

DEFINIZIONE TIPOLOGIA DI STUDIO

- indagine GLOBALE o PARZIALE o CARATTERISTICA
- studi TRASVERSALI o LONGITUDINALI } 1 UNITA' STATISTICA
- " OCCASIONALI o RIPETUTI } + UNITA' STATISTICHE
- " OSSERVAZIONALI o SPERIMENTALI } MODALITA' DI INTERVENIRE
- " DESCRITTIVI o ANALITICI } FINALITA' DI RICERCA

DEFINIZIONE DELLE VARIABILI

- | | |
|-----------------------------------|-------------------------------------|
| QUALITATIVE
CATEGORIALI | QUANTITATIVE
NUMERICHE |
| • ORDINABILI | • CONTINUE <small>QUALSIASI</small> |
| • NOMINALI | • DISCRETE <small>NATURALI</small> |

TIPOLOGIE GRAFICI

- | | |
|---------------------|---|
| DIAGRAMMA A BARRE | VARIABILI QUALITATIVE
O QUANTITATIVE DISCRETE
O PER VARIABILI MISURATE
CON SCALA NOMINALE O ORDINALE |
| DIAGRAMMA A TORTA | |
| DIAGRAMMA DI PARETO | |
| ISTOGRAMMA | VARIABILI QUANTITATIVE
CONTINUE |
| POLIGONO | |
| POLIGONO CUMULATIVO | |

ANALISI DESCRITTIVA

- tecniche: studio:
- FREQUENZE → FORMA delle DISTRIBUZIONI
 - VALORI MEDI → TENDENZA CENTRALE ...
 - MISURE DI VARIABILITA' → VARIABILITA' ...

VARIABILE QUALITATIVA NOMINALE → MODA

- esempio tipo di dissetto
SCALA NOMINALE {I, FC, IC...}
- TIPO DI GRAFICO: DIAGRAMMA A BARRE
- forma DISTRIBUZIONE:
- TENDENZIALMENTE UNIFORME
 - TENDENZIALMENTE NON UNIFORME
 - UNIMODALE
 - PURIMODALE

VARIABILE QUALITATIVA ORDINATA → MODA

- esempio grado di urbanizzazione in 50 comuni
SCALA ORDINALE {BASSO, MEDIO/BASSO, MEDIO...}
- TIPO DI GRAFICO: DIAGRAMMA A BARRE
- forma DISTRIBUZIONE:
- TENDENZIALMENTE UNIFORME → QUANTILI
 - TENDENZIALMENTE NON UNIFORME
 - UNIMODALE TENDENZIALMENTE
 - PURIMODALE
 - FORMA A J
 - FORMA A U
 - NEGATIVA → MEDIANA
 - SIMMETRICA → MEDIANA
 - POSITIVA → MEDIANA

VARIABILE QUANTITATIVA DISCRETA → MODA

- esempio n° difetti in 50 barre di acciaio
SCALA NUMERICA: {0, 1, 2, 3, ...}
- TIPO DI GRAFICO: DIAGRAMMA A BARRE
DIAGRAMMA A BASTONCINI
DOT PLOT NUMEROSITA' < 30
- forma DISTRIBUZIONE:
- TENDENZIALMENTE UNIFORME → QUANTILI
 - TENDENZIALMENTE NON UNIFORME
 - UNIMODALE TENDENZIALMENTE
 - PURIMODALE
 - FORMA A J
 - FORMA A U
 - NEGATIVA → MEDIANA
 - SIMMETRICA → MEDIA
 - POSITIVA → MEDIANA

VARIABILE QUANTITATIVA CONTINUA → MODA

- esempio resistenza a compressione in MPa
SCALA NUMERICA CONTINUA
- TIPO DI GRAFICO: ISTOGRAMMA
POLIGONO
- forma DISTRIBUZIONE:
- TENDENZIALMENTE UNIFORME → QUANTILI
 - TENDENZIALMENTE NON UNIFORME
 - UNIMODALE TENDENZIALMENTE
 - PURIMODALE
 - FORMA A J
 - FORMA A U
 - NEGATIVA → MEDIANA
 - SIMMETRICA → MEDIA
 - POSITIVA → MEDIANA

VALORI MEDI

- INTENSITA' MEDIA: MEDIA
- POSIZIONE CENTRALE: MEDIANA $\frac{n+1}{2}$
- POSIZIONI CENTRALI E NON: QUANTILI
- MODALITA' PIU' FREQUENTI: MODA
- SE PARI $\frac{n}{2}$ e $\frac{n+1}{2}$ E MEDIA DEI 2 VALORI!!!
- $1^o \rightarrow \frac{(n+1)}{4}$
- $2^o \rightarrow \frac{(n+1)}{2}$
- $3^o \rightarrow \frac{3(n+1)}{4}$

MISURE DI VARIABILITA'

- RANGE
 - RANGE INTERQUARTILE
 - VARIANZA
 - DEVIATIONE STANDARD
 - COEFFICIENTE DI VARIAZIONE
 - FREQUENZA RELATIVA DELLA MODA
- per variabili QUANTITATIVE:
- CONFRONTO TRA MEDIA E MEDIANA
 - BOXPLOT

MISURA VARIABILITA'

- UNIFORMI: RANGE, IQR
- SIMMETRICHE: DS, RANGE, IQR
- ASIMMETRICHE: RANGE, IQR
- PURIMODALI: RANGE, IQR

VARIANZA:
$$\frac{\sum_{j=1}^n (x_j - \bar{x})^2}{n}$$

$$\sigma^2 = \frac{1}{N} \sum_{j=1}^N (x_j - \bar{x})^2$$

$$S^2 = \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2$$

DEVIAZIONE STANDARD:

$$\sigma = \sqrt{\frac{1}{N} \sum_{j=1}^N (x_j - \bar{x})^2}$$

$$S = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2}$$

↳ distanza media dei valori dalla loro media

$$S = \sqrt{\frac{\sum_{i=1}^k (x_i - \bar{x})^2 \times n_i}{n-1}}$$
 NON si dividono in classi

COEFFICIENTE di VARIAZIONE:

indipendente dall'ordine di grandezza della variabile

esempio trovare quale di due prodotti con vendite diiferenti presenta un prezzo maggiormente variabile:

A) media = 50 euro dev. stand. = 5 euro $CV_A\% = \left(\frac{5}{50}\right) \times 100 = 10\%$

B) media = 100 \$ dev. stand. = 5 \$ $CV_B\% = \left(\frac{5}{100}\right) \times 100 = 5\%$

