



Appunti universitari

Tesi di laurea

Cartoleria e cancelleria

Stampa file e fotocopie

Print on demand

Rilegature

NUMERO: 2238A

ANNO: 2017

A P P U N T I

STUDENTE: Tosti Michela

MATERIA: Statistica - Teoria - Prof. Vicario

Il presente lavoro nasce dall'impegno dell'autore ed è distribuito in accordo con il Centro Appunti.

Tutti i diritti sono riservati. È vietata qualsiasi riproduzione, copia totale o parziale, dei contenuti inseriti nel presente volume, ivi inclusa la memorizzazione, rielaborazione, diffusione o distribuzione dei contenuti stessi mediante qualunque supporto magnetico o cartaceo, piattaforma tecnologica o rete telematica, senza previa autorizzazione scritta dell'autore.

**ATTENZIONE: QUESTI APPUNTI SONO FATTI DA STUDENTIE NON SONO STATI VISIONATI DAL DOCENTE.
IL NOME DEL PROFESSORE, SERVE SOLO PER IDENTIFICARE IL CORSO.**

STATISTICA

CAP 1

Definizione di Probabilità:

o classica $P[E] = \frac{S}{n}$
 n° risultati favorevoli
 n° risultati possibili

Attenzione!
 o risultati casualmente possibili
 o incompatibili

o $0 \leq P \leq 1$
 ↑ evento certo
 ↓ evento impossibile

o frequenzista
 frequenza $f_E = \frac{NE}{N}$
 ↑ evento
 Relativa

$\frac{\text{n° volte in cui si presenta l'evento}}{\text{n° esperimenti}}$

ovvero come valori della P di un evento E, il valore limite a cui tende la frequenza relativa di quell'evento al tendere del n° di prove all'∞

$$\lim_{N \rightarrow \infty} P[f_E - P[E]] = 0$$

o sovrattiva
 VENTE SULLA FIDUCIA con cui il soggetto ritiene possibile il verificarsi dell'evento

o ASSIOMATICA

PREMESSE:
 o Fenomeno Casuale (Aleatorio): un fenomeno empirico il cui risultato non è prevedibile a priori → non ha quindi REGOLA
 DETERMINISTICA NE STATISTICO.

INDIPENDENTI:
 $P[A \cap B] = P[A] \cdot P[B]$
 $P[A \cup B] = P[A] + P[B] + P[A \cap B]$
 ESCLUDENTESI:
 $P[A \cap B] = \emptyset$
 $P[A \cup B] = P[A] + P[B] - P[A \cap B]$
 DIPENDENTI
 $P[A \cap B] = \rightarrow \infty$

Per trovare $P[A \cup B] \rightarrow$ uso Formula:
 $P[A] = P[A \cap B] + P[A \bar{B}]$
 $P[B] = P[A \cap B] + P[\bar{A} \cap B]$

INDIPENDENTI: $P[A] = P[A/B]$
 $P[B] = P[B/A]$

detto E un qualsiasi evento costituito da #E punti campione si può definire:

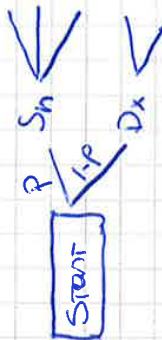
$$P[E] = \frac{\#E}{\#S}$$

↳ se ho k componenti: $\rightarrow \frac{\#A_k}{\#S} = \frac{\binom{n}{k} k^k (n-k)^{n-k}}{\binom{n}{n}}$

da cui:

$$P[A_k] = \frac{\binom{k}{k} (n-k)^{n-k}}{\binom{n}{n}}$$

diagramma ad albero



Probabilità Condizionata

Probabilità che ha un evento A verificarsi condizionatamente al verificarsi dell'evento B

$$P[A|B]$$

e si calcola:

$$P[A|B] = \frac{\#AB}{\#B} = \frac{P[AB]}{P[B]}$$

scoprisse tutti e 3 gli orsioni di R.

Probabilità totali

teorema:

Sia E_1, E_2, \dots, E_n una collezione di eventi incompatibili tali che $S = \bigcup_{i=1}^n E_i$ (eventi esaustivi) e $P[E_i] \neq 0 \forall i=1 \dots n$. Qualsiasi sia $F \subseteq S$ si ha:

$$P[F] = \sum_{i=1}^n P[F|E_i]P[E_i]$$

Formule di Bayes

teorema:

Sia E_1, E_2, \dots, E_n una collezione di eventi incompatibili ed esaustivi e sia $P[E_i] \neq 0 \forall i=1 \dots n$. Qualsiasi sia $F \subseteq S$ si ha:

$$P[E_k|F] = \frac{P[F|E_k]P[E_k]}{\sum_{i=1}^n P[F|E_i]P[E_i]}$$

$\forall k=1 \dots n$

↓ se $n=2$

$$P[E|F] = \frac{P[F|E]P[E]}{P[F|E]P[E] + P[F|\bar{E}]P[\bar{E}]}$$

Regole per il calcolo delle P

moltiplicazione:

Se E_1, E_2, \dots, E_n sono eventi appartenenti allo stesso spazio degli eventi e $P[E_1, E_2, \dots, E_n] \neq 0$, si ha:

$$P[E_1 E_2 E_3 \dots E_m] = P[E_1] \times P[E_2|E_1] \times P[E_3|E_1 E_2] \dots \times P[E_n|E_1 E_2 \dots E_{n-1}]$$

La funzione a distribuzione, spiega che sono distribuiti i valori della variabile casuale, il termine cumulativo, talvolta oneroso, sta a ricordare che i valori di tale funzione vengono forniti in forma cumulata.

Definizione: Tale funzione è univocamente definita per ciascuna variabile casuale.



Caratteristiche di tale funzione sono:

- non negativa
 - non decrescente
- $\lim_{x \rightarrow -\infty} F_X(x) = 0$ $\lim_{x \rightarrow +\infty} F_X(x) = 1$
 $F_X(x) \leq F_X(x_2)$

- Continua e $dx \rightarrow \lim_{h \rightarrow 0} F_X(x+h) = F_X(x)$

Funzione di Densità

Definizione: Data una variabile casuale discreta X con coordinate $\{x_1, x_2, \dots\} \subset \mathbb{R}$ e la funzione $f_X(x)$ definita da:

$$f_X(x) = \begin{cases} P\{X=x_j\} & x=x_j, \text{ con } j=1, \dots, n \\ 0 & x \neq x_j \end{cases}$$

è una funzione da \mathbb{R} nell'intervallo chiuso $[0, 1]$ de, per come è definita, gode delle seguenti proprietà:

- i) $f_X(x) \geq 0 \quad \forall x \in \mathbb{R}$. (Altre volte si scrive $f_X(x) \geq 0 \quad \forall x = x_j$ e $f_X(x) = 0 \quad \forall x \neq x_j$)
- ii) $\sum_{j=1}^n f_X(x_j) = 1$ dove la sommatoria è estesa a tutti i x nome

teorema: la funzione di densità discreta $f_X(x)$ è legata alla funzione di distribuzione cumulativa $F_X(x)$ della stessa variabile casuale X delle seguenti relazioni:

$$F_X(x) = \sum_{j: x_j \leq x} f_X(x_j)$$

$$f_X(x) = \begin{cases} F_X(x_j) - \lim_{h \rightarrow 0^+} F_X(x_j - h) & \text{se } x = x_j \\ 0 & \forall x \neq x_j \end{cases}$$

Funzione di Densità di Probabilità (Continua)

È la variabile X continua la funzione $f_X(x)$ tale che

$$F_X(x) = \int_{-\infty}^x f_X(x) dx$$

è detta funzione di densità di probabilità di X

Proprietà: (analoghe alla discreta:)

- i) $f_X(x) \geq 0, \quad \forall x \in \mathbb{R}$
- ii) $\int_{-\infty}^{+\infty} f_X(x) dx = 1$

NOTA: se la funzione di densità discreta ha come coordinate l'intervallo $[0, 1]$ la funzione di densità di \mathbb{R} $f_X(x)$ ha come coordinate $[0, +\infty)$

NOTA2:

$$f_X(x) = \frac{dF_X(x)}{dx}$$

Disuguaglianza di Tchebyche F.

Se X è una variabile casuale e $g(x)$ una funzione non negativa per $\forall x \in \mathbb{R}$, allora:

$$P[g(x) \geq k] \leq \frac{E[g(x)]}{k} \quad \downarrow k > 0$$

condizione delle disuguaglianze di T.

$$P[|X - \mu_x| \geq t\sigma_x] = P[(X - \mu_x)^2 \geq t^2\sigma_x^2] \leq \frac{1}{t^2}$$

in forma equivalente:

$$P[|X - \mu_x| < t\sigma_x] \geq 1 - \frac{1}{t^2}$$

e anche:

$$P[|X - \mu_x| < X < \mu_x + t\sigma_x] \geq 1 - \frac{1}{t^2}$$

è la probabilità che la somma di una qualsiasi variabile casuale X differa da meno di 2 deviazioni standard dalla sua media.

