

1

# Formulario statistica descrittiva

solo per caratteri QUANTITATIVI!

**MEDIA**

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{x} = \frac{1}{n} \sum_{i=1}^k x_i \cdot n_i, \quad \bar{x} = \frac{1}{N} \sum_{i=1}^k x_i g_i$$

FR. ASSOLUTA                      FR. RELATIVA

**MEDIA ARITMETICA PONDERATA**

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

**MEDIA GEOMETRICA**

$$\bar{x}_g = \sqrt[n]{\prod_{i=1}^n x_i} \rightarrow \text{prodottoria} = x_1 \cdot x_2 \cdot x_3 \dots$$

$$\bar{x}_g = \sqrt[n]{(x_1)^{n_1} \cdot \dots \cdot (x_k)^{n_k}} \quad \text{FR. ASSOLUTA}$$

$$\bar{x}_g = \sqrt[n]{(x_1)^{g_1} \cdot \dots \cdot (x_k)^{g_k}} \quad \text{FR. RELATIVA}$$

**MODA** → modalità che presenta la MASSIMA FREQUENZA

**MEDIANA** → SOLO PER CARATTERI ORDINABILI!

↓  
corrisponde alla modalità per cui  $FX = 0,5$

**VARIANZA**

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

$$\sigma^2 = E[(x - \bar{x})^2]$$

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 \cdot n_i - \bar{x}^2$$

**VARIANZA DI UNA TRASFORMAZIONE LINEARE** →  $Y = \alpha X + \beta$

$$\sigma_y^2 = \alpha^2 \sigma_x^2 \rightarrow v(y) = \alpha^2 v(x)$$

$$v(\alpha x) = \alpha^2 v(x)$$

$$v(aX + bY) = a^2 v(X) + b^2 v(Y) + 2ab \text{cov}(X, Y) \rightarrow \text{se } x \text{ e } y \text{ sono indipendenti } v(x \pm y) = v(x) + v(y)$$

**DEVIATIONE STANDARD** o scarto quadratico medio

$$\sigma = \sqrt{\sigma^2}$$

**COEFFICIENTE DI VARIAZIONE**

$$CV = \frac{\sigma}{|\bar{x}|} \cdot 100$$

**STANDARDIZZAZIONE**  $\Rightarrow$  rende qualsiasi variabile con.

$$z = \frac{x_i - \bar{x}}{s}$$

MEDIA NULLA

VARIANZA UNITARIA

**CAMPO DI VARIAZIONE**

$R = X(n) - X(1)$   $\rightarrow$  differenza tra il valore max assunto dalla variabile e il valore min

**DIFFERENZA INTERQUARTILICA**

$W = Q_3 - Q_1$   $\rightarrow$  include il 50% delle osservazioni

**MEDIA CONDIZIONATA**

$$\bar{y}_{x=x_i} = \frac{1}{n_i} \sum_{j=1}^R y_j \cdot n_{ij} \longleftrightarrow \bar{y} = \frac{1}{n} \sum_{i=1}^H \bar{y}_{x=x_i} \cdot n \cdot \bar{i}$$

medie totale della popolazione.

**VARIANZA CONDIZIONATA**

$$s^2_{Y|X=x_i} = \frac{1}{n_i} \sum_{j=1}^H (y_j - \bar{y}_{x=x_i})^2 \cdot n_{ij}$$

**BARICENTRO DI UNA DISTRIBUZIONE DOPPIA**

punto che ha coordinate  $(\bar{x}, \bar{y})$

# Indipendenza statistica

## REGOLA DI FATTORIZZAZIONE

$$n_{ij}^1 = \frac{n_{i.} \cdot n_{.j}}{n} \rightarrow \text{frequenza teorica che troverei nella tabella } X/Y \text{ nel caso di indipendenza}$$

## Dipendenza o interdipendenza perfetta

X \ Y	y <sub>1</sub>	y <sub>2</sub>	y <sub>3</sub>
X <sub>1</sub>	0	0	•
X <sub>2</sub>	0	•	0
X <sub>3</sub>	•	0	0
X <sub>4</sub>	•	0	0

$$H \neq R$$

X \ Y	y <sub>1</sub>	y <sub>2</sub>
X <sub>1</sub>	•	0
X <sub>2</sub>	0	•

$$H = R$$

relazione  
BIDIREZIONALE

## CONTINGENZE

$$c_{ij} = n_{ij} - n_{ij}^1 \rightarrow \text{frequenza teorica in caso di indipendenza}$$