PREZZO B MELO VARIABILE dell'A

$CV\% = \left(\frac{5}{100}\right) \times 100$

SCALA NOMINALE → freq. moda

ORDINALE → freq. moda, range, range inter

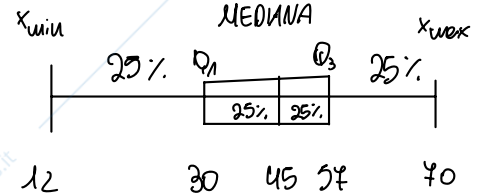
NUMERICA → freq. moda, range, range inter deviazione standard

BOX PLOT

- X_{min} (minimo)
- PRIMO QUARTILE (Q_1)
- MEDIANA
- TERZO QUARTILE (Q_3)
- X_{max} (massimo)

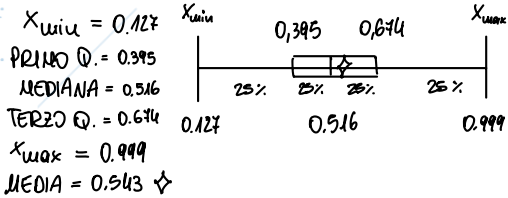
esempio

$X_{min} = 12$
PRIMO = 30
MED = 45
TERZO = 57
 $X_{max} = 70$



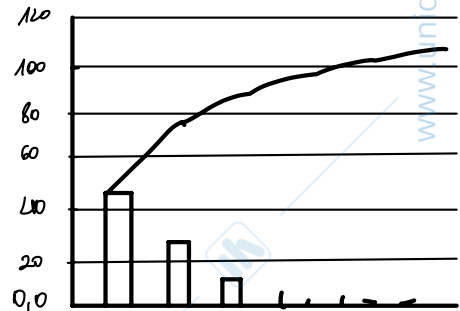
ESERCITAZIONE es 3:

CLASSE IV_0-10



es 5:

	freq. ASS	freq. rel. %	freq. % cum.
60:69 CIVILE	87	$\frac{87}{903} \times 100 = 42.9$	42.9
ISTRUZIONE SANITA'	58	$\frac{58}{208} \times 100 = 28.6$	71.5
REUBIODO	21	$\frac{21}{203} \times 100 = 10.3$	81.9
MILITARE	10	4.9	86.8
NON IDENT.	10	4.9	91.7
MOB e TRASP.	5	2.5	94.2
TECNOLOGICO	1	2.5	100.0
TOTALE	203	100	



FORMA	VALORE MEDIO	MISURA VARIABILITA'
UNIFORME	QUANTILI	RANGE RANGE IQR
SIMMETRICHE	MEDIA	D.S RANGE, IQR
ASIMMETRICHE	MODA, MEDIANA QUANTILI	RANGE, IQR
PURIMODALI	QUANTILI	RANGE, IQR

LA VARIABILITÀ, INCERTEZZA E PROBABILITÀ

PROVA, EVENTO, PROBABILITÀ

3 APPROCCI: IMPOSTAZIONE CLASSICA, FREQUENTISTA, SOGGETTIVISTA

- esperimenti con:
- 1) incertezza multipla
 - 2) ripetibili nelle stesse condizioni
 - 3) con risultati equiprobabili
- con:
- 1)
 - 2)
- eventi unici senza esperienza consolidata

E = EVENTO
 $f_n(E)$ = numero di volte che E si è verificato
 Allora $P(E)$ = PROBABILITÀ

tra 0 e 1 compresi
 0 → evento impossibile
 1 → evento certo

VARIABILI CASUALI

MISURE	PARAMETRI DI POPOLAZIONE	STATISTICHE CAMPIONARIE
MEDIA	μ	\bar{x}
VARIANZA	σ^2	S^2
DEVIAZIONE STANDARD	σ	S

MODELLI PER VARIABILI CASUALI CONTINUE

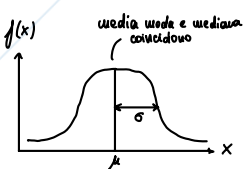
se la variabile può assumere tutti i valori compresi in un intervallo reale
 $f(x) = Pr(X \in A)$

DISTRIBUZIONE NORMALE O GAUSSIANA

diversi fenomeni continui sembrano seguire una distribuzione normale; può approssimare numerose distribuzioni discrete; teorema del limite centrale.

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

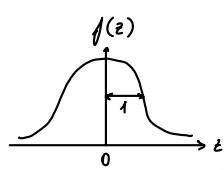
- e = 2.718
- π = 3.14
- μ = media della pp.
- σ = dev. stand.
- x = valori assunti dalla variabile X



DISTRIBUZIONE NORMALE STANDARDIZZATA

consente di ottenere una variabile casuale normale standardizzata con valore atteso pari a 0 e deviazione standard pari a 1. La standardizzazione è data dalla trasformazione di X in Z:

$$Z = \frac{X - \mu}{\sigma}$$

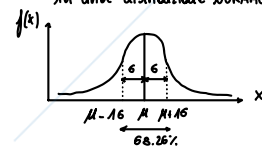


+ TAVOLE DISTRIBUZIONE NORMALE STANDARDIZZATA

- riga: parte intera di Z
- colonna: seconda cifra decimale di Z
- incrocio: PROBABILITÀ di trovare un valore $\leq z$

se l'incrocio dei possibili valori (o supposto) è un insieme discreto di numeri reali

in una distribuzione NORMALE:



- $\mu \pm 1\sigma$ include circa il 68.26% dei valori
- $\mu \pm 2\sigma$ include circa il 95.4% dei valori
- $\mu \pm 3\sigma$ include circa il 99.7% dei valori

tre discordanze REEDEL EMPIRICHE: (n=30)
 $\bar{x} \pm 1s \dots 68\% \dots$
 $\bar{x} \pm 2s \dots 95.4\% \dots$
 $\bar{x} \pm 3s \dots 99.7\% \dots$

MODELLI PER VARIABILI CASUALI DISCRETE

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

MODELLO di BERNOULLI

Una prova può dar luogo a due risultati: successo di E e INSUCCESSO di E

$x = 0, 1$

$Pr(x=1) = p$
 $Pr(x=0) = (1-p)$ PROBABILITÀ

$X \sim Ber(p)$

RISULTATI DICOTOMICI (binomio)

MODELLO BINOMIALE

esperimento casuale ripetuto n volte nelle medesime condizioni, in cui in ogni singola prova può dar luogo a due possibili risultati in un dato esito: evento verificato (successo) ed evento non verificato (insuccesso)

$(x = 0, 1, \dots, n)$
 $Pr(E) = p$ data p costante

$X \sim Bin(n, p)$

MODELLO di POISSON

in un certo ambito circoscritto, un evento (E) successo può verificarsi più volte e in modo indipendente

$x = 0, 1, 2, \dots$ → supposto discreto infinito numerabile

$X \sim Po(\lambda)$

$\lambda = E(x)$ = n° medio di successi nell'ambito considerato

MODELLO GEOMETRICO

esperimento ipotetico: successive di sottoprove bernoulliane indipendenti. ogni sottoprova: SUCCESSO (E) e INSUCCESSO (E)

$x = 1, 2, \dots$ → supposto discreto infinito numerabile

$X \sim Geo(p)$

$Pr(E) = p$
 $E(x) = 1/p$ = n° medio di sottoprove fino al successo

MODELLI PER VARIABILI CASUALI CONTINUE

DISTRIBUZIONE ESPOENZIALE

introdotta per fenomeni di cui interessa la DURATA di vita, di resistenza...