Altre misure di dispersione; MOMENTI

Quantile: si definisce quantile q -esimo di X o di una variabile casuale continua X il più piccolo valore $X \in \mathbb{R}$ tale che $F_X(X_q) = q$

In caso equivalente, il quantile è il + piccolo valore $X \in \mathbb{R}$ tale che l'area sotto alla curva di densità di F a sx di tale valore sia uguale a q .



per variabile casuale discreta si intende che il quantile X_q sia il più piccolo N^o per cui si ha $F_X(X_q) \geq q$

Mediana: valore che soddisfa:

$$P[X \leq \text{med}(X)] \geq \frac{1}{2} \quad \text{e} \quad P[X \geq \text{med}(X)] \geq \frac{1}{2}$$

$$X \text{ continua: } \int_{-\infty}^{\text{med}(x)} f_X(x) dx = \int_{\text{med}(x)}^{+\infty} f_X(x) dx$$

Nota: valore in cui $f_X(x)$ ha il suo max (solo per caso)

Escursione (range): differenza tra il valore max e il valore min che X può assumere

escursione interquartile: differenza $Q_3 - Q_1$, tra il valore del 3° e del 1° quartile delle $u.c. X$

Momento di ordine n di X : valore atteso X^n se esiste, cioè:

$$\mu_n = E[X^n]$$

Momento centrale di ordine n : μ_n valore atteso di $(X - \mu_x)^n$ cioè:

$$\mu_n = E[(X - \mu_x)^n]$$

NOTA: $\mu_1 = E[X - \mu_x] = 0$

$\mu_2 = E[(X - \mu_x)^2] = \text{var}(X)$

Se la funzione di densità $f_X(x)$ è simmetrica rispetto a μ_x , allora tutti i momenti centrali dispari sono nulli.

Introduciamo ora la v.c. $V/n = (X - np)/n$, essa rappresenta lo scarto relativo delle v.c. X rispetto al suo valore medio

$$E[V/n] = E[X/n - p] = \frac{1}{n} E[X] - p = 0$$

$$\text{var}[V/n] = \frac{1}{n^2} \text{var}[X] = \frac{pq}{n}$$

cioè al crescere del numero delle prove diventano sempre più probabili i valori di piccoli scarti relativi, in quanto il valore medio dello scarto relativo è 0 e la sua dispersione decresce al crescere di n .

di considerare ora il coefficiente delle disuguaglianze di Tchebycheff, esprimendo come v.c. X la frequenza X/n in condizioni che $E[X/n] = p$ e $\text{var}[X/n] = \frac{pq}{n}$ si ha:

$$P[|X/n - p| < \sqrt{\frac{pq}{n}}] \geq 1 - \frac{1}{n^2}$$

da cui, ponendo $\epsilon = \sqrt{\frac{pq}{n}}$ si ha:

$$P[|X/n - p| < \epsilon] \geq 1 - \frac{1}{n\epsilon^2}$$

che vuol dire che, fissata una quantità piccola e fissato $\epsilon > 0$, il n° delle prove tende ad infinito $P[|X/n - p| < \epsilon]$ tende all'unità.

↳ Risultato è il Teorema di Bernoulli → enunciato è

In una serie di n prove, in ciascuna delle quali un evento ha probabilità p di manifestarsi, la probabilità che lo affermano tra le n prove con cui l'evento si è manifestato è p sia, in valore assoluto, inferendo ad una quantità assegnata e inferendo limitare e tende alla certezza col crescere del numero n delle prove.

Così S è un insieme di n -uple in cui l' i -esimo elemento indica il risultato della i -esima prova

$$(P[S] = p, P[A] = 1-p)$$

de P che si verificano determinate n -uple costituite da X successi, e quindi $n-X$ insuccessi, è data da $P^X q^{n-X}$ a causa dell'indipendenza delle prove e dal fatto che la stessa P si applica ad ogni prova

Teorema di Bernoulli

si consideri una variabile casuale X Binomiale di parametri n e p , cioè X rappresenta il numero delle volte che in n prove indipendenti e ripetute si verifica un determinato evento E che ha probabilità p di presentarsi in ciascuna prova.

Assoluta da variabile casuale $V = X - np$ rappresenta lo scarto assoluto della v.c. X dal suo valore medio np quanto

$$E[X] = np$$

$$E[V] = E[X - np] = 0$$

$$\text{var}[V] = \text{var}[X - np] = \text{var}[X] = npq$$

Così V è una v.c. che ha la sua distribuzione con pmf meno $-np, -np+1, \dots, np$ con medio $= 0$ e dispersione che cresce al crescere del n° delle prove n . Il fatto che $\text{var}[V]$ aumenti al crescere di n significa che all'aumentare del n° delle prove cresce la probabilità di trovare per valori delle variabile X sempre più distanti dal valore medio.

Consideriamo ora la variabile casuale X/n cioè il rapporto tra il n° di successi X in n prove ed il n° stesso delle prove (frequenza relativa). Applicando le proprietà dei valori attesi:

$$E[X/n] = \frac{1}{n} E[X] = p$$

$$\text{var}[X/n] = \frac{1}{n^2} \text{var}[X] = \frac{pq}{n} = \frac{p(1-p)}{n}$$

Perciò la frequenza ha una distribuzione con E nome $0, 1/n, 2/n, \dots, 1$ valore medio p e dispersione che tende a 0 per $n \rightarrow \infty$

code: $P[N, \text{eventi in } (t_1, t_2)] = n_2 \text{ eventi in } (t_2, t_1) = P[(t_1, t_2)]$

se tali hp sono soddisfatte si ha il Teorema:

la variabile casuale $N(t)$ che indica il N° di volte che un evento (per il quale sono soddisfatte le 3 hp) si verifica in un intervallo qualsiasi di lunghezza t segue una distribuzione di Poisson con parametro $\lambda = \alpha t$ cioè:

$$P[N(t) = x] = e^{-\alpha t} \frac{(\alpha t)^x}{x!} \quad \text{per } x = 0, 1, \dots, n$$

Distribuzione Geometrica

Def: la variabile casuale X ha una distribuzione geometrica se la densità discreta di X è data da:

$$f_X(x) = f_X(x; p) = \begin{cases} p(1-p)^{x-1} & x = 1, 2, \dots, n \\ 0 & \text{altrimenti} \end{cases}$$

con $0 < p \leq 1$

Come sempre: è una funzione di densità perché:

1) $f_X(x, p) > 0$

2) $\sum_{x=1}^{\infty} f_X(x, p) = 1$

si ha: $E[X] = \frac{1-p}{p}$ var $[X] = \frac{1-p}{p^2}$

Teorema: se X è una v.c. con densità geometrica con parametro p si ha:

$$P[k \geq i + j | \dots | i] = P[X \geq j]$$

con $i, j = 0, 1, \dots, n$

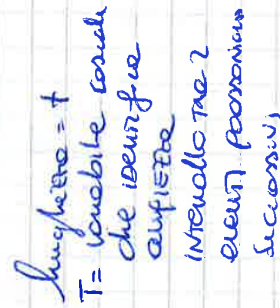
Processi stocastici es: - code

- temperatura e richiesta energia elettrica
- sistemi di telecomunicazione

sviluppo secondo le leggi probabilistiche

Processi di Poisson

sono per esempio N° di difetti per unità di superficie e per unità di lunghezza o di un certo prodotto proveniente da una produzione continua
 N° di radioattività emessa per unità di tempo da sostanze radioattive



hp: 1) esiste quantità $\alpha > 0$ tale che la P che si verifici esattamente un evento in un piccolo intervallo di lunghezza Δt sia \approx uguale a $\alpha \Delta t$

$$P(\text{esattamente 1 evento in } \Delta t) = \alpha \Delta t + o(\Delta t)$$

con $\lim_{\Delta t \rightarrow 0} \frac{o(\Delta t)}{\Delta t} = 0$

2) la P che nell'intervallo Δt si verifichi più di un evento è trascurabile rispetto alle P che se ne verifichi esattamente 1:

$$P[2 \text{ o } + \text{ eventi in } \Delta t] = o(\Delta t)$$

3) l'evento E_1 rappresenta il verificarsi di n_1 eventi in qualsiasi intervallo a tempo (t_1, t_2) ed è indipendente da E_2 rappresentante il verificarsi di n_2 eventi in un intervallo (t_3, t_4) non contenuto in (t_1, t_2)

→ qualunque sia la v.c. X che essa appartenga ad un intervallo simmetrico rispetto al valore medio e di deviazione 2 volte la deviazione standard è circa 95%
 ↳ inoltre la \mathbb{P} che una variabile casuale normale appartenga ad un intervallo centrato sul valore medio e di deviazione tre volte le σ è 99,73%.

