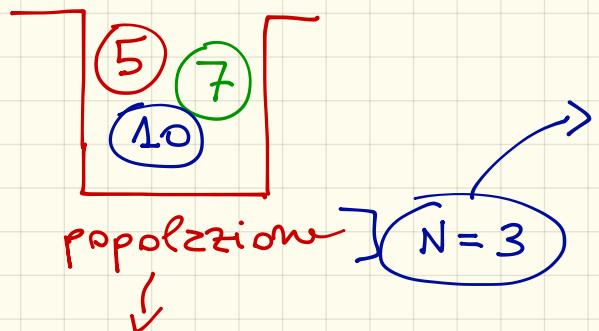
$\bar{J}_m = \hat{\pi}$  è una trasformazione di una r.c. binomiale

$\xrightarrow{\quad} \bar{J}_m = \frac{\sum X_i}{m} \xrightarrow{\quad} \text{bin}(m, \pi) \quad ]$  è possibile fare riferimento al  
modello binomiale per ottenere la  
distribuzione esatta di  $J_m$

Essendo  $\bar{J}_m$  una combinazione lineare di v.c. i.i.d.  $\sim \text{Ber}(\pi)$  è  
possibile sfruttare il teorema del limite centrale  $\hat{\pi} \xrightarrow{m \rightarrow \infty} N(\pi, \frac{\pi(1-\pi)}{m})$

DISTRIBUZIONE  
ASINTOTICA DI  $\bar{J}_m$

# CASO STUDIO (ci permette di introdurre i concetti)



$X$  m. errori per  
pagina siti tra  
olitografi

$X$	$p(x)$
5	$1/3$
7	$1/3$
10	$1/3$
	1

modelli che  
descrivono la  
popolazione

$$E(X) = 7,33$$

$$\sqrt{v}(X) = 4,22$$

due "parametri" caratteristici  
delle popolazioni

Ipotizziamo di estrarre un campione di numerosità

$$n = 2$$

?) quanti sono i possibili campioni che ?  
potrai osservare

$$3 \times 2 = 6$$

estrazione senza  
rimessa

$$3 \times 3 = 9$$

estraz. con  
rimessa

$X_1 \rightsquigarrow$  risultati che potrai osservare alle prime  
estrazioni

$X_2 \rightsquigarrow$  risultati che potrai osservare alle seconde  
estrazioni

$X_1$	$X_2$	[questi tre campioni si hanno nel caso di estrazione con rimessa]	
		$X_1 \neq X_2$	$X_1 = X_2$
5	7		
5	10		
7	5		
7	10		
10	5		
10	7		
5	5		
7	7		
10	10		

$X_1 \neq X_2$	senza rimessa	con rimessa
5	$\frac{2}{6} = \frac{1}{3}$	$\frac{3}{9} = \frac{1}{3}$
7	$\frac{2}{6} = \frac{1}{3}$	$\frac{3}{9} = \frac{1}{3}$
10	$\frac{2}{6} = \frac{1}{3}$	$\frac{3}{9} = \frac{1}{3}$

le singole  $X_i$  (osservazione campionariose) si distribuisce allo stesso modo della  $X$  che definisce la popolazione

$X_i$  è i.o. come la  $X$

nel caso di estrazioni con rimessa le  $X_i$  saranno anche indipendenti:  $X_i$  è i.i.o. come la  $X$

$X_1$	$X_2$
5	7
5	10
7	5
7	10
10	5
10	7
5	5
7	7
10	10

parametro o rettore  
di parametri che caratterizzano  $f$

$$X \sim f(x, \theta)$$

fenomeno di interesse  
popolazione

nell'es. è il m. di errori  
commesso dalle plattigrafe

modello di probabilità che si  
pensa possa descrivere adeguatamente  
il fenomeno

tutto quello che vedremo & seguirà  
si basa sull'ipotesi che la  $f$  sia  
adeguata

$X_1$	$X_2$	$\bar{X}_m$
5	7	6
5	10	7,5
7	5	6
7	10	8,5
10	5	7,5
10	7	8,5
5	5	5
7	7	7
10	10	10

principio di  
naturalità

↓  
applico la stessa  
statistica che aspetterei  
sulla popolazione all'unico  
campione disponibile

→ distribuzione campionaria  
delle statistiche: comportamento di  $\bar{X}_m$   
su tutti i possibili campioni di  
numerosità  $n = 2$

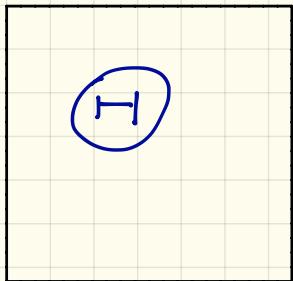
Immaginiamo di essere interessati a studiare le  
medie della popolazione

	Task Description	Due Date	Status
	I problemi dell'inferenza statistica: stima e verifica di ipotesi		
	Proprietà degli stimatori: la non distorsione.		
	Lo stimatore s'è per le varianze		
	Efficienza relativa ed efficienza assoluta		
	Criterio dell'errore quadratico medio come misura congiunta della posizione e delle variabilità della distribuzione campionaria		
	Proprietà finite ed asintotiche di uno stimatore		
	Consistenza di uno stimatore		
	<b>RIEPILOGO DISTRIBUZIONI CAMPIONARIE CASO DI INFERENZA SU <u>UNA SOLA</u> POPOLAZIONE</b>		
	Inferenza sulla <u>media</u> : caso di varianza nota, di varianza incognita e rimozione dell'ipotesi di normalità		
	Inferenza sulla <u>varianza</u> di una popolazione normale		
	Inferenza sulla <u>proporzione</u> di una popolazione bernoulliana		

# GLI "SPAZI" DELL'INFERENZA



SPAZIO  
PARAMETRICO



lo spazio di interesse

rispetto cui devo "decidere"

di solito non si lavora direttamente su  $X_m$  ma su una sua opportuna sintesi

STATISTICA

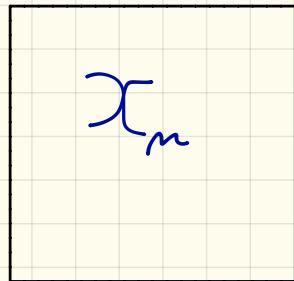
SPAZIO DELLA  
STATISTICA



$T_m$

in molti problemi è di dimensione univariata

SPAZIO  
CAMPIONARIO



lo spazio di lavoro

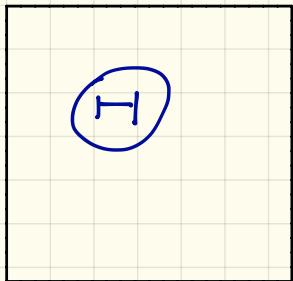
sul quale effettua la "decisione"

