



Appunti universitari
Tesi di laurea
Cartoleria e cancelleria
Stampa file e fotocopie
Print on demand
Rilegature

NUMERO: 2140A-

ANNO: 2017

A P P U N T I

STUDENTE: Pinna Eleonora

MATERIA: Statistica - Prof. Pellerey e 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.**

METODO DEI MINIMI QUADRATI

DOMANDE D'ESAME UVALE

• Il metodo dei quadrati è una tecnica di ottimizzazione che permette di trovare una funzione rappresentata da una curva di regressione che si avvicini ad un insieme di dati.

La funzione trovata deve essere quella che minimizza la somma dei quadrati delle distanze tra i dati osservati e quelli della curva che rappresenta la funzione stessa.

oppure

• È una tecnica di ottimizzazione grazie alla quale è possibile determinare una stima $\hat{y} = b_0 + b_1 x$ della retta di regressione $y = \beta_0 + \beta_1 x + \epsilon$ minimizzando la somma dei quadrati degli errori casuali $\epsilon_i = y_i - \hat{y}_i$ tra i valori delle variabile casuale y e i valori previsti con la retta in corrispondenza del valore della variabile indipendente x .

$$\begin{cases} b_0 = \bar{y} - b_1 \bar{x} \\ b_1 = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2} \end{cases}$$

RESIDUI

Sono la realizzazione della nostra variabile casuale, erro casuale, se è davvero casuale.

errore casuale = ϵ . Se manca di tutti effetti di errore, di tutte le variabili non considerate nel modello quindi la componente casuale omiva a rappresentare altri errori di variabile.

$$\text{residui} : \sum_{i=1}^n (y_i - \hat{y}_i) = 0$$

Se elevo al quadrato:

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 = \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}_{\text{var. dovuta alla regressione}} + \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}_{\text{variabile errore residua}}$$

\downarrow SS_{tot} \downarrow SS_{reg} \downarrow SS_{res}

Differenza tra valore osservato e il valore predetto dal modello

Formula di Bayes /

Se usa tutte e volte in cui abbiamo un esperimento o fatto successive, viene usata per calcolare la P di una causa che ha scatenato l'evento verificato.

Sia E_1, E_2, \dots, E_n una collezione di eventi incompatibili ed esaurivi e sia $P[E_i] \neq 0$ per ogni i .

Qualora sia $F \subseteq S$ si ha:

$$P[A_i | B] = \frac{P(A_i) \cdot P(B | A_i)}{\sum_{j \in I} P(A_j) \cdot P(B | A_j)}$$

Definizione Probabilità classica /

La probabilità di un evento E è definita come il rapporto tra il numero n_E di risultati favorevoli (cioè il numero dei risultati che determinano E) ed il numero dei risultati possibili:

$$P[E] = \frac{n_E}{n} \rightarrow \text{perché tutti i risultati sono ugualmente possibili e tra loro incompatibili.}$$

Definizione Probabilità frequentista /

Si definisce come frequenza assoluta n_E di un evento E il numero delle volte in cui si è ripetuto l'esperimento nelle medesime condizioni, il rapporto

$$P(A) = \lim_{n \rightarrow \infty} \frac{n_A}{n} \rightarrow \text{valida nell'ipotesi di ripetere esp. all'infinito}$$

presentato l'evento favorevole della N il numero delle volte in cui si è ripetuto...

Probabilità assiomatica e Kolmogoroff /

Una funzione di probabilità P è una funzione di insieme che ha come dominio lo spazio degli eventi e come codominio l'intervallo $[0, 1]$ e che soddisfa i seguenti assiomi:

- $P(A) \geq 0$ sempre $\rightarrow 0 \leq P(A) \leq 1$
- $P(S) = 1$, S = spazio campione
- Per la successione di eventi A_1, A_2, \dots, A_n , tra loro mutuamente esclusivi (incompatibili), se l'unione infinita è allo spazio di eventi A , allora:

$$P[\cup_{i=1}^{\infty} E_i] = \sum_{i=1}^{\infty} P[E_i]$$

CONSEQUENZE

- $P[\bar{E}] = 1 - P[E] \rightarrow P[\emptyset] = 0$
- $P[E_1 \cup E_2 \cup \dots \cup E_n] = \sum_{i=1}^n P[E_i]$
- $P[E \cup F] = P[E] + P[F] - P[E \cap F]$

ERRORI PRIMA E SECONDA SPECIE

Si commette errore di prima specie quando rifiuto l'ipotesi nulla quando questa è invece corretta.

Si commette errore di seconda specie quando si accetta l'ipotesi nulla, nel caso in cui questa sia invece falsa.

L'errore di prima specie è identificato con α , quello di seconda specie con β .

PROBABILITÀ CONDIZIONATA

Sia dato lo spazio di probabilità (S, \mathcal{A}, P) . Sia $B \in \mathcal{A}$, tal per cui $P(B) > 0$.

Dato l'evento $A \in \mathcal{A}$, è detta "probabilità di A condizionata a B ":

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

→ esprime la probabilità che si verifichi A condizionata al fatto che si verifichi B .

- Non sempre una prob. condiz. diminuisce.

DISTRIBUZIONE POISSON (APPROSS. BINOMIALE)

È il modello adatto per modellare il numero di difetti o non conformità che si trovano in una unità di prodotto.

È una distribuzione discreta che esprime la probabilità per il n° di eventi che si verificano successivamente ed indipendentemente in un dato intervallo di tempo, sapendo che mediamente se ne verificano λ .

$$p_X(x) = p_X(x; \lambda) = \begin{cases} \frac{e^{-\lambda} \cdot \lambda^x}{x!} & \text{per } x=0, \dots, \infty \\ 0 & \text{altrove} \end{cases}$$

È una funzione sempre > 0

$E[X] = \lambda$ e $VC[X] = \lambda$, significa che la dist. di Poisson è tanto più dispersa quanto più è grande la sua media.

\sim Bin → Se consideri una parametri n e p , se il n° delle prove tende all'infinito e la prob. di successo tende a 0, in modo che np rimanga costante, si ha:

$$\binom{n}{x} p^x (1-p)^{n-x} \rightarrow \frac{e^{-\lambda} \cdot \lambda^x}{x!}$$

IPERGEOM \sim BIN

Quando abbiamo un campionamento con rimpiazzamento, la distribuzione che modella il n° di successi è la Bin.

Se il campionamento è senza rimpiazzamento usiamo la dist. ipergeo.

→ se per π e K grandi e se $\frac{K}{\pi}$ è costante, posto $\frac{K}{\pi} = p$

Torino \sim Bin.

TASSATI E TAVITI

Supponiamo che x sia una variabile casuale relativa all'esperimento con funzione di ripartizione F e densità f

Leviamo n repliche indipendenti dell'esperimento semplice e otteniamo un campione (x_1, \dots, x_n) di v.c. indipendenti.

Utilizzo le statistiche d'ordine per mettere in ordine i dati

$$Z = \min\{x_1, \dots, x_n\}$$

$$Y = \max\{x_1, \dots, x_n\}$$

$$F_Y(t) = P[Y \leq t] = P[x_1 \leq t \dots x_n \leq t] = \prod_{i=1}^n F_i(t)$$

- in particolare se $F_i = F \rightarrow F_Y(t) = (F(t))^n$
 $f_Y(t) = n F^{n-1}(t) \cdot f(t)$

$$F_Z(t) = P[\min\{x_1, \dots, x_n\} \leq t] = \dots = 1 - \prod_{i=1}^n (1 - F_i(t))$$

- in particolare se $F_i = F \rightarrow 1 - F_Z(t) = (1 - F(t))^n$

~~CAMPIONE~~

~~Sistema di n variabili casuali indipendenti, con stessa distribuzione~~

~~$$\{x_1, \dots, x_n\}$$~~

CONTRASTO

Si definisce contrasto una combinazione lineare delle medie dei trattamenti, in modo che i coefficienti sommati diano 0.

$C \rightarrow$ lo uso come media comp.

$V[C] \rightarrow$ faccio tabella con repliche e poi moltiplico

Poi standardizzo

Legge dei grandi numeri!

Se la numerosità n delle prove tende all'infinito, la probabilità che \bar{x}_n assuma valori di fuori di $(\mu - \epsilon, \mu + \epsilon)$ tende a 0, qualunque sia l'ampiezza ϵ dell'intervallo. In formula

$$\lim_{n \rightarrow +\infty} P[|\bar{x}_n - \mu| > \epsilon] = 0$$

Equivalentemente:

Se n tende ad infinito, la probabilità che la media campionaria \bar{x}_n converga alla media comune μ delle x_i è uguale a 1, ossia per ogni $\epsilon > 0$,

$$\lim_{n \rightarrow +\infty} P[|\bar{x}_n - \mu| < \epsilon] = 1$$

Formula probabilità totale:

etale che $U_{A_i} = S$ e $A_i \neq \emptyset$

Sono A_1, A_2, \dots, A_n eventi incompatibili. \forall B un altro evento qualsiasi, $B \subseteq S$

$$\text{Vale } P[B] = \sum_{i=1}^n P(A_i) \cdot P(B|A_i)$$

Ipotesi di Poisson:

Affinché un fenomeno in cui è in atto un conteggio possa essere modellizzato mediante la dist. di Poisson, devono essere rispettate 3 ipotesi:

1) Esiste quantità λ positiva tale che la probabilità che si verifichi esattamente un evento in un piccolo intervallo di lunghezza Δt sia approssimativamente uguale a $\lambda \Delta t$:

$$P[\text{esattamente 1 evento nell'intervallo di tempo } \Delta t] = \lambda \Delta t + o(\Delta t)$$

2) La probabilità che nell'intervallo Δt si verifichino più di un evento è trascurabile (è più alta) rispetto alla probabilità che non ne verifichi esattamente 1:

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

$$P[0 \text{ eventi in } \Delta t] = 1 - P[1 \text{ evento}] - P[2 \text{ o più eventi}] = 1 - (\lambda \Delta t + o(\Delta t) - o(\Delta t))$$

3) L'evento E_1 rappresente il verificarsi di n_1 eventi in qualche s intervallo (t_1, t_2) , è indipendente dall'evento E_2 rappresente il verificarsi di n_2 eventi in un intervallo (t_3, t_4) non contenuto in (t_1, t_2) :

$$P[n_1 \text{ eventi in } (t_1, t_2) \text{ e } n_2 \text{ eventi in } (t_3, t_4)] = P[E_1] \cdot P[E_2]$$

R^2 FORMULA E SPIEGAZIONE

$$SS_{TC} = SS_{reg} + \text{variaz. residue}$$

$$R^2 = \frac{SS_{reg}}{SS_{TC}} \rightarrow \text{indice di determinazione multiplo}$$

È sempre compreso $0 \leq R^2 \leq 1$, esprime la frazione della variabilità totale dei dati sperimentali rispetto al loro valore medio spiegato dal modello di regressione

- $R^2 = 0 \rightarrow$ la variabilità dovuta alla retta di regressione è nulla, perciò la retta di regressione è parallela all'asse delle ascisse. Non ha dunque nessuno scopo interpretativo.
- $R^2 = 1 \rightarrow$ la variabilità dovuta alla regressione coincide con quella totale, mentre quella residua è nulla. In questo caso tutti i punti sperimentali sono collocati sulla retta di regressione e quindi tutti allineati.
 ↓
 situazione ideale.
- $R^2 = \frac{0}{0} \rightarrow$ tutti i punti osservati sono allineati parallelamente all'asse x e la retta di regressione è $y = y_i = \bar{y}$

INDICE DI CORRELAZIONE LINEARE

Si usa per analizzare la relazione tra due variabili:

$$r_{xy} = \frac{\text{Cov}[x, y]}{\sqrt{\text{Var}[x] \cdot \text{Var}[y]}} \rightarrow -1 \leq r_{xy} \leq 1$$

Se:

- $r_{xy} > 0 \rightarrow x$ e y sono correlate positivamente.
- $r_{xy} = 0 \rightarrow x$ e y sono incorrelate. (sono indipendenti)
- $r_{xy} < 0 \rightarrow x$ e y sono correlate negativamente.