$\lambda (x \sim E(\lambda))$

valore medio $1/\lambda$ è la DURATA media del fenomeno riferimento agli esperimenti di Poisson

$S(x) = 1 - F(x) = Pr(X > x)$

DISTRIBUZIONE di WEIBULL

Esponenziale generalizzata con $\lambda > 0$ e $\beta > 0$ introdotta per fenomeni di cui interessa la DURATA di vita, di resistenza...

- $\beta = 1$, funzione di rischio costante (Esp.)
- $\beta > 1$, " " " " è funzione crescente del tempo
- $\beta < 1$, " " " " è funzione decrescente del tempo

DISTRIBUZIONE GENERALIZZATA DEI VALORI ESTREMI

usata come approssimazione per modellare i massimi o i minimi di lunghe (finite) sequenze di variabili aleatorie

DISTRIBUZIONE LOGNORMALE

$Y \sim N(\mu, \sigma)$ dove N NORMALE
 $Y = \ln X$ v.c. Lognormale X

Viene usata al posto della NORMALE quando il fenomeno di interesse non può assumere valori negativi, ma è UNIMODALE e con forma approssimativamente campanulare e simmetrica

Esponenziale di Weibull GEV
 Variabili casuali connesse alla Normale:
 Lognormale
 Chi-quadrato di Fisher di Student

di Bernoulli
 Binomiale
 Geometrica di Poisson

INTRODUZIONE ALLE TECNICHE DI INFERENZA STATISTICA

GENERALIZZARE LE INFORMAZIONI TRATTE DAL CAMPIONE ALLA TOTALITÀ DELLA POPOLAZIONE

UNICITÀ = $2z_{\alpha/2} (\sigma/\sqrt{n})$
 $z_{\alpha/2}$ = VALORE CRITICO
 probabilità che μ si trovi in uno dei due estremi dell'intervallo
 (50%) probabilità di confidenza

- SELEZIONARE CAMPIONE
 - VALUTARE GRADO DI ATTENDIBILITÀ
 - INFERENZA STATISTICA
 - ESTENDERE RISULTATI ALLA POPOLAZIONE
 - CONTROLLARE GRADO DI INCERTEZZA
- MISURE INTRODOTTE AL FINE DELL'ANALISI DESCRITTIVA DIVENTANO OGGETTO DI INFERENZA

STIMA } PUNTUALE } QUANTITÀ INDICATA AD UN UNICO VALORE. È UNA SINGOLA STATISTICA. MEDIA CAMPIONARIA È STIMATORE PUNTUALE DELLA MEDIA DELLA POPOLAZIONE μ
 PER INTERVALLO } QUANTITÀ INDICATA AD UN INSIEME DI VALORI COSTITUENDO UN INTERVALLO ENTRO CUI CI POTREBBE ESSERE IL VALORE VERO DELLA QUANTITÀ INDICATA

INTERVALLI DI CONFIDENZA PER LA MEDIA:

- 1) VARIANZA NOTA } $P(\bar{X} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}) = 1 - \alpha$
- 2) VARIANZA NON NOTA } $(\bar{X} - t_{n-1, \alpha/2} \cdot \frac{s}{\sqrt{n}}, \bar{X} + t_{n-1, \alpha/2} \cdot \frac{s}{\sqrt{n}})$
- 3) POPOLAZIONE NON NORMALE } $(\bar{X} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}})$ $(\bar{X} - t_{n-1, \alpha/2} \cdot \frac{s}{\sqrt{n}}, \bar{X} + t_{n-1, \alpha/2} \cdot \frac{s}{\sqrt{n}})$

TEOREMA DEL LIMITE CENTRALE

PER LA MAGGIOR PARTE DELLE POPOLAZIONI, LA DISTRIBUZIONE DELLA MEDIA CAMPIONARIA È APPROSSIMATIVAMENTE NORMALE, PURTANTO SI CONSIDERANO CAMPIONI DI AMPIEZZA 30 OSSERVAZIONI

SE DISTRIBUZIONE POPOLAZIONE È ABBASTANZA SIMMETRICA, LA DISTRIBUZIONE DELLA MEDIA → DISTRIBUZIONE NORMALE PURTANTO CI BASTANO ANCHE 5-15 OSSERVAZIONI

SE POPOLAZIONE HA DISTRIBUZIONE NORMALE, LA MEDIA CAMPIONARIA È DISTRIBUITA SECONDO LA LEGGE NORMALE, INDIPENDENTEMENTE DA AMPIEZZA

$$\bar{X} \sim N(\mu, \sigma/\sqrt{n})$$

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0,1)$$

$$P(\bar{X} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}) = 1 - \alpha$$

VARIANZA σ NON NOTA: $t = \frac{\bar{X} - \mu}{s/\sqrt{n}}$

$$\bar{X} - t_{n-1, \alpha/2} \cdot \frac{s}{\sqrt{n}} \leq \mu \leq \bar{X} + t_{n-1, \alpha/2} \cdot \frac{s}{\sqrt{n}}$$

INTERVALLO DI CONFIDENZA PER LA PROPORZIONE: $p - z_{\alpha/2} \cdot \sqrt{p(1-p)/n} \leq \pi \leq p + z_{\alpha/2} \cdot \sqrt{p(1-p)/n}$

VERIFICA DI IPOTESI: TEST SU UN CAMPIONE

PROCEDURA INFERENZIALE, SI FA UN'IPOTESI SU QUANTITÀ INDICATA DELLA POPOLAZIONE (PARAMETRO) E DECIDE SE CONFERMA IL TEST O SE ESSA È ACCETTABILE O NO.