# GLI "SPAZI" DELL'INFERENZA

MEMO

la distribuzione di probabilità di  $T_m$  si chiama distribuzione campionaria

SPAZIO PARAMETRICO



SPAZIO DELLA STATISTICA



SPAZIO CAMPIONARIO



e che quindi può essere calcolata

qualunque funzione del campione reale che non dipende da parametri incogniti

$$T_m = T(X_1, X_2, \dots, X_i, \dots, X_n)$$

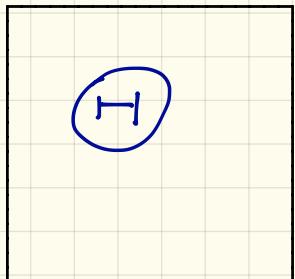
è una r.c. (funzione di m r.c.)

$$t_m = T(y_1, y_2, \dots, y_i, \dots, y_n)$$

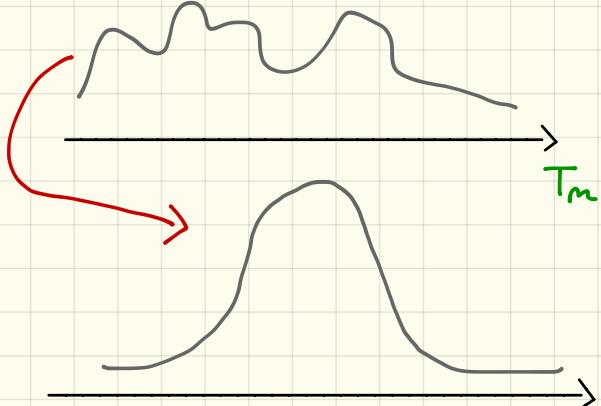
valore assunto dalla statistica in corrispondenza del campione osservato

# MEMO

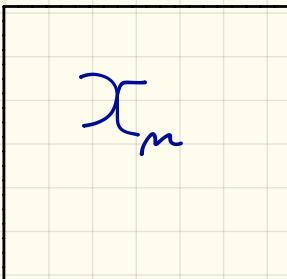
SPAZIO  
PARAMETRICO



SPAZIO DELLA  
STATISTICA



SPAZIO  
CAMPIONARIO



in molti problemi di interesse è possibile fare riferimento a statistiche caratterizzate da distribuzioni campionarie normali

$$X \sim \text{Bin}(\pi)$$

$$\bar{x}_m = \hat{\pi} = \frac{\sum x_i}{m} \rightarrow N\left(\pi, \frac{\pi(1-\pi)}{m}\right)$$

distribuzione approssimata  
(TLC)

$$X \sim N(\mu, \sigma^2)$$

$$\bar{X}_m = \frac{\sum x_i}{m} \sim N\left(\mu, \frac{\sigma^2}{m}\right)$$

distrib. eratta (proprietà riproduttiva normale)

$$X \sim f(x; \theta)$$

$$\bar{X}_m \xrightarrow{D} N\left(\mu, \frac{\sigma^2}{m}\right)$$

distribuzione approssimata (TLC)

# SULLA SCELTA DELLA "MIGLIORE" STATISTICA $T_m$

MEMO

↳

esistono una serie di proprietà formali che consentono di scegliere fra statistiche alternative

( tali proprietà dipendono anche dal particolare utilizzo che si fa di  $T_m$  )

↳ problemi di stima  $\rightarrow T_m$  si chiama stimatore

↳ verifica di ipotesi  $\rightarrow T_m$  si chiama statistica test

due requisiti essenziali:

- NATURALITÀ applico la stessa statistica che aspetterei sulla popolazione all'unico campione disponibile
- SUFFICIENZA  $\rightarrow$  il requisito essenziale è che  $T_m$  sia in grado di catturare tutte le informazioni utili su  $\theta$  contenute nel campione

una volta calcolata la statistica  $T_m$  non è necessario conservare il campione

## INFERENZA (parametrica)

problemi di stima



si vuole stimare il parametro  $\theta$   
a partire da osservazioni campionarie



la statistica  $T_m$ , sintesi delle  
osservazioni campionarie, si chiama  
STIMATORE

STIMATORE



verifica di ipotesi



si vuole verificare le  
valigilità di un'ipotesi  
su  $\theta$  a partire dall'esperienza  
campionaria



la statistica  $T_m$  si chiama  
in questo caso  
STATISTICA TEST

in entrambi i casi le scelte della migliore statistica  $T_m$  da usare  
si basa su una serie di proprietà riconducibili al comportamento  
della distribuzione campionaria di  $T_m$  (ovvero del modello di proba-  
bilità che può essere usato per  $T_m$ )

# PROBLEMI DI STIMA



proprietà degli stimatori  $\rightarrow$  riconducibili a proprietà della distribuzione campionaria di  $T_m$

POSIZIONE



Non distorsione



riassume il comportamento  
medio delle distribuzioni  
campionarie di  $T_m$

VARIABILITÀ



efficienza



riassume la distribuzione  
campionaria di  $T_m$  in  
termini di dispersione rispetto  
alla media

FORMA



"NORMALITÀ"



molti stimatori possono essere  
descritti o approssimati da una  
legge normale

## PROBLEMI DI STIMA

Stima di  $\mu$



POPOLAZIONE

$$X \sim N(\mu, \sigma^2)$$

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

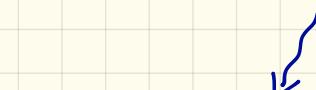
$$X \sim f$$

$$\text{stimate} \bar{X}_m = \frac{\sum X_i}{m}$$

$$\Rightarrow \bar{X}_m \xrightarrow{m \geq 30} N\left(\mu, \frac{\sigma^2}{m}\right)$$