COVARIANZA:

Indice che permette di verificare una relazione lineare tra due variabili statistiche, x e y .

$$\text{Cov}[x, y] = \begin{cases} \text{continue} \rightarrow \iint (x - E[x])(y - E[y]) - E_{xy}(x, y) \, dxdy \\ \text{discrete} \rightarrow \sum_{x_i} (x_i - E[x])(y_i - E[y]) - P_{xy}(x, y) \end{cases}$$

- $\text{Cov} > 0 =$ correlazione positiva
- $\text{Cov} < 0 =$ correlazione negativa
- $\text{Cov} = 0 \rightarrow$ incorrelate

• Def. Assiomatrica → \mathcal{A} = Spazio degli eventi = $\{A_i, A_i \subseteq S\}$
o algebra

La terna (S, \mathcal{A}, P)
 che è lo spazio
 di probabilità

se:

- 1) $S \in \mathcal{A}$
- 2) $A \in \mathcal{A} \rightarrow \bar{A} \in \mathcal{A}$
- 3) $\{A_n, n \in \mathbb{N}, A_n \in \mathcal{A}\} \rightarrow \bigcup_{n=1}^{\infty} A_n \in \mathcal{A}$

dove S = insieme di risultati possibili:
 o, meglio, spazio campione

Legge di Morgan : $\overline{A \cap B} = \bar{A} \cup \bar{B}$
 $\overline{A \cup B} = \bar{A} \cap \bar{B}$

DEFINIZIONE PROBABILITÀ Kolmogorov

Dato un esperimento, dato S e dato \mathcal{A} , diciamo
 probabile una funzione

$$P: \mathcal{A} \rightarrow \mathbb{R}$$

$$P: A \in \mathcal{A} \mapsto P(A) \in \mathbb{R}$$

tale che:

- i) $0 \leq P(A) \leq 1$
- ii) $P(S) = 1$
- iii) Presa la successione di eventi A_1, A_2, \dots, A_n , tra di loro
 mutuamente esclusivi, vale $P(\bigcup_{n=1}^{\infty} A_n) = \sum_{n=1}^{\infty} P(A_n)$

PROPRIETÀ ELEMENTARI DELLA PROBABILITÀ:

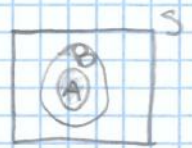
• $A \in \mathcal{A}$ $P(\bar{A}) = 1 - P(A)$ → infatti $A \cup \bar{A} = S$ e
 $P(A \cup \bar{A}) = P(S) = 1$



• $P(\emptyset) = 0$ perché $P(\bar{A}) = 1 - P(A)$
 $P(\bar{S}) = 1 - P(S) = 1 - 1 = 0$ \emptyset = insieme vuoto

• $A, B \in \mathcal{A}$ $A \subseteq B \rightarrow P(A) \leq P(B)$

Dimostrazione: $P(B) = P(A \cup [B-A]) = P(A) + P(B-A) \geq P(A)$
 perché $P(B-A) > 0$



• Per $A, B \in \mathcal{A}$ vale: $P(A \cup B) = P(A) + P(B) - P(A \cap B)$



va tolto perché
 altrimenti andrebbe
 calcolato due volte

ELEMENTI STOCASTICAMENTE INDIPENDENTI

Quando la probabilità $P(A|B) = P(A)$ allora A e B sono detti stocasticamente indipendenti.

- Vera solo se $P(B) \neq 0$

- lo sono quando B in realtà non condiziona A

Per verificare che sono stoc. indipendenti:

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

PER ESERCIZI: se sono stoc. indipendenti, posso ricavare l'intersezione da formula sopra. se non lo sono mi deve essere dato nel testo.

NOTA! indipendente $\not\leftrightarrow$ incompatibilità

esempio $\rightarrow A = \{pari\}$

$B = \{ \leq 4 \}$

sono indipendenti ma non incompatibili

$A \cap B = \{2, 4\}$

generalmente se sono incompatibili non sono indipendenti, sono sempre entrambi quando $P(A) \text{ o } P(B) = 0$ e quindi $A \text{ o } B \text{ e } \emptyset$

INDIPENDENZA PER 3 o più elementi

A, B, C

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

ii) $\begin{cases} P(A \cap B) = P(A) \cdot P(B) \\ P(A \cap C) = P(A) \cdot P(C) \\ P(B \cap C) = P(B) \cdot P(C) \end{cases}$

i $\not\leftrightarrow$ ii \rightarrow Dimostrazione:

$P(A \cap B \cap C) = P(\emptyset) = 0 = P(A) \cdot P(B) \cdot P(C)$

Pero: $P(A \cap B) = \frac{1}{2} \neq P(A) \cdot P(B) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$

oppure

$P(A \cap B) = P(\{b\}) = \frac{1}{4}$ e $P(A) \cdot P(B) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$
 $P(A \cap C) = P(\{a\}) = \frac{1}{4}$ e $P(A) \cdot P(C) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$

$S = \{a, b, c, d\}$
 $A = \{a, b\}$ $B = \{a, b\}$
 $C = \emptyset$

$S = \{a, b, c, d\}$
 $A = \{a, b\}$
 $B = \{b, c\}$

ma: $P(A \cap B \cap C) = P(\emptyset) = 0 \not\leftrightarrow P(A) \cdot P(C) \cdot P(B) = \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{8}$ $C = \{a, c\}$

FORMULA DI BAYES

Sia dato (S, \mathcal{G}, P) . Sia data la partizione di S

$$P(A_j | B) = \frac{P(A_j) \cdot P(B | A_j)}{\sum_{j \in I} P(A_j) \cdot P(B | A_j)}$$

Si usa tutte le volte in cui abbiamo un esperimento a fase successive. (es. cosa è successo a meta?)

Dim:

$P(E_1 \cap E_2) = P(E_1) \cdot P(E_2 | E_1)$ oppure $P(E_1 \cap E_2) = P(E_2) \cdot P(E_1 | E_2)$
ossia che entrambi hanno primo membro uguale, posso quindi eguagliare 1° e 2°

$$P(E_2) \cdot P(E_1 | E_2) = P(E_1) \cdot P(E_2 | E_1)$$

da cui ottengo

$$P(E_1 | E_2) = \frac{P(E_1) \cdot P(E_2 | E_1)}{P(E_2)} \quad \text{oppure} \quad P(E_2 | E_1) = \frac{P(E_2) \cdot P(E_1 | E_2)}{P(E_1)}$$

- Usato per calcolare la P di una causa che ha scatenato e' evento verificato.

Funzioni MISURABILI

$$(s_1, a_1) \xrightarrow{x} (s_2, a_2)$$

se hanno 2 esperimenti

$$x: S_1 \rightarrow S_2$$

x grande è una funzione definita tra S_1 e S_2

osservo che x sia misurabile se preso un qualsiasi evento $A_2 \in \mathcal{A}_2$ la sua controimmagine $x^{-1}(A_2)$ è un evento nello spazio \mathcal{A}_1 , cioè $x^{-1}(A_2) \in \mathcal{A}_1$

esempio sul quale

Definizione:

Dato (S, \mathcal{A}, P) e dato (R, \mathcal{B}, \dots) è detta variabile casuale una funzione $x: S \rightarrow R$ $x: S \rightarrow x(s) \in R$ se tale funzione è misurabile.

$\mathcal{B} \rightarrow$ algebra di Borel (sottinsieme dei reali) è il più piccolo spazio degli eventi contenente tutti gli intervalli (a,b) , $[a,b]$, $[a,b)$, $(a,b]$ e le loro unioni ed intersezioni infinite.

Dato la variabile casuale x è poi possibile definire una probabilità P_x così:

$$\forall B \in \mathcal{B} \quad P_x(B) = P[x^{-1}(B)] = P[x \in B]$$

è σ ben definita grazie alla misurabilità

È detta variabile casuale:

$$x: S \rightarrow R$$

$$x: s \in S \rightarrow x(s) \in R \quad \text{tale che sia misurabile.}$$

Come si descrivono le variabili casuali?

Notazione con lettere minuscole sono usate nei momenti di esperimenti PRIMA di essere effettuati.

es. x = esperimento lancio del dado

y = tempo vita auto dopo l'esperimento. Diventano costanti denotate con le maiuscole X, Y, Z .

$$X = x \quad Y = y$$

Ha come dominio lo spazio dei campioni S e come codominio lo retto reale

• $\lim_{t \rightarrow -\infty} F_x(t) = 0$

Inoltre:

$\lim_{t \rightarrow -\infty} F_x(t) = \lim_{t \rightarrow -\infty} P[X \in (-\infty, t)] = P[X \in (-\infty, -\infty)] = F_x(-\infty) = 0$

• F_x sono continue da destra, ma non da sinistra, o almeno non per parte.

$\lim_{t \rightarrow t_0^+} F_x(t) = F_x(t_0)$

$\lim_{t \rightarrow t_0^+} F_x(t) = \lim_{t \rightarrow t_0^+} P[X \in (-\infty, t)] = P[\lim_{t \rightarrow t_0^+} X \in (-\infty, t)] = P[X \in (-\infty, t_0)] = F_x(t_0)$

VARIABILI DISCRETE O CONTINUE

due variabili possono essere continue o discrete:

- x è detta continua se F_x è derivabile (tranne al più in un numero finito di punti)
- x è detta discreta se può assumere valori in una insieme di cardinalità finita o al più numerabile.