Distribuzione Esponenziale

Def: la v.c. X che ha una distribuzione esponenziale (negativa) se la sua funzione di densità $f(x)$ è:

$$f(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & x < 0 \end{cases} \quad \text{con } \lambda \in \mathbb{R}^+$$

Funzione di densità probabilitaria come primo, perché:

- 1) $f(x) \geq 0$
- 2) $\int_{-\infty}^{\infty} f(x; \lambda) dx = 1$

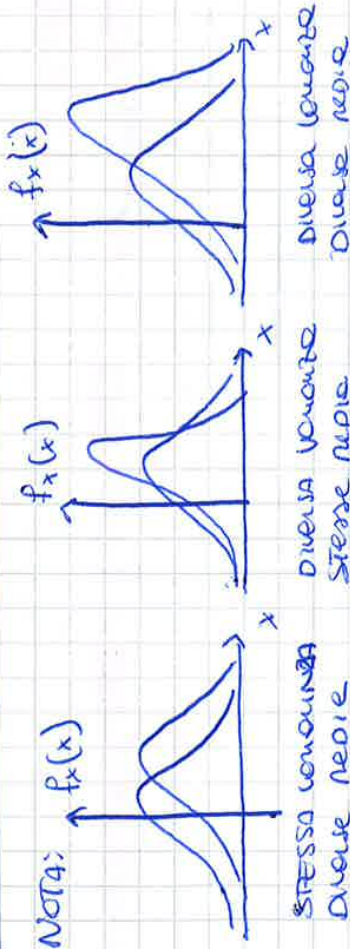
Si ha: $E[X] = \frac{1}{\lambda}$ $Var[X] = \frac{1}{\lambda^2}$

Teorema: data X una variabile casuale esponenziale con parametro λ si ha:

$$\mathbb{P}(X > s + t | X > s) = \mathbb{P}(X > t) \quad \text{con } s, t \in \mathbb{R}$$

→ si può anche definire una v.c. T che rappresenta la lunghezza dell'intervallo temporale che intercorre tra 2 manifestazioni successive dell'evento

$$\mathbb{P}(T > t) = \mathbb{P}[\text{nessuna manifestazione in un intervallo di lunghezza } t]$$



ATTENZIONE → $\mu = 0$ e $Var = 1$ Distribuzione STANDARDIZZATA

↳ Si trasforma la v.c. X in Z con $Z = \frac{(X - \mu)}{\sigma}$

→ la normale è simmetrica rispetto a $X = \mu$ (p.e. la standardizzata è simmetrica rispetto a $X = 0$)

→ sic che $x \rightarrow +\infty$ la $f(x)$ tende a 0

la funzione di distribuzione cumulativa normale è:

$$\Phi_{\mu, \sigma^2}(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(v-\mu)^2}{2\sigma^2}} dv$$

a parte di quanto: la ordinata con cumulate e fu ordinate ed ha parametro λ .

si ha: $E[X] = \mu$ e $Var[X] = \sigma^2$

se normale standardizzata

$$E[Z] = 0 \quad Var[Z] = 1$$

Campione: (sottoinsieme della popolazione), è un insieme di n unità statistiche, scelte tra quelle che costituiscono la popolazione, in base a criteri opportuni.

Campionamento sistematico: si tratta di scegliere come elemento che amovici a formare il campione, ad esempio ogni 30-esimo pezzo prodotto da una macchina, ogni 20-esimo nome di una lista di nomi, si ha che il **PASSO**

DI CAMPIONAMENTO $K = \lfloor N/n \rfloor$, il più grande intero contenuto in N/n

Campionamento stratificato: si divide la popolazione in un numero prestabilito di sottopopolazioni e strat che quindi estraggono delle unità che andranno a costituire il campione totale. Ad esempio, si viene a formare un campione di un certo tipo di componenti per rilevare le loro difettosità, una stratificazione possibile potrebbe essere quella di dividere la popolazione costituita dai componenti, in 3 strati costituiti da 3 tipi di materiali. La variabilità interna allo strato sarà $<$ della variabilità delle caratteristiche sotto indagine.

Conteggio: aumento efficiente perché si vuole l'errore di campionamento sempre aumentare la numerosità - possibilità di precisione più grande stratificati con precise previsioni, aumentano la numerosità del campione negli strati in cui si vuole una confidenza più approfondita. **Oversampling** è diminuzione incoerente **undersampling**

Campionamento con il rinvio delle quote
Diviso la popolazione in gruppi sulle basi delle caratteristiche (oggetto di studio) per i quali sono noti i pesi % di ciascuno nei confronti della popolazione

a questo punto vengono definite le quote, cioè il n° di elementi da prelevare da ciascun gruppo. Il campione, sono l'insieme costituito da tutte le unità estratte.

Campionamento a grappoli: la popolazione viene vista come un insieme di grappoli (cluster) non sovrapposti, in quanto si trovano nelle difficoltà di estrazione un campione dall'intero popolazione, e differenza dei precisi metodi di campionamento, in questo caso sono i grappoli ad essere scelti in modo casuale. In questi confronti la popolazione fidele le unità appartenenti ai grappoli scelti sono a pari parte del campione. Va fatto osservare che in tal caso la numerosità del campione, elemento determinante per l'infertilità statistica, è nota tra il termine del campionamento, in quanto la numerosità dei grappoli, in genere non è uguale e non è nota a priori.

Campione fisso

Insieme di n elementi per ognuno dei quali vengono rilevate determinate caratteristiche (individui, oggetti, ...)

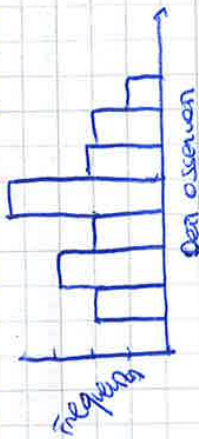
Campione influito

Insieme corrispondente alle note altre caratteristiche del campione fisso

MODALITÀ: termine usato per indicare uno dei possibili modi con cui le caratteristiche viene rilevate sull'unità statistica

SPUNTA: modo sintetico di trascrivere una serie di dati consistente nell'indicare ripetutamente un simbolo grafico il n° di ripetizioni osservate per ogni delle note altre unità

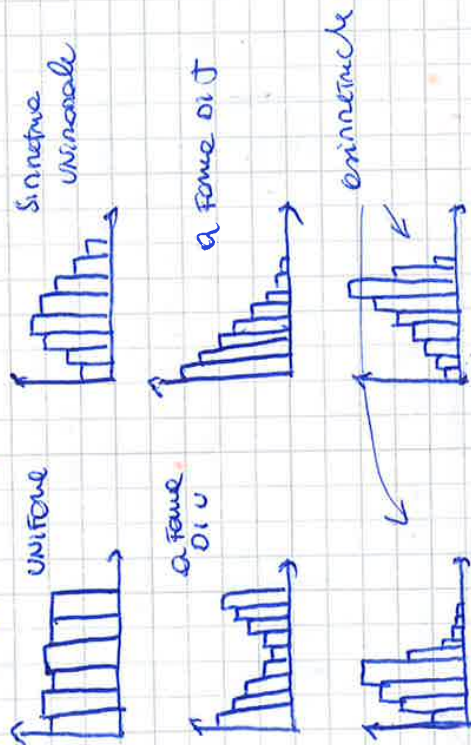
come componiamo. Interpretazione visiva più semplice, se le classi hanno ampiezza uguale in quanto in tal caso l'altezza è proporzionale alla frequenza.



Se combino 2 classi, la frequenza diventa la somma delle frequenze, ed è uguale all'area del rettangolo di cui sono le due più basse e lungo.

Frequenza specificata: si ottiene dividendo ogni frequenza di classe per la propria ampiezza.

Gli istogrammi: delineano indistintamente la forma delle distribuzioni delle popolazioni delle quali proviamo il campione esaminato.