Stima di  $\pi$  ~~~  $X \sim \text{Bee}(\pi)$  ~~~ la distribuzione esatta di  $\hat{\pi}$  può essere ottenuta facendo riferimento ad una binomiale

$$\text{stimate} \hat{\pi}_m = \frac{1}{m} = \bar{X}_m = \frac{\sum X_i}{m} \sim \text{bin}(m, \pi)$$



sfruttando il TLC, se  $m \geq 30$  si ha:

$$\hat{\pi} \rightarrow N\left(\pi, \frac{\pi(1-\pi)}{m}\right)$$

Stima di  $\sigma^2$  → lo stimatore "maturo" sarebbe:

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^m (x_i - \mu)^2}{m}$$

questo richiederebbe la conoscenza di  $\mu$  (poco plausibile)

Si può stimare  $\mu$  usandolo  $\bar{x}_m$

STIMATORE  
VARIANZA  
CAMPIONARIA

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^m (X_i - \bar{X}_m)^2}{m}$$

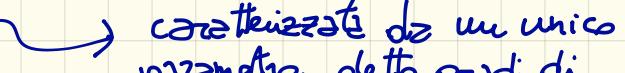


per correggere un (piccolo) problema di distorsione di  $\hat{\sigma}^2$  si divide per  $m-1$  (vedi pagine più avanti per i dettagli tecnici)

$$s^2 = \frac{\sum (X_i - \bar{X}_m)^2}{m-1}$$

STIMATORE  
VARIANZA  
CAMPIONARIA  
CORRETTA

La distribuzione campionaria di  $s^2$  richiede l'introduzione di un nuovo

modello di probabilità: La v.c.  $X^2$  

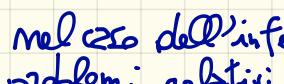


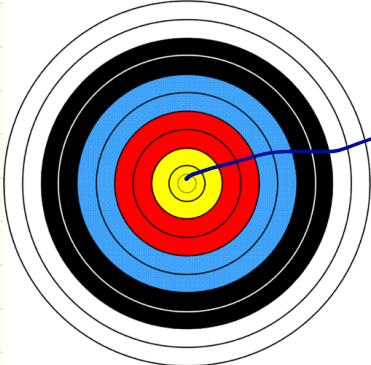
caratterizzata da un unico  
parametro, detto gradi di  
libertà (g.d.l.)

è definita dalla

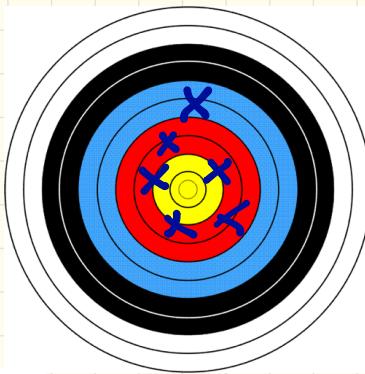
somma di  $m$  v.c. normali  
standardizzate al quadrato



 nel caso dell'inferenza su varianze è necessaria l'ipotesi di normalità (rispetto ai problemi relativi alle medie o alle proporzioni, l'inferenza sulle varianze è più sensibile a violazioni dell'ipotesi di normalità)



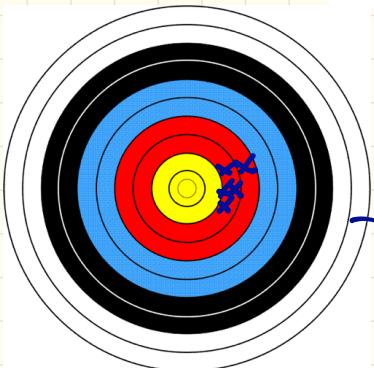
il centro del bersaglio è  $\theta$   
( l'obiettivo finale )



→ comportamento  
medi

$$E(T_m) = \theta$$

$T_m$  è non distorto per  $\theta$



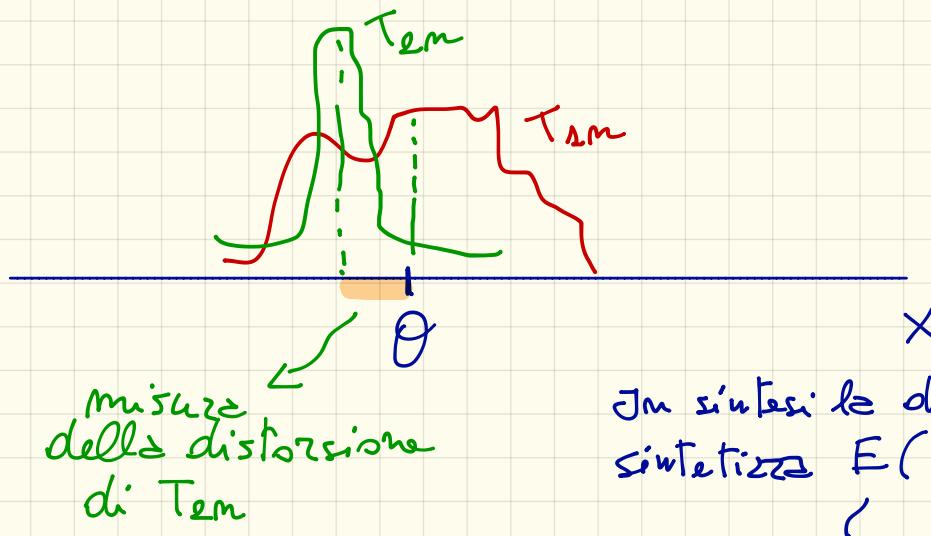
$$E(T_m) \neq \theta$$

qui  $T_m$  è distorto:  $E(T_m) - \theta =$   
 $= d(T_m) = b(T_m)$

**distorsione (bias)**

Se  $b(T_m) < 0 \rightarrow$  sottostimato  $\theta$

$b(T_m) > 0 \rightarrow$  sovrestimato  $\theta$



Im sintesi la distorsione  
sintetizza  $E(T_m)$

posizione delle v.c. stimatore

$$\bar{X}_m = \frac{\sum X_i}{n}$$



MEMO: veoli lezione sulle combinazioni lineari di v.c.