## INDICE DI ASSOCIAZIONE CHI-QUADRO DI PEARSON

$$\chi^2 = \sum_{i=1}^H \sum_{j=1}^R \frac{c_{ij}^2}{n_{ij}^1} \rightarrow \text{in caso di indipendenza } \chi^2 = 0$$

$$\chi^2 = \sum_{i=1}^H \sum_{j=1}^R \frac{n_{ij}^2}{n_{ij}^1} - n$$

$$\chi^2 = n \left( \sum_{i=1}^H \sum_{j=1}^R \frac{n_{ij}}{n_{i.} \cdot n_{.j}} - 1 \right)$$

sottraggo 1 una volta sola

## INDICE NORMALIZZATO

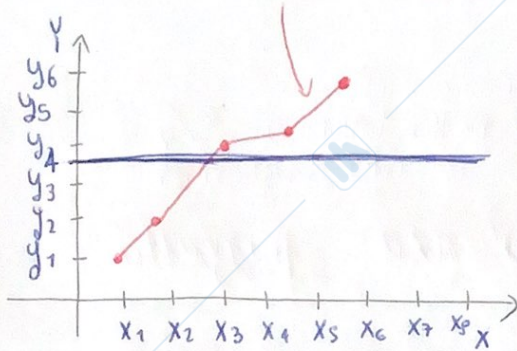
$$\tilde{\chi}^2 = \frac{\chi^2}{n \cdot \min\{H-1, R-1\}}$$

il suo campo di variazione è

indipend. statistica  $[0, 1]$  dipendenza o interdipendenza perfetta

# Indipendenza in media

## SPEZZATA DI REGRESSIONE



$\bar{y}$  → non caso di indipendenza la spezzata assume questa forma

## SCOMPOSIZIONE DELLA VARIANZA

$$\sigma^2 Y = \sigma^2 \text{media}(Y|X) + \text{media}(\sigma^2 Y|X)$$

### VARIANZA SPIEGATA

$$\sigma^2 \text{media}(Y|X) = \frac{1}{n} \sum_{i=1}^H (\bar{y}_{x=x_i} - \bar{y})^2 \cdot n_i \rightarrow \text{varianza delle medie di } Y \text{ condizionate ad } X$$

### VARIANZA RESIDUA

$$\text{media}(\sigma^2 Y|X) = \frac{1}{n} \sum_{i=1}^H \sigma^2 Y|X=x_i \cdot n_i \rightarrow \text{media delle varianze condizionate della } Y \text{ date le } X$$

## RAPPORTO DI CORRELAZIONE

$$r^2_{Y|X} = \frac{\sigma^2 \text{media}(Y|X)}{\sigma^2 Y} \rightarrow \text{il suo campo di variazione e'}$$

$$r^2_{Y|X} = 1 - \frac{\text{media}(\sigma^2 Y|X)}{\sigma^2 Y}$$

[0, 1]
indipendenza statistica
dipendenza statistica

# Incorrelazione

COVARIANZA

misura

la forza di legame lineare che  $\exists$  tra le variabili

concordanza

discordanza

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

$$s_{XY} = E[(X - \bar{x})(Y - \bar{y})]$$

il suo campo di variazione oscilla tra

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

$$[-s_x s_y] + s_x s_y$$

$s_{XY} = E[XY] - E[X] \cdot E[Y]$      $s_{(a+bX)cY} = bc s_{XY}$

(prodotto degli scarti quadratici med.)

## TRASFORMAZIONE LINEARE DELLA COVARIANZA

se  $x \rightarrow ax+b$

$y \rightarrow cy+d \Rightarrow s_{x'y'} = a \cdot c \cdot s_{xy}$

## COEFFICIENTE DI CORRELAZIONE LINEARE DI BRAVAIS E PEARSON

indica l'addensamento dei dati (clustering) nello scatter plot

$$r_{XY} = \frac{s_{XY}}{s_x \cdot s_y}$$

il suo campo di variazione oscilla tra

$$[-1; 1]$$

perfetto legame lineare con caratteri DISCORDI

perfetto legame lineare con caratteri CONCORDI

$r_{XY} = 0$

è possibile che ci sia indipendenza statistica ma non certo, in quanto questo indice non è abbastanza forte da poterlo determinare

se x e y sono correlate allora non è detto che x sia la causa di y o viceversa

potrebbe esistere un terzo fattore z che è la causa di entrambe

2

# Probabilità

$0 \leq P(\cdot) \leq 1$   
 ↙ evento nullo  
 ↘ evento certo ( $P(\Omega)$ )

↳ somma logica  
 $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

$P(A \cap B) = P(A) + P(B) - P(A \cup B)$   
 ↓  
 prodotto logico

## DEFINIZIONE DI PROBABILITA'

$P(E) = \frac{n \text{ di casi favorevoli}}{n \text{ di casi possibili}}$  → questi devono essere EQUIPROBABILI!!!