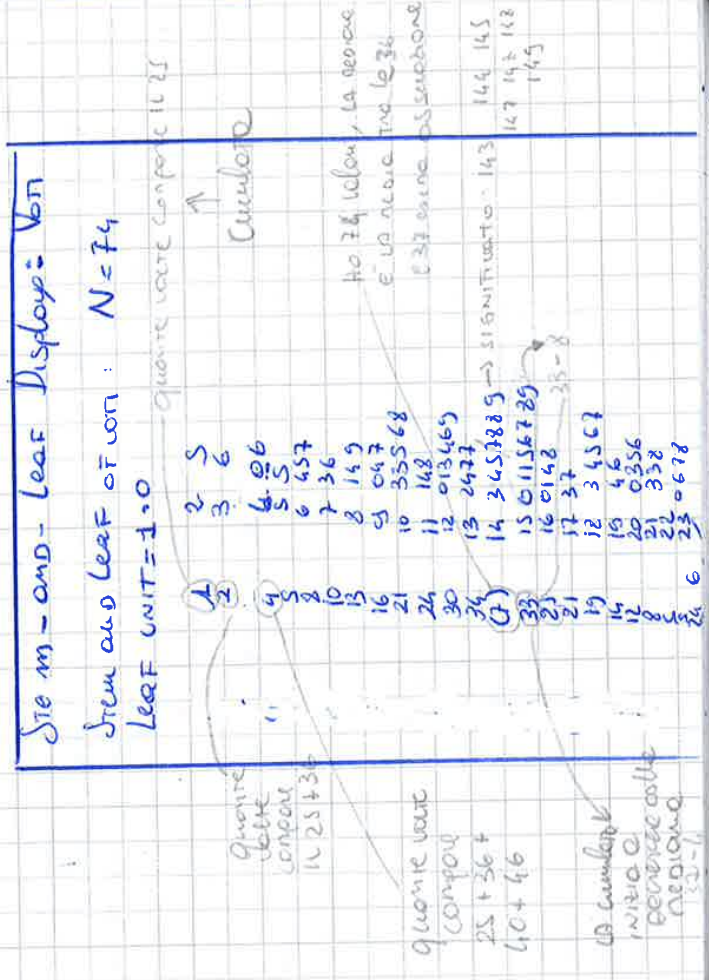
STEM AND LEAF plot (Diagrammi a Foglie e Ramo) rappresentano una locazione rispetto agli istogrammi, non si perde però info grazie all'insieme di coefficienti numerici per costruire le barre. La suddivisione in classi (di eguali ampiezze), elencate in ordine, è di regola su base decimale. (Vedi esempio).

L'aspetto corrisponde a quello di un istogramma, e differisce di quello, però, consente di indicare ad ogni singolo valore dell'insieme esaminato.

→ la 1ª colonna indica le frequenze considerate partendo dai 2 estremi e vi è indicato il ramo parentesi la frequenza corrispondente al ramo centrale e la mediana.

→ la 2ª colonna (ramo) corrisponde alle decine e le foglie (Foglie) ai centesimi.

La suddivisione in rami o classi (di eguale ampiezza) elencate nelle 2 colonne è di regola su base decimale e il n° Max di rami si determina sempre con la formula del logaritmo.



Indipendenza e Connessione

Frequenze congiunte (Assolute) di una combinazione (x_i, y_j) di due variabili, il n° di volte che compare le 2 caratteristiche X e Y secondo le modalità x_i, y_j e lo si indicherà con N_{ij} .

$$\sum_{i=1}^h \sum_{j=1}^k N_{ij} = \sum_{i=1}^h n'_{i0} = n$$

Frequenze congiunte relative: f_{ij} si ottiene dividendo la frequenza assoluta per il totale

$$f_{ij} = N_{ij}/n$$

Si ottiene lo stesso risultato tutte le frequenze congiunte conosciute a quelle conosciute di una delle variabili che consideriamo fissate.

$$n_{i0} = \sum_{j=1}^k N_{ij} \quad n_{0j} = \sum_{i=1}^h N_{ij}$$

La Freq. assoluta di una delle caratteristiche osservate in un fenomeno casuale nel quale vengono rilevate contemporaneamente 2 o 3 caratteristiche è detta:

FREQUENZA MARGINALE (Assoluta)

freq. marginale relativa si ottiene dividendo quelle assolute x il tot.

$$f_{i0} = n_{i0} / n \quad f_{0j} = n_{0j} / n$$

n_{i0} è la frequenza marginale osservata a una delle 2 variabili, che riguardano più caratteristiche.

osservazione) dicono le frequenze marginali e le condizionate sono distribuzioni di frequenze, si potranno definire le due misure descrittive (μ, σ) e secondo delle tipologie di caratteristiche.

INDIPENDENZA E CONNESSIONE

→ se non vi è alcun legame si parla di **INDIPENDENZA** tra le variabili che rappresentano le caratteristiche sotto indagine.

→ se si hanno 2 o 3 variabili, esse sono indipendenti se la distribuzione di frequenza delle variabile Y condizionata a $X=x_i$ non cambia di valore di X_i con $i=1, \dots, h$ (idem X condiz. a $Y=y_j$) dette le sinistre delle frequenze di indipendenza.

$$\frac{N_{ij}}{n_{i0}} = \frac{n_{0j}}{n} \quad \text{opp.} \quad \frac{N_{ij}}{n_{0j}} = \frac{n_{i0}}{n}$$

quando se le due variabili sono indipendenti, la frequenza relativa congiunta è uguale al prodotto delle 2 frequenze relative marginali.

→ Anche in presenza di indipendenza si osservano scostamenti tra le frequenze congiunte osservate e quelle espresse dalle formule sopra indicate dette frequenze teoriche, tali scostamenti vengono chiamati **CONTINGENZE (Assolute)**

$$C_{ij} = N_{ij} - \frac{n_{i0} \cdot n_{0j}}{n} \quad \begin{matrix} i=1, \dots, h \\ j=1, \dots, k \end{matrix}$$

ovvero definite come le differenze tra le frequenze assolute osservate e quelle teoriche.

→ di visuale in cui che sono medie assolute (indice di connessione di Pearson) e quadratiche (indice di normalizzati).

→ **INDICI DI CONNESSIONE** in quanto si indicano la presenza di un legame tra le 2 caratteristiche sotto indagine.

funzione interpolante. \rightarrow spesso $\hat{=}$ una retta
 $Y = a + bX$,
 $a, b \in \mathbb{R}$

Il metodo dei minimi quadrati richiede che sia

$$\min_{a,b} Q(a,b) = \sum_{i=1}^n [y_i - (a + bx_i)]^2$$



Per ottenere una dimostrazione si ottiene:

$$a = \bar{y} - b\bar{x}$$

$$b = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{S_{xy}}{S_x^2}$$

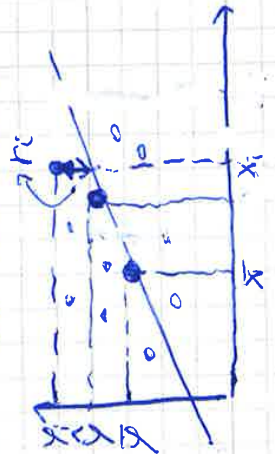
Però l'equazione della retta interpolante è la funzione di regressione $\hat{=}$:

$$Y = \bar{y} + \frac{S_{xy}}{S_x^2} (x - \bar{x})$$

\rightarrow coefficiente b strettamente legato alla covarianza tra le 2 variabili, perciò la retta è crescente se $cov > 0$, decrescente se $cov < 0$

RESIDUO $y_i - (a + bx_i)$ $(i=1 \dots n)$ e

le differenze tra i valori osservati e i valori in output della funzione interpolante



\rightarrow la somma dei residui è sempre $= 0$

$$\sum_{i=1}^n (y_i - \hat{y}_i) = \sum_{i=1}^n (y_i - \bar{y}) - b \sum_{i=1}^n (x_i - \bar{x}) = 0$$

con

$$r = \frac{S_{xy}}{S_x S_y}$$

COEFFICIENTE DI CORRELAZIONE LINEARE

la covarianza misura perciò una eventualità correlazione tra le 2 variabili sotto indagine

$cov = 0 \rightarrow$ assenza di relazione lineare \rightarrow VARIABILI NON CORRELATE

assente di legami lineari tra le 2 variabili che NON implica l'assenza di ogni tipo di relazione

\rightarrow INDIPENDENZA IMPLICA NON CORRELAZIONE
 \rightarrow NON CORRELATE NON IMPLICA INDIPENDENZA

$$r^2 = 1 - \frac{\text{var}[y - (a + bx)]}{\text{var}[y]} = 1 - \frac{\text{varianza residua di interpolazione}}{\text{varianza tot } y}$$

COEFFICIENTE DI DETERMINAZIONE LINEARE

\rightarrow SPIEGA LA PROPORZIONE DI VARIANZA DELLA VARIANZA TOT, DAVUTA ALLA RETTA INTERPOLANTE

$$0 \leq r^2 \leq 1 \quad \text{e} \quad -1 \leq r \leq 1$$

se $r^2 = r = S_{xy} = 0$ la retta interpolante non spiega alcuna varianza della variabile risale \rightarrow RETTA INTERP. $Y = \bar{y}$ e $X = \bar{x}$ non correlata

se $r^2 = 1$ ($r = \pm 1$) la varianza residua è nulla avendo tutti i punti (x_i, y_i) giaccono sulla retta interpolante \rightarrow TRE X E Y CORRELAZIONE PERFETTA

DEFINIZIONE 3: Dato un campione $\{x_1, x_2, \dots, x_n\}$ proveniente da una popolazione con densità $f(x)$ si definisce **media campionaria** la statistica

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

è concisa con la media estrinseca delle variabili che costituiscono il campione.

