



TECNICHE DI ANALISI DEI DATI

AA 2017/2018

PROF. V.P. SENESE

Questi materiali sono disponibili per tutti gli studenti al seguente indirizzo:

<https://goo.gl/hxL9zG>

Seconda Università di Napoli (SUN) – Dipartimento di Psicologia – TECNICHE DI ANALISI DEI DATI – © Prof. V.P. Senese

INFERENZA STATISTICA

Una volta definito un certo **piano sperimentale**, potremmo non limitarci a voler considerare una singola specifica manifestazione di tale fenomeno (e l'insieme particolare di numeri-misure che in tale singola esecuzione si ottiene), ma fare invece riferimento alla **pluralità delle possibili esecuzioni del medesimo**.

Sotto questa prospettiva, che diremo probabilistico-inferenziale, **ogni singola misura è una variabile** poiché, in esecuzioni ripetute del medesimo piano sperimentale, può ottenere differenti risposte.

INFERENZA STATISTICA

Ogni grandezza variabile (ossia, con possibili variazioni da esecuzione ad esecuzione) è supposta governata da una certa legge di probabilità (detta anche **distribuzione di probabilità**).

Per sottolineare questo presupposto caratteristico, la grandezza variabile è detta **variabile casuale** o **aleatoria** ossia variabile la cui instabilità (di prova in prova) è qualificata da certe regolarità di tipo casuale.

LE DISTRIBUZIONI TEORICHE

Per poter analizzare e interpretare un fenomeno osservato è necessario riferirlo alla sua **probabilità di accadimento**. A tale scopo in statistica vengono impiegate le **distribuzioni teoriche di probabilità** il cui vantaggio consiste nella possibilità di stimare per ciascun evento la sua probabilità di accadimento.

In Psicologia molte sono le distribuzioni teoriche utilizzate:

Normale
Binomiale
Chi quadro
F di fisher
t di student
...

DISTRIBUZIONE NORMALE

È la funzione di probabilità che viene utilizzata per descrivere le **variabili casuali continue**.



Gauss

È DEFINITA DA:

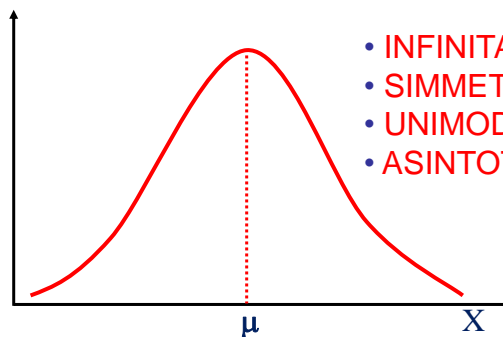
$\bar{X} = \mu$ = media della popolazione

$\sigma = ds$ della popolazione

HA LE SEGUENTI CARATTERISTICHE:

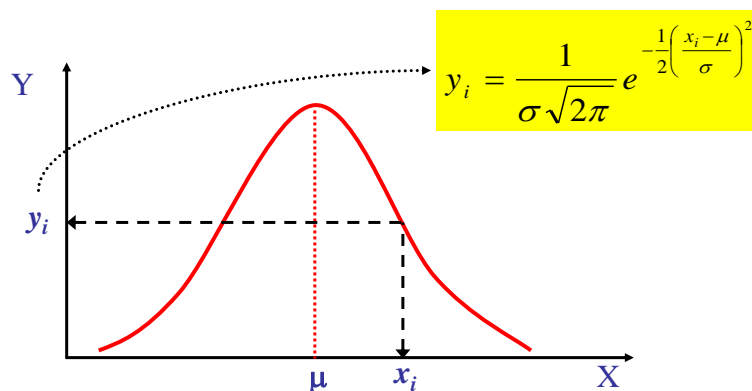
- INFINITA
- SIMMETRICA
- UNIMODALE
- ASINTOTICA

L'area sottesa alla curva è pari ad 1.

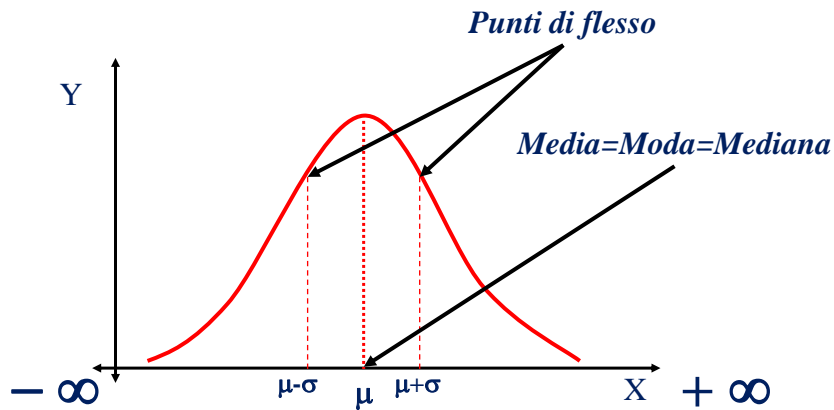


DISTRIBUZIONE NORMALE

Per qualsiasi valore x_i che la variabile X può assumere, attraverso una funzione si calcola la y_i corrispondente (probabilità associata).