ASSOLUTAMENTE CONTINUE

è così definita se $\exists \frac{dF(t)}{dt} = f_x(t)$

$f_x(t) =$ DENSITÀ

$P[X \in A] = \int_A f_x(t) dt \rightarrow F_x(t) = P[X \leq t] = P[X \in (-\infty, t]] = \int_{-\infty}^t f_x(s) ds$

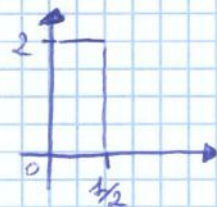
Proprietà delle f_x :

• $f_x(t) \geq 0$

• $\int_{-\infty}^{\infty} f_x(t) dt = 1$

Nota: $0 \leq f_x(t) \leq 1$

$f_x(t) = \begin{cases} 0 & t \in [0, \frac{1}{2}] \\ 2 & t \in (\frac{1}{2}, 1] \end{cases}$



Densità congiunte

(x_1, x_2) è discreta se i valori assumibili sono di cardinalità finita o numerabile

→ $P_{\vec{x}}(t_1, t_2) = P_{(x_1, x_2)}(t_1, t_2) = P[x_1=t_1, x_2=t_2]$

(x_1, x_2) è ass. continua se $\exists f_{(x_1, x_2)}(t_1, t_2)$ tale che $P[(x_1, x_2) \in A] = \iint_A f_{(x_1, x_2)}(t_1, t_2) dt_1 dt_2$, $A \subset \mathbb{R}^2$

DISTRIBUZIONI MARGINALI

Fanno riferimento alla sola x_1 o alla sola x_2

$P_{x_1}(t_1) = \sum_{t_2} P_{(x_1, x_2)}(t_1, t_2)$ (in colonna nella tabella)

ESAME : Definizione : x_1 e x_2 sono dette stocastiche indipendenti se $P_{(x_1, x_2)}(t_1, t_2) = P_{x_1}(t_1) \cdot P_{x_2}(t_2)$
 $\forall t_1, t_2 \in \mathbb{R}$

Nota: INDIPENDENZA: $P_{(x_1, x_2)}(t_1, t_2) \rightarrow P_{x_1}(t_1), P_{x_2}(t_2)$
~~←~~ No (a meno che non siano indipendenti)

ti vede questo argomento

Nel caso delle ass. continue

$f_{x_1}(t_1) = \int_{\mathbb{R}} f_{(x_1, x_2)}(t_1, t_2) dt_2$

INDIPENDENTI se $f_{(x_1, x_2)}(t_1, t_2) = f_{x_1}(t_1) \cdot f_{x_2}(t_2) \forall t_1, t_2$

Infatti:

$P[x \in A, y \in B] = P[x \in A] \cdot P[y \in B]$

$\forall B, A \in \mathbb{R}$

$F_{xy}(a, b) = F_x(a) \cdot F_y(b)$

Definizione:

Dati $A_1, A_2, A_3, \dots, A_n$ diciamo che sono indipendenti se preso un loro sottoinsieme qualsiasi, ~~vale~~

$\{A_2, A_5, A_{13}\}$ vale $P(A_2 \cap A_5 \cap A_{13}) = P(A_2) \cdot P(A_5) \cdot P(A_{13})$

Nota

Se A, B sono indipendenti, allora lo sono anche i complementari \bar{A}, \bar{B}

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

$I_s: P(\bar{A} \cap \bar{B}) = P(\bar{A}) \cdot P(\bar{B})$

$$P(\bar{A} \cap \bar{B}) = P(\overline{A \cup B}) = 1 - P(A \cup B) = 1 - [P(A) + P(B) - P(A \cap B)] =$$

$$= 1 - P(A) - P(B) + P(A) \cdot P(B) = 1 \cdot (1 - P(A)) - P(B) \cdot (1 - P(A)) =$$

$$= (1 - P(B)) \cdot (1 - P(A)) = P(\bar{B}) \cdot P(\bar{A})$$

Formula prodotto:

(S, \mathcal{A}, P) sia data una famiglia di eventi A_1, A_2, \dots, A_n
vale $P(A_1 \cap A_2 \dots \cap A_n) = P(A_1) \cdot P(A_2 | A_1) \cdot P(A_3 | A_1 \cap A_2) \dots P(A_n | A_1 \cap A_2 \dots A_{n-1})$

VARIANZA

$Var[X] = E[(X - E[X])^2]$ = varianza di x , cresce al crescere della dispersione della v.c. x

Al posto della varianza può essere usata la DEVIAZIONE STANDARD $\sigma = \sqrt{Var}$

OSSERVAZIONE:

$$V[X] = \begin{cases} \sum_{x_i} (x_i - E[X])^2 \cdot P_X(x_i) & \text{DISCRETA} \\ \int_{\mathbb{R}} (x - E[X])^2 \cdot f_X(x) dx & \text{ASS. CONTINUA} \end{cases}$$

Nota: togliendo il quadrato

$$\int_{\mathbb{R}} (x - E[X]) dx \cdot f_X(x) dx = \int_{\mathbb{R}} x \cdot f_X(x) dx - E[X] \int_{\mathbb{R}} f_X(x) dx = E[X] - E[X] \cdot 1 = 0$$

Perché non con il valore assoluto?

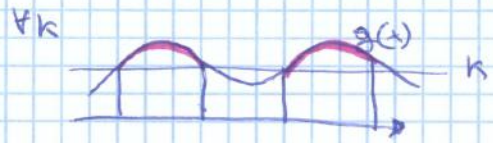
$E[|X - E[X]|]$ = scarto medio assoluto

DISTRIBUZIONE DI TCHEBY CHEFF

Teorema

Sia x v.c. qualsiasi, sia $g: \mathbb{R} \rightarrow \mathbb{R}^+$

Allora: $P[g(x) \geq k] \leq \frac{E[g(x)]}{k}$



Demonstrazione (ass. continue)

$$E[g(x)] = \int_{\{x: g(x) < k\}} g(x) \cdot dx(x) \, dx + \int_{\{x: g(x) \geq k\}} g(x) \cdot dx(x) \, dx \geq \int_{\{x: g(x) \geq k\}} g(x) \cdot dx(x) \, dx$$

$$\geq \int_{\{x: g(x) \geq k\}} k \cdot dx(x) \, dx$$

$$E[g(x)] \geq \dots \geq \int_{\{x: g(x) \geq k\}} k \cdot dx(x) \, dx = k \int_{x: g(x) \geq k} dx(x) \, dx = k \cdot P[g(x) \geq k]$$

$$P[g(x) \geq k] \leq \frac{E[g(x)]}{k} \quad \rightsquigarrow \quad g(x) = (x - E[x])^2$$

Corollario di Tchebycheff

$P[|x - \mu_x| < t\sigma_x] \geq 1 - \frac{1}{t^2}$
 ↳ minore parte o ampie impione

Dato y , $P[|y - E[y]| \geq a \cdot \sigma_y] \leq \frac{1}{a^2}$

$$y = \bar{x}_n \quad \varepsilon = a\sigma_y \leftrightarrow a = \frac{\varepsilon}{\sigma_y} \leftrightarrow a^2 = \frac{\varepsilon^2}{\sigma_y^2} = \frac{\varepsilon^2}{\frac{P(1-P)}{n}} \rightarrow \frac{1}{a^2} = \frac{P(1-P)}{n\varepsilon^2}$$

$P[|\bar{x}_n - \mu| \geq \varepsilon] \leq \frac{P(1-P)}{n\varepsilon^2} \rightarrow 0$
 $n \rightarrow \infty$

LEGGI DEI GRANDI NUMERI

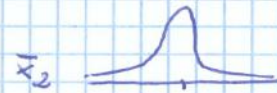
x_i indipendenti, stessa distribuzione

$$E[x] = \mu < \infty$$

$$V[x] = \sigma^2 < \infty$$

↓
deviazione standard

$$\bar{x}_n = \frac{x_1 + x_2 + \dots + x_n}{n}$$



Man mano che aumenti gli esperimenti (\bar{x}), la media diventa sempre più stretta.

DISTRIBUZIONE DI POISSON

È il modello adatto per modellare il numero di difetti o non conformità che si trovano in una unità di prodotto.

$x \sim \text{Poisson}(\lambda)$ con $\lambda > 0$ e $\lambda \in \mathbb{R}^+$, $\lambda = n \cdot p$

se può assumere valori in $N = \{0, 1, \dots\}$ con probabilità

$$P[x=k] = \frac{\lambda^k}{k!} \cdot e^{-\lambda} \quad \forall k \in \{0, 1, \dots\}$$

- Nota: $\sum_{k=0}^{\infty} \frac{\lambda^k}{k!} e^{-\lambda} = e^{-\lambda} \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} = e^{-\lambda} \cdot e^{\lambda} = 1 \rightarrow$ quindi questa è una distribuzione.

- È un caso particolare della binomiale:

Se $x \sim \text{Bin}(n, p)$ con $n \rightarrow \infty$, $p \rightarrow 0$ sotto il vincolo che $n \cdot p = \lambda$ e resto costante

$$\begin{aligned} P[x=k] &= \binom{n}{k} p^k (1-p)^{n-k} \rightarrow \frac{n!}{k!(n-k)!} \cdot \frac{\lambda^k}{n^k} \cdot \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-k} = \\ &= \frac{\lambda^k}{k!} \cdot \frac{n(n-1)\dots(n-k+1) \cdot (n-k)!}{n^k \cdot (n-k)!} \cdot \left(1 - \frac{\lambda}{n}\right)^k \cdot \left(1 - \frac{\lambda}{n}\right)^{-k} = \frac{\lambda^k}{k!} \cdot 1 \end{aligned}$$

$\downarrow \lambda, n \rightarrow \infty$ $\downarrow 1, n \rightarrow \infty$

ovvero $\frac{\lambda^k}{k!} \cdot e^{-\lambda} \rightarrow$ distribuzione di Poisson quando n è molto grande e p molto piccolo.

$n \gg 100$ e $p \leq 0,05$

Media: $E[x] = \lambda = n \cdot p$ varianza = $n \cdot p(1-p) = np - p(np) = 1 - p \lambda$
in caso di $p \rightarrow 0$

Proprietà geometrica:

Sia $x \sim \text{Geo}(p)$, allora $P[x > k] = (1-p)^k$

Dimostrazione:

$$\begin{aligned}
 P[x > k] &= \sum_{i=k}^{\infty} P[x=i] = \sum_{i=k}^{\infty} (1-p)^i \cdot p = p \sum_{i=k}^{\infty} (1-p)^i = p \sum_{i=k}^{\infty} (1-p)^{i+k} \\
 &= p \sum_{i=0}^{\infty} (1-p)^i \cdot (1-p)^k = (1-p)^k \cdot p \cdot \underbrace{\sum_{i=0}^{\infty} (1-p)^i}_{=1} = (1-p)^k
 \end{aligned}$$

perché $p \cdot \sum_{i=0}^{\infty} (1-p)^i = p \cdot \frac{1}{1-(1-p)} = 1$

PROPRIETÀ DI NON-INTERFERENZA

Facendo la ripetizione della stessa cosa mille volte, questo non cambia la probabilità.