## PROBABILITA' CONDIZIONATA DI A DATO B

$P(A|B) = \frac{n \text{ di casi favorevoli ad } (A \cap B)}{n \text{ di casi favorevoli a } B} = \frac{P(A \cap B)}{P(B)}$

## PRINCIPIO DELLE PROBABILITA' COMPOSTE

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

## INDIPENDENZA TRA EVENTI

$P(B|A) = P(B)$        $P(A|B) = P(A)$

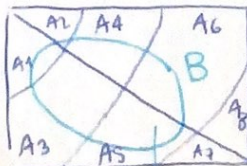
$P(A \cap B) = P(A) \cdot P(B)$

nel caso di distribuzione di probabilità congiunta

↓  
 SI HA INDIPENDENZA SE LE PROBABILITA' CONDIZIONALI SONO UGUALI!

## TEOREMA DI BAYES → utile per problemi a più casi

$P(A_i|B) = \frac{P(A_i \cap B)}{P(B)} = \frac{P(A_i) \cdot P(B|A_i)}{P(B)}$



$P(B) = \sum_{i=1}^n P(B \cap A_i)$   
 ↳ teorema delle probabilità totali

$= \frac{P(A_i) \cdot P(B|A_i)}{P(B)}$       ↳ verosimiglianze  
 $\sum_{i=1}^n P(A_i) \cdot P(B|A_i)$

3

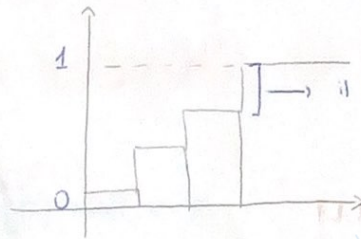
## Variabili casuali discrete

### FUNZIONE DI RIPARTIZIONE

↳ si basa sulle probabilità cumulate

$$F(x) = P(X \leq x) = \sum_{w \leq x} P(X=w) \quad x \in \mathbb{R} \quad \begin{cases} \text{dominio } \mathbb{R} \\ \text{codominio } [0,1] \end{cases}$$

↓  
funzione a gradini



il singolo gradino corrisponde all'ultima probabilità che viene aggiunta

### FUNZIONE DI PROBABILITÀ

associa ad ogni possibile valore di  $x$  la probabilità che questo si verifichi  $\rightarrow P(X=x_i)$

VALORE ATTESO  $\rightarrow$  misura dove i valori tendono a concentrarsi

$$E(X) = \sum_i x_i \cdot P(x_i)$$

### VARIANZA

$$V(X) = \sum_i [x_i - E(X)]^2 \cdot P(x_i)$$

$$V(X) = E[(X - E(X))^2]$$

$$V(X) = E(X^2) - [E(X)]^2$$

momento secondo  $\Rightarrow E(X^2) = \sum_i x_i^2 \cdot P(x_i)$

### FUNZIONE LINEARE DI UNA VARIABILE CASUALE

$$Y = a + bX$$

$$E(Y) = a + b(E(X))$$

$$V(Y) = b^2(V(X))$$

$$SD(Y) = \sqrt{b^2 V(X)}$$

## DISTRIBUZIONE UNIFORME DISCRETA

$X \sim Ud(a, s) \rightarrow$  il numero dei possibili valori  
valore più piccolo assumibile

funzione di probabilità

$P(X=x) = \frac{1}{s} \rightarrow$  l'intervallo in cui questa funzione è definita è  $[a, a+s-1]$

$$E(X) = a + \frac{s-1}{2}$$

$$V(X) = \frac{s^2-1}{12}$$

## DISTRIBUZIONE DI BERNOULLI

$X \sim Ber(\pi) \rightarrow$  probabilità di successo

si riferisce ad osservazioni BINARIE

X	0	1
P(X)	$1-\pi$	$\pi$

$\Leftarrow$

$X=0$   
quando l'evento si VERIFICA

$X=1$   
quando l'evento non si verifica

$$P(X=x) = \pi^x (1-\pi)^{1-x}$$

$$E(X) = \pi$$

$$V(X) = \pi(1-\pi)$$

## DISTRIBUZIONE BINOMIALE $\rightarrow$ Successione di n prove di tipo Bernoulliano

$X \sim Bin(n, \pi)$  numero di prove effettuate

$$P(X=x) = \binom{n}{x} \pi^x (1-\pi)^{n-x} \rightarrow \text{numero di Successi ottenuti}$$

BINOMIO DI NEWTON

$$\binom{n}{x} = \frac{n!}{x!(n-x)!}$$

$$E(X) = n\pi$$

$$V(X) = n\pi(1-\pi)$$

**DISTRIBUZIONE DI POISSON** → applicazione nel conteggio di un evento casuale in uno specifico  $\Delta t$

$$X \sim \text{Poi}(\lambda)$$

$$P(X=x) = \frac{\lambda^x}{x!} \cdot e^{-\lambda}$$

$$\left. \begin{array}{l} E(X) = \lambda \\ V(X) = \lambda \end{array} \right\} \text{genovano} \\ \text{EQUIDISPERSO}$$