TEOREMA: Dato un campione $\{x_1, \dots, x_n\}$ proveniente da una popolazione avente densità $f(x)$, posto $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n x_i$ si ha:

$$E[\bar{X}_n] = \mu \quad \text{var}[\bar{X}_n] = \frac{\sigma^2}{n}$$

→ media e varianza di $f(x)$ sono delle

ci dice che, qualunque sia la distribuzione della media campionaria, esse $\mu \equiv \mu_{popolazione}$ e $\sigma^2 \equiv \sigma_{popolazione}^2$ sono sempre proporzionali alla $\frac{1}{n}$. Questo significa che al crescere della numerosità del campione, i valori delle componenti medie campionarie tendono a concentrarsi sempre di più attorno al loro valore medio che è la media della popolazione.

⇒ Da Tchebycheff si ha: **Legge dei Grandi Numeri**

$$P\left[|\bar{X}_n - \mu| < t \sqrt{\frac{\sigma^2}{n}}\right] \geq 1 - \frac{1}{t^2} \Rightarrow \lim_{n \rightarrow \infty} P\left[|\bar{X}_n - \mu| < \epsilon\right] = 1$$

DEFINIZIONE 4: **Varianza campionaria**

$$s_n^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{X}_n)^2$$

TEOREMA: Dato un campione $\{x_1, \dots, x_n\}$ proveniente da una popolazione con $f(x)$, e posto s_n^2 si ha: $E[s_n^2] = \sigma^2$ $\text{var}[s_n^2] = \frac{1}{n} \left(\mu_4 - \frac{n-3}{n-1} \sigma^4 \right)$

con μ_4 e σ^2 rispettivamente il momento di ordine 4 e la varianza di $f(x)$

DISTRIBUZIONE MEDIA CAMPIONARIA

TEOREMA DEL LIMITE CENTRALE Sia data una popolazione distribuita con densità $f(x)$ avente media μ e varianza σ^2 finite. Dato \bar{X}_n la media di un campione casuale di dimensioni n estratto da esse, allora la media campionaria segue una distribuzione normale, al tendere di $n \rightarrow \infty$

Si può dire in modo equivalente che, detto Z la standardizzazione di \bar{X}_n , segue la distribuzione normale standardizzata:

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

$$\text{con } E[\bar{X}_n] = \mu \quad \text{var}[\bar{X}_n] = \frac{\sigma^2}{n}$$

Si vede graficamente che se il campione cresce ($n \rightarrow \infty$) la forma della distribuzione delle medie campionarie tende a quella di una normale.

Osservazione

Se la distribuzione della popolazione è normale, la distribuzione delle medie campionarie è normale, con qualsiasi dimensione campionaria; infatti, la media campionaria è una particolare combinazione lineare e qualsiasi combinazione lineare di v.c. normali, indipendenti, segue una normale.

Distribuzione di Fisher

DEF: siano $U \in V_2$ e $V \in V_m$ due indipendenti variabili entrambe distribuite χ^2 con m e n GDL la V.C.:

$$F = \frac{U/m}{V/m}$$

h_1 : una distribuzione F di FISHER con m e n GDL ovvero $F(m, n)$

Sto sopra i campioni S_1 e S_2 nelle 1^a e 2^a colonne $F(m, n)$ n 2^{da} che m 1^{da} colonna

Teorema: Se $\{x_1, \dots, x_n\}$ è un campione casuale di dimensione n estratto da una popolazione distribuita normalmente con media μ_x e varianza σ^2 e $\{y_1, \dots, y_m\}$ è un campione casuale di dimensione m estratto da una popolazione distribuita normalmente con media μ_y e varianza σ^2 e se i 2 campioni sono indipendenti, allora la V.C.

$$\sum_{i=1}^m (x_i - \bar{x}_m)^2 / (m-1)$$

segue una distribuzione

$$F_{m-1, n-1}$$

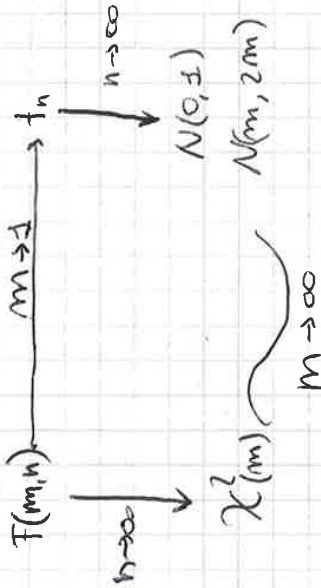
GDL \uparrow
GDL

$$\sum_{i=1}^n (y_i - \bar{y}_m)^2 / (n-1)$$

→ se voglio i percentuali ... allora si applica

$$f_{m, n, \alpha} = \frac{1}{f_{m, m, \alpha}}$$

Schemi mensurativo



STIMATORI PUNTUALI

METODO DEI MOMENTI Data una popolazione distribuita secondo la funzione di densità $f(\cdot, \theta_1, \dots, \theta_k)$ se i momenti $M'_k = E[X^k] = \int x^k f(\cdot, \theta_1, \dots, \theta_k)$ delle variabili casuali $X \sim f(\cdot, \theta_1, \dots, \theta_k)$, esistono finiti in vece di quello dei parametri θ_j che si desidera stimare, allora il metodo dei momenti consiste nel risolvere rispetto a $\theta_1, \dots, \theta_k$ le equazioni:

$$\begin{cases} M'_1(\theta_1, \dots, \theta_k) = M'_1 \\ M'_2(\theta_1, \dots, \theta_k) = M'_2 \\ \vdots \\ M'_k(\theta_1, \dots, \theta_k) = M'_k \end{cases}$$

contenute nelle prime k equazioni che si ottengono eguagliando i momenti campionari

$$M'_r = \frac{1}{n} \sum_{i=1}^n X_i^r$$

generati dal campione $\{x_1, \dots, x_n\}$ provenienti da $f(\cdot, \theta_1, \dots, \theta_k)$ ai momenti M'_r

loro varianza, e quindi, se preferisci lo stimatore con varianza minore perciò:

$$E\{T_1 | T_2\} = \frac{var(T_1)}{var(T_2)}$$

DEFINIZIONE: uno stimatore $T = (X_1, \dots, X_n)$ del parametro θ si dice **CONSISTENTE IN MEDIO** e **CONVOLGENTE** per il parametro θ se:

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

per esempio la media campionaria \bar{X}_n è uno stimatore corretto e consistente in medio quando della media μ della popolazione da cui proviene il campione poiché $var(\bar{X}_n) = \sigma^2/n \rightarrow 0$ per $n \rightarrow \infty$

DEFINIZIONE: Uno stimatore $T = (X_1, \dots, X_n)$ del parametro θ si dice **CONSISTENTE IN PROBABILITÀ** se con la legge in probabilità del parametro θ , cioè se per $\forall \epsilon > 0$ si ha

$$\lim_{n \rightarrow \infty} P\{|T - \theta| < \epsilon\} = 1$$

Dove la dipendenza dell'argomento del limite da n deriva dalla dipendenza dello stimatore dalla numerosità del campione

osservazione 1 la consistenza in medio equivale a dire che lo stimatore è consistente in probabilità;

osservazione 2 la media campionaria \bar{X}_n nella classe degli stimatori lineari è il migliore stimatore per la media della popolazione

osservazione 3 la varianza campionaria S_n^2 , oltre ad essere uno stimatore corretto della varianza σ^2 della popolazione da cui proviene il campione è anche consistente

Stima per intervalli

DEFINIZIONE: Si definisce intervallo di fiducia per il parametro θ un intervallo con cui è inclusa il vero valore del parametro con una prefissata probabilità **livello di fiducia**