TEST DI IPOTESI } PROCESSO CONSUNTIVO O NON
 } IPOTESI DI RICERCA

PRESENZA DI DIFFERENZA PARAMETRO DA VALORE INTERESSE
 DIFFERENZE TRA PARAMETRI DI GRUPPI DIVERSI
 RELAZIONI TRA VARIABILI

VERIFICA DI IPOTESI PARAMETRICA COMPOSTA DA

- IPOTESI NULLA H_0 : ASSENZA DIFFERENZA DEL PARAMETRO DA VALORE ASSENZA DIFFERENZE TRA PARAMETRI O ASSENZA DI RELAZIONI TRA VARIABILI
- IPOTESI ALTERNATIVA H_1 : NUOVA IPOTESI DEL RISPETTIVO

1) FASI TEST STATISTICO: TEST Z PER LA MEDIA (σ NOTO)

1° PASSO: FORMULAZIONE IPOTESI

- $H_0: \mu = 368 \text{ gr}$
 - $H_1: \mu \neq 368 \text{ gr}$
- TEST DI IPOTESI BIDIREZIONALE O A DUE CODE

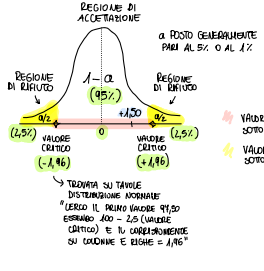
2° PASSO: SCELTA DI UNA STATISTICA TEST

STATISTICA CAMPIONARIA LA CUI DISTRIBUZIONE CAMPIONARIA DEVE ESSERE NOTA E COMPLETAMENTE SPECIFICATA SOTTO L'IPOTESI NULLA ($\mu = 6$)
 SE X È NORMALE POSSIAMO APPLICARE IL TEOREMA DEL LIMITE CENTRALE:

$$\bar{X} \sim N(\mu, \sigma/\sqrt{n}) \quad Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0,1)$$

3° PASSO: REGIONE DI ACCETTAZIONE E DI RIFIUTO

IPOTIZZIAMO VERA L'IPOTESI NULLA
 AD ESEMPIO: $\bar{X} \sim N(368, 15/\sqrt{25})$ e $Z = \frac{\bar{X} - 368}{15/\sqrt{25}}$



VALORI ALTAMENTE PROBABILI SOTTO L'IPOTESI NULLA
 VALORI POCO PROBABILI SOTTO L'IPOTESI NULLA

UN ALTRO MODO PER IDENTIFICARE IL TEST È QUELLO DI CALCOLARE IL P-VALUE:
 - SE IL P-VALUE È > α , IPOTESI NULLA NON È RIFIUTATA
 - SE IL P-VALUE È < α , IPOTESI NULLA È RIFIUTATA

NELL'ESEMPPIO LA MEDIA CAMPIONARIA DIVENTA 372,5

$$Z = \frac{372,5 - 368}{15/\sqrt{25}} = +1,50$$

CHE RISULTA ESSERE ALL'INTERNO DELLA REGIONE DI ACCETTAZIONE. IL TEST NON CI PERMETTE DI RIFIUTARE L'IPOTESI NULLA H_0 LA DIFFERENZA TRA LA MEDIA CAMPIONARIA OSSERVATA (372,5) E IL VALORE IPOTIZZATO REL LA MEDIA DELLA POPOLAZIONE (368) NON È STATISTICAMENTE SIGNIFICATIVA. IL PROCESSO È SOTTO CONTROLLO

ESEMPIO: SE I VALORI CRITICI = $\pm 1,50$ e $\pm 1,50$ GIUSTA LA PROBABILITÀ CHE Z ASSUMA UN VALORE MINORE DI $+1,50$ È DI $0,9332$. PER CUI LA PROBABILITÀ CHE ASSUMA UN VALORE SUPERIORE A $+1,50$ È: $1 - 0,9332 = 0,0668$. ALLO STESSO MODO LA PROBABILITÀ CHE ASSUMA UN VALORE MINORE A $-1,50$ È $0,0668$. PERTANTO IL P-VALUE ASSUME UN VALORE DI: $0,0668 + 0,0668 = 0,1336$ QUINDI $0,1336 > \alpha (0,05) \rightarrow$ NON RIFIUTIAMO H_0

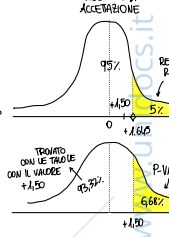
I TEST AD UNA CODA

IPOTESI NULLA $H_0: \mu = 368$
 IPOTESI ALTERNATIVA: $H_1: \mu > 368$

FISSANDO $\alpha = 0,05 (5\%)$ LA REGIONE CRITICA È CONTENUTA NELLA CODA DESTRA DELLA DISTRIBUZIONE NORMALE STANDARDIZZATA E CORRISPONDE AL 5%

IL VALORE CRITICO LASCIA ALLA SUA SINISTRA UNA MASSA DI PROBABILITÀ PARI A $0,95 (95\%)$ QUINDI, DALLE TABELLE IL VALORE CRITICO È $+1,645$ (QUANDO IL VALORE CHE NOI SI AMMOVA A OLTREDO IL VALORE CRITICO È LA MEDIA TRA DUE VALORI E PRELIEVO I CORRISPETTIVI VALORI DELLA CODA SINISTRA E DESTRA = $1,645$ O $0,05$)

SECONDO IL P-VALUE, DOBBIAMO CALCOLARE LA PROBABILITÀ CHE Z SIA $> +1,50$.
 P-VALUE = $1 - 0,9332 = 0,0668 (6,68\%)$
 $0,0668 > 0,05 \rightarrow$ L'IPOTESI NULLA NON PUÒ ESSERE RIFIUTATA



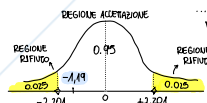
2) TEST DI IPOTESI T PER LA MEDIA (σ NON NOTO)

$$t = \frac{\bar{X} - \mu}{s/\sqrt{n}}$$

SE NORMALE O NON T-Student
 $n-1$ GRADI DI LIBERTÀ

ESEMPIO:
 $H_0: \mu = 120$
 $H_1: \mu \neq 120$
 $n = 12$
 $\bar{X} = 112,85 \text{ m}^2$
 $s = 20,80 \text{ m}^2$

TEST A DUE CODE BASTA DETERMINARE I VALORI CRITICI CON LE TABELLE DELLA DISTRIBUZIONE T DI STUDENT. GRADO DI LIBERTÀ = $n - 1$ AREA CODA DI DESTRA = $0,05$
 $\alpha/2 = 0,025 \rightarrow 0,015$
 VALORE CRITICO = $\pm 2,204$



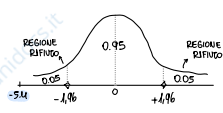
3) TEST DI IPOTESI Z PER LA PROPORZIONE

PROPORZIONE CAMPIONARIA X/n
 SE NUMERO SUCCESSI X E DI INSUCCESSI $(n-X)$ SONO ENTRAMBI MAGGIORI DI 5, LA DISTRIBUZIONE DELLA PROPORZIONE DI SUCCESSI PUÒ ESSERE APPROSSIMATA DALLA DISTRIBUZIONE NORMALE E QUINDI:

$$Z = \frac{p - \pi}{\sqrt{\pi(1-\pi)/n}}$$

$\alpha = 0,05$ REGIONE ACCETTAZIONE = $0,95$
 $\frac{0,05 - 0,05}{0,5(1-0,05)\sqrt{64}} = -5/4$ CADE NELLA REGIONE DI RIFIUTO QUINDI RIFIUTIAMO L'IPOTESI NULLA