$$E(\bar{X}_m) = \text{combinat. lineare di} \underset{n \text{ r.c. i.i.d.}}{X_i} = \mu$$



la media campionaria  
è non distorta per  $\theta = \mu$

$$\hat{\pi}_m = \bar{X}_m \quad \text{particolare media di} \underset{\text{osservazioni 0,1}}{\text{estrazioni}}$$



$$E(\hat{\pi}_m) = \mu = \pi$$

estrazioni da una  
popolazione  $X \sim \text{Ber}(\pi)$

$\hat{\pi}_m$  è uno stimatore non distorto per  $\pi$  ]  $\hat{\pi}_m = \hat{\pi}$  (altro possibile nome)

Immaginiamo di essere interessati a stimare  $\sigma^2$



il nostro  $\theta$

Sulla popolazione

$$\sigma^2 = \frac{\sum_{i=1}^N (\bar{x}_i - \mu)^2}{N}$$

| utilizzando il principio di naturalità

$$\frac{\sum_{i=1}^m (\bar{x}_i - \mu)^2}{m}$$

$\mu$  è un parametro di  
disturbo (non è il reale  
obiettivo dell'inferenza)

è poco plausibile che  
io abbia informazioni  
sue  $\mu$

Varianza  
campionaria

$$\hat{\sigma}^2 = \frac{\sum (\bar{x}_i - \bar{x}_m)^2}{m}$$

utilizzo  
la stima  
di  $\mu$

$$E(\hat{\sigma}^2) \neq \sigma^2$$

esempio di stimatore obiettivo

$$\begin{aligned} E[\hat{\sigma}^2] &= E\left[\frac{\sum (X_i - \bar{X}_m)^2}{m}\right] = E\left[\frac{\sum X_i^2}{m} - \bar{X}_m^2\right] = \\ &= E\left[\frac{\sum X_i^2}{m}\right] - E(\bar{X}_m^2) = \end{aligned}$$

$$\frac{\sum E(X_i^2)}{m}$$

$$\hookrightarrow \text{NOTA: } \text{Var}(\bar{X}_m) = E[(\bar{X}_m - E(\bar{X}_m))^2]$$

$$E(\bar{X}_m^2) - [E(\bar{X}_m)]^2$$



$$E(\bar{X}_m^2) = \text{Var}(\bar{X}_m) + [E(\bar{X}_m)]^2 =$$

$$= \frac{\sigma^2}{m} + \mu^2$$

$$E(\hat{\sigma}^2) = \frac{\sum_i E(X_i^2)}{m} - \left( \frac{\sigma^2}{m} + \mu^2 \right) =$$

$$= \frac{\sum_i E(X_i^2)}{m} - \mu^2 - \frac{\sigma^2}{m} = \sigma^2 - \frac{\sigma^2}{m} =$$

$$= \frac{m\sigma^2 - \sigma^2}{m} =$$

$$= \frac{m-1}{m} \sigma^2$$

$$\text{Var}(\bar{x}_i) = \text{Var}(x)$$

$\hat{\sigma}^2$  è distorto per campioni finiti ma asintoticamente non distorto

← [la distorsione diminuisce al crescere di  $m$ ]

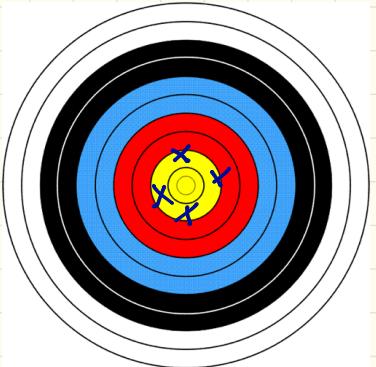
$$m \rightarrow \infty \text{ allora } \frac{m-1}{m} \rightarrow 0$$

In questo caso è possibile correggere la distorsione:

$$T_m = \frac{m}{m-1} \hat{\sigma}^2 = \frac{\cancel{m}}{m-1} \frac{\sum (x_i - \bar{x}_m)^2}{\cancel{m}} = S^2$$

Varianza  
campionaria  
corretta

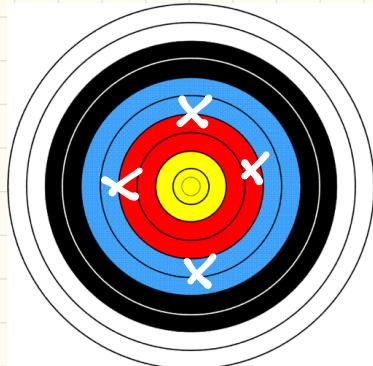
$$\begin{aligned} E(S^2) &= E\left(\frac{m}{m-1} \hat{\sigma}^2\right) = \frac{m}{m-1} E(\hat{\sigma}^2) = \\ &= \frac{m}{m-1} \frac{m-1}{m} \sigma^2 = \sigma^2 \end{aligned}$$



$T_{sm}$  e  $T_{2m}$  sono  
due stimatori non  
distorti per  $\theta$



$$E(T_{sm}) = E(T_{2m}) = \theta$$



$$\text{Var}(T_{sm}) = E[(T_{sm} - E(T_{sm}))^2] = \\ = E[T_{sm} - \theta]^2$$

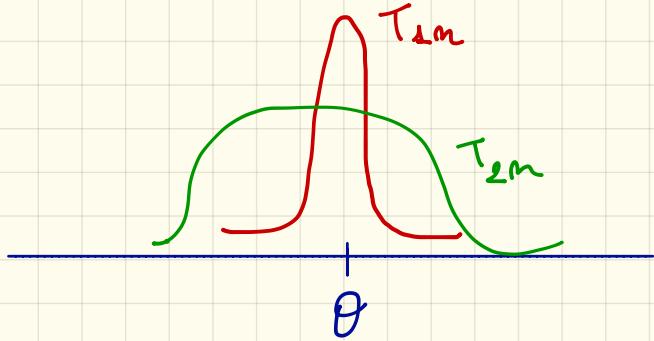
$$\text{Var}(T_{2m}) = E[T_{2m} - \theta]^2$$