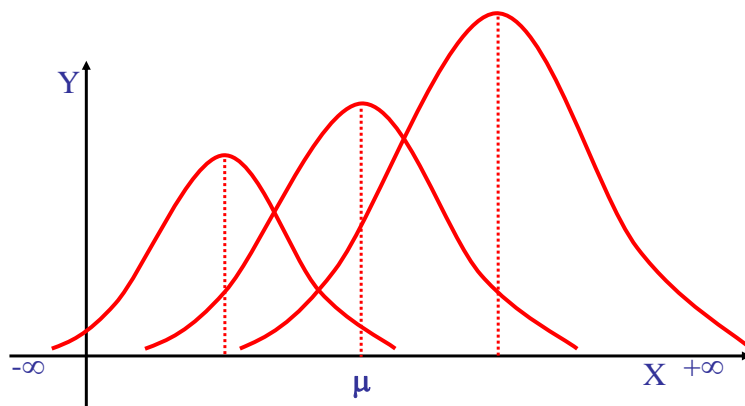


DISTRIBUZIONE NORMALE



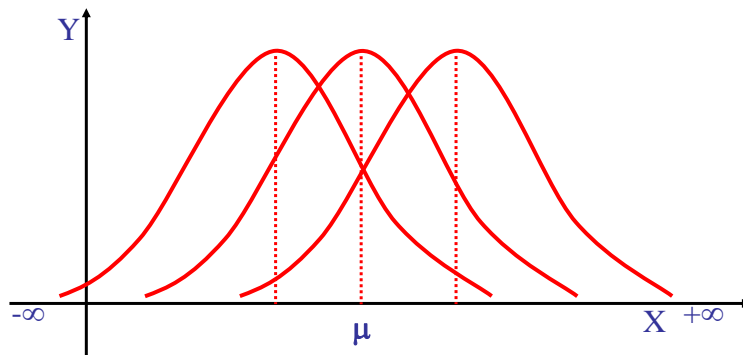
DISTRIBUZIONE NORMALE

La curva **NORMALE** è definita dai parametri μ e σ . Abbiamo un'ampia famiglia di distribuzioni normali con medie e deviazioni standard diverse...



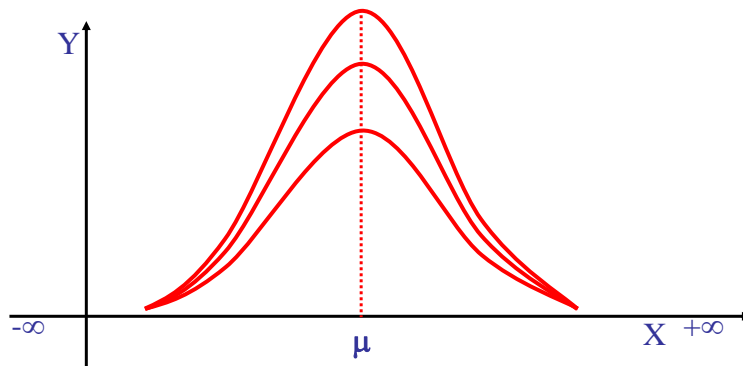
DISTRIBUZIONE NORMALE

...con uguali ds ma diverse medie...



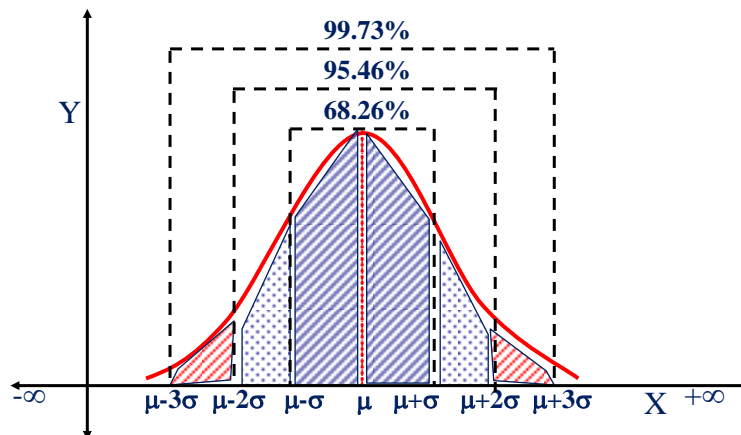
DISTRIBUZIONE NORMALE

...oppure con uguali medie ma diverse ds .



DISTRIBUZIONE NORMALE

Qualsiasi siano i parametri μ e σ , l'area della porzione di curva delimitata dalla media e un'ascissa espressa in termini di deviazioni standard è **costante**.



DISTRIBUZIONE NORMALE

Per verificare se la distribuzione dei valori di una variabile è **normale**, è possibile utilizzare due indicatori statistici:

- la **SIMMETRIA** (*skewness*) \Rightarrow sx; dx
- la **CURTOSI** (*kurtosis*) \Rightarrow bassa; alta

Variano tra $-\infty$ e $+\infty$ e quando assumono valore **0** indicano una distribuzione **perfettamente normale**.

Nella **prassi psicologica** si considerano normalmente distribuite variabili con valori compresi tra ± 1 (\Rightarrow max 2)

DISTRIBUZIONE NORMALE

Nell'ambito della famiglia delle distribuzioni normali la possibilità di **trasformare i valori di una variabile continua in valori standardizzati** consente agli studiosi di far riferimento ad un'unica distribuzione di probabilità che risulta **indipendente dalla specifica variabile oggetto di studio**:

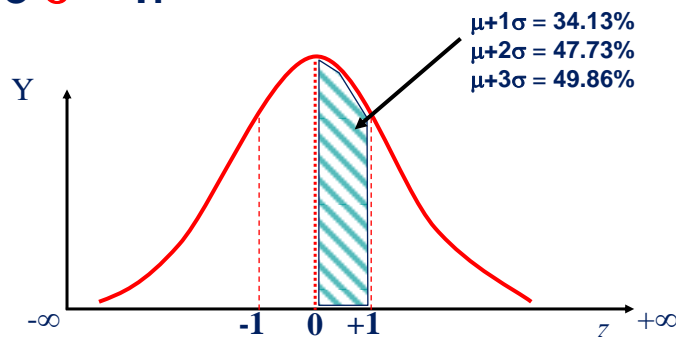
i punti z: la distribuzione normale standard.

$$z_i = \frac{x_i - \mu}{\sigma}$$

DISTRIBUZIONE NORMALE *STANDARD*

Si tratta di una distribuzione di probabilità di forma normale associata ai valori di una variabile standardizzata i cui parametri caratteristici sono:

$\mu = 0$ e $\sigma = 1$.