## Variabili casuali continue

FUNZIONE DI DENSITA'

$$P(X) = \int_a^b f(x) dx$$

$P(X=x) = 0$   
non è possibile calcolare la probabilità di un singolo evento

FUNZIONE DI RIPARTIZIONE

$$F(x) = \int_{-\infty}^x f(w) dw$$

VALORE ATTESO

$$E(X) = \int_a^b x \cdot f(x) dx$$

VARIANZA

$$V(X) = \int_a^b (x - E(X))^2 f(x) dx$$

$$V(X) = E(X^2) - [E(X)]^2 = \int_{-\infty}^{\infty} x^2 f(x) dx - [E(X)]^2$$

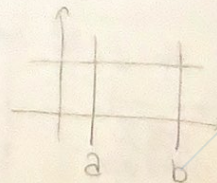
**DISTRIBUZIONE UNIFORME CONTINUA**

$X \sim U(a, b)$ , valore più grande  
valore più piccolo

$$f(x) = \frac{1}{b-a}$$

$$E(X) = \frac{a+b}{2}$$

$$V(X) = \frac{(a-b)^2}{12}$$



**DISTRIBUZIONE NORMALE** → simmetrica rispetto a  $\mu$

$$X \sim \mathcal{N}(\mu; \sigma^2)$$

**FUNZIONE DI DENSITA'**

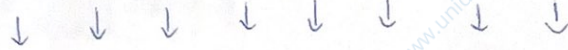
$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \cdot e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$$

**VALORE ATESO** → PARAMETRO DI LOCAZIONE  
determina la posizione della  
curva su  $X$

$$E(X) = \mu$$

**VARIANZA** → PARAMETRO DI SCALA  
più  $\sigma^2$  aumenta più la  
curva è schiacciata

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



la sua standardizzazione è

$$\frac{x-\mu}{\sigma} = z \sim \mathcal{N}(0,1) \rightarrow \text{distribuzione che uso per le tavole}$$

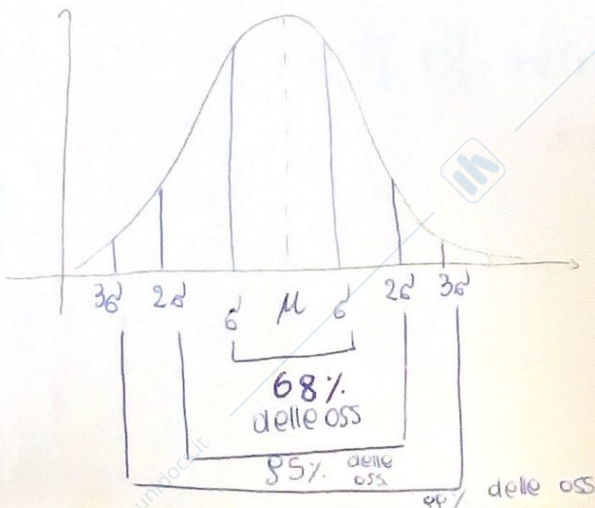
**TRASFORMAZIONE LINEARE DI  $z \sim \mathcal{N}(0,1)$**

$$Y = a + bX \sim \mathcal{N}(a, b^2) \rightarrow \text{quando } X \sim \mathcal{N}(0,1)$$

**SOMMA DI DUE NORMALI**

$$(X_1 + X_2) \sim \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$$

**LEGGE DEI TRE SIGMA**



## teorema del limite centrale

I.I.D.

DATA UNA SUCCESSIONE DI VARIABILI CASUALI TUTTE UGUALI  
NE CONSIDERO LA MEDIA CAMPIONARIA.

SE  $n \rightarrow \infty$  LA MEDIA CAMPIONARIA, A PRESCINDEERE  
DELLA DISTRIBUZIONE DELLE VARIABILI, QUANDO E'  
STANDARDIZZATA, SI DISTRIBUISCE COME UNA  
NORMALE STANDARD  $N(0,1)$

$$\text{se } n \rightarrow \infty \rightarrow z_n = \frac{(\bar{x}_n - \mu) \sqrt{n}}{\sigma} \sim N(0,1)$$

$\downarrow$   
 media  
 campionaria  
 standardizzata

- SE  $\bar{x}$  NON FOSSE STANDARDIZZATA

$$\text{se } n \rightarrow \infty \quad \bar{x}_n \sim N\left(\mu; \frac{\sigma^2}{n}\right)$$

- SE CONSIDERASSIMO LA SOMMA DI  $n$  V.C. I.I.D

$$\text{se } n \rightarrow \infty \quad S_n \sim N(n\mu; n\sigma^2)$$

## CURTOSI

$$\gamma = E \left( \frac{X - E(X)}{SD(X)} \right)^4 = \frac{1}{n} \sum_{i=1}^n \left( \frac{X_i - \bar{x}}{s} \right)^4$$

se  $\gamma < 3$   
la distribuzione è  
IPERNORMALE,  
ovvero ha la coda  
più spessa  
rispetto alla normale