$\text{Var}(T_{sm}) < \text{Var}(T_{2m})$

$T_{sm}$  è preferibile

$T_{sm}$  è più efficiente  
di  $T_{2m}$  (effic. in  
senso relativo rispetto  
a  $T_{2m}$ )

sono equivalenti



Efficienza in senso relativo consiste in un confronto tra le misure di variabilità delle distribuzioni campionarie dei due stimatori

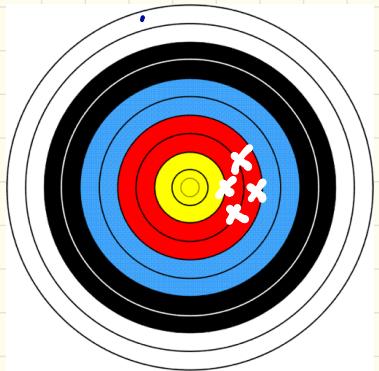
Efficienza in senso assoluto  $\rightarrow$  permette di stabilire un valore limite per la varianza di uno stimatore per un dato parametro al di sotto del quale non si può scendere

limite di Cramér - Rao

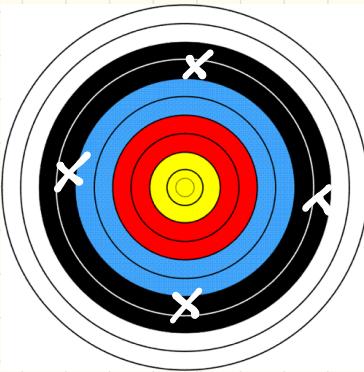
$\downarrow$

←

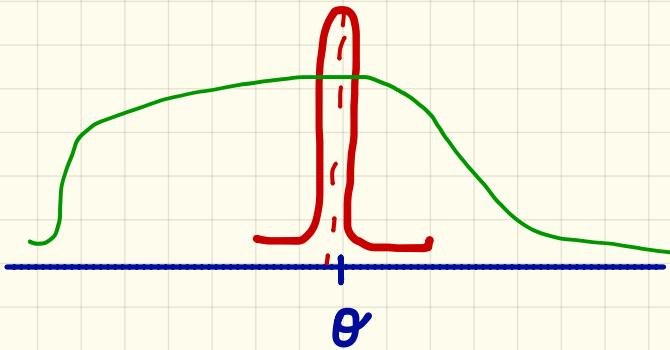
Se  $\text{Var}(T_m) = \text{valore limite}$  ]  $T_m$  è efficiente in senso assoluto



$T_m$



$T_m$



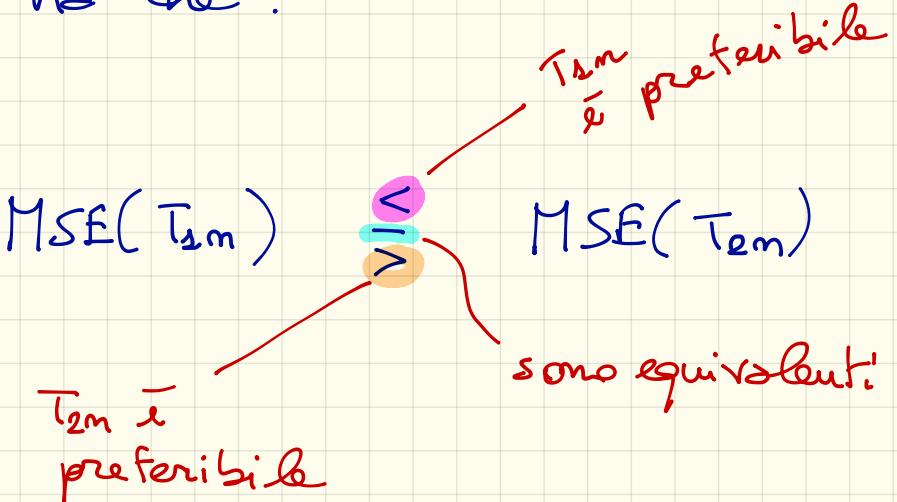
Si introduce un indice : errore quadratico medio EQM  
( mean squared error MSE )

che tiene conto contemporaneamente delle due proprietà

$$MSE(T_m) = E[(T_m - \theta)^2] = \text{Var}(T_m) + b^2(T_m)$$

$$\begin{aligned}
 \text{MSE}(\bar{T}_m) &= E[(\bar{T}_m - \theta)^2] = \\
 &= E[(\bar{T}_m - E(\bar{T}_m) + E(\bar{T}_m) - \theta)^2] = \\
 &= E\left\{ [(\bar{T}_m - E(\bar{T}_m)]^2 + [E(\bar{T}_m) - \theta]^2 \right\} = \\
 &= E[\underbrace{\bar{T}_m - E(\bar{T}_m)}_{\text{Varia}}]^2 + E[\underbrace{E(\bar{T}_m) - \theta}_{\text{costante}}]^2 + 2E[(\bar{T}_m - E(\bar{T}_m)) \underbrace{[E(\bar{T}_m) - \theta]}_{\text{costante}}] \\
 &\quad \downarrow \\
 &\quad E(\text{costante}) = \text{costante} \\
 &\quad \downarrow \\
 &\quad [\underbrace{E(\bar{T}_m) - \theta}_{\text{Varia}}]^2 \\
 &\quad \downarrow \\
 &\quad b^2(\bar{T}_m)
 \end{aligned}$$

Dati due stimatori  $\bar{T}_{1m}$  e  $\bar{T}_{2m}$  (non necessariamente non distorti) si ha che :



Esempio di scelta tra due stimatori alternativi

$$X \sim N(\mu, \sigma^2)$$

popolazione

obiettivo stima

simmetrica

$$\mu = \bar{\mu}_e$$

$\bar{X}_m$  media  
campionaria

$\hat{\mu}_e$  mediana  
campionaria

$$\text{non distorti} \rightarrow E(\bar{X}_m) = E(\hat{\mu}_e) = \mu$$

$$\text{Var}(\bar{X}_m) = \frac{\sigma^2}{n} < \frac{\pi \sigma^2 / \varrho_m}{\varrho^2 n} = \text{Var}(\hat{\mu}_e)$$