In altri parole si tratta di trovare un prefissato grado di probabilità che il valore stimato del parametro si discosti dal suo valore vero per meno di una certa quantità, e in modo equivalente:

$$P[L_i < \theta \leq L_s] = 1 - \alpha$$

d. L_i = limite inferiore
 L_s = limite superiore
 $1 - \alpha$ = livello di fiducia
 ↑ intervallo di fiducia
 ↓ livello $1 - \alpha$
 - Più è basso **UNILATERALE** (inferiore superiore)
 - Bilaterale (simmetrico)

Stima per intervalli: il medio

La **popolazione** nota

Si è $\{X_1, \dots, X_n\}$ campione estratto da una popolazione avente medie μ incognite e varianza σ^2 nota

Stimatore più efficace per la media:

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

Si è visto che $E[\bar{X}_n] = \mu$ con $var[\bar{X}_n] = \frac{\sigma^2}{n}$, inoltre \bar{X}_n segue una distribuzione normale con parametri μ e σ^2/n

$$Z = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}$$

segue la distribuzione normale standardizzata con parametri $\mu = 0$ e $\sigma^2 = 1$

Possiamo scrivere che

$$P[-z_{1-\alpha/2} < Z \leq z_{1-\alpha/2}] = 1 - \alpha$$

Da cui dopo opportuni calcoli si ottiene

$$P\left[\bar{X}_n - z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}} < \mu \leq \bar{X}_n + z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}\right] = 1 - \alpha$$

Supponiamo $\{X_{1i}, \dots, X_{n_1}\}$ Campione casuale di n_1 dimensioni
 n_1 estratto da una popolazione distribuita normalmente
 con media μ_1 e varianza σ_1^2 e che $\{X_{2i}, X_{22}, \dots, X_{2n_2}\}$
 sia anch'esso (valevole)
 la miglior stimate x la differenza delle medie $\mu_1 - \mu_2$ è:

$$\bar{X}_1 - \bar{X}_2 = \frac{1}{n_1} \sum_{i=1}^{n_1} X_{1i} - \frac{1}{n_2} \sum_{i=1}^{n_2} X_{2i}$$

stimatori lineari

Per le varianze:

a) caso con varianze note

$$\sigma_{\Delta}^2 = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} \quad \rightarrow \quad Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sigma_{\Delta}}$$

$$P[-z_{1-\alpha/2} < \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sigma_{\Delta}} < z_{1-\alpha/2}] \quad \text{si ottiene}$$

$$(L_1, L_2) = \left[(\bar{X}_1 - \bar{X}_2) - z_{1-\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, (\bar{X}_1 - \bar{X}_2) + z_{1-\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right]$$

b) caso con varianze non note non utilizzabili =

Ognuno dei 2 campioni può essere usato per
 ottenere una stima per la varianza comune σ^2
 Il stimatori comuni per σ_1^2 e σ_2^2 sono:

$$S_1^2 = \frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (X_{1i} - \bar{X}_1)^2 \quad \text{per st. non. } \sigma_1^2$$

per st. non. σ_2^2

Siccome S_1^2 e S_2^2 sono entrambi stimatori comuni, possiamo
 usarli entrambi per ottenere la stima per σ^2 data da:

$$S_{\text{pool}}^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2} = \frac{\sum_{i=1}^{n_1} (X_{1i} - \bar{X}_1)^2 + \sum_{i=1}^{n_2} (X_{2i} - \bar{X}_2)^2}{n_1 + n_2 - 2}$$

$$s_{\text{pool}}^2 = \left(\frac{1}{n_1} + \frac{1}{n_2} \right)$$

Ne consegue che la v.c.

$$T = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{s_{\text{pool}} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \quad \rightarrow \quad n_1 + n_2 - 2 \text{ GDL}$$

segue la distribuzione t di Student con $n_1 + n_2 - 2$
 GDL
 quindi si può affermare che:

Media, Varianza incognite

Considero campione casuale $\{X_1, \dots, X_n\}$ di dimensioni n proveniente da una distribuzione normale, volendo sottoporre il test $H_0: \mu = \mu_0$ e σ^2 non nota con media μ_0/μ_0 e σ^2 non nota:

$$T = \frac{\bar{X}_n - \mu_0}{\frac{S}{\sqrt{n}}} = \frac{\bar{X}_n - \mu_0}{\frac{S}{\sqrt{n}}}$$

T STUDENT
(n-1) G.D.L.

Tabella

H_0	H_a	Rifutro H_0	Scheda grafica
$H_0: \mu = \mu_0$	$H_a: \mu \neq \mu_0$	$T > t_{n-1, 1-\alpha/2}$ o $T < -t_{n-1, 1-\alpha/2}$	
$H_0: \mu = \mu_0$	$H_a: \mu > \mu_0$	$T > t_{n-1, 1-\alpha}$	
$H_0: \mu = \mu_0$	$H_a: \mu < \mu_0$	$T < -t_{n-1, 1-\alpha}$	

CURVA CARATTERISTICA OPERATIVA E POTENTE DEL TEST

Procedimento per il calcolo del rischio di errore di II specie (ipotesione in torto l'ipotesi nulla)

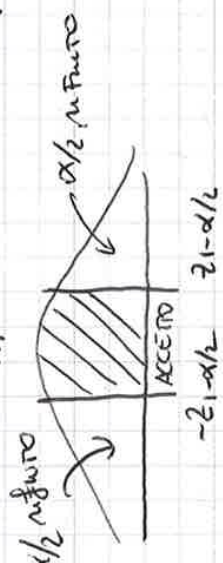
$$\beta = P(\bar{X}_n < \bar{X}_n \in \bar{X}_{n\alpha} | H_a: \mu)$$

Importante su di un diagramma cartesiano i valori β (SI) e $1-\beta$ (NO) in funzione dei valori dell' H_a si ottiene il grafico della curva caratteristica operativa

Il problema è quello di testare $H_0: \mu = \mu_0$ contro l' H_a alternativa $H_a: \mu \neq \mu_0$ (Test bilaterale). L'IP nulla H_0 è da rifiutare se il valore osservato di Z è $>$ di $z_{1-\alpha/2}$ o minore di $-z_{1-\alpha/2}$ quantile tale \times cui è data la potenza della curva caratteristica. ella sua dx e poi ed $\alpha/2$

REGIONE DI RIFIUTO:

$$dub: \left(\mu_0 - z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}} \right) < \mu < \left(\mu_0 + z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}} \right)$$



NOTA: oltre al test bilaterale esiste anche l'unilaterale superiore e inferiore

rispetto al bilaterale sono ventoseggiosi x gli errori di II specie

Tab monotone

H_0	H_a	Rifutro H_0	Scheda grafica
$H_0: \mu = \mu_0$ Bilaterale	$H_a: \mu \neq \mu_0$	$Z > z_{1-\alpha/2}$ $Z < -z_{1-\alpha/2}$	
$H_0: \mu = \mu_0$ monotone superiore	$H_a: \mu > \mu_0$	$Z > z_{1-\alpha}$	
$H_0: \mu = \mu_0$ monotone inferiore	$H_a: \mu < \mu_0$	$Z < -z_{1-\alpha}$	

TEST D'IP sulle proporzioni (ASINTOTICO)

- Campione Bernoulliano vs se campioni elaborato (p. referiamo approssimare)
 $H_0: p = p_0$ con livello di significatività α .

STATISTICA: $Z = \frac{p - p_0}{\sqrt{\frac{p(1-p)}{n}}}$ che segue dist. normale standard

se $|z_{calc}| > z_{1-\alpha/2}$ si rifiuta H_0 a favore di H_1

Se sottoposto il test: $H_0: p_1 = p_2$ (con n_1 e n_2 due campioni indipendenti di 2 popol.)

Si costruiscono 2 stime di max verosimiglianza p_1 e p_2 . lo stimatore differente $p = p_1 - p_2$ essendo le dist. di 2 v.c. indipendenti popol. ed indipendenti segue ancora una distribuzione normale con media $p_1 - p_2$ e

varianza $\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}$

Se H_0 vera ($p_1 = p_2 = p$) si ha:

$$P_{H_0}(0, \sqrt{p(1-p)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)})$$

si rifiuta H_0 se $|z_{calc}| > z_{1-\alpha/2}$

Regression
OUTPUT

Results for: data - A

Regression Analysis: Y1 versus X1
 ↳ X_1 INDIPENDENT VARIABLE

The regression equation is:

$$Y1 = 9,235 + 9932 X1$$

PREDICTION CONSTANT	COEF	SE COEF	T	P
X_1	9,23471	0,08288	2,83	0,006
	9,93224	0,02758	33,80	0,000

$S = 9,300090$
 ↳ $R-Sq = 93,2\%$
 ↳ $R-Sq(Adj) = 93,1\%$

Analysis of variance:

Source	DF	SS	MS	F	P
Regression	1	102,89	102,89	114,50	0,000
Residual Error	84	7,56	0,09		
Lack of Fit	79	7,51	0,10	9,16	0,010
Pure Error	5	0,05	0,01		
Total	85	110,45			

77 rows with no replicates

over zero per entrare nel coefficienti significativi
 ↳ $H_0: \beta_1 = 0$
 ↳ $H_1: \beta_1 \neq 0$

$R^2 = 93,2\%$ OTTIMO PERCENTUALE DI VARIABILITÀ SPiegata dal modello che ha come variabile indipendente X_1

$H_0: \beta_1 = 0$
 ↳ $F_{0,1,84} = 2,00$
 ↳ $F_{0,1,84} = 9,16 > 2,00$
 ↳ H_0 è rifiutata

↳ $F_{0,1,84} = 9,16$

conclusione: si può dire che il modello è significativo

↳ se otteniamo $S.F.$ con il livello di significatività, rifiutiamo l'ipotesi nulla

Ma non ci dice quale è il modo
più bello di costruire.