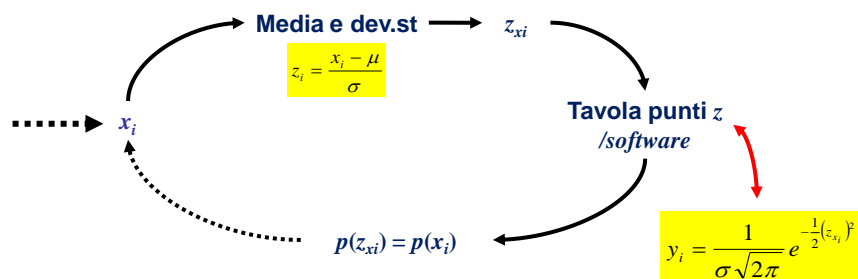
DISTRIBUZIONE NORMALE STANDARD

z	0	1	2	3
.0	.0000	.0040	.0080	.0120
.1	.0398	.0438	.0478	.0517
.2	.0793	.0832	.0871	.0910
.3	.1179	.1217	.1255	.1293
.4	.1554	.1591	.1628	.1664
.5	.1915	.1950	.1985	.2019
.6	.2257	.2291	.2324	.2357
.7	.2580	.2611	.2642	.2673
.8	.2881	.2910	.2939	.2967
.9	.3159	.3186	.3212	.3238

$$z = 0.52 \Rightarrow p = .1985$$

DISTRIBUZIONE NORMALE STANDARD

In pratica, mediante i **punti z**, conoscendo la *media* e la *deviazione standard* della distribuzione, per qualsiasi valore osservato x_i , è possibile calcolare la probabilità di accadimento. Basta fare riferimento alle tavole dei **punti z**, o utilizzare un *software*.



DISTRIBUZIONE NORMALE **STANDARD**

Riassumendo:

1. i **punti z** consentono di **calcolare in modo immediato la probabilità di accadimento di un dato valore** in una distribuzione di cui se ne conoscano i parametri;
2. inoltre, i **punti z** , che come altre scale standardizzate (es. percentili, decili, ecc.), consentono di **trasformare su una metrica comune** punteggi che appartengono a distribuzioni differenti, **possono anche essere utilizzati per confrontare i diversi punteggi**.

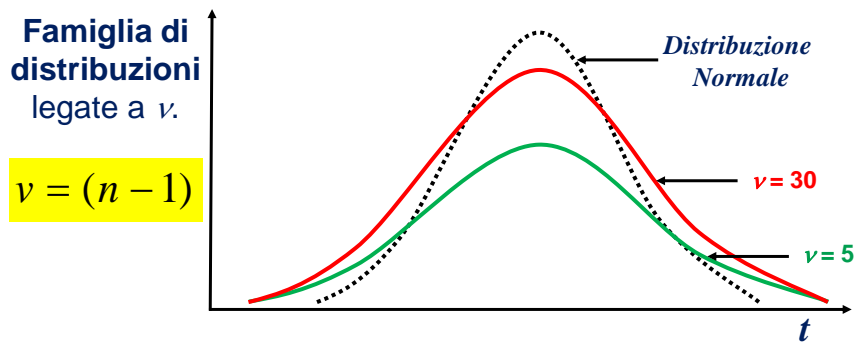
ALTRI PUNTEGGI **STANDARDIZZATI**

Nome	Caratteristica	Formula per il calcolo a partire dai punti z
QI di Wechsler	Media = 100 Deviazione standard = 15	$100 + z \times 15$
QI Stanford-Binet	Media = 100 Deviazione standard = 16	$100 + z \times 16$
Stanine (<i>STANDARD NINE</i>)	Media = 5 Deviazione standard = 2	$5 + z \times 2$
Sten (<i>Standard TEN</i>)	Media = 5,5 Deviazione standard = 2	$5,5 + z \times 2$
Scaled Scores	Media = 10 Deviazione standard = 3	$10 + z \times 3$

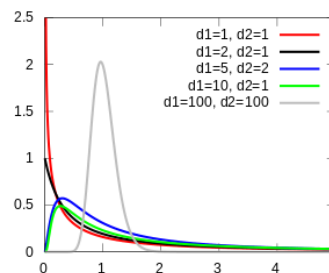
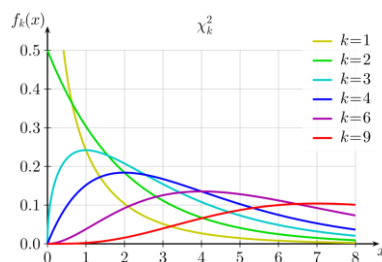
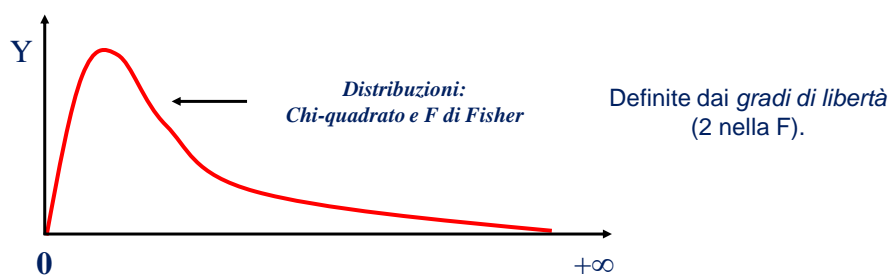
t di Student

La forma della distribuzione t varia secondo la numerosità delle osservazioni (n). Ciascuna distribuzione t è definita da tre parametri: μ , σ e ν ($ni - \text{gradi di libertà}$).

PROPRIETÀ:
 • INFINITA
 • SIMMETRICA
 • UNIMODALE
 • ASINTOTICA



F FISHER E CHI-QUADRATO



DISTRIBUZIONI TEORICHE

Mediante le **distribuzioni teoriche** è possibile confrontare e interpretare qualsiasi valore osservato. In questo modo non soltanto è possibile confrontare un valore (campione) con la popolazione, ma anche le popolazioni tra loro.

LA DISTRIBUZIONE CAMPIONARIA DELLA MEDIA (DCM)

La **DISTRIBUZIONE CAMPIONARIA DELLA MEDIA (DCM)** è una distribuzione teorica di primaria importanza nella metodologia della ricerca psicologica.