$\gamma = 3$   
la distribuzione  
è NORMALE

se  $\gamma > 3$   
la distribuzione è  
IPONORMALE  
ovvero i casi estremi  
sono ancora più  
estremi

## DISTRIBUZIONE CHI-QUADRO

$$f(x) = \frac{1}{2^{\frac{g}{2}} \Gamma(\frac{g}{2})} \cdot x^{\frac{g}{2}-1} \cdot e^{-\frac{x}{2}}$$

VALORE ATTESO

$$E(X) = g$$

VARIANZA

$$V(X) = 2g$$

→ corrisponde al quadrato  
della distribuzione normale  
standard

$$\chi^2_g \sim \sum_{i=1}^g X_i^2 \rightarrow \mathcal{N}(0,1)^2$$

↳ somma g volte delle  
normali standard  
al quadrato

## DISTRIBUZIONE T DI STUDENT

$$X = \frac{z}{\chi^2_g}$$

## DISTRIBUZIONE F di FISHER

$$F_{g_1, g_2} = \frac{\chi^2_{g_1}}{\chi^2_{g_2}}$$

# 4 Inferenza statistica

## FRAZIONE DI CAMPIONAMENTO

$n$  → dimensione campionaria

$N$  → numerosità della popolazione

## DISTRIBUZIONE DELLA MEDIA CAMPIONARIA → $\bar{x}$ NELLE POPOLAZIONI INFINITE

se  $x \sim N(\mu; \sigma^2)$  allora  $\bar{x} \sim N(\mu; \frac{\sigma^2}{n})$

se  $x \sim \text{Ber}(\pi)$  allora  $\bar{x} \sim \frac{1}{n} \text{Bin}(n; \pi)$

$n > 30$

per T.L.C.  $\bar{x} \sim N(\mu; \frac{\sigma^2}{n})$

## DISTRIBUZIONE DELLA $\bar{x}$ NELLE POPOLAZIONI FINITE

se  $x \sim N(\mu; \sigma^2)$  allora  $\bar{x}_n \sim N(\mu; \frac{(N-n)}{(n-1)} \frac{\sigma^2}{n})$

## PROPRIETA' DEGLI STIMATORI

↓  
QUANTITA' TEORICA

= v.e. utilizzata per stimare una determinata caratteristica  $\theta$  della popolazione

FATTORE DI CORREZIONE

- **CORRETTEZZA** → la media della distr. che ipotizzo è centrata nella media reale

$$E(T) = \theta \quad \text{f.e. possibile}$$

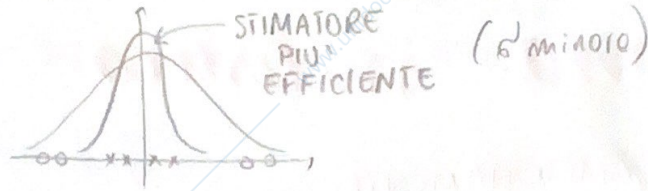
↓

se  $E(T) \neq \theta \Rightarrow T$  sarebbe uno stimatore DISTORTO

↓  
DISTORSIONE

$$B(T) = E(T) - \theta$$

• EFFICIENZA



ERRORE QUADRATICO MEDIO

$$MSE(T) = E[(T - \theta)^2]$$

$$MSE(T) = \underbrace{V(T)}_{\text{varianza}} + \underbrace{[B(T)]^2}_{\text{distorsione}}$$

$$V(T) = E[(T - E(T))^2]$$

nella scelta tra due stimatori, prediligilo quello che presenta un MSE minore  
 ↓  
 i suoi valori saranno più prossimi a  $\theta$

• CONSISTENZA ( $\infty$ )

↓  
 all'aumentare della  $n$  cresce la precisione dello stimatore

LEGGI DEI GRANDI NUMERI  
 uno stimatore è consistente se la probabilità che una sua realizzazione cada in un intervallo che contiene il vero valore nella popolazione tende ad uno a mano a mano che la numerosità del campione aumenta

$$\lim_{n \rightarrow \infty} MSE(T_n) = \lim_{n \rightarrow \infty} E[(T_n - \theta)^2] = 0$$

→ UNO STIMATORE SI DICE CONSISTENTE IN MEDIA QUADRATICA SE  $V(T_n)$  E  $B(T_n)$  TENDONO A ZERO

• CORRETTEZZA ASINTOTICA ( $\infty$ )

$$\lim_{n \rightarrow \infty} B(T) = \lim_{n \rightarrow \infty} E(T_n) - \theta = 0$$