Sia $x \sim \text{Geo}(p)$. Siano $m, k \in \mathbb{N}$, allora

$$P[x = k+m | x > k] = P[x = m]$$

Dimostrazione:

come mai non è moltiplicato?

$$\begin{aligned}
 P[x = k+m | x > k] &= \frac{P[x = k+m, x > k]}{P[x > k]} = \frac{P[x = k+m]}{P[x > k]} = \\
 &= \frac{(1-p)^{k+m} \cdot p}{(1-p)^k} = (1-p)^m \cdot p = P[x = m]
 \end{aligned}$$

Somma di 2 geometriche indipendenti non è più indipendente

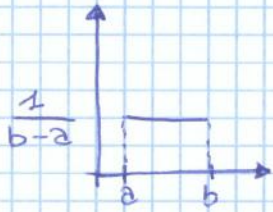
DISTRIBUZIONI PER FUNZIONI CONTINUE

DISTRIBUZIONE UNIFORME

$x \sim U[a, b]$

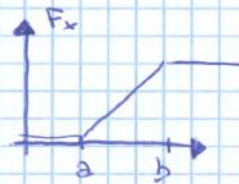
$a < b$

se $f_x(t) = \begin{cases} \frac{1}{b-a} & t \in [a, b] \\ 0 & \text{altrimenti} \end{cases}$



Funzione di distribuzione cumulativa:

$F_x(t) = \begin{cases} 0 & t < a \\ \frac{t-a}{b-a} & t \in [a, b] \\ 1 & t > b \end{cases}$



MEANA: $E[x] = \frac{a+b}{2}$

Dimostrazione:

$$E[x] = \int_{\mathbb{R}} t \cdot f_x(t) dt = \int_a^b t \cdot \frac{1}{b-a} dt = \frac{1}{b-a} \left[\frac{t^2}{2} \right]_a^b = \frac{1}{b-a} \cdot \frac{1}{2} (b^2 - a^2) = \frac{(a+b)(b-a)}{2(b-a)} = \frac{a+b}{2}$$

VARIANZA: $V[x] = \frac{(b-a)^2}{12}$

Dimostrazione:

$$V[x] = E[x^2] - (E[x])^2 = \int_a^b t^2 \cdot \frac{1}{b-a} dt - \frac{(a+b)^2}{4} = \frac{1}{b-a} \left[\frac{t^3}{3} \right]_a^b - \frac{(a+b)^2}{4}$$

\downarrow non integrabile

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

Proprietà esponenziale

- Se $x \sim \text{Exp}(\lambda)$, allora $r_x(t) = \lambda$ costante

Infatti,

$$r_x(t) = \frac{f_x(t)}{\bar{F}_x(t)} = \frac{\lambda \cdot e^{-\lambda t}}{e^{-\lambda t}} = \lambda$$

- Se sono x_1, x_2, \dots, x_n

$x_i \sim \text{Exp}(\lambda_i)$ indipendenti. Allora $Y = \min\{x_1, x_2, \dots, x_n\}$ è ancora esponenziale di parametro $\lambda = \lambda_1 + \lambda_2 + \dots + \lambda_n$

MEAN: $E[x] = \frac{1}{\lambda}$

Dimostrazione:

$$E[x] = \int_0^{\infty} \lambda t \cdot f_x(t) dt = \int_0^{\infty} \lambda t \cdot e^{-\lambda t} dt = \dots = \frac{1}{\lambda}$$

VARIANZA: $V[x] = \frac{1}{\lambda^2}$

Dimostrazione:

$$V[x] = \int_0^{\infty} \lambda t^2 \cdot e^{-\lambda t} dt - (E[x])^2 = \dots = \frac{1}{\lambda^2}$$

WEIBULL

$x \sim \text{Exp}(\lambda)$

$r_x(t) = \lambda$

Consideriamo un tasso

$r_x(t) = \beta \cdot t^{\beta-1} \cdot \alpha^\beta$

$\alpha, \beta > 0$

$\bar{F}_x(t) = e^{-(\alpha t)^\beta}$

$\rightarrow x \sim \text{Weibull}(\alpha, \beta)$

N.B.: se $\beta = 1 \rightarrow$ distribuzione esponenziale

Varianza della popolazione non nota:

σ^2 sconosciute $\rightarrow S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2$

$T = \frac{\bar{x}_n - \mu}{S_n / \sqrt{n}}$

$T \sim$ variabile casuale di Student con $n-1$ grad. di libertà

$P[-t_{n-1, 1-\frac{\alpha}{2}} < T \leq t_{n-1, 1-\frac{\alpha}{2}}] = 1-\alpha$

Intervallo di fiducia:

$I = (L_i, L_s) = (\bar{x}_n - t_{n-1, 1-\frac{\alpha}{2}} \cdot \frac{S}{\sqrt{n}}, \bar{x}_n + t_{n-1, 1-\frac{\alpha}{2}} \cdot \frac{S}{\sqrt{n}})$

utilizzabile solo con popolazioni distribuite solo con una normale $\sim N(\mu, \sigma^2)$

DETERMINAZIONE DELLA NUMEROSITÀ CAMPIONARIA:

$I_d =$ semi ampiezza

$I_d = z_{1-\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}$ oppure $I_d = t_{n-1, 1-\frac{\alpha}{2}} \cdot \frac{S}{\sqrt{n}}$

$n \gg \left(\frac{z_{1-\frac{\alpha}{2}} \cdot \sigma}{I_d} \right)^2$

TEST IPOTESI

$H_0: \mu = \mu_0$ livello di fiducia $1-\alpha$, $x \sim N(\mu, \sigma^2)$
 $H_A: \mu \neq \mu_0$ $n = \text{num. campioni}$, σ^2 nota

$$\bar{z} = \frac{\bar{x}_n - \mu_0}{\sigma/\sqrt{n}}$$
 $\rightarrow \begin{cases} \text{se } \bar{z}_{\text{calc}} \text{ è "vicina" a zero, accetto } H_0 \\ \text{se } \bar{z}_{\text{calc}} \text{ è "lontano" da zero, rifiuto } H_0 \end{cases}$

Regione accettazione: $R_{acc} \left(-z_{1-\frac{\alpha}{2}}, z_{1-\frac{\alpha}{2}} \right)$

Regione rifiuto: $R_{rifi} \left(-\infty, -z_{1-\frac{\alpha}{2}} \right) \cup \left(z_{1-\frac{\alpha}{2}}, +\infty \right)$

$\alpha = \text{prob. errore del primo tipo}$

Tre tipi di test:

TEST BILATERALE:

$H_0: \mu = \mu_0$ $\bar{z}_{\text{calc}} > z_{1-\frac{\alpha}{2}}$
 $H_A: \mu \neq \mu_0$ $\bar{z}_{\text{calc}} < -z_{1-\frac{\alpha}{2}}$

TEST UNILATERALE: (\bar{x}_{max})

$H_0: \mu = \mu_0$ $\bar{z}_{\text{calc}} > z_{1-\alpha}$
 $H_A: \mu > \mu_0$

TEST UNILATERALE: (\bar{x}_{min})

$H_0: \mu = \mu_0$ $\bar{z}_{\text{calc}} < -z_{1-\alpha}$
 $H_A: \mu < \mu_0$

TEOREMA LIMITE CENTRALE

Sono x_1, x_2, \dots, x_n delle variabili ^{casuali} indipendenti con stessa distribuzione, con $E[x_i] = \mu$ e $V[x_i] = \sigma^2$

Sia $Y_n = \frac{x_1 + x_2 + \dots + x_n - n\mu}{\sqrt{n} \cdot \sigma}$, allora per "n grande" vale

è approssimazione $Y_n \sim N(0, 1)$

per grande n in un v.o. $n \sim 20$ simmetrica

se però potessero anche avere $n \rightarrow \infty$

Sia data una popolazione distribuita con densità $f(x)$ e anche avente media μ e varianza σ^2 finite. Dato \bar{x}_n la media di un campione casuale di dimensione n estratto da essa, posso allora la media campionaria seguire una dist. normale e tendere di n ad infinito.

RICORDA INFATTI

Sono x_1, \dots, x_n indipendenti e $\sim N(\mu, \sigma^2)$

$\bar{x}_n = \frac{S}{n}$ dove $S = x_1 + x_2 + \dots + x_n \sim N(n \cdot \mu, n \cdot \sigma^2)$

da cui ottengo

$$\bar{x}_n = \frac{x_1 + \dots + x_n}{n} \sim N\left(\mu, \frac{\sigma^2}{n}\right) \rightarrow \begin{cases} E[\bar{x}_n] = \frac{n \cdot \mu}{n} = \mu \\ V[\bar{x}_n] = V\left[\frac{x_1 + \dots + x_n}{n}\right] = \frac{1}{n^2} \cdot n \sigma^2 = \frac{\sigma^2}{n} \end{cases}$$

Conferma la legge dei grandi numeri

Se $n \rightarrow \infty$, $\bar{x}_n \sim N(\mu, 0)$

Standard di \bar{x}_n :

$$\frac{\bar{x}_n - \mu}{\sqrt{\frac{\sigma^2}{n}}} = \frac{x_1 + x_2 + \dots + x_n}{n} - \frac{n \cdot \mu}{n} = \frac{x_1 + x_2 + \dots + x_n - n \mu}{n} = \frac{x_1 + x_2 + \dots + x_n - n \mu}{\sqrt{\frac{\sigma^2}{n}}} \sim N(0, 1) \text{ vale per } x_i \sim N(\mu, \sigma^2)$$

Sea x una v.c., sea $Y = \phi(x)$, dove $\phi: \mathbb{R} \rightarrow \mathbb{R}$. Conosciamo F_x ed f_x
 Come sono dette F_y e f_y ?