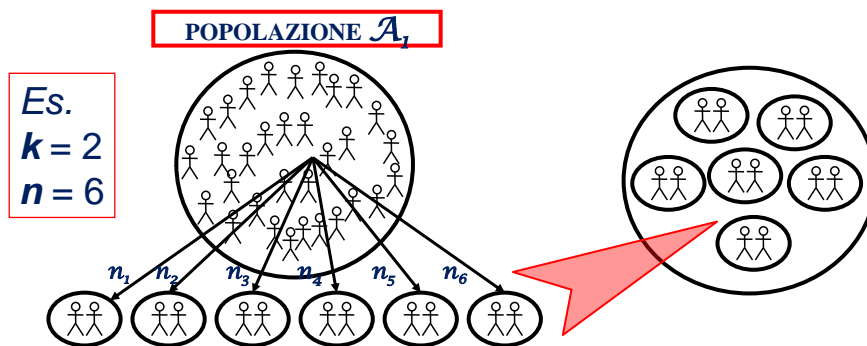
Consente di svolgere due operazioni importanti:

descrivere probabilisticamente le caratteristiche di un campione;

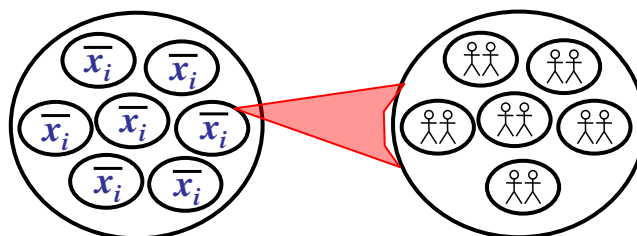
stimare la probabilità associata ai parametri della popolazione di provenienza.

LA DISTRIBUZIONE CAMPIONARIA DELLA MEDIA (DCM)

La DCM è una **distribuzione di valori medi**. Può essere definita come la distribuzione di tutte le medie ottenibili da n campioni casuali di ampiezza k , estratti a partire da una popolazione.



LA DISTRIBUZIONE CAMPIONARIA DELLA MEDIA (DCM)

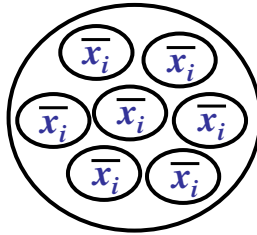


media

$$\mu_{DCM} = \frac{\sum_{i=1}^N \bar{x}_i}{N}$$

\bar{x}_i = media del campione i -esimo
 N = numerosità totale dei campioni

LA DISTRIBUZIONE CAMPIONARIA DELLA MEDIA (DCM)



In base al **teorema del limite centrale** sappiamo che tale distribuzione è approssimata alla **normale**.

LA DISTRIBUZIONE CAMPIONARIA DELLA MEDIA (DCM)

Se la popolazione è **infinita** o se il campionamento è **con reinserimento**:

la **media** della distribuzione campionaria è **uguale** alla media della popolazione.

$$\mu_{DCM} = \mu$$

l'**errore standard** è uguale alla deviazione standard della popolazione fratto la radice di n :

$$\sigma_{DCM} = \frac{\sigma}{\sqrt{n}}$$

LA DISTRIBUZIONE CAMPIONARIA DELLA MEDIA (DCM)

Se la popolazione è **finita** (N) o se il campionamento è **senza reinserimento**:

la **media** della distribuzione campionaria è **uguale** alla media della popolazione:

$$\mu_{DCM} = \mu$$

l'**errore standard** è uguale alla deviazione standard della popolazione fratto la radice di n moltiplicato un fattore di correzione:

$$\sigma_{DCM} = \frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}$$

LA DISTRIBUZIONE CAMPIONARIA DELLA MEDIA (DCM)

Se **non è nota la varianza** (o la ds) della popolazione:

la **media** della distribuzione campionaria è **uguale** alla media della popolazione:

$$\mu_{DCM} = \mu$$

l'**errore standard** è stimabile a partire dalla **varianza del campione** e quindi:

$$\hat{s} = \sigma = \sqrt{\frac{s^2}{n-1}}$$

$$\sigma_{DCM} = \frac{\hat{s}}{\sqrt{n}} = \frac{\sqrt{\frac{s^2}{n-1}}}{\sqrt{n}}$$

LA DISTRIBUZIONE CAMPIONARIA DELLA MEDIA (DCM)

Usando la **DCM**, mediante la trasformazione in **punti z** è possibile stimare la probabilità di osservare una data media (\bar{x}_i) calcolata su di **un campione**.

$$z_{\bar{x}} = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

\bar{x} = media del campione
 μ = media della popolazione e
 σ = ds della popolazione e
 n = ampiezza del campione

LA DISTRIBUZIONE CAMPIONARIA DELLA MEDIA (DCM)

La trasformazione in **punti z** consente di esprimere lo scarto tra il **VALORE ATTESO** (media della popolazione) e il **VALORE OSSERVATO** (media del campione) in una **UNITÀ DI MISURA STANDARD**. In questo modo è possibile conoscere e interpretare la **probabilità** di accadimento del valore osservato.

LA DISTRIBUZIONE CAMPIONARIA DELLA MEDIA (DCM)

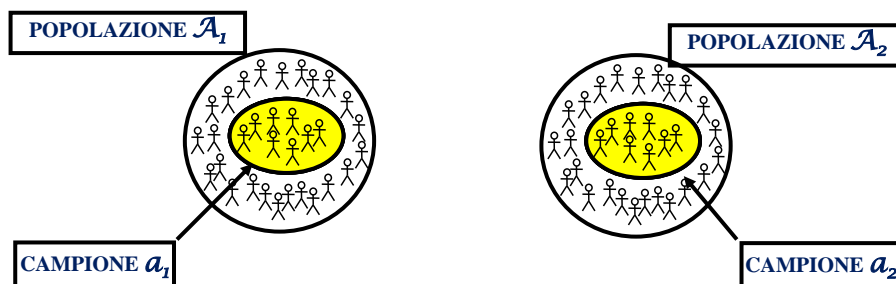
Quando il campione è **piccolo** ($n < 30$) la distribuzione delle medie dei campioni che si usa non è quella dei punti z ma è quella ***t di Student***.

$$t = \frac{\bar{x} - \mu}{s} \sqrt{n - 1}$$

\bar{x} = media del campione
 μ = media della popolazione e
 s = ds del campione
 n = ampiezza del campione

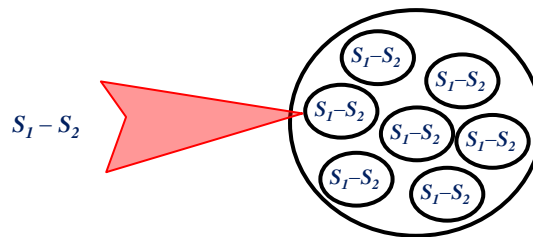
LA DISTRIBUZIONE CAMPIONARIA DELLA DIFFERENZA TRA MEDIE

Supponiamo che siano date **due popolazioni** (N_1 e N_2) dalle quali vengano estratti due sottocampioni (n_1 e n_2).