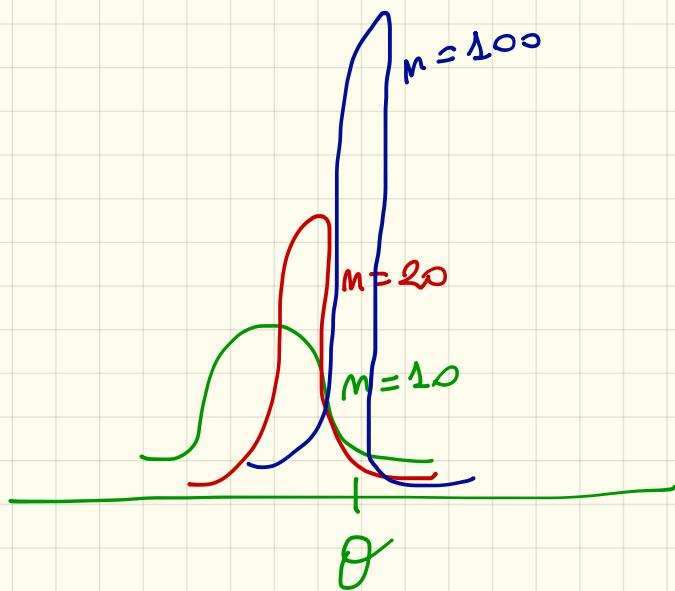
$$\frac{\text{Var}(\hat{\mu}_e)}{\text{Var}(\bar{X}_m)} = \frac{\pi \sigma^2 / \varrho_m}{\sigma^2 / n} = \frac{\pi}{\varrho} \approx 1.57$$

Le proprietà degli stimatori possono essere studiate sia su campioni finiti ( PROPRIETA' FINITE) che al crescere delle numerosità campionarie ( PROPRIETA' ASINTOTICHE)

$E(\hat{\sigma}^2) \neq \sigma^2$  distorto  $\rightarrow$  diverso non distorto asintoticamente

Tra le differenti proprietà asintotiche (che studiano il comportamento della distribuzione campionaria quando  $n \rightarrow +\infty$ ) vale la pena richiamare le proprietà della consistenza in media quadratica

$$\hookrightarrow \lim [E(T_n - \theta)^2] = 0$$



La consistenza in  
media quadratica vale  
se e solo se :

$$\textcircled{1} \quad \lim_{n \rightarrow \infty} E(T_m) = \theta$$

$T_m$  è assintoticamente  
non distorto

$$\textcircled{2} \quad \lim_{n \rightarrow \infty} \text{Var}(T_m) = 0$$

Una proprietà utile relative alle forme delle distribuzioni campionarie di uno stimatore  $T_m$  è la normalità  
assintotica :

$$\lim_{n \rightarrow \infty} \frac{T_m - E(T_m)}{\sqrt{\text{Var}(T_m)}} \rightarrow N(0,1)$$

RIEPILOGO DISTRIBUZIONI CAMPIONARIE  
CASO DI INFERNENZA SU UNA SOLA POPOLAZIONE

# Schemi di sintesi (lessico chi base)

$$X \sim f(x; \theta)$$

modello  
popolazione

parametro, o vettore di parametri

ci limitiamo allo studio  
di un parametro alla volta

$$X_1, X_2, \dots, X_m$$

campione casuale  
n-pz di r.c.

$$x_1, x_2, \dots, x_m$$

le  $x_i$  sono i.i.d.  $\sim X$

ipotesi di  
campionamento con  
ripetizione

campione  
osservato



(unico)

realizzazione disponibile  
del campione casuale

l'inferenza lavora su questo unico  
campione

Le scelte dei metodi da usare si basa sul

## PRINCIPIO DEL CAMPIONAMENTO RIPETUTO

ovvero lo studio delle statistiche di sintesi campionarie su tutti i "possibili" campioni  $\hookrightarrow$  le scelte della statistica (stimatore) da usare si basa sulle caratteristiche della corrispondente distribuzione campionaria

$$T(X_1, X_2, \dots, X_m) = T_m \quad \text{stimatore (r.c.)}$$
$$T(x_1, x_2, \dots, x_m) = t_m \quad \text{stima}$$

The diagram illustrates the relationship between the population parameters and their sample counterparts. It shows four variables: \$X\_1, X\_2, \dots, X\_m\$ above the equation \$T(X\_1, X\_2, \dots, X\_m) = T\_m\$, and \$x\_1, x\_2, \dots, x\_m\$ below the equation \$T(x\_1, x\_2, \dots, x\_m) = t\_m\$. Four arrows point downwards from each \$X\_i\$ to its corresponding \$x\_i\$. Another arrow points downwards from \$T\_m\$ to \$t\_m\$. To the right of the equations, the word "stima" is written inside a green oval, indicating the final estimated value.

NOTA:  $\sqrt{\text{var}(T_m)}$  (lo scarto quadratico medio) viene talvolta chiamato STANDARD ERROR, ERRORE STANDARD o DEVIAZIONE STANDARD

# INFERENZA SULLA MEDIA $\mu$ DI UNA POPOLAZIONE

①  $X \sim N(\mu, \sigma^2)$

$\theta$

mot $\Delta$

$$\bar{X}_m$$

$E(\bar{X}_m) = \mu$

$\sqrt{\text{var}}(\bar{X}_m) = \frac{\sigma}{\sqrt{m}}$

consistent

raggiunge il  
limite di Cramér-Rao  
(efficiente in senso  
assoluto)

① come si distribuisce ②

proprietà riproduttiva  $\rightarrow$  Se  $X \sim N$

allora  $\bar{X}_m \sim N$

distribuzione  
esalta

$$\frac{\bar{X}_m - \mu}{\sigma/\sqrt{m}} \sim N(0, 1)$$

(2)

$$X \sim N(\mu, \sigma^2)$$

$\theta$



incognito

[ PARAMETRO DI  
DISTURBO ]