STIMA PUNTUALE DELLA MEDIA DI UNA POPOLAZIONE

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

stimatore naturale della media

MEDIA CAMPIONARIA  $\bar{x} / E(\bar{x}) = \mu$

stimatore corretto / consistente

STIMA PUNTUALE DELLA PROPORZIONE IN UNA POPOLAZIONE

$$\hat{\pi} = \bar{x}$$

$$E(\bar{x}) = \pi$$

$$V(\bar{x}) = \frac{\pi(1-\pi)}{n}$$

stimatore corretto / consistente

## STIMA PUNTUALE DELLA VARIANZA DELLA POPOLAZIONE

$\hat{\sigma}^2$   
(var. campionaria) → STIMATORE DISTORTO

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \rightarrow \text{VARIANZA CAMPIONI CORRETTA}$$

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

↓  
equivale a

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2$$

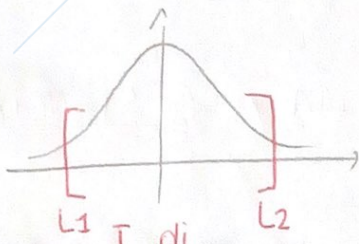
MEDIA QUADRATICA

stimatore non corretto di  $\sigma^2$

→ per passare da uno stimatore non corretto  $\hat{\sigma}^2$  allo stimatore corretto  $S^2$

$$S^2 = \hat{\sigma}^2 \cdot \frac{n}{n-1}$$

## STIMA PER INTERVALLO



I di confidenza di livello  $1-\alpha$  → arbitrario per il parametro  $\theta$

La sua realizzazione in corrispondenza del valore osservato e' detto

INTERVALLO DI CONFIDENZA STIMATO

NON abbiamo la certezza che il nostro intervallo stimato attraversa i dati del campione contenga  $\theta$

$1-\alpha\%$  dei campioni possibili producono un I che contiene il vero valore di  $\theta$

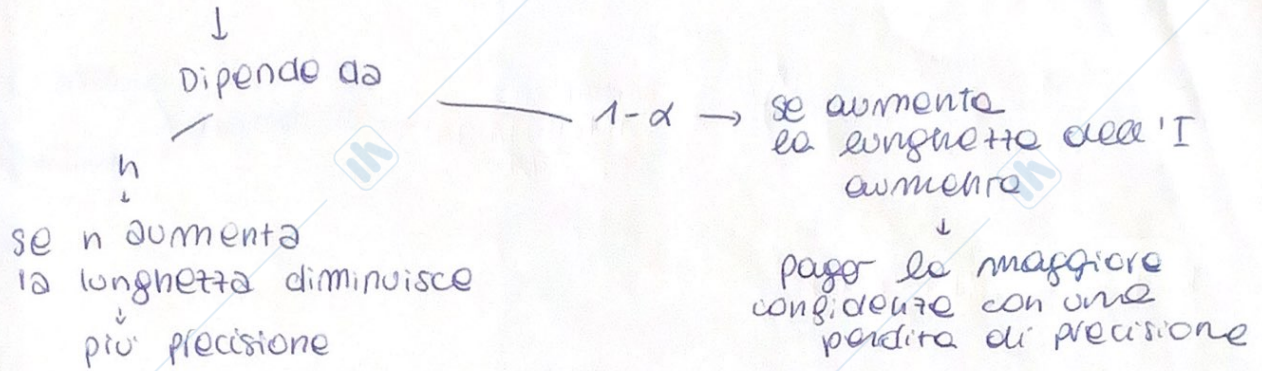
$\alpha\%$  dei campioni casuali che posso estrarre dalla popolazione producono un I che contiene  $\theta$

## I DI CONFIDENZA PER LA MEDIA $\mu$ ( $\sigma^2$ NOTA)

$$\left[ \bar{x} - z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} ; \bar{x} + z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \right]$$

## LUNGHEZZA DELL'INTERVALLO DI CONFIDENZA

$$\text{lunghezza} = 2 \cdot z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}$$



## INTERVALLO DI CONFIDENZA PER LA MEDIA $\mu$ ( $\sigma^2$ IGNOTA)

SD campionaria e corretta

$$\left[ \bar{x} - t_{\frac{\alpha}{2}} \cdot \frac{S}{\sqrt{n}} ; \bar{x} + t_{\frac{\alpha}{2}} \cdot \frac{S}{\sqrt{n}} \right]$$

t di student con n-1 g.d.e.

MAX DIFFERENZA ASSOLUTA CHE CI ASPETTIAMO TRA LA NOSTRA STIMA E IL VERO PARAMETRO

**MARGINE D'ERRORE**  
↳ semilunghezza dell'I di confidenza

## INTERVALLO DI CONFIDENZA PER LA PROPORZIONE (campioni elevati)

CONDIZIONE PER L'USO DEL T.L.O.