CASO DISCRETO

Sea ϕ invertibile (e sea ϕ^{-1} la sua inversa), conviene lavorare con ϕ e ϕ^{-1} .

$$P_y(t) = P[Y=t] = P[\phi(x)=t] = P[x=\phi^{-1}(t)] = P_x[\phi^{-1}(t)]$$

Se invece ϕ non fosse invertibile

$$P_y(t) = P[Y=t] = P[\phi(x)=t] = P[\{x_i : \phi(x_i)=t\}]$$

CASO CONTINUO

Conviene usare la dist. cumulata, $x, Y = \phi(x)$

Se ϕ invertibile $\rightarrow \phi^{-1}$ è l'inversa

[ϕ monotona crescente]

$$F_x \rightarrow F_y: F_y(t) = P[Y \leq t] = P[\phi(x) \leq t] = P[x \leq \phi^{-1}(t)] = F_x(\phi^{-1}(t))$$

$$f_y(t) = \frac{dF_y(t)}{dt} = \frac{dF_x(\phi^{-1}(t))}{dt} = f_x(\phi^{-1}(t)) \cdot \frac{1}{\phi'(\phi^{-1}(t))}$$

per definizione

RELATIONI TRA INDIPENDENZA E INCORRELATIONS

(x, y) indipendenti \rightarrow (x, y) sono anche incorrelate

\downarrow
non è sempre vero il contrario

Problema

Analizzo la relazione tra x ed y .

$Cov[x, y] = 5$ 5 è "grande" oppure è "vicino" a zero

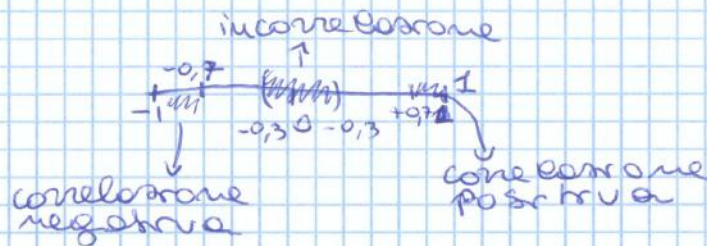
\rightarrow indice adimensionale \rightarrow INDICE DI CORRELAZIONE LINEARE

$$r_{x,y} = \frac{Cov[x, y]}{\sqrt{V[x]} \cdot \sqrt{V[y]}}$$

$x \sim$ anni
 $y \sim$ euro

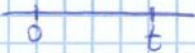
Teorema

$$-1 \leq r_{xy} \leq 1$$



PROCESSO DI POISSON

Sono una particolare classe dei processi di conteggio, si applicano quando valgono le ipotesi sottostanti:



Ipotesi 1:

Esiste un certo λ positiva tale che la probabilità che si verifichi esattamente un evento in un piccolo intervallo di lunghezza Δt sia approssimativamente uguale a $\lambda \Delta t$, cioè:

$$P[\text{esattamente un evento nell'intervallo di tempo } \Delta t] = \lambda \Delta t + o(\Delta t)$$

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

Ipotesi 2:

La probabilità che nell'intervallo Δt si verifichino più di un evento è trascurabile (e più alta) rispetto alla probabilità che se ne verifichino esattamente uno, cioè:

$$P[\text{due o più eventi nell'intervallo } \Delta t] = o(\Delta t)$$

$$P[\text{no eventi in } \Delta t] = 1 - P[1 \text{ evento}] - P[2 \text{ o più eventi}] = 1 - (\lambda \cdot \Delta t + o(\Delta t)) - o(\Delta t)$$

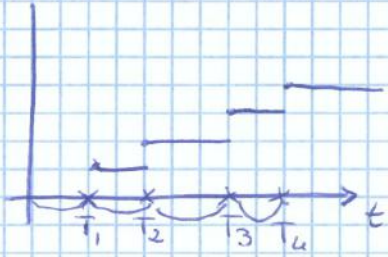
Ipotesi 3:

L'evento E_1 rappresenta il verificarsi di n_1 eventi in qualsiasi intervallo di tempo (t_1, t_2) e indipendente dall'evento E_2 rappresentante il verificarsi di n_2 eventi in un intervallo (t_3, t_4) non contenuta (parzialmente o del tutto) in (t_1, t_2) , cioè:

$$P[n_1 \text{ eventi in } (t_1, t_2) \text{ e } n_2 \text{ eventi in } (t_3, t_4)] = P[E_1] \cdot P[E_2]$$

Un processo di conteggio che soddisfa queste 3 ipotesi è detto **PROCESSO DI POISSON D'INTENSITÀ λ**

TEOREMA



T_i = intervalli tra eventi

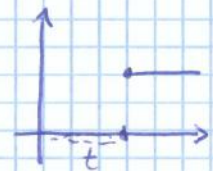
Sia $N = \{N(t), t \geq 0\}$ un processo di Poisson d'intensità λ ($\lambda > 0$).

Allora le variabili T_i sono variabili con distribuzione esponenziale di parametro λ , indipendenti. È viceversa.

Dimostrazione:

$$\bar{F}_{T_1}(t) = P[T_1 > t] = P[N(t) = 0] = \frac{(\lambda t)^0}{0!} \cdot e^{-\lambda t} = e^{-\lambda t}$$

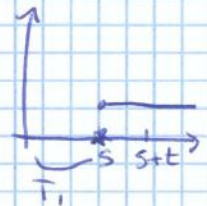
e questa è proprio la



$$\begin{aligned} \bar{F}_{T_2}(t) &= P[T_2 > t] = \int_0^{\infty} P[T_2 > t \mid T_1 = s] \cdot d_{T_1}(s) ds = \\ &= \int_0^{\infty} P[N(t+s) - N(s) = 0 \mid T_1 = s] \cdot d_{T_1}(s) ds \cdot T_1 \end{aligned}$$

↓
per indipendenza

$$\begin{aligned} &= \int_0^{\infty} P[N(t) = 0] \cdot d_{T_1}(s) ds = P[N(t) = 0] \cdot \int_0^{\infty} d_{T_1}(s) ds = \\ &= P[N(t) = 0] = \frac{(\lambda t)^0}{0!} \cdot e^{-\lambda t} \end{aligned}$$



Campione:

Un insieme di n variabili casuali indipendenti $(x_1, x_2, x_3, \dots, x_n)$ e ciascuna x_i ha la stessa distribuzione: $x_i \sim f(\cdot)$

- Attenzione: sono variabili casuali, non numeri!

Dopo l'estrazione si osservano i dati \rightarrow realizzazione del campione

STATISTICA:

$g(x_1, x_2, \dots, x_n)$ = una funzione di variabili casuali (e quindi a sua volta una variabile casuale) che non contiene parametri da stimare.

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

$$\max(x_1, x_2, \dots, x_n) - \min(x_1, x_2, \dots, x_n)$$

$$x_1 + 2x_2 - 4x_3 + \dots - x_n$$

\vdots

$x_n - \mu$ non è una stima

STIMATORI:

Statistiche che vengono usate per stimare un parametro θ o una qualche sua funzione

Esistono stimatori migliori di altri?

PROPRIETA' STIMATORI

- correttezza:

$T = g(x_1, x_2, \dots, x_n)$ stimatore di θ

$E[T] = \theta \rightarrow$ T stimatore corretto o non deviato

$E[T] - \theta = d$ è detta distorsione

È corretto se $\bar{x}_n = \frac{x_1 + x_2 + x_3 + \dots + x_n}{\text{numero}} = \mu$

- consistenza:

T stimatore corretto, si dice uno stimatore consistente se $\text{Var}[T] \rightarrow 0$ quando n di n del campione $\rightarrow +\infty$:

consistente se $\lim_{n \rightarrow \infty} \text{Var}[T_n] = 0$

Dimostrazione:

$$\begin{aligned}
 E[S_n^2] &= E\left[\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2\right] = \frac{1}{n-1} \cdot E\left[\sum_{i=1}^n \{x_i^2 - 2x_i\bar{x}_n + \bar{x}_n^2\}\right] = \\
 &= \frac{1}{n-1} \cdot E\left[\sum_{i=1}^n x_i^2 - 2\sum_{i=1}^n x_i \bar{x}_n + \sum_{i=1}^n \bar{x}_n^2\right] = \\
 &= \frac{1}{n-1} \cdot \left\{ \sum_{i=1}^n E[x_i^2] - 2E\left[\sum_{i=1}^n x_i \cdot \bar{x}_n\right] + \sum_{i=1}^n E[\bar{x}_n^2] \right\} = \\
 &= \frac{1}{n-1} \cdot \left\{ n \cdot (\sigma^2 + \mu^2) - 2 \cdot n \cdot \bar{x}_n E\left[\frac{(x_1 + x_2 + \dots + x_n)}{n} \cdot \bar{x}_n\right] + n E[\bar{x}_n^2] \right\} = \\
 &= \frac{1}{n-1} \cdot \left\{ n\sigma^2 + n\mu^2 - n\left(\frac{\sigma^2}{n} + \mu\right) \right\} = \frac{1}{n-1} \left\{ n\sigma^2 + n\mu^2 - \sigma^2 - n\mu^2 \right\} = \\
 &= \frac{\sigma^2(n-1)}{(n-1)} = \sigma^2
 \end{aligned}$$

Invece la sua varianza: ???

$$V[S_n^2] = \frac{1}{n} \left(E[x_i^4] - \frac{n-3}{n-1} \cdot E[(x_i - E[x_i])^4] \right), \text{ per } n \rightarrow +\infty = 0$$

ERRORE QUADRATICO MEDIO (MSE)

Si usa per valutare l'efficienza di stimatori non correlati
 $J = E[T] - \theta \rightarrow$ dispersione

$$\begin{aligned}
 MSE_T &= E[(T - \theta)^2] = E[(T - E[T]) + (E[T] - \theta)^2] = \\
 &= E[(T - E[T])^2 + 2(T - E[T]) \cdot (E[T] - \theta) + (E[T] - \theta)^2] = \\
 &= E[(T - E[T])^2] + 2E[(T - E[T]) \cdot (E[T] - \theta)] + E[S^2] = \\
 &= V[T] + 2(E[T] - \theta) \cdot (E[T] - E[E[T]]) + \sigma^2 = \\
 &= V[T] + 2(E[T] - \theta) \cdot \{E[T] - E[E[T]]\} + \sigma^2 = \\
 &= V[T] + \sigma^2
 \end{aligned}$$

poiché $E[E[T]] = E[T]$ perché il valore atteso è un n° e quindi il suo valore atteso è il n° stesso.
 Da qui ottengo che quella parte della formula si annulla.

Teorema

$E[S_n^2] = \sigma^2$

$V[S_n^2] = \frac{1}{n} \left(\mu_4 - \frac{n-3}{n-1} \sigma^4 \right)$

Distribuzione S_n^2 ?

STIME INTERVAZZARI

Stime rappresentate da un intervallo di valori (invece che un singolo valore) in cui cade il parametro da stimare con data probabilità.