LA DISTRIBUZIONE CAMPIONARIA DELLA DIFFERENZA TRA MEDIE

Se per ciascun campione estratto calcoliamo la statistica S (S_1 e S_2) e calcoliamo lo scarto tra S_1 e S_2 . Otteniamo una distribuzione delle differenze che è detta **distribuzione delle differenze delle statistiche campionarie**.



LA DISTRIBUZIONE CAMPIONARIA DELLA DIFFERENZA TRA MEDIE

Se S_1 e S_2 sono le **medie campionarie** delle due popolazioni allora la distribuzione delle differenze avrà come **media**:

$$\mu_{\bar{X}_1 - \bar{X}_2} = \mu_{\bar{X}_1} - \mu_{\bar{X}_2} = \mu_1 - \mu_2$$

mentre come **deviazione standard**:

$$\sigma_{\bar{X}_1 - \bar{X}_2} = \sqrt{\sigma_{\bar{X}_1}^2 + \sigma_{\bar{X}_2}^2} = \sqrt{\frac{\sigma_1^2}{N_1} + \frac{\sigma_2^2}{N_2}}$$



TECNICHE DI ANALISI DEI DATI

AA 2017/2018

PROF. V.P. SENESE

Questi materiali sono disponibili per tutti gli studenti al seguente indirizzo:

<https://goo.gl/hxL9zG>

Seconda Università di Napoli (SUN) – Dipartimento di Psicologia – TECNICHE DI ANALISI DEI DATI – © Prof. V.P. Senese

INFERENZA STATISTICA

Sul campione si calcolano le **STATISTICHE** (ad es., \bar{x} , s , ...) per conoscere i **PARAMETRI** (ad es., μ , σ , ...) della popolazione.

Questo passaggio dalle **statistiche** ai **parametri** si basa sulla conoscenza delle proprietà delle **DISTRIBUZIONI CAMPIONARIE** dei parametri e viene detto **INFERENZA STATISTICA**.

INFERENZA STATISTICA

Insieme dei **metodi** basati sulla **teoria della probabilità** che consentono di formulare delle conclusioni sulla **variabile casuale** associata ad una certa caratteristica di una popolazione, riferendosi all'osservazione di un **campione** di osservazioni.

Le conclusioni a cui si può giungere si distinguono in due classi:

- ① **verifica di ipotesi**
- ② **stima dei parametri**

LA VERIFICA DI IPOTESI

Nella ricerca, una volta definito un **problema** di interesse si definisce l'**ipotesi generale**. Essa rappresenta un "grappolo di problemi" che non sono mai direttamente osservabili a meno che non vengano concretizzati in un'**ipotesi scientificamente osservabile**: le **ipotesi di ricerca**.

Un **ipotesi di ricerca** è un'affermazione su uno o più **parametri** relativi alla distribuzione di probabilità di una popolazione.

Pertanto la verifica delle ipotesi si basa sulle **distribuzioni teoriche** di volta in volta chiamate in causa.

LA VERIFICA DI IPOTESI

Nelle ipotesi di ricerca si suppone una relazione tra una variabile **Y** con una o più variabili **X**.

L'ipotesi di ricerca tipicamente è formulata come:

- se **X**, allora ne consegue che **Y**
- se **X**, allora probabilmente **Y**

Generalmente **X** è una *variabile indipendente* mentre **Y** è una *variabile dipendente*.

LA VERIFICA DI IPOTESI

L'ipotesi di ricerca:

- (1) deve essere formulata in termini operativi
- (2) deve considerare tutti i possibili risultati
- (3) deve essere empiricamente verificabile

Un'ipotesi si ritiene **verificata** quando è stata sottoposta alla **prova dei fatti**. Al termine della verifica l'ipotesi può essere:

- (a) **confermata**
- (b) **respinta**

LA VERIFICA DI IPOTESI

La prima tappa del processo di verifica consiste nel definire l'**ipotesi nulla** (H_0). L'**ipotesi nulla** è una ipotesi di **non effetto** e dovrebbe rappresentare la negazione dell'ipotesi che si vuole verificare.

Quando si respinge l'ipotesi nulla si supporta l'**ipotesi alternativa** (H_1). L'**ipotesi alternativa** è l'espressione dell'**ipotesi di ricerca** da cui parte lo sperimentatore

ESEMPIO

Supponiamo che una teoria socio-psicologica porti a prevedere che due gruppi di persone (es., uomini e donne) differiscono tra loro per il tempo dedicato alla cura dei figli.

Questa supposizione potrebbe trasformarsi in un'**ipotesi di ricerca**. Più specificamente l'ipotesi di ricerca potrebbe essere: **se** è vero che le due popolazioni (di cui i due gruppi sono rappresentativi) sono diverse, **allora** dovrebbero differire nel **tempo medio giornaliero dedicato a giocare con i figli** (μ_G).

Le ipotesi allora potrebbero essere:

$$H_0 \Rightarrow \mu_{G_1} = \mu_{G_2}$$

$$H_1 \Rightarrow \mu_{G_1} \neq \mu_{G_2}$$

$$H_1 \Rightarrow \mu_{G_1} > \mu_{G_2} \Rightarrow \mu_{G_1} < \mu_{G_2}$$

LA VERIFICA DI IPOTESI

L'ipotesi alternativa (H_1) o sperimentale si definisce:

SEMPLICE

si fissa un unico valore
del parametro

$$H_1 \Rightarrow \mu_{G_1} = 50$$

COMPOSTA

- BIDIREZIONALE
- MONODIREZIONALE

si fissano diversi valori
del parametro

$$H_1 \Rightarrow \mu_{G_1} \neq \mu_{G_2}$$

$$H_1 \Rightarrow \mu_{G_1} > \mu_{G_2} \Rightarrow \mu_{G_1} < \mu_{G_2}$$

LA VERIFICA DI IPOTESI

Accettare l'ipotesi nulla o l'ipotesi alternativa non è mai una decisione assoluta. La decisione è **sempre soggetta ad errore**.