- $n \cdot \hat{\pi} \geq 5$
- $n(1 - \hat{\pi}) \geq 5$

↳ distr. reale  
 $\hat{\pi}/\bar{x} \sim \frac{1}{n} \text{Bin}(n; \pi)$

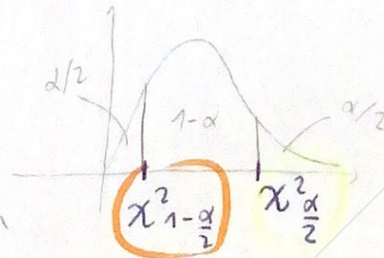
$$\left[ \bar{x} - z_{\frac{\alpha}{2}} \cdot \sqrt{\frac{\bar{x}(1-\bar{x})}{n}} ; \bar{x} + z_{\frac{\alpha}{2}} \cdot \sqrt{\frac{\bar{x}(1-\bar{x})}{n}} \right]$$

distrib. della varianza campionaria corretta  $S^2$   
↓  
ignota!  
↳ trasformazione in  $\chi^2$  per il calcolo del p.v. →  
 $\frac{(n-1) \cdot S^2}{\sigma^2} \sim \chi^2_{n-1}$

## INTERVALLO DI CONFIDENZA PER LA VARIANZA $\sigma^2$

$$\left[ \frac{(n-1)S^2}{\chi^2_{\frac{\alpha}{2}}}; \frac{(n-1)S^2}{\chi^2_{1-\frac{\alpha}{2}}} \right]$$

chi-quadro con n-1 g.d.e.



NON È SIMMETRICA  
↓  
dovrò identificare separatamente entrambi i valori.

→ se dovess. stimare  $\sigma$  dovrò moltiplicare tutto sotto  $\sqrt{\quad}$

5

# 5 passi da seguire per verificare l'ipotesi

Hp. sulla  
quale si basa  
la verifica d'ipotesi

1. che il campione  
sia estratto casualmente (ogni xi e' indipendente dalle altre) SOSPETTO SUL PARAMETRO  
2. che le variabili siano tra loro **INDIPENDENTI**

1. Formulare il sistema d'ipotesi,  $H_0$  e  $H_1$
2. Determinare la statistica-test più appropriata e determinare la distribuzione campionaria sotto  $H_0$

es. se sto facendo un test sulla proporzione o sulla media  $\mu$  la statistica test sarà basata sulla media campionaria  
 $\bar{x}$

3. scelgo il livello di significabilità  $\alpha$  e la numerosità campionaria  $n$  ecc. } → connesso all'errore di I tipo

se e' importante che questo errore sia minimo basterà  $\alpha$  piccolo

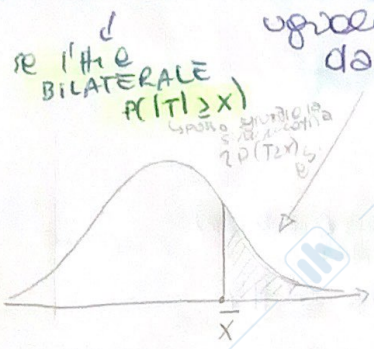
4. Determino la regione di rifiuto del test, calcolandone i valori critici  $c$ , a seconda dell'alternativa

5. Estraggo il campione e svolgo tutte le analisi
6. calcolo il valore campionario dello statistico-test
7. vedo se questo valore cade nella regione di accettazione o di rifiuto del test

8. prendo una decisione (statistica) per poterne descrivere le conseguenze nei termini del problema reale che si sta affrontando.

# P-VALUE

→ Prob. di osservare un valore della statistica Test uguale o più estremo del valore ottenuto dal campione sotto l'Hp nulla



se  $P(T \geq \bar{x})$  è MOLTO GRANDE non rifiuto l'Ho  $\alpha > 0,1$

se  $P(T \geq \bar{x})$  è PICCOLO ( $P < \alpha$ ) rifiuto l'Ho

# ERRORI DI PRIMO E SECONDO TIPO

(es. falsi negativi)  
RIFIUTARE  
L'Hp NULLA QUANDO  
È VERA → + imp!!