Sia x popolazione, μ è da stimare, σ^2 nota e $\sigma^2 = 400$.
 Per stimare μ uso \bar{x}_n , $n = 64 \rightarrow \hat{\mu}$ [Stima di μ attraverso F_{64}]

$P[|\bar{x}_n - \mu| < 5]$

prob. che il valore con promesso σ discosti dalla media di non più di 5

$\bar{x}_n \sim N\left(\mu, \frac{\sigma^2}{64}\right)$

$= P[-5 < \bar{x}_n - \mu < 5] = P\left[\frac{-5}{\sigma/\sqrt{64}} < \frac{\bar{x}_n - \mu}{\sigma/\sqrt{64}} < \frac{5}{\sigma/\sqrt{64}}\right] = P\left[\frac{-5}{20/8} < z < \frac{5}{20/8}\right] =$
 $= P[-2 < z < 2] = 0,96$



Cercando la tavola

$$U = \sum_{i=1}^2 \left(\frac{x_i - \bar{x}_m}{\frac{1}{\sigma}} \right)^2 = \left(x_1 - \frac{x_1+x_2}{2} \right)^2 + \left(x_2 - \frac{x_1+x_2}{2} \right)^2 = \left(\frac{x_1-x_2}{2} \right)^2 + \left(\frac{x_2-x_1}{2} \right)^2 =$$

$$= 2 \left(\frac{x_1-x_2}{2} \right)^2 = \left(\frac{x_1-x_2}{\sqrt{2}} \right)^2 = z^2 \sim \chi_1^2$$

↓
Una normale standard
al quadrato

ci serve per:

Dato (x_1, x_2, \dots, x_n)

$$U = \frac{(n-1) S^2}{\sigma^2} = \frac{(n-1) \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_m)^2}{\sigma^2} = \sum_{i=1}^n \left(\frac{x_i - \bar{x}_m}{\sigma} \right)^2 \sim \chi_{n-1}^2$$

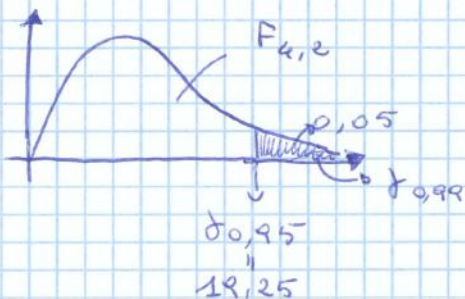
$$\rightarrow S^2 = \frac{\sigma^2}{n-1} \cdot U \quad \text{dove } U \sim \chi_{n-1}^2$$

DISTRIBUZIONE F (di Fisher)

U e V due variabili casuali indipendenti con distribuzione χ^2 rispettivamente con m e n gradi di libertà

$$F = \frac{U/m}{V/n} \sim F_{\frac{m}{n}} \quad [F_{\frac{m}{n}} \neq F_{\frac{n}{m}}]$$

Si dice di avere una distribuzione F con m e n gradi di libertà



Esistono due tavole diverse: - quantile 95
- quantile 99
di Fisher

$$d_{m,n,1-d} = \frac{1}{d_{m,n,d}} \rightarrow \text{inverti } m \text{ e } n$$

Quando si usa la t student?

$$\{x_1, x_2, \dots, x_n\} \quad x_i \sim N(\mu, \sigma^2)$$

$$\bar{x}_n \sim N\left(\mu, \frac{\sigma^2}{n}\right) \quad \text{e} \quad \frac{\bar{x}_n - \mu}{\sqrt{\frac{\sigma^2}{n}}} \sim N(0, 1)$$

Supponiamo che σ^2 non sia nota \rightarrow la sostituisco con una stima

$$\frac{\bar{x}_n - \mu}{\sqrt{\frac{s_n^2}{n}}} = \frac{\bar{x}_n - \mu}{\frac{s}{\sqrt{n}}} = \frac{\bar{x}_n - \mu}{\frac{s}{\sqrt{n}}} = \frac{z}{\sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x}_n)^2}{n-1}}}$$

Nota! la distribuzione di t è più larga di quella standard.

STIMA PUNTUALE (Stima per intervalli)

Ad ogni stima del rapporto σ si associa un intervallo di valori per il quale possiamo affermare che, con un certo grado di fiducia, contiene il valore vero del parametro.

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

lim. inferiore di fiducia \uparrow \uparrow lim. superiore di fiducia

Stima per intervalli:

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

livello di fiducia

$$\text{tipicamente } \alpha = \begin{cases} 0,01 \\ 0,05 \\ 0,1 \end{cases}$$

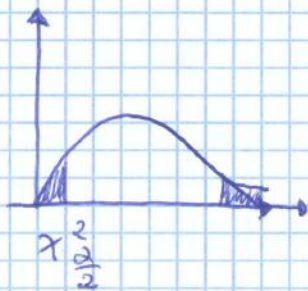
$$1 - \alpha = \begin{cases} 0,99 \\ 0,95 \\ 0,9 \end{cases}$$

6) Test varianza

$H_0: \sigma^2 = \sigma_0^2$ α, n fissati

$V = \frac{(n-1)S^2}{\sigma_0^2} \sim \chi^2_{n-1}$ (se è vero $\sigma^2 = \sigma_0^2$)

la decisione \rightarrow $\begin{cases} \text{se assume valori compatibili con una } \chi^2_{n-1} \rightarrow \text{accetto} \\ \text{altrimenti rifiuto} \end{cases}$



7) Test uguaglianza tra 2 varianze

$X_1 \sim n_1$ $X_2 \sim n_2$

$H_0: \sigma_1^2 = \sigma_2^2 \rightarrow \frac{\sigma_1^2}{\sigma_2^2} = 1$

$F = \frac{S_1^2 / \sigma_1^2}{S_2^2 / \sigma_2^2} \sim F(n_1-1, n_2-1)$



È vero se $\frac{\sigma_1^2}{\sigma_2^2} = 1 \rightarrow F = \frac{S_1^2}{S_2^2} \sim F(n_1-1, n_2-1)$

PRIMA DELL'ANOVA - SCOPRO STRONZE DELLA VARIANZA

Esempio:

2 popolazioni, 3 osservazioni per ogni popolazione. $n=6$

$$\left. \begin{array}{l} A \rightarrow x_1, x_2, x_3 \rightarrow \mu_x, \sigma_x^2 \\ B \rightarrow y_1, y_2, y_3 \rightarrow \mu_y, \sigma_y^2 \end{array} \right\} \rightarrow \mu = \sigma^2 \text{ (totali)}$$

Valore:

$$\mu = \frac{1}{6} \left[\sum_{i=1}^3 x_i + \sum_{i=1}^3 y_i \right] = \frac{1}{2} \left[\frac{\sum x_i}{3} + \frac{\sum y_i}{3} \right] = \frac{\mu_x + \mu_y}{3}$$

La media totale è la media delle medie

$$\sigma^2 = \frac{1}{6} \left[\sum_{i=1}^3 (x_i - \mu)^2 + \sum_{i=1}^3 (y_i - \mu)^2 \right] = \frac{1}{6} \left[\sum (x_i - \mu_x + \mu_x - \mu)^2 + \sum (y_i - \mu_y + \mu_y - \mu)^2 \right]$$

$$= \frac{1}{6} \left[\sum (x_i - \mu_x)^2 + \sum (\mu_x - \mu)^2 + 2 \sum (x_i - \mu_x) (\mu_x - \mu) + 2 \sum (x_i - \mu_y) (\mu_y - \mu) \right]$$

$$= \frac{1}{6} \left[3\sigma_x^2 + 3(\mu_x - \mu)^2 + \dots \right]$$

moltiplico e divido per 3 e mi viene il quadrato della varianza

$$\begin{aligned} & 2(\mu_x - \mu) \cdot \sum (x_i - \mu) = \\ & = 2 \cdot (\mu_x - \mu) \left[\frac{3\sum x_i}{3} - 3\mu \right] = 0 \end{aligned}$$

stesso per y

Da cui ottengo:

$$\sigma^2 = \frac{\sigma_x^2 + \sigma_y^2}{2} + \frac{(\mu_x - \mu)^2 + (\mu_y - \mu)^2}{2}$$

Media delle varianze

Varianza delle medie

→ se molto alta indica che le medie sono molto diverse

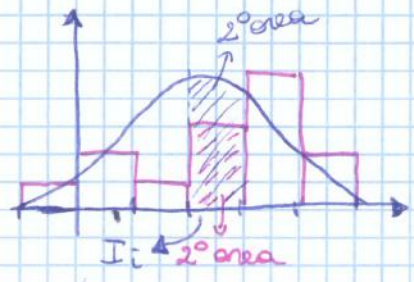
VARIANZA INTERNA AI GRUPPI (WITHIN)

VARIANZA TRA GRUPPI (BETWEEN)

TEST χ^2 PER BONTÀ ADATTAMENTO

$H_0: \mu_x = \mu_0$

$H_A: \mu_x \neq \mu_0$



- μ_0
 - μ empirico (istogramma dati campione)

Devo misurare la distanza

\Rightarrow μ teorico di osservazione su I_i

μ = ipotesi "teorica" di individui che avrei dovuto osservare su I_i (se H_0 vera)

$$\sum_{i=1}^N \frac{(n_i - n p_i)^2}{n p_i} = W$$

n_i = n° osservazioni del campione
 p_i = prob. di cadere I_i
 n = tot. osservazioni campione
 N = tot di sotto intervalli

Teorema:

Se vale H_0 ($\mu_x = \mu_0$) allora $W \sim \chi^2$

g.l. = $\begin{cases} n-1 & \text{se nelle } \mu_0 \text{ teorica non ci sono} \\ & \text{parametri stimati} \\ n-k-1 & \text{se nelle } \mu_0 \text{ teorica ci sono } k \\ & \text{parametri stimati con campione} \end{cases}$

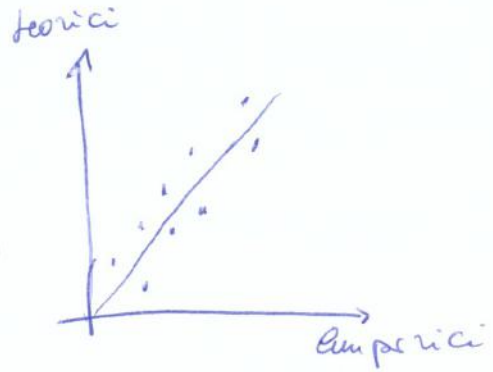
$R_{Acc} = [0, \chi^2_{1-\alpha}]$

Prezzo: può essere utilizzato anche per discrete

Defetto: campioni molto numerosi (almeno 5 osservazioni in ogni sottinsieme)

fit - plot

n° oss	osser	quant. empirico	quant. teorico
1	x_1	x_1	$F_0^{-1}(x_1)$
2	x_2	x_2	$F_0^{-1}(x_2)$
3	x_3	x_3	$F_0^{-1}(x_3)$
⋮	⋮		⋮