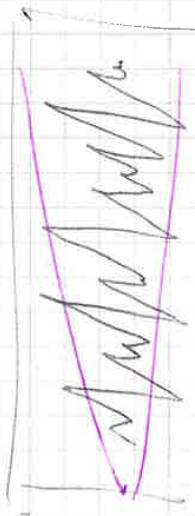
↳ Il nuovo modo di costruire è:

↳ sono tutte in termini del 2° ordine
e nei grafici si evidenzia un

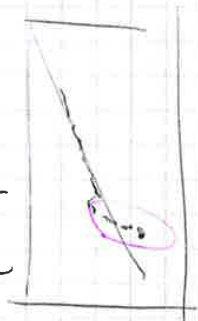
↳ stame del 2° ordine nel verso PTS



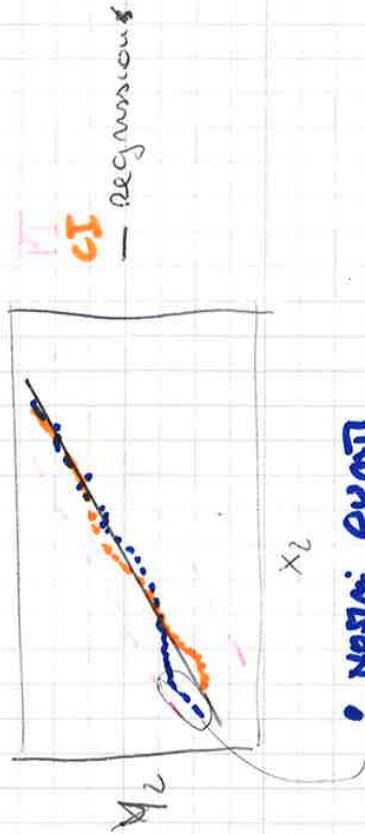
e anche nell'or versus energia, si
vede che conobilità non costante,
aumentare.



e anche l'ossimazione di rosobato
con po' farbare



Fitted line Plot, si evidenzia l'intervento di
previsione e l'intervento di confidenza



• **nostri punti**

Si vede che la retta di regressione non coglie
i nostri PT, ci sono fuori da PT le parti, ~~la parte~~
(e il 95% area storica!)
non stavo anche quasi fuori dall'intervento
di previsione

e poi o nono abbiamo come 50% sotto e 50% sopra
 di PT (UNITAS & COLLEGA con i spettro)
 ↳ ad occhio vedo due t e - il 50% sta sotto
 e 50% sta sopra quindi ok, ne la simmetria
 e due: questa parte si sono centron ne
 un po deviat verso il basso, qui invece:
 ↳ se ne centron ne deviat verso l'alto e
 PSR ha coordinate una nota ≠

OUT LAYER
 vanno tutti sotto controllo
 → possono essere eventi
 di ballatura
 → possono essere errori di
 campionamento

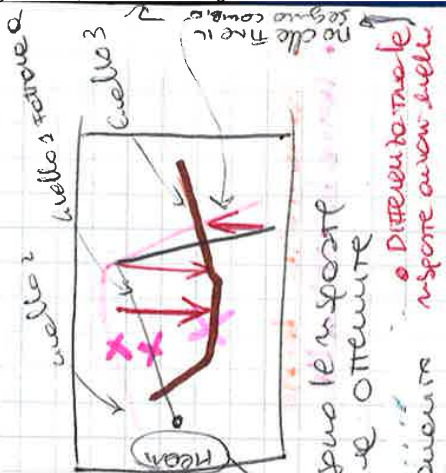
NOTRE SI INTRODUCE UN
 TROUPO (PARTE TUTT SU E POI TUTT GIU) → NON È
 una bella cosa perché i residui sono determinati
 delle variabili casuali e proprio perché in
 U.C. non devono seguire alcun trend

NOTRE, anche la variabile non è costante
 di residui

QUESTO GRAFICO X QNT DELLO NON MI PIACE
 MOLTO, non mi sembra molto tracion di una
 variabile casuale. Ma anche con residuo = 0
 e variabile costante!

GRAFICO INTERAZIONI

3 Spettate che corrispondono
 ai 3 livelli del fattore 2
 al variare dei livelli del
 fattore 1



questo sono le risposte
 MMEORE OTTENUTE
 ↳ Differente tra le
 risposte con variabili
 ↳ Spettate del
 livello 3 che tutti altro
 ↳ SIMILIA
 ↳ Spettate del
 livello 3 che tutti altro

DF	SS	MS	F	P
2	843,6	421,8	2,57	0,097
3	451,7	150,57	9,25	0,000
6	470,5	78,75	4,78	0,002
24	3936,3	164,0		
35	1403,1			

$R^2 = 71,95\%$
 $R^2_{adj} = 59,10\%$
 $S = 12,81$

Se la risposta

è la /
 di variabile casuale
 nelle risposte y
 spiegate dai 2
 fattori e della
 loro interazione

↳ quello che rimane
 da attribuire
 all'errore

$SS_{tot} = SS_{a} + SS_{b} + SS_{interazione}$
 SS_{tot}

ne ho tante
 quanti sono i livelli
 del fattore colonne
 e quanti del
 fattore righe

GRAFICA

1° grafico
 rappresenta
 come le interazioni
 degli eventi casuali. E i

↳ se sono vere le interazioni
 fra tutte, questi residui
 dovrebbero essere
 randomizzati in una U.C.
 che segue
 una distribut
 normale

↳ errore
 ↳ aspetto
 dell'ordine di
 osservazione

↳ verso righe +
 ↳ osservazione ordine



no in realtà bisogna che non sia concettuale solo con \neq .

6 e questo è importante l'alta significatività a questo fattore

14) quando si verifica la questa situazione uso i

CONTRASTI

per indagare meglio cose capite

o CONTRASTO μ_i è un parametro ordinato

$T = \sum_{i=1}^k c_i \mu_i$ cambiamento livello delle note tale per cui $\sum_{i=1}^k c_i = 0$

nel no output \times vedere il problema su b

STATORIA μ_i

$$\frac{\mu_1 + \mu_2 + \mu_3}{2} = \mu_4$$

$$\frac{1}{3} + \frac{1}{3} - 1 = 0 \text{ è un CONTRASTO}$$

le distri. binome dello stator C_i è una variabile μ_i QNT è una combinatori di μ_i normali

$$C \sim N(\mu_c, \text{var}[c])$$

$\mu_c = E[\sum_{i=1}^k c_i \mu_i] = T$ e uno stator con μ_i

$$\text{var}[c] = \text{var}[\sum_{i=1}^k c_i \mu_i] = \sum_{i=1}^k c_i^2 \frac{\sigma^2}{n b}$$

esad²

ovvero $\mu_i \in O$ POSSO FARE TESTING e INTERPRETARE DI CONGIUNTE

NOTA
 μ_i colonne
 μ_i valori nelle celle
 il primo DNE che lo faccio \times tutti le note

grafici: 60

COME AGIRE

~~Analisi della varianza:~~

~~ANOVA ONE-WAY~~

N° GDL DEL LACK OF FIT

$$n-2 = n-k + 1$$

$$\hookrightarrow n-2 - n-k \Rightarrow k-2$$

$$\frac{SS_{LOF}}{k-2} = MS_{LOF} = \hat{\sigma}^2$$

stima covariate di σ^2

reg. solo sotto l'HP

nulle che il

modello sia adeguato

$\sigma^2 + bias$

attinenti

e

ho 2 stime di σ^2 una senza covariate dovuta al modello inadeguato e una con le covariate

le confronti

sono i modelli adeguati, altrimenti NO

H₀: modello adeguato
H₁: non "

$$\frac{MS_{LOF}}{MS_{pe}} \sim 1 \text{ se modello adeguato}$$

F_{k-2, n-k}

$$y_i - \hat{y}_i = \underbrace{\{ (y_i - \hat{y}_i) - E(y_i - \hat{y}_i) \}}_{R_i} + E(y_i - \hat{y}_i)$$

errore sistematico

$$MS_{pe} = \hat{\sigma}^2$$

stima covariate di σ^2 (se il modello è adeguato, altrimenti non)

$$SS_{LOF} = \sum_{i=1}^k (y_i - \hat{y}_i)^2$$

$$SS_{GDL} = \sum_{i=1}^k (n_i - 1) = (n-k)$$

no tot punti

$$SS_{res} = SS_{pe} + SS_{LOF}$$

BACK OF FIT \rightarrow non conto di σ^2 per costruzione
Senza distorsione di stime del modello vero

i DATI \rightarrow R^2 ADJUSTED che per un fenomeno è $\approx R^2$ ma per più fenomeni è migliore perché tiene conto delle interazioni delle variabili usate