$\alpha$  = prob. di commettere l'errore di I tipo

$1 - \alpha$  = COEFF. DI CONFIDENZA DEL TEST

(es. falsi positivi)  
NON RIFIUTARE l'Hp NULLA QUANDO È FALSA

$\beta$  = prob. di commettere l'errore di II tipo

$1 - \beta$  = POTENZA DEL TEST

$\alpha$  e  $\beta$  sono INVERSAMENTE PROPORZIONALI

NON possono essere minimizzati entrambi

è molto più imp. minimizzare gli errori di I tipo, anche a discapito di quelli di II tipo

## Test per $\mu_1 - \mu_2$ ( $\sigma^2$ note)

$$X_1 \sim N(\mu_1, \sigma_1^2)$$

$$X_2 \sim N(\mu_2, \sigma_2^2)$$

$$Z = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0, 1)$$

↓  
se pongo  $\mu_D = \mu_1 - \mu_2$   
allora

$$\begin{cases} H_0: \mu_D = 0 \\ H_1: \mu_D \neq 0 \end{cases}$$

→ mi riconduco  
al test per le  
medie classiche

$$X_1 \sim N(\mu_1, \sigma_1^2)$$

$$X_2 \sim N(\mu_2, \sigma_2^2)$$

• SOMMA DI NORMALI

$$X_1 + X_2 \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$$

• DIFFERENZA TRA  
NORMALI

$$X_1 - X_2 \sim N(\mu_1 - \mu_2, \sigma_1^2 + \sigma_2^2)$$

completo  
somma delle  $\sigma^2$

## Test $\mu_1 - \mu_2$ ( $\sigma^2$ ignote ma UGUALI)

$$T = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{S_p^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}} \stackrel{H_0}{\sim} t_{(n_1 + n_2 - 2)}$$

↓  
STIMATORE CONGIUNTO  
(pooled)

corrisponde alle medie  
ponderate di  $s_1$  e  $s_2$   
con pesi proporzionali alle  
loro numerosità.

$$S_p^2 = \frac{(n_1 - 1) s_1^2 + (n_2 - 1) s_2^2}{n_1 + n_2 - 2}$$

## Test per la media (σ² nota)

→  $X \sim N(\mu, \sigma^2)$   
oppure in tutti i casi  
in cui  $n \approx 120$

Se  $\mu_1 = \mu_2$   
sotto  $H_0$ :  $Z = \frac{\bar{X} - \mu_0}{s/\sigma/\sqrt{n}} \sim N(0, 1)$

## Test per la media (σ² ignota)

sotto  $H_0$ :  $T = \frac{\bar{X} - \mu_0}{s/\sqrt{n}} \sim t_{n-1}$

stimatore  
della varianza  
 $s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2$

## Test per la proporzione

→  $X \sim \text{Ber}(\pi)$

↓  
devo poter applicare il T.C.C.

↳ verifico che  $|n\pi_0| \geq 5$ ;  $|n(1-\pi_0)| \geq 5$

$H_0: \pi = \pi_0$

$$Z = \frac{\bar{X} - \pi_0}{\sqrt{\frac{\pi_0(1-\pi_0)}{n}}} \stackrel{H_0}{\sim} N(0, 1)$$

## Test per la variante

$$\chi = (n-1) \frac{s^2}{\sigma_0^2} \stackrel{H_0}{\sim} \chi^2_{n-1}$$

→ la chi-quadrato  
NON  
è simmetrica!!!

## Test per il rapporto tra due varianze

$X_1$  e  $X_2$  devono distribuirsi come  $N(\mu, \sigma^2)$

Sotto H<sub>0</sub>:  $\sigma_1^2 = \sigma_2^2$

varianze costanti

$$F = \frac{S_1^2}{S_2^2} \sim F_{(n_1-1), (n_2-1)}$$

F di Fisher (non simm.)

DECIDO QUANTO METTERE SOPRA O SOTTO IN BASE A QUALE MI PERMETTE DI AVERE IL RAPPORTO > 4  
→ meglio sopra la più grande

## Test per la differenza tra due proporzioni

devo sapere un campione GRANDE ⇒ T.L.C.

$$X_1 \sim \text{Ber}(\pi_1)$$

$$X_2 \sim \text{Ber}(\pi_2)$$

Sotto H<sub>0</sub>:  
 $\hat{\pi}_1 = \hat{\pi}_2$

$$Z = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\bar{X}_p (1 - \bar{X}_p) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}} \sim N(0, 1)$$

STIMATORE CONGIUNTO DI  $\pi$

$$\bar{X}_p = \frac{n_1 \cdot \bar{X}_1 + n_2 \cdot \bar{X}_2}{n_1 + n_2}$$

↓  
 bisogna sempre verificare che

$$\begin{cases} n \hat{\pi}_A \geq 5 \\ n(1 - \hat{\pi}_A) \geq 5 \end{cases} \text{ e } \begin{cases} n \hat{\pi}_B \geq 5 \\ n(1 - \hat{\pi}_B) \geq 5 \end{cases}$$

Se un test passa dall'essere bilaterale ( $\neq$ )  
all'essere unilaterale ( $\leq$ )

↓  
il valore critico passa da

$$z_{\frac{\alpha}{2}} / t_{\frac{\alpha}{2}} \dots \rightarrow t_{\alpha} / t_{\alpha} \rightarrow \text{E' PIU' PICCOLO}$$

la regione di rifiuto

AUMENTA