$$\frac{\bar{X}_m - \mu}{\sigma/\sqrt{m}} \sim N(0, 1)$$

$s$

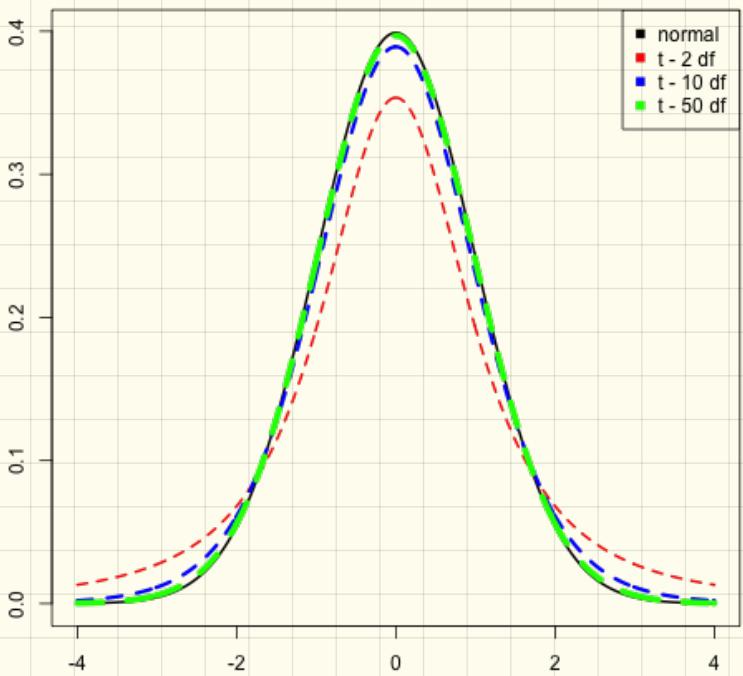
uso il campione per stimare  
anche il parametro di  
disturbo

$$t_{m-1}$$

gradi di libertà  
unico parametro t  
di Student

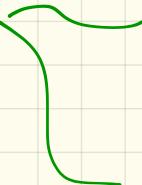
al crescere dei g. d. l., la v.c. t converge  
ad un normale standardizzata

esistono delle tavole che contengono alcuni  
percentili di più comune utilizzo di questa  
distribuzione



(3)

$$X \sim f$$



posso conoscere e anche non è  
normale

non la conosco

Sono comunque interessato a stimarne la  $\mu$  sulla base  
delle osservazioni campionarie

$$\frac{\bar{X}_m - \mu}{\sigma/\sqrt{m}} \xrightarrow{\text{converge}} N(0,1)$$

distribuzione asintotica

$$\frac{\bar{X}_m - \mu}{s/\sqrt{m}} \rightarrow t_{m-1} \rightarrow N(0,1)$$

se m cresce

# INFERNZA SULLA VARIANZA $\sigma^2$

$$X \sim N(\mu, \sigma^2)$$

di solito non ho informazioni su  $\mu$ : parametro di disturbo

l'inferenza su  $\sigma^2$  è meno robusta rispetto alle riduzioni dell'ipotesi di normalità

Sotto queste ipotesi si può mostrare

$$\frac{(n-1) s^2}{\sigma^2} \sim \chi_{m-1}^2$$

unico parametro  
(gradi di libertà)

## INFERNZA SULLA PROPORTIONE $\pi$

$$X \sim \text{Bin}(\pi)$$

↳ in questo caso la distribuzione esatta di  $\hat{\pi}_m = \bar{X}_m = \hat{\pi}$   
può essere ricavata sfruttando la distribuzione binomiale



se  $m$  è "abbastanza" grande si può anche in questo caso sfruttare  
il teorema del limite centrale per ricavare la distribuzione asintotica

$$\frac{\hat{\pi} - E(\hat{\pi})}{\sqrt{\text{Var}(\hat{\pi})}} = \frac{\hat{\pi} - \pi}{\sqrt{\frac{\hat{\pi}(1-\hat{\pi})}{m}}} \xrightarrow{\text{TLC}} N(0, 1)$$



In molti casi si può essere interessati a confrontare i parametri di due popolazioni. In particolare:

- confronto fra le due medie (caso di varianze note e caso di varianze incognite) di due popolazioni normali

↓

questo schema, sfruttando il TLC, può essere esteso anche al caso di due generiche popolazioni se le numerosità campionarie sono adeguate

- confronto fra le varianze di due popolazioni normali
- confronto fra le proporzioni di due popolazioni di Bernoulli

## CONFRONTO TRA LE MEDIE DI DUE POPOLAZIONI NORMALI (caso di varianze note)

Siano  $X$  e  $Y$  due v.c. normali ed indipendenti

$$X \sim N(\mu_x, \sigma_x^2)$$

↓  
note a priori



$$Y \sim N(\mu_y, \sigma_y^2)$$

↓  
note a priori



$$\bar{X}_m \sim N\left(\mu_x, \frac{\sigma_x^2}{m_x}\right) \rightsquigarrow \text{per le proprietà riproduttive della normale} \leftarrow \bar{Y}_m \sim N\left(\mu_y, \frac{\sigma_y^2}{m_y}\right)$$



se si considera la v.c.  $\bar{X}_m - \bar{Y}_m$  (teorema sulle v.c. normale)

$$\bar{X}_m - \bar{Y}_m \sim N\left(\mu_x - \mu_y, \frac{\sigma_x^2}{m} + \frac{\sigma_y^2}{m}\right)$$

qualsiasi combinazione lineare di v.c. normali è ancora normale

$$\frac{(\bar{X}_m - \bar{Y}_m) - (\mu_x - \mu_y)}{\sqrt{\frac{\sigma_x^2}{m_x} + \frac{\sigma_y^2}{m_y}}} \sim N(0, 1)$$

## CONFRONTO TRA LE MEDIE DI DUE POPOLAZIONI NORMALI (caso di varianze incognite)

Se nel precedente schema teorico le due varianze

$\sigma_x^2$  e  $\sigma_y^2$  sono incognite si parla di problema di Behrens-Fisher



come nel caso dell'inferenza su una sola popolazione, quando si sostituisce alla varianza la sua stima corretta, la distribuzione campionaria dello stimatore passa da una distribuzione normale ad una t di Student



le formule per il calcolo dei g.d.l. della t di Student risultante sono molto complesse



il caso più semplice da gestire si basa sull'ipotesi di OMOSCEDASTICITÀ: le due popolazioni hanno varianza incognita che però si può assumere uguale  $\rightarrow \sigma_x^2 = \sigma_y^2 = \sigma^2$