PROPORZIONE CAMPIONARIA = $369/644 = 0,573$
 VALORE CRITICO CON TABELLE: $0,9750 \rightarrow 1,96$



4) LEGAME TRA INTERVALLI DI CONFIDENZA E VERIFICA DI IPOTESI

$n = 25$
 $\bar{X} = 372,5$ grammi
 $\sigma = 15$ grammi
 INTERVALLO DI CONFIDENZA = 95%
 $\alpha = 0,05 (5\%)$
 VALORE CRITICO = $1,96$

AVENDO QUINDI IL SEGUENTE INTERVALLO DI CONFIDENZA:
 $\bar{X} \pm 1,96 \cdot \frac{15}{\sqrt{25}} = (366,6; 378,4)$
 AVENDO IPOTIZZATO NELL'IPOTESI NULLA 368 ED ESSENDO CHE $368 \notin (366,6; 378,4)$ NON RIFIUTIAMO L'IPOTESI NULLA

VERIFICA DI IPOTESI: TEST SU DUE CAMPIONI

CONFRONTO TRA MEDIE DI DUE POPOLAZIONI INDIPENDENTI

(LA SELEZIONE DI UN
CAMPIONE NON HA EFFETTO
SULLA SECONDA)

CONSIDERIAMO DUE POPOLAZIONI E SCELTIAMO DI ESTRARRE
UN CAMPIONE DI AMPIEZZA n_1 DALLA 1^a POP E UN CAMPIONE
 n_2 DALLA 2^a POP, ENTRAMBI OSSERVATI IN RELAZIONE AD UNA
VARIABLE CASUALE QUANTITATIVA X

μ_1 e μ_2 LE MEDIE
 σ_1 e σ_2 LE DEV. STANDARD

POPOLAZIONI CON DISTRIBUZIONE NORMALE OGNUNE ($n_1 \geq 30$
NON NORMALE MA CON AMPIEZZE ABBASTANZA GRANDI ($n_1 \geq 30$
PER POTER APPLICARE IL TEOREMA DEL LIMITE CENTRALE)

VUOLIAMO VERIFICARE L'IPOTESI NULLA CHE LE MEDIE DELLE
DUE POPOLAZIONI SIANO UGUALI TRA LORO:

$$H_0: \mu_1 = \mu_2 \quad (\mu_1 - \mu_2 = 0)$$

CONTRO UN'IPOTESI ALTERNATIVA BIDIREZIONALE O UNIDIREZIONALE:

$$H_1: \mu_1 \neq \mu_2 \quad (\mu_1 - \mu_2 \neq 0) \quad H_1: \mu_1 > \mu_2 \quad (\mu_1 - \mu_2 > 0)$$

$$H_1: \mu_1 < \mu_2 \quad (\mu_1 - \mu_2 < 0)$$

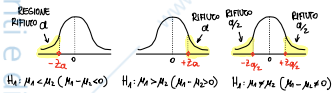
LA STATISTICA TEST È UNA FUNZIONE DELLA DIFFERENZA TRA LE
MEDIE CAMPIONARE DEI DUE CAMPIONI: $\bar{X}_1 - \bar{X}_2$

PER IL CONFRONTO TRA LE MEDIE DI 2 POP. INDIPENDENTI POSSIAMO:

- 1) TEST Z SULLA DIFFERENZA DI DUE MEDIE, VARIANZE σ_1 E σ_2 NOTE;
- 2) TEST t SULLA DIFFERENZA DI DUE MEDIE, VARIANZE OMOGENEE ($\sigma_1 = \sigma_2 = \sigma$) NON NOTE;
- 3) TEST t SULLA DIFFERENZA DI DUE MEDIE, VARIANZE NON OMOGENEE E NON NOTE.

$$1) Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

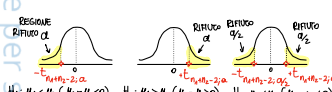
NORMALE $N(0,1)$; $\alpha(5\% \text{ o } 1\%)$



$H_1: \mu_1 < \mu_2$ ($\mu_1 - \mu_2 < 0$) $H_1: \mu_1 > \mu_2$ ($\mu_1 - \mu_2 > 0$) $H_1: \mu_1 \neq \mu_2$ ($\mu_1 - \mu_2 \neq 0$)

$$2) t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{s_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

DOVE $s_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{(n_1-1) + (n_2-1)}$; t DI STUDENT (n_1+n_2-2) GRADI DI LIBERTÀ; $\alpha(5\% \text{ o } 1\%)$



$H_1: \mu_1 < \mu_2$ ($\mu_1 - \mu_2 < 0$) $H_1: \mu_1 > \mu_2$ ($\mu_1 - \mu_2 > 0$) $H_1: \mu_1 \neq \mu_2$ ($\mu_1 - \mu_2 \neq 0$)

INTERVALLO DI CONFIDENZA PER LA DIFFERENZA TRA LE MEDIE DI
DUE POPOLAZIONI INDIPENDENTI ($\mu_1 - \mu_2$):

$$(\bar{X}_1 - \bar{X}_2) \pm t_{n_1+n_2-2, \alpha/2} \sqrt{s_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$$

VALORE CRITICO
AREA ACCETTAZIONE: $1 - \alpha/2$

CONFRONTO TRA MEDIE DI DUE POPOLAZIONI NON INDIPENDENTI

(CAMPIONI APPAIATI
CAMPIONI RELATIVE)

L'ANALISI NON CONSIDERA LE LE MEDIE DEI SINGOLI CAMPIONI
MA LA MEDIA DELLE DIFFERENZE TRA I DUE CAMPIONI

VALORI	1 CAMPIONE	2 CAMPIONE	DIFFERENZA
1	X_{11}	X_{21}	$D_1 = X_{11} - X_{21}$
...
j	X_{1j}	X_{2j}	$D_j = X_{1j} - X_{2j}$
...
n	X_{1n}	X_{2n}	$D_n = X_{1n} - X_{2n}$