Come vedremo i dati permettono solo di fare **affermazioni probabilistiche** riguardo alle ipotesi (basate sulle distribuzioni di probabilità di volta in volta utilizzate).

L' **ipotesi nulla** (H_0) serve a specificare la distribuzione campionaria sulla quale si basa la verifica.

LA VERIFICA DI IPOTESI

Dopo aver formulato le **ipotesi di ricerca** la fase successiva è relativa alla **scelta del test statistico** da utilizzare per la verifica delle ipotesi.

Per operare una giusta scelta è necessario conoscere i “**principi**” e le **proprietà** dei diversi test statistici utili alla loro applicabilità.

LA VERIFICA DI IPOTESI

Definito il test appropriato il passo successivo è la scelta del **criterio di significatività** (α e β) e della **dimensione campionaria** (N).

Il **livello di significatività** rappresenta il **criterio** in base al quale decido se **confermare** l'ipotesi nulla (H_0) o **rifiutarla** in favore dell'**ipotesi alternativa** (H_1).

ERRORE DI I° TIPO

Si definisce **errore di I° tipo** la probabilità di **rifiutare** l'ipotesi nulla (H_0) quando è vera. La probabilità di questo errore è simboleggiata con la lettera greca α .

Nella verifica delle ipotesi i valori di α e N si fissano anticipatamente.

ERRORE DI I° TIPO

La scelta del livello di significatività α implica la definizione del **grado di sicurezza** con il quale si vuole prendere la decisione.

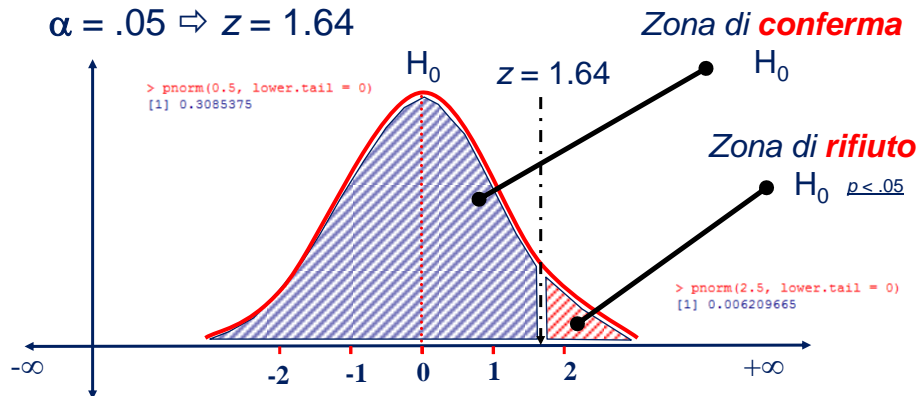
Per **convenzione** si adottano come livelli critici per la decisione:

- $\alpha = .05$ \Rightarrow rischio di **sbagliare** rifiutando H_0 quando essa è vera **5** volte su **100**
- $\alpha = .01$ \Rightarrow rischio di **sbagliare** rifiutando H_0 quando essa è vera **1** volta su **100**
- $\alpha = .001$ \Rightarrow rischio di **sbagliare** rifiutando H_0 quando essa è vera **1** volta su **1000**

LA VERIFICA DI IPOTESI

Supponiamo di avere stimato una **statistica**, che i valori di questa statistica siano regolati dalla **distribuzione normale**, di avere un'ipotesi **monodirezionale** e di aver scelto come livello critico $\alpha = .05$.

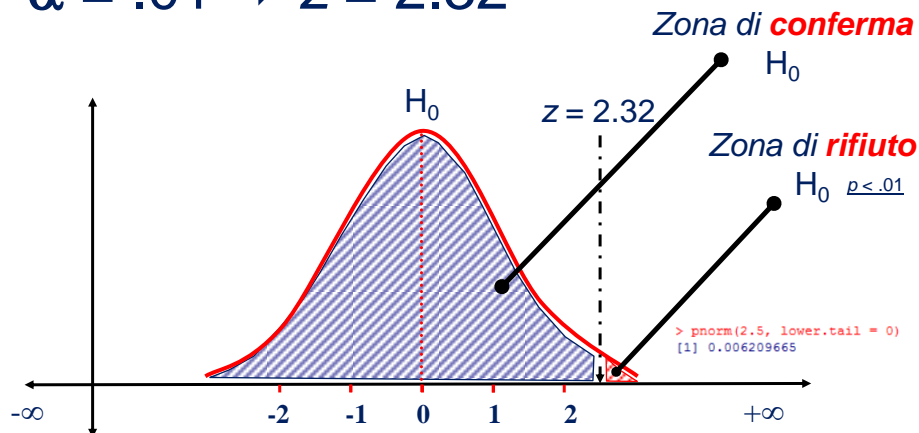
Trasformando la statistica in **punti z** è possibile conoscere il livello di probabilità associato con il valore osservato e confrontarlo con il livello critico:



LA VERIFICA DI IPOTESI

Ipotesi monodirezionale, $\alpha = .01$

$\alpha = .01 \Rightarrow z = 2.32$

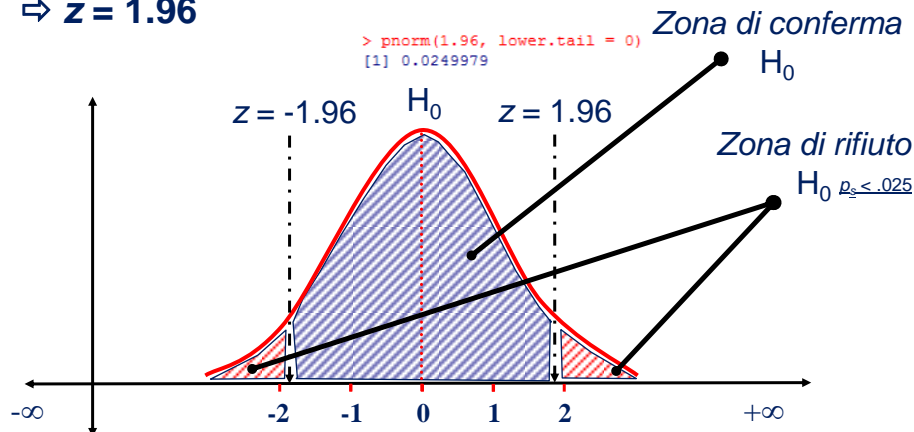


LA VERIFICA DI IPOTESI

Ipotesi bidirezionale, $\alpha = .05$

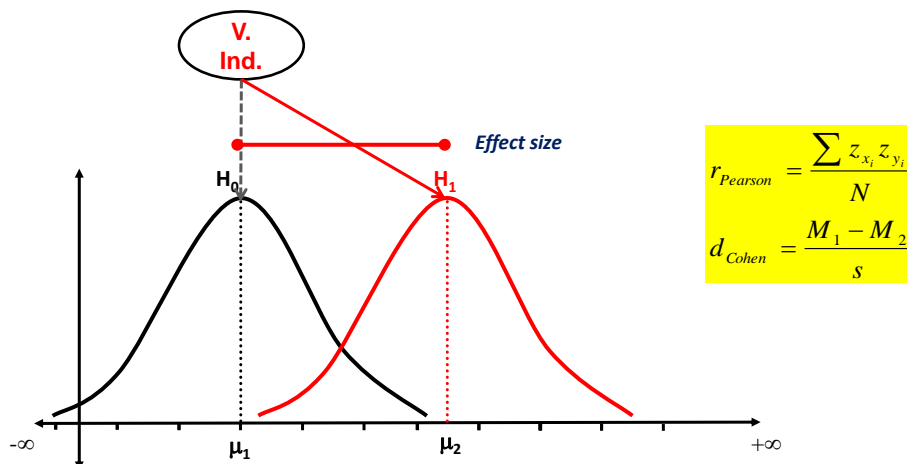
$\alpha = .05 \Rightarrow \alpha = .05/2 = .025$

$\Rightarrow z = 1.96$



EFFECT SIZE

Si definisce **grandezza dell'effetto** (*effect size*) la forza della relazione tra due variabili.



ERRORE DI II° TIPO E POTENZA

Si definisce **errore di II° tipo** la probabilità di **accettare l'ipotesi nulla** (H_0) quando è **falsa**. La probabilità di questo errore è simboleggiata con la lettera greca β .

Nella verifica delle ipotesi il valore di β **non viene scelto arbitrariamente** ma **si calcola** in base ai parametri α , n e alla **grandezza dell'effetto**.

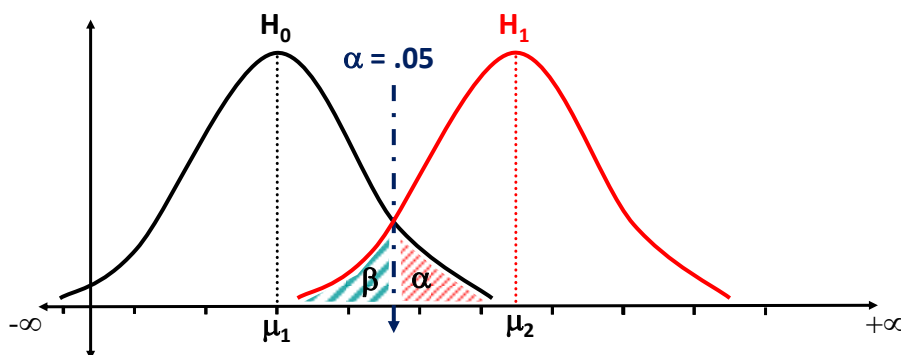
La **potenza statistica** (*power*) è collegata all'errore di **II° tipo**: è la probabilità di prendere la decisione giusta. Vale a dire:

rifiutare l'ipotesi nulla (H_0) quando è falsa

ERRORE DI I° e II° TIPO

In ogni **inferenza statistica** esiste il rischio di commettere uno dei due tipi di **errori alternativi**. Se α diminuisce β aumenta.

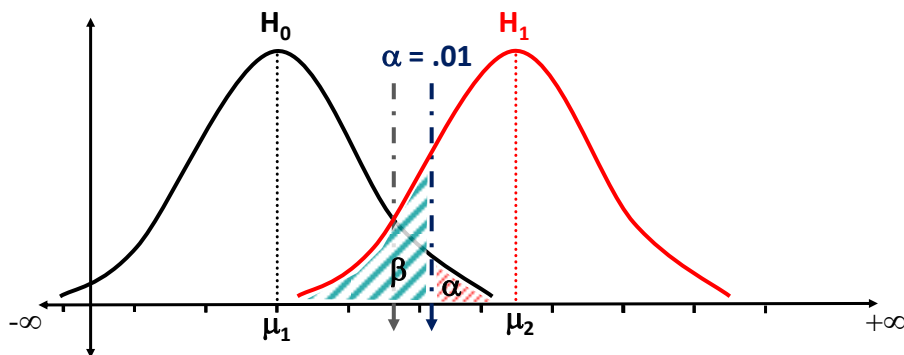
Evitare errori di **I° tipo** può portare ad una elevata probabilità di commettere errori di **II° tipo**.