Si ha allora:

$$X \sim N(\mu_x, \sigma^2)$$



$$\bar{X}_m \sim N\left(\mu_x, \frac{\sigma^2}{m_x}\right)$$

$$Y \sim N(\mu_y, \sigma^2)$$



$$\bar{Y}_m \sim N\left(\mu_y, \frac{\sigma^2}{m_y}\right)$$

$$\frac{(\bar{X}_m - \bar{Y}_m) - (\mu_x - \mu_y)}{\sqrt{\frac{\sigma^2}{m_x} + \frac{\sigma^2}{m_y}}} = \frac{(\bar{X}_m - \bar{Y}_m) - (\mu_x - \mu_y)}{\sqrt{\sigma^2 \left( \frac{1}{m_x} + \frac{1}{m_y} \right)}}$$



sotto l'ipotesi di omoschedasticità ho un solo parametro di disturbo

posso stimarlo uscendo lo  $\rightarrow \frac{(\bar{X}_m - \bar{Y}_m) - (\mu_x - \mu_y)}{\sqrt{s_p^2 \left( \frac{1}{m_x} + \frac{1}{m_y} \right)}} \sim t_{m_x + m_y - 2}$

$s_p^2$

p sta per pooled (congiunto)

$$s_p^e = \frac{(m_x - 1)s_x^e + (m_y - 1)s_y^e}{m_x + m_y - 2}$$

lo stimatore congiunto  
della varianza comune  
è una media ponderata  
delle due varianze campionarie  
con pesi date dalle rispettive  
numerosità

esplorando le formule per  $s_x^e$  e  $s_y^e$ , si ottiene:

$$s_p^e = \frac{\frac{(m_x - 1) \sum_i (x_i - \bar{x}_{m_x})^2}{(m_x - 1)} + \frac{(m_y - 1) \sum_i (y_i - \bar{y}_{m_y})^2}{(m_y - 1)}}{(m_x + m_y - 2)} =$$

$$= \frac{\sum_i (x_i - \bar{x}_{m_x})^2 + \sum_i (y_i - \bar{y}_{m_y})^2}{m_x + m_y - 2} = \frac{\text{Dev}(X) + \text{Dev}(Y)}{m_x + m_y - 2}$$

lo stimatore congiunto di  $\sigma^2$   
è la somma delle due devianze  
campionarie

$X \sim f$

rimuoviamo l'ipotesi di  
normalità delle due popolazioni

$Y \sim f$



$$\bar{X}_m \xrightarrow{(m_x \geq 30)} N\left(\mu_x, \frac{\sigma_x^2}{m_x}\right)$$



$$\bar{Y}_m \xrightarrow{(m_y \geq 30)} N\left(\mu_y, \frac{\sigma_y^2}{m_y}\right)$$

$$\bar{X}_m - \bar{Y}_m \sim \cancel{N}\left(\mu_x - \mu_y, \frac{\sigma_x^2}{m_x} + \frac{\sigma_y^2}{m_y}\right)$$



posso usare ancora le stesse  
statistiche

robuste rispetto all'ipotesi di  
violatione di normalità

anche in questo caso, nel caso di varianze  
incognite, per semplicità lavoriamo sotto  
l'ipotesi di omoschedasticità

## CONFRONTO TRA LE VARIANZE DI DUE POPOLAZIONI NORMALI

$$X \sim N(\mu_x, \sigma_x^2)$$



$$\frac{(m_x - 1) s_x^2}{\sigma_x^2} \sim \chi_{m_x - 1}^2$$

$$Y \sim N(\mu_y, \sigma_y^2)$$



$$\frac{(m_y - 1) s_y^2}{\sigma_y^2} \sim \chi_{m_y - 1}^2$$

$$\frac{\frac{(m_x - 1) s_x^2}{\sigma_x^2}}{\frac{(m_y - 1) s_y^2}{\sigma_y^2}} \sim F_{m_x - 1, m_y - 1}$$

||

NOTA: il rapporto fra due v.c.  $\chi^2$ , ciascuna divisa per i corrispondenti galli è una v.c. F di Fisher

anche in questo caso  
è disponibile una  
tavola con i percentili  
di più comune utilizzo

$$\frac{s_x^2 / \sigma_x^2}{s_y^2 / \sigma_y^2} \sim F_{m_x - 1, m_y - 1}$$

due parametri

gall del numeratore      gall del denominatore

NOTA TECNICA: mentre nel caso dell'inferenza sul confronto fra le due medie, nel caso di numerose campionarie adeguate si può utilizzare comunque lo stesso schema anche se le due popolazioni non sono normali, nel caso di inferenza sulle due varianze la procedura risulta abbastanza sensibile a violazioni dell'ipotesi di normalità

# CONFRONTO TRA LE PROPORZIONI DI DUE POPOLAZIONI BERNOULLIANE

$$X \sim \text{Ber}(\pi_x)$$

$$Y \sim \text{Ber}(\pi_y)$$

la distribuzione esatta di  $\hat{\pi}_x = \bar{X}_m = \hat{\pi}_x$   
è ricavabile a partire dalla v.c. binomiale

lo stesso per  
 $\hat{\pi}_y = \bar{Y}_m = \hat{\pi}_y$

se i due campioni sono sufficientemente grandi si può sfruttare il TLC per ottenere le due distribuzioni campionarie

$$\hookrightarrow (\hat{\pi}_x - \hat{\pi}_y) \xrightarrow{\text{TLC}} N \left( \pi_x - \pi_y, \frac{\pi_x(1-\pi_x)}{m_x} + \frac{\pi_y(1-\pi_y)}{m_y} \right)$$

la varianza congiunta può essere  
stimata usando lo stimatore  $\hat{\pi}_p$

$$\hat{\pi}_p = \frac{\sum X_i + \sum Y_i}{m_x + m_y}$$

proporzione di successi nei  
due campioni

$$\frac{(\hat{\pi}_x - \hat{\pi}_y) - (\pi_x - \pi_y)}{\sqrt{\hat{\pi}_p(1-\hat{\pi}_p)\left(\frac{1}{m_x} + \frac{1}{m_y}\right)}} \xrightarrow{TLC} N(0,1)$$

where:  $\hat{\pi}_p = \frac{\sum_{i=1}^{m_x} X_i + \sum_{i=1}^{m_y} Y_i}{m_x + m_y}$