SI VUOLE VERIFICARE L'IPOTESI NULLA CHE LA MEDIA DELLE DIFFERENZE
TRA LE DUE POPOLAZIONI UGUALE A 0:

$$H_0: \mu_D = 0$$

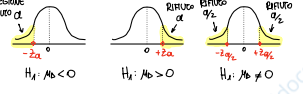
CONTRO UN'IPOTESI ALTERNATIVA BIDIREZIONALE O UNIDIREZIONALE:

$$H_1: \mu_D \neq 0 \quad H_1: \mu_D < 0 \quad H_1: \mu_D > 0$$

TEST Z SULLA MEDIA DELLE DIFFERENZE, VARIANZA DELLE DIFFERENZE σ_D NOTA

$$Z = \frac{D - \mu_D}{\frac{\sigma_D}{\sqrt{n}}}$$

DOVE $D = \frac{\sum D_j}{n}$; $\sigma_D = \sqrt{\frac{\sum (D_j - D)^2}{n-1}}$; NORMALE $N(0,1)$; $\alpha(5\% \text{ o } 1\%)$

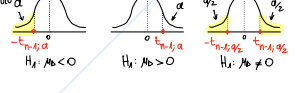


$H_1: \mu_D < 0$ $H_1: \mu_D > 0$ $H_1: \mu_D \neq 0$

TEST t SULLA MEDIA DELLE DIFFERENZE, VARIANZA DELLE DIFFERENZE σ_D NON NOTA

$$t = \frac{D - \mu_D}{\frac{s_D}{\sqrt{n}}}$$

DOVE $s_D = \sqrt{\frac{\sum (D_j - D)^2}{n-1}}$; t DI STUDENT ($n-1$) GRADI DI LIBERTÀ; $\alpha(5\% \text{ o } 1\%)$



$H_1: \mu_D < 0$ $H_1: \mu_D > 0$ $H_1: \mu_D \neq 0$

INTERVALLO DI CONFIDENZA PER LA MEDIA DELLE DIFFERENZE PER
DUE POPOLAZIONI NON INDIPENDENTI (μ_D):

$$D \pm t_{n-1, \alpha/2} \frac{s_D}{\sqrt{n}}$$

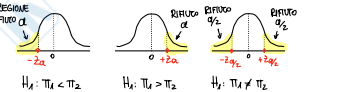
AREA ACCETTAZIONE: $1 - \alpha/2$

CONFRONTO TRA LE PROPORZIONI DI DUE POPOLAZIONI
 π_1 e π_2 PROPORZIONI DEI SUCCESSI NELLA POPOLAZIONE
 p_1 e p_2 PROPORZIONI CAMPIONARE DEI SUCCESSI (APPROSSIMAZIONE)
IPOTESI NULLA $H_0: \pi_1 = \pi_2$
IPOTESI ALTERNATIVA $H_1: \pi_1 < \pi_2$; $H_1: \pi_1 > \pi_2$; $H_1: \pi_1 \neq \pi_2$

TEST Z SULLA DIFFERENZA TRA DUE PROPORZIONI
DOVE X_1 E X_2 SONO IL NUMERO DI SUCCESSI NEI DUE CAMPIONI

$$Z = \frac{(p_1 - p_2) - (\pi_1 - \pi_2)}{\sqrt{p(1-p) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

DOVE $p_1 = \frac{X_1}{n_1}$; $p_2 = \frac{X_2}{n_2}$; $p = \frac{X_1 + X_2}{n_1 + n_2}$
NORMALE $N(0,1)$; $\alpha(5\% \text{ o } 1\%)$



$H_1: \pi_1 < \pi_2$ $H_1: \pi_1 > \pi_2$ $H_1: \pi_1 \neq \pi_2$

INTERVALLO DI CONFIDENZA PER LA DIFFERENZA TRA DUE PROPORZIONI

$$\pi_1 - \pi_2 : (p_1 - p_2) \pm z_{\alpha/2} \sqrt{\frac{p(1-p)}{n_1} + \frac{p(1-p)}{n_2}}$$

AREA ACCETTAZIONE: $1 - \alpha/2$

CONFRONTO TRA VARIANZE DI DUE POPOLAZIONI INDIPENDENTI

IPOTESI NULLA $H_0: \sigma_1^2 = \sigma_2^2$
IPOTESI ALTERNATIVA $H_1: \sigma_1^2 < \sigma_2^2$; $H_1: \sigma_1^2 > \sigma_2^2$; $H_1: \sigma_1^2 \neq \sigma_2^2$

TEST F SULLA DIFFERENZA TRA DUE VARIANZE CAMPIONARE:

$$F = \frac{s_1^2}{s_2^2}$$

DISTRIBUZIONE F DI FISCHER CON (n_1-1) E (n_2-1) GRADI DI LIBERTÀ
 $\alpha(5\% \text{ o } 1\%)$

VALORE CRITICO SUPERIORE F_U SI TROVA CON LE TAVOLE DI FISCHER;
MENTRE VALORE CRITICO INFERIORE F_L CON $F_L = 1/F_U^*$

DOVE F_U^* È IL VALORE CRITICO SUPERIORE CON GRADI DI LIBERTÀ INVERTITI:
 (n_2-1) E (n_1-1)

AREA ACCETTAZIONE

ANALISI DELLA VARIANZA (ANOVA) A UNA VÍA

GRUPPI DEFINITI SULLA BASE DI 1 SINGOLO FATTORE

CONSENTE DI CONFRONTARE
MEDIE DI PIÙ DI 2 GRUPPI (POPOLAZIONI)
INDIVIDUATI SULLA BASE
DI 1 O PIÙ FATTORE DI
INTERESSE E SONO FORMATI
SECONDO I LIVELLI (MODALITÀ)

ANCHE SE SI PARLA DI ANALISI DELLA VARIANZA IN REALTÀ
L'OBIETTO DI INTERESSE SONO LE DIFFERENZE TRA MEDIE
NEI DIVERSI GRUPPI E TRAMITE L'ANALISI DELLA VARIANZA

CONSIDERIAMO UN FATTORE DI INTERESSE X CARATTERIZZATO DA C LIVELLI (TRATTAMENTI)
E UNA VARIABLE CASUALE QUANTITATIVA Y (RISPOSTA O RISULTATO)
I C GRUPPI RAPPRESENTANO POPOLAZIONI CON DISTRIBUZIONE NORMALE,
CON TUTTI LA STESSA VARIANZA E LE OSSERVAZIONI ESPRIME CASUALMENTE

GRUPPO 1: LIVELLO A_1 $X_{11}, X_{12}, \dots, X_{1n_1}$
GRUPPO 2: LIVELLO A_2 $X_{21}, X_{22}, \dots, X_{2n_2}$
GRUPPO C: LIVELLO A_C $X_{c1}, X_{c2}, \dots, X_{cn_c}$



IPOTESI NULLA $H_0: \mu_1 = \mu_2 = \dots = \mu_C$

IPOTESI ALTERNATIVA H_1 : NON TUTTE LE μ_j SONO UGUALI
(CON $j = 1, 2, \dots, C$)