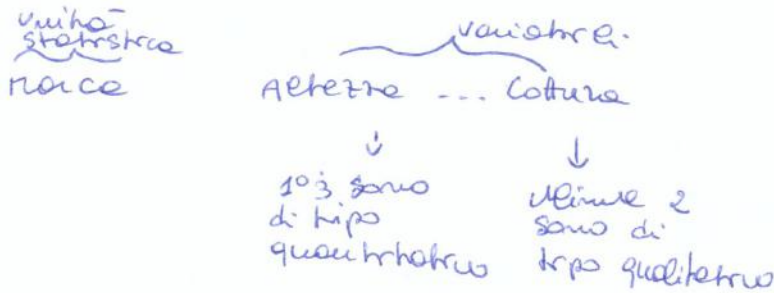
↓
 controimmagine
 di dist. empirica
 ($\frac{1}{n}u$)

Qui

LEZIONE VICARIO - SUDE

Matrice → n° righe > n° colonne

se tutte le variabili (variabile) sono qualitative, otterremo una matrice di tipo matematico



ottimo: qualitativa ordinale → ho un grado di giudizio

Riduzione dimensionale spaziale tra le colonne

REGRESSIONE

2m 0/10

variabile stat. - continua: età: non è discreta perché in realtà se dico ho "25 anni" non sono mai esatto, quindi è continua

Plot correlazione (tabella)

p-value → è stato fatto il test di ipotesi

$$H_0: \rho_{ij} = 0 \quad \text{contro} \quad H_a: \rho_{ij} \neq 0$$

p-value
se

accetto se ho un ~~errore~~ 10% in più,
 ↓
 non posso rifiutare che
 ci sia un legame
 lineare

MATRICE BLOCCO: vengono rapp. valori delle variabili a due a due

Come è vuoi rappresentare:

- matrice completa
- triang. superiore
- triang. infer. → lei preferisce queste

se invece dei punti è lungo e sottile, ~~però~~ ho
 relazione di tipo lineare.
 p-value e grafici non soddisfano mai da soli,

Test:

lett. greche = ? (parametri)
 maius = stimatori
 minusc. = stime

$H_0: \rho_{x,y} = 0$

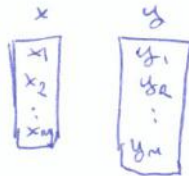
$H_A: \rho_{x,y} \neq 0$

Se può dimostrare che se H_0 è vera, la dist. campionaria del coef. di correlazione è tale che la statistica:

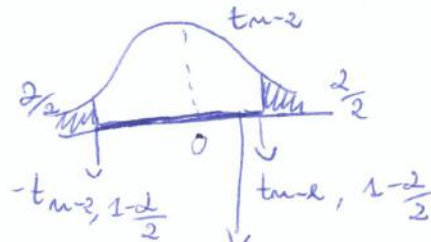
gradi di libertà

$$T_r = \frac{R}{\sqrt{1-R^2}} \sqrt{n-2} \sim t_{n-2}$$

$\{(x_1, y_1), (x_2, y_2) \dots (x_n, y_n)\}$



$$t_{calc} = \frac{r}{\sqrt{1-r^2}} \cdot \sqrt{n-2}$$



regione di accettazione

Es. su slide:



$\bar{x}_{min}, \bar{x}_{max}$
 $\bar{x}_s, \min \bar{x}_s, \max$

un campione 0,80
 non mi dice niente
 sul fatto che siano
 correlate o no.
 (riverde)

$$t_{n-2, 1-d/2} = \frac{r_{calc}}{\sqrt{1-r_{calc}^2}} \sqrt{n-2}$$

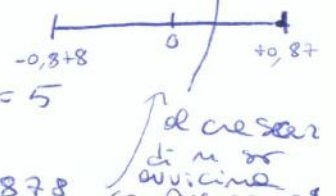
$$t_{n-2, 1-d/2}^2 = \frac{r_c^2}{1-r_c^2} (n-2) \rightarrow t_{n-2, 1-d/2}^2 = r_c^2 t_{n-2, 1-d/2}^2 = r_c^2 (n-2) =$$

$$= r_c^2 (t_{n-2, 1-d/2}^2 + n-2) = t_{n-2, 1-d/2}^2$$

$$|r_c| = \sqrt{\frac{t_{n-2, 1-d/2}^2}{t_{n-2, 1-d/2}^2 + n-2}}$$

 facendo limite $\frac{t_{n-2, 1-d/2}^2}{t_{n-2, 1-d/2}^2 + n-2}$

$\rightarrow \alpha = 5\% \Rightarrow n = 5$
 $-0,878 \quad +0,878$



CONTRASTI (contrasts)

Tramite ANOVA (e relativi output) è possibile effettuare test relativi a contrasti non nulli e negati delle medie (differenza fissata, diversi livelli)

ES:

ANOVA $H_0: \mu_0 = \mu_1 = \dots = \mu_k$

supponete test \rightarrow rifiutare H_0

allora testate $H_0: C_0\mu_0 + C_1\mu_1 + \dots + C_k\mu_k = d$

ES: rifiuto $\mu_1 = \mu_2 = \mu_3$, non zero che $\mu_2 = \frac{\mu_2 + \mu_3}{2}$

$$2\mu_1 - \mu_2 - \mu_3 = 0$$

$$H_0: \sum_{i=1}^k C_i \mu_i = d$$

Possiamo essere fatti $\sum_{i=1}^k C_i = 0$

In questo caso si parla di contrasti

Esempio:

1) $\mu_1 = 2\mu_2 - \mu_3 \rightarrow \mu_1 - 2\mu_2 + \mu_3 = 0$

$$\bar{c} = (1, -2, 1) \rightarrow \text{è un contrasto}$$

2) $\mu_1 = \mu_2 + \mu_3 \rightarrow \mu_1 - \mu_2 - \mu_3 = 0$

$$\bar{c} = (1, -1, -1) \Rightarrow \sum \bar{c} \neq 0$$

non è un contrasto

$$C \sim N(\sum C_i \mu_i, \sum C_i^2 \frac{\sigma^2}{n})$$

$H_0: \sum_{i=1}^k C_i \mu_i = d$ usare $C = \sum_{i=1}^k C_i \bar{y}_{\cdot i}$ come stimatore

$$E[\bar{c}] = E\left[\sum_{i=1}^k C_i \bar{y}_{\cdot i}\right] = \sum_{i=1}^k C_i \cdot E[\bar{y}_{\cdot i}] = \sum_{i=1}^k C_i \cdot \mu_i \rightarrow \text{corretto}$$

$$V[\bar{c}] = V\left[\sum_{i=1}^k C_i \bar{y}_{\cdot i}\right] = \sum_{i=1}^k C_i^2 \cdot V[\bar{y}_{\cdot i}] = \sum_{i=1}^k C_i^2 \cdot \frac{\sigma^2}{n}$$

dove $n = n^0$ misura di ogni livello
 $\sigma^2 =$ varianza errore sperimentale

2 fattori $\begin{cases} A \sim K \\ B \sim M \end{cases}$ più repliche (p repliche per ogni coppia)

$$C = \begin{cases} \sum C_i \bar{y}_{.i} & \text{per fatt A} \\ \sum C_j \bar{y}_{.j} & \text{per fatt B} \end{cases}$$

$$\sigma_C^2 = \underbrace{\sigma^2 \sum_{i=1}^K C_i^2}_{\text{per fatt A}} = \underbrace{\sigma^2 \sum_{j=1}^M C_j^2}_{\text{per fattore B}}$$

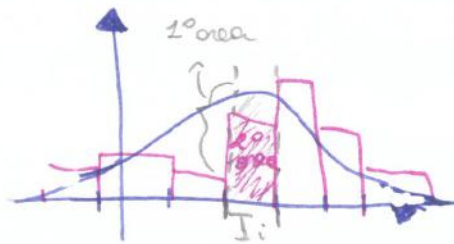
gradi libertà? \rightarrow quelli dell'errore sperimentale

prego e difetto di questo $n \dots$ 3
 difetto \rightarrow ultratrotori solo per F_0 ass. continue
 prego \rightarrow ultratrotori anche con campioni di numerosità molto piccola ($n=5,6 \dots$)

TEST χ^2 per BONTÀ ADATTAMENTO

$H_0: f_x = f_0$

$H_1: f_x \neq f_0$



f_0
 f empirico (istogramma del campione)

Devo misurare la "distanza"

$n_i =$ numero di osservazioni su I_i

$n_{i1} =$ "teorica" di individui che avrei dovuto osservare su I_i (se H_0 vera)

$$\sum_{i=1}^N \frac{(n_i - n p_i)^2}{n p_i} = W$$

$n_i =$ n° osservazioni dal campione

$p_i =$ prob di cadere I_i

$n =$ tot. oss. campione

$N =$ tot di sottintervalli

Teorema

se vale $H_0 (f_x = f_0)$ allora $W \sim \chi^2$

$g.l. =$ $\frac{\quad}{\quad} N-1$ se nella f_0 teorica non ci sono parametri stimati con campione

\searrow $N-k-1$ se nella f_0 teorica ci sono k parametri stimati con campione

Raccogliamone $[0, \chi^2_{1-\alpha}]$

ES (contrast) Anova, 2 fattori sotto controllo (A-4 livelli, B-3 livelli)
2 repliche corroni

	DF	SS	MS	F	P-value
A	3	50	15,6		
B	2	210	105	~48	0,0001
Error	18	40	2,2		
Tot	23	300			

$H_0: \mu_1 = \mu_2 = \mu_3$

→ rifiuto H_0

$\hat{\mu}_1 = 3.5 \quad \hat{\mu}_2 = 2.1 \quad \hat{\mu}_3 = 2.4$

$H_0: \mu_1 = \frac{1}{3}\mu_2 + \frac{2}{3}\mu_3 + 1$ → $\mu_1 - \frac{1}{3}\mu_2 - \frac{2}{3}\mu_3 = 1$

$C \left(1, -\frac{1}{3}, -\frac{2}{3} \right)$ $\sum c_i = 0$ *giusto* *struttura* *corretto*

$\sum_{i=1}^3 c_i \hat{\mu}_i = 1 \cdot \hat{\mu}_1 - \frac{1}{3} \hat{\mu}_2 - \frac{2}{3} \hat{\mu}_3 = 1 \cdot 3,5 - \frac{1}{3} \cdot 2,1 - \frac{2}{3} \cdot 2,4 = 1,2$

se fosse vera H_0 , $E[c] = 1$

$V[c] = \sigma^2 \cdot \sum c_i^2 = \frac{\sigma^2}{8} \left(1 + \frac{1}{9} + \frac{4}{9} \right) = \frac{\sigma^2}{8} \left(\frac{14}{9} \right) = \frac{2,2}{8} \left(\frac{14}{9} \right) = 0,4$