ERRORE DI I° e II° TIPO

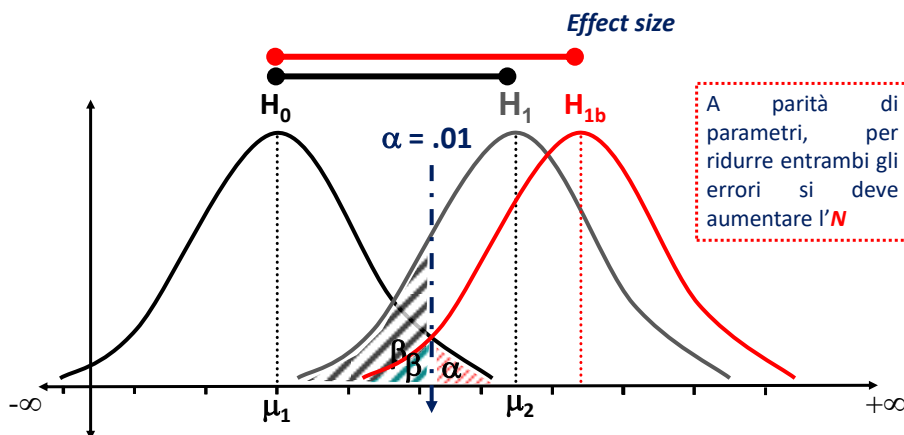
In ogni **inferenza statistica** esiste il rischio di commettere uno dei due tipi di **errori alternativi**. Se α diminuisce β aumenta.

Evitare errori di **I° tipo** può portare ad una elevata probabilità di commettere errori di **II° tipo**.



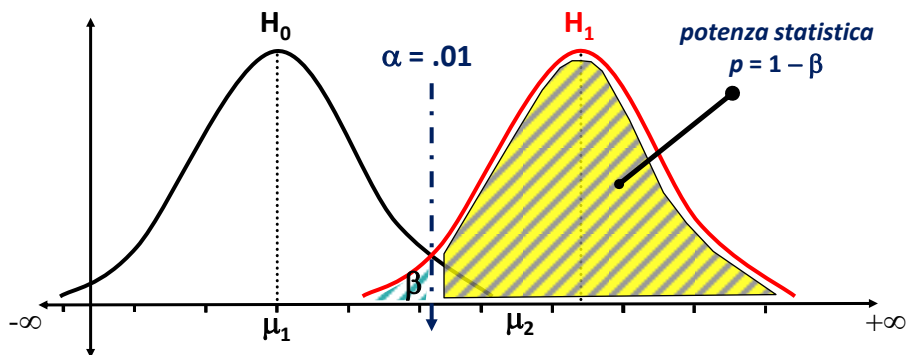
ERRORE DI I° e II° TIPO

Anche l'**effect size** influenza l'errore di **II tipo**. **Maggiore** è il suo valore **minore** è la probabilità di incorrere in un errore di **II tipo**.

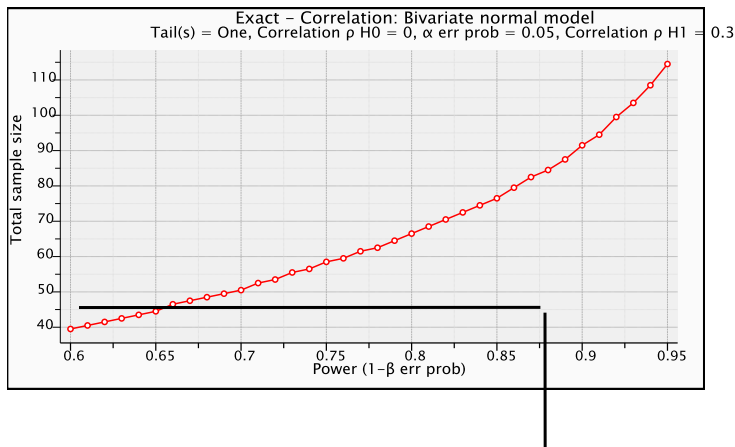


ERRORE DI I° e II° TIPO

Per poter interpretare correttamente i risultati di un'analisi, un indice necessario è la **potenza statistica** che esprime la capacità di trovare un effetto quando quest'effetto esiste realmente. La potenza dovrebbe essere **almeno uguale a .80**.



Es. POTENZA



ERRORE DI I° e II° TIPO

		LA REALTÀ NELLA POPOLAZIONE	
		H ₀ <u>VERA</u>	H ₀ <u>FALSA</u>
LA DECISIONE STATISTICA	H ₀ <u>CONFERMATA</u> H ₁ RIFIUTATA	Decisione corretta nessun errore Prob.: $1 - \alpha$	Decisione errata errore di II° tipo Prob.: β
	H ₀ <u>RIFIUTATA</u> H ₁ <u>CONFERMATA</u>	Decisione errata errore di I° tipo Prob.: α	Decisione corretta nessun errore Prob.: $1 - \beta$

STIMA DEI PARAMETRI

Si definisce **corretta** la statistica il cui valore stimato sul campione corrisponde al rispettivo parametro della popolazione. In caso contrario la statistica si definisce **stimatore distorto**.

Es.1 La media campionaria è uno **stimatore corretto** della media della popolazione, essendo:

$$E(X) = \mu$$

Es.2 La deviazione standard campionaria è uno **stimatore distorto** della ds della popolazione, essendo:

$$E(s) \neq \sigma$$

STIMA DEI PARAMETRI

Una stima è tanto più **efficiente** quanto **minore** è la **varianza** (o la **ds**) della propria distribuzione campionaria.

Una stima è detta **PUNTUALE** quando la stima del parametro è data da **un valore unico**.

$$\mu = 20$$

Una stima è detta **INTERVALLARE** quando viene **definito un intervallo** entro il quale è compreso il valore del parametro reale.

$$18 \leq \mu \leq 22$$

STIMA DEI PARAMETRI

Nella stima intervallare è **possibile definire la precisione della stima**. Questo le rende preferibili alle stime puntuali.

Es.1 se dopo aver somministrato un **test di intelligenza** ($\mu = 100$; $\sigma = 15$) affermiamo che il soggetto i ha un **QI** pari a **110** produciamo una **stima puntuale**.

Es.2 se affermiamo che allo **stesso test** il medesimo soggetto i ha un **QI** compreso tra **95** e **125**, produciamo una **stima intervallare**, che (*in base alla distribuzione normale*) è vera al **68.3%** ($x_i \pm 1\sigma$).

$$IC[68.3\%] = x_i \pm 1\sigma = 110 \pm 1 \cdot 15$$