$\mu_1 = \mu_2 = \mu_3$ $\mu_1 \neq \mu_2 \neq \mu_3$

PER VERIFICARE L'IPOTESI NULLA, LA VARIABILITÀ TOTALE (SST)

VIENE SCONIATA IN DUE COMPONENTI:

- VARIABILITÀ TRA I GRUPPI (SSA)

- VARIABILITÀ ENTRO I GRUPPI (SSW)

$$SST = SSA + SSW$$

S.L. n-1 G.L. n-1 G.L. n-1

VARIABILITÀ TOTALE: $SST = \sum_{j=1}^C \sum_{i=1}^{n_j} (X_{ij} - \bar{X})^2$

DOVE $\bar{X} = \frac{\sum_{j=1}^C \sum_{i=1}^{n_j} X_{ij}}{n}$ = MEDIA COMPLESSIVA, $n = \sum_{j=1}^C n_j$

VARIABILITÀ TRA GRUPPI: $SSA = \sum_{j=1}^C n_j (\bar{X}_j - \bar{X})^2$

DOVE $\bar{X}_j = \frac{\sum_{i=1}^{n_j} X_{ij}}{n_j}$ = MEDIA CAMPIONARE NEL j-ESIMO CAMPIONE

VARIABILITÀ ENTRO GRUPPI: $SSW = \sum_{j=1}^C \sum_{i=1}^{n_j} (X_{ij} - \bar{X}_j)^2$

DOVE $\bar{X}_j = \frac{\sum_{i=1}^{n_j} X_{ij}}{n_j}$ = MEDIA CAMPIONARE NEL j-ESIMO CAMPIONE

LA STATISTICA F HA DISTRIBUZIONE F DI FISCHER CON
(C-1) GRADI DI LIBERTÀ AL NUMERATORE E (n-C) GRADI
DI LIBERTÀ AL DENOMINATORE "TABELLE DI FISCHER" -> F₀

REGIONE DI RIFIUTO: α $\alpha = F_0$

REGIONE ACCETTAZIONE: $1 - \alpha$

SE $F \geq F_{\alpha}$ L'IPOTESI NULLA VIENE RIFIUTATA

CALCOLATO TAVOLE = VALORE CRITICO

$$SST = \sum_{j=1}^C \sum_{i=1}^{n_j} (X_{ij} - \bar{X})^2$$

$$SST = \sum_{j=1}^C \sum_{i=1}^{n_j} X_{ij}^2 - C$$

$$SSA = \sum_{j=1}^C n_j (\bar{X}_j - \bar{X})^2$$

$$SSA = \sum_{j=1}^C \frac{T_j^2}{n_j} - C$$

$$SSW = \sum_{j=1}^C \sum_{i=1}^{n_j} (X_{ij} - \bar{X}_j)^2$$

$$SSW = SST - SSA$$

DIVIDENDO CIASCUNA SOMMA DEI QUADRATI PER I RESPECTIVI
GRADI DI LIBERTÀ SI OTTENGONO TRE VARIANZE,
O MEDIE DEI QUADRATI: MST, MSA, MSW

$$MST = \frac{SST}{n-1} \quad MSA = \frac{SSA}{C-1} \quad MSW = \frac{SSW}{n-C}$$

INOLTRE VIENE RIPORTATA LA STATISTICA TEST F PER
VERIFICARE L'IPOTESI NULLA:

$$F = \frac{MSA / (C-1)}{MSW / (n-C)} = \frac{MSA}{SSW / (n-C)}$$

SE L'IPOTESI NULLA È VERA
IL VALORE DEDURREBBE ESSERE CIRCA 1

SE L'IPOTESI NULLA È FALSA
IL VALORE DEDURREBBE ESSERE MOLTO
SUPERIORE ALL'UNITÀ

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

DUE MEDIE SONO SEGNALETE COME SIGNIFICANTEMENTE
DIVERSE TRA LORO SE LA DIFFERENZA IN VALORE
ASSOLUTO DI TALI MEDIE È MAGGIORE DEL
RANGE CRITICO = $Q_{\alpha} \sqrt{\frac{MSW}{2} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$

NELL'ANOVA CI SONO ALMENO 2 MEDIE
SIGNIFICANTEMENTE DIVERSE, PER TROVARE
USANDO DEI METODI CHIAMATI CONFRONTI MULTILI
-> PROCEDURA DI TUKEY-CRAMEL $C(C-1)/2$

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

DOVE Q_{α} È IL VALORE CRITICO CON:
C GRADI AL NUMERATORE E n-C GRADI
AL DENOMINATORE

ANALISI DELLA VARIANZA (ANOVA) A DUE VIE

IPOTESI:

- EFFETTO PRINCIPALE DEL FATTORE A
- EFFETTO PRINCIPALE DEL FATTORE B
- INTERAZIONE TRA A E B

Origine della variabilità	Somma dei quadrati	Gradi di libertà	Media dei Quadrati	F ₀
A trattamenti	SS _A	a-1	MS _A = SS _A / (a-1)	MS _A / MS _E
B trattamenti	SS _B	b-1	MS _B = SS _B / (b-1)	MS _B / MS _E
Interazione	SS _{AB}	(a-1)(b-1)	MS _{AB} = SS _{AB} / ((a-1)(b-1))	MS _{AB} / MS _E
Errore	SS _E	ab(n-1)	MS _E = SS _E / (ab(n-1))	
Totale	SS _T	abn-1		

- GRAFICO DI PROBABILITÀ NORMALE DEI RESIDUI
- GRAFICO DEI RESIDUI RISPETTO AI VALORI STIMATI
- GRAFICO DEI RESIDUI NELL'ORDINE TEMPORALE DELLA RACCOLTA DEI DATI
- GRAFICO DEI RESIDUI RISPETTO AD ALTRE VARIABILI

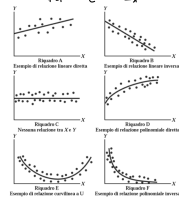
CONNESSIONE E CORRELAZIONE

SE VI È ASSENZA DI RELAZIONI TRA DUE O PIÙ VARIABILI
 → INDIPENDENTI (QUANDO IL VALORE DI UNA NON MODIFICA L'ALTRA)

MISURE DI CONNESSIONE SE VARIABILI CATEGORICALI (4M, ASS E REL)
 MISURE DI CORRELAZIONE SE VARIABILI NUMERICHE (LEGARE STATISTICO)