S \rightarrow non usare una F-test è sempre coerente lo $\sqrt{D1}$ HS Error

CI \rightarrow intervalli di confidenza per le medie

POOLED ST DEV \rightarrow media delle ST dev dei singoli Fattori

TWO WAY ANOVA **Parte 1**

\rightarrow una caratteristica o misura 2 Fattori

STESSA CONDIZIONE:

$H_0: \mu_1 = \mu_2 = \dots = \mu_k$
 $H_1: \mu_i \neq \mu_k$

DF \rightarrow Fattori: $1^{\circ} \rightarrow (K-1) \rightarrow x \rightarrow$ n° di celle
 $2^{\circ} \rightarrow (n-1)$

TOT: valore \rightarrow valore \rightarrow n° dati
 se \neq vuol dire che ho più dati in ogni cella (REPLICAZIONI P)

Problema \rightarrow se qualche Fattori ha qualche valore $\leq 10\%$ è abbastanza significativo che vuol dire che le differenze sulle risposte che si ottengono cambiano i livelli del Fattore in questione, non è un dato significativo quanto visto

lo si può confrontare anche con SS che in quel Fattori sono "accidentalmente" < due regni altri \rightarrow viene ed il p-value $\leq 5\%$ \rightarrow ottenere significato

Problema \rightarrow INTERACTION \rightarrow se ottenere significativo vuol dire anche avere l'INTERACTION plot per verificare in che modo influisce sulle risposte

$R^2 \rightarrow$ (stessa def) = $\frac{SS_{Fattori1} + SS_{Fattori2} + SS_{Interazione}}{SS_{TOT}}$
 $R^2_{ADJ} \rightarrow$ ("")
 Nota: $SS_{TOT} = SS_{Fattori1} + SS_{Fattori2} + SS_{Fattori3} + SS_{Error}$

S \rightarrow $\sqrt{MS_{F1}}$

Parte 2: Intervalli di FIDUCIA

per verificare al livello di FIDUCIA indicato sopra ogni intervallo se essi sono tra loro concorrenti

se sono tutti conc. $H_0: \mu_1 = \mu_2 = \dots$ NON può essere rifiutare a tale livello di FIDUCIA
 se non tutti concor. \rightarrow uso i confronti per capire (F)

Problema \rightarrow anche per altre cose si distaccano (per vedere pg 255)

MAIN EFFECTS PLOT:

Distingue le reazioni delle risposte ottenute rispetto ai livelli del 1° e del 2° fattore separatamente
 ↳ in presenza di interazione è utile vedere in che modo ottenere le risposte nominali (es: se il pt più alto del grafico corrisponde ad un livello 200 del fattore A come dire che posso ottenere il max della risposta con quel fattore)
 ↳ se presente interazione tra i fattori ed esse è significativa, bisogna prima controllare l'INTERACTION PLOT



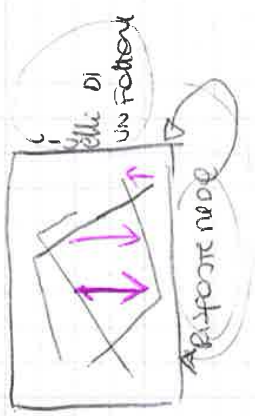
INTERACTION PLOT

spiega che corrispondono ai valori del fattore e al valore del fattore b (o viceversa)

↳ se linee ~~non~~ c'è interazione
 se linee non // c'è interazione

Per vedere come ottenere le risposte max con un certo livello del fattore in base a bente vedere quella delle linee assegnate solo + in alto.

• Differente le risposte ai vari livelli, si comporta PT segue cambio, interazione forte



REGRESSIONE

(lineare semplice)

→ ha equazione della retta con 1 variabile dipendente e una (o più) indipendente (o +)
 → sono ho in colonne i vari predictor e viene x estronabili restere Ho: $\beta_0 = 0$ Ha: $\beta_0 \neq 0$
 e Ho: $\beta_1 = 0$ Ha: $\beta_1 \neq 0$
 → se p-value altamente significante rifetto H_0 (e qas)

↳ se $COEF = \sqrt{Var(\hat{\beta}_i)}$ → standard error
 abbiamo sempre J, R-Sq, R-sq (ADJ)
 → posso avere: F = T-Test (che mi conferma che mi interesso fare un test di H₀)

Nelle detone dell'analisi della varianza: (STUDIO ^{di base} ^{di base} ^{di base})
 P-value del LACK OF FIT → se non significativo
 modello buono (è un indicatore delle bente del modello → non posso rifiutare H₀: modello adeguato)

↳ confrontare con R-Sq

NOTA: più ovvio → lo posso calcolare LACK OF FIT → Ss: lo ottengo con differenza tra Residual error e pure error

NOTA: se ho il T-Test il p-value è: $2P(T_{n-2} > |t_{0.05}|)$
 ↳ sistema $b_i - H_0 / \sqrt{Var(\hat{\beta}_i)}$

NOTA: P-value regression altamente significativo indica che non rifiutare H₀ con un certo livello di confidenza

Residui: $\sum_{i=1}^n (y_i - \hat{y}_i) = 0$ NOTA
 $E_i = \{ (y_i - \hat{y}_i) - E[y_i - \hat{y}_i] \}$ + $E[y_i - \hat{y}_i]$

DIFFICILMENTE
 valori osservati
 delle risposte e
 valori previsti del
 modello

costante
 ↳ errore sistematico
 introdotto dall'
 utilizzo di un modello
 non adatto

$E_i = R_i + d_i$

variabile
 consideri in quanto
 differenza tra
 V.C. $(y_i - \hat{y}_i)$ e il suo
 valore atteso

inoltre $\sum_{i=1}^n R_i^2$ segue una distribuzione χ^2_{n-2}
 con valore atteso $(n-2)\sigma^2$

elencos el quoziente tutti i numeri si ha

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n R_i^2 + \sum_{i=1}^n d_i^2$$

DA cui: $E[\sum_{i=1}^n (y_i - \hat{y}_i)^2 / (n-2)] = \sigma^2 \pm \sum_{i=1}^n d_i^2$

↳ LA somma dei quadrati componenti alla
 varianza totale (SMA DI σ^2) / INDEPENDENT
 delle regressioni \Rightarrow PURE ERROR \Rightarrow SI C'è STIMA
 CONVERTE OIGAR
 SPER

$MSE = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2}$ MEAN SQUARE

$MS_{PE} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-k}$ MEAN SQUARE

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n [(y_i - \bar{y}) - (\hat{y}_i - \bar{y})]^2$$

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

misura
 della variazione
 totale

misura
 della variazione
 delle risposte
 previste nel
 modello
 rispetto alle
 risposte
 osservate

SSreg
 SSreg

SSres
 SSres

SS Tot
 SS Tot

SS Tot = SSreg + SSres

SSreg \rightarrow scrivere C
 e tutti
 i dati
 con
 sulle
 righe
 scrivere
 la
 colonna
 non occupata

$R^2 =$ frazione delle
 variazioni reali di data spiegata
 rispetto al loro valore medio
 spiegato dal modello di regressione

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{SSreg}{SS Tot} = R^2 \leq 1$$

$R^2_{ADJ} = 1 - \frac{MSE}{SS_{tot}/n-1}$

$R^2 = r^2_{xy}$

Correlazione con il quoziente
 del COEF di correlazione lineare
 tra y e \hat{y} (non osservabile e
 variabile INDIP)

previsi con un
 modello
 con la regressione
 con i valori
 osservati

$\hat{y} = b_0 + b_1 x$