PER DUE VARIABILI QUANTITATIVE SI PUÒ USARE UN DIAGRAMMA DI DISPERSIONE → SERVE A RILEVARE ESSENTIALI TENDENZE

COEFFICIENTE DI CORRELAZIONE TRA +1 E -1



COEFFICIENTE DI CORRELAZIONE:

$$r = \frac{\text{COV}(X,Y)}{S_x S_y}$$

$$\text{COV}(X,Y) = \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})(y_j - \bar{y})$$

$$S_x = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2}$$

$$S_y = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (y_j - \bar{y})^2}$$

OGNI n VALORE DI x → SOMMA VALORI X / NUMERO DELLE X
 SOMMA VALORI Y / NUMERO DELLE Y

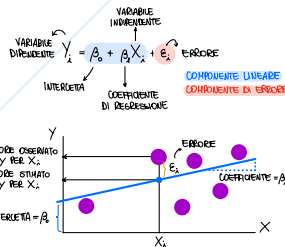
CORRELAZIONE E REGRESSIONE

STUDIO DELLA ASSOCIAZIONE TRA VARIABILI QUANTITATIVE

MODELLO ATTUALE CON PREVEDERE IL VALORE DI UNA VARIABILE DIPENDENTE O ESPONATA A PARTIRE DAL VALORE DI UNA O PIÙ VARIABILI INDIPENDENTI O ESPLICATIVE

- LINEARE SEMPLICE

LINEARE O CURVILINEA
 RELAZIONE FORTE, DEBOLE
 O NESSUNA RELAZIONE



STIMA O VALORE PREDETTO DI Y PER L'UNITÀ X-ESIMA
 $\hat{Y}_x = b_0 + b_1 X_x$
 STIMA DELL'INTERCETTA b_0
 STIMA DEL COEFFICIENTE DI REGRESSIONE b_1
 IL SEGNO DI b_1 INDICA SE LA RELAZIONE LINEARE È POSITIVA O NEGATIVA

METODO DEI MINIMI QUADRATI
 - DETERMINARE b_0 e b_1

$$\min \sum (Y_i - \hat{Y}_i)^2 = \min \sum (Y_i - (b_0 + b_1 X_i))^2$$

$$b_1 = \frac{SQXY}{SQX}; \quad b_0 = \bar{Y} - b_1 \bar{X}$$

$$\bar{Y} = \frac{\sum Y_i}{n}; \quad \bar{X} = \frac{\sum X_i}{n}$$

$$SQXY = \sum (X_i - \bar{X})(Y_i - \bar{Y}) = \sum X_i Y_i - \frac{(\sum X_i)(\sum Y_i)}{n}$$

$$SQX = \sum (X_i - \bar{X})^2 = \sum X_i^2 - \frac{(\sum X_i)^2}{n}$$

INTERPOLAZIONE ED ESTRAPOLAZIONE

PREVISIONE DI Y CONSIDERANDO AD UN VALORE DI X INTERIO ALL'INTERVALLO

PREVISIONE DI Y CONSIDERANDO AD UN VALORE DI X CHE NON CADE NELL'INTERVALLO

VARIABILITÀ TOTALE DELLA Y IN DUE QUADRANTI

$$SST = SSR + SSE$$

SST: SOMMA TOTALE DEI QUADRATI

$$SST = \sum (Y_i - \bar{Y})^2$$

MISURA LA DIMENSIONE DEI VALORI DELLA Y ATTORNO ALLA RICERCA MEDIA

SSR: SOMMA DEI QUADRATI DELLA REGRESSIONE

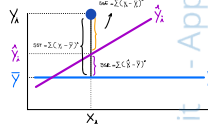
$$SSR = \sum (\hat{Y}_i - \bar{Y})^2$$

MISURA LA VARIABILITÀ DELLA Y ATTRIBUIBILE ALLA VARIABILE X

SSE: SOMMA DEI QUADRATI DEGLI ERRORI

$$SSE = \sum (Y_i - \hat{Y}_i)^2$$

MISURA LA VARIABILITÀ DELLA Y NON ATTRIBUIBILE ALLA X



COEFFICIENTE DI DETERMINAZIONE

$$r^2 = \frac{SSR}{SST} \quad 0 \leq r^2 \leq 1$$

NESSUNA RELAZIONE TRA X E Y

CORRELAZIONE LINEARE DIRETTA (+) O INVERSA (-) PERFETTA

ERRORE STANDARD DELLA STIMA

$$S_{YX} = \sqrt{\frac{SSE}{n-2}} = \sqrt{\frac{\sum (Y_i - \hat{Y}_i)^2}{n-2}}$$

MINORE S_{YX} → RELAZIONE FORTE

MAGIORE S_{YX} → RELAZIONE DEBOLE

ANALISI DEI RESIDUI e_i

$e_i = Y_i - \hat{Y}_i$ RESIDUO

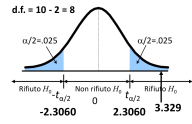
- IPOTESI DI LINEARITA' (retta orizzontale)
- DISTRIBUZIONE NORMALE DEGLI ERRORI
- OMOSCHEDASTICITA' (costante)
- INDIPENDENZA DAGLI ERRORI (autocorrelazione)

INFERENZA NEL MODELLO DI REGRESSIONE

$H_0: \beta_0 = 0$ (NESSUNA RELAZIONE LINEARE)

$H_1: \beta_0 \neq 0$ (ESISTENZA RELAZIONE LINEARE)

$t_{STAT} = \frac{b_0 - \beta_0}{S_{b_0}} = \frac{0.10927 - 0}{0.033284} = 3.32828$



d.f. = 10 - 2 = 8

$F_{STAT} = \frac{MSR}{MSE} = \frac{SSR}{SSE / (n-2)}$

RIFIUTO H_0 SE $F > F_\alpha$
 CON F_α VALORE CRITICO CHE
 LASCIA A DESTRA PROBABILITA'
 PARI AD α .

$b_0 \pm t_{\alpha/2} \cdot S_{b_0}$ G.L. = n-2

APPLICAZIONE CORRETTA DELL'ANALISI DI REGRESSIONE

Sintetizziamo i passi per una corretta applicazione dell'analisi di regressione:

1. iniziare l'analisi con un'attenta osservazione del diagramma di dispersione, per cogliere l'eventuale relazione tra X e Y;
2. dopo la stima del modello e prima di passare a impiegarne i risultati, verificare se le assunzioni alla base del modello sono soddisfatte:
 - a) grafico dei residui rispetto alla variabile X per verificare l'assunzione di omoschedasticità delle varianze;
 - b) istogramma dei residui e normal probability plot, per verificare l'assunzione di normalità degli errori;
 - c) grafico dei residui rispetto all'ordine temporale di raccolta dei dati, per verificare l'assunzione di indipendenza degli errori.
3. se le ipotesi alla base del modello non risultano violate, procedere alle inferenze sul modello (test di ipotesi sui coefficienti, intervalli di confidenza per la risposta media e la previsione).

www.unidocs.it - Appunti e dispense per superare i tuoi esami universitari

www.unidocs.it - Appunti e dispense per superare i tuoi esami universitari