\downarrow
n° di osservazioni: 4 livelli del fattore A e 30 preso a misurazione

$V[\sum c_i \hat{\mu}_i] =$

$\sigma^2 \rightarrow 2,2$

$T_{calc} = \frac{C_{calc} - E[c]}{\sqrt{V[c]}} = \frac{1,2 - 1,0}{\sqrt{0,42}} = 0,31$

$\alpha = 0,05$

Regione Accettazione? = $(t_{18, 1-0,025}, t_{18, 0,975})$
 \downarrow \downarrow
errore in tabella \uparrow da output

se poi $x_i \sim t$

$(F_i \equiv F)$

MAX + MIN

$F_z(t) = 1 - (1 - F(t))^n$

$f_z(t) = n(1 - F(t))^{n-1} \cdot f(t) \rightarrow$ inversa

$F(t) \rightarrow$ dist cumulata

$F(t) = P[x \leq t]$

$\bar{F}(t) \rightarrow$ funzione di sopravvivenza

$\bar{F}(t) = P[x > t]$

$\bar{F}(t) = 1 - F(t)$

↓
distribuzione di

$\rightarrow 1 - F_z(t) = (1 - F(t))^n$

$\bar{F}_z(t) = (\bar{F}(t))^n$

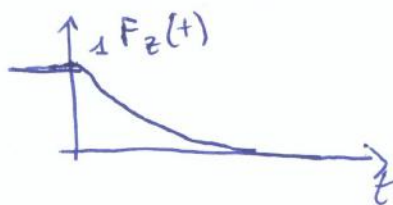
Es 1)

Si sono x_1, x_2, \dots, x_n indipendenti e taliche $x_i \sim \text{Exp}(\lambda_i)$

$Z = \min(x_i) \rightarrow E[Z]$

$Y = \max(x_i) \rightarrow E[Y]$

$F_z(t) = \bar{F}_1(t) \cdot \bar{F}_2(t) \cdot \dots \cdot \bar{F}_n(t) = e^{-\lambda_1 t} \cdot e^{-\lambda_2 t} \cdot \dots \cdot e^{-\lambda_n t} = e^{-(\lambda_1 + \lambda_2 + \dots + \lambda_n)t}$



$\rightarrow Z \sim \text{Exp}(\sum_{i=1}^n \lambda_i)$

$E[Z] = \frac{1}{\lambda_1 + \lambda_2 + \dots + \lambda_n}$

↓
valore atteso del minimo ha valore atteso inferiore a quello degli altri

$$Y = \max(x_1, \dots, x_n)$$

$$F_Y(t) = F_1(t) \cdot F_2(t) \cdot \dots \cdot F_n(t) \quad \forall t \geq 0$$

$$= (1 - e^{-\lambda_1 t}) (1 - e^{-\lambda_2 t}) \dots (1 - e^{-\lambda_n t}) = \prod_{i=1}^n (1 - e^{-\lambda_i t})$$

$$E[Y] = ? \rightarrow \text{difficile?}$$

Ci limitiamo al caso $n=2$, $\lambda_1 = \lambda_2 = \lambda$

$$x_1, x_2 \sim \text{Exp}(\lambda) \quad \text{ind.} \quad Y = \max(x_1, x_2)$$

$$F_Y = (1 - e^{-\lambda t})^2 = 1 - 2e^{-\lambda t} + e^{-2\lambda t} \quad (t \geq 0)$$

$$f_Y(t) = 2\lambda e^{-\lambda t} - 2\lambda e^{-2\lambda t} \quad (t \geq 0)$$

↓
densità

$$E[Y] = \int_0^{\infty} t \cdot (2\lambda e^{-\lambda t} - 2\lambda e^{-2\lambda t}) dt = 2 \int_0^{\infty} t \cdot \lambda e^{-\lambda t} dt - \int_0^{\infty} t \cdot 2\lambda e^{-2\lambda t} dt =$$

$$= 2 \cdot \frac{1}{\lambda} - \frac{1}{2\lambda} = \left(2 - \frac{1}{2}\right) \cdot \frac{1}{\lambda} = \frac{3}{2} \cdot \frac{1}{\lambda}$$

Analisi della varianza

Un contadino vuole valutare l'efficacia di un fertilizzante per coltivazioni di legumi. Sul mercato sono disponibili 4 fertilizzanti prodotti da altrettante ditte. Allo stato attuale, il contadino è convinto che i 4 fertilizzanti sono sostanzialmente equivalenti. **Come fare per accettare o confutare la sua convinzione??**



Si scelgono 5 appezzamenti di terreno diversi, vengono divisi in 4 parti uguali e ciascuna delle quattro parti viene *trattata* con un fertilizzante diverso. Si seminano poi i legumi ed alla fine della stagione i raccolti risultanti sono:

Campo	Fertiliz. 1	Fertiliz. 2	Fertiliz. 3	Fertiliz. 4	Tot. riga
I	12.6	13.1	12.7	12.9	51.3
II	6.2	7.0	6.7	7.0	26.9
III	4.8	5.1	4.8	5.2	19.9
IV	3.5	3.8	3.6	3.7	14.6
V	2.9	3.0	2.7	2.9	11.5
Tot. colon.	30	32	30.5	31.7	124.2

Rese (kg/m²) ottenute negli appezzamenti trattati con fertilizzanti diversi, per ognuno di cinque campi (*)

(*) sono anche indicati i totali per riga, per colonna e generale

se un concetto su una colonna, poiché ho lo stesso numero di osservazioni la distribuzione è appross. ad una normale

Requisiti:

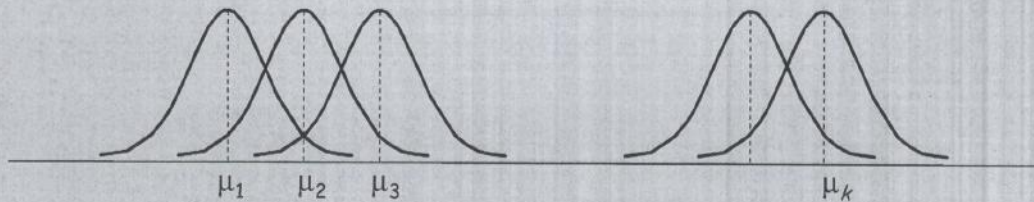
ognuno dei k gruppi di osservazioni è un campione casuale che proviene da una popolazione $\sim N(\mu_i, \sigma^2)$, cioè la varianza è costante tra i gruppi



$H_0: \mu_1 = \mu_2 = \dots = \mu_k + \sigma^2$ costante equivale a:

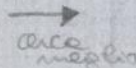
costante per tutti detto ERRORE SPERIMENTALE

I campioni provengono tutti dalla stessa popolazione



Come fare???

è possibile



Calcolo 2 varianze, se è soddisfatta la condizione che è nulla, allora è verificata

Per ogni colonna $j, j = 1, 2, \dots, k$, si calcoli:

$$s_j^2 = \frac{1}{n_j - 1} \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{\cdot j})^2$$

Ogni $s_j^2, j = 1, 2, \dots, k$, è una stima (corretta) della varianza comune σ^2 . Ma allora:

numero di \uparrow *varianza j-esima*

$$\hat{\sigma}^2 = s_{pool}^2 \equiv \frac{1}{\sum_{j=1}^k (n_j - 1)} \sum_{j=1}^k (n_j - 1) \left(\frac{1}{n_j - 1} \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{\cdot j})^2 \right)$$

stima varianza comune

se i campioni hanno lo stesso numero di osservazioni

$$\hat{\sigma}^2 = s_{pool}^2 \equiv \frac{1}{k} \sum_{j=1}^k s_j^2 = \frac{1}{k(n-1)} \sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{\cdot j})^2$$

1° STIMA:

$$\hat{\sigma}^2 = \frac{1}{k(n-1)} \sum_{j=1}^k \sum_{i=1}^n (v_{ij} - \bar{y}_j)^2 = \frac{SS_W}{k(n-1)} = MS_W$$

averò n° delle colonne della tabella
stima errore sperimentale

$SS_W =$ Sum of Squares Within Samples
campioni
 $MS_W =$ Mean Squares Within Samples
media
 $k(n-1) =$ gradi di libertà (degree of freedom)

I Mean Squares Within Samples (MS_W) sono una stima (sempre corretta, sia o non sia sotto H_0) della varianza σ^2 (incognita) comune alle k popolazioni

$$\frac{SS_W}{\sigma^2} \sim \chi^2_{k(n-1)}$$

09/11/15 Analisi della varianza 9

2° STIMA:

Ma si può anche calcolare:

stima delle varianze delle medie compresse
media compresse
 supponiamo ora vera l'ipotesi nulla (tutte le medie uguali) tra di loro a un certo valore μ
 $H_0: \mu_1 = \mu_2 = \dots = \mu_k$
 $\bar{y}_{.j} \sim N(\mu, \frac{\sigma^2}{n})$

$$\hat{\sigma}_y^2 \equiv S_y^2 = \frac{\sum_{j=1}^k (\bar{y}_{.j} - \bar{y}_{..})^2}{k-1}$$

E' una stima corretta solo sotto H_0 della varianza σ_y^2 delle k medie delle popolazioni rispetto alla media generale

Ma: $\sigma_y^2 = \frac{\sigma^2}{n}$ *ogni singolo individuo* $\rightarrow \sigma_y^2 = n \cdot \hat{\sigma}_y^2$ *cioè n volte e moltiplicando della formula vista sopra*
 (nel caso dei fertilizzanti: $s_y^2 = 0.0364$)

$$S_y^2 = n S_{\bar{y}}^2 = \frac{n \sum_{j=1}^k (\bar{y}_{.j} - \bar{y}_{..})^2}{k-1} = \frac{SS_B}{k-1} = MS_B$$

$$\frac{SS_B}{\sigma^2} \sim \chi^2_{k-1}$$

gradi di libertà, perché il campione delle medie = k

09/11/15 Analisi della varianza 10

Cosa fare???

$$\frac{MS_B}{MS_W} = \frac{\frac{SS_B}{k-1}}{\frac{SS_W}{k(n-1)}} = \frac{\frac{\sum x^2_{k-1}}{k-1}}{\frac{\sum x^2_{k(n-1)}}{k(n-1)}} = F_{k-1, k(n-1)}$$

Facciamo un test sull'uguaglianza dei due quadrati medi MS_B e MS_W ...
 (le due stime)

diventa grande quando la differenza tra le med sono grandi - se fosse vero e l'Ho non avremmo grandi

09/11/15 Analisi della varianza 13

Costruiamo una tavola dell'analisi della varianza (ANOVA) con un fattore sotto controllo ...

Tabella dell'Analisi della Varianza (ANOVA) con un solo fattore sotto controllo

Origine della variazione	Gradi di libertà	Somme dei quadrati	Quadrati medi	F_{cal}
Tra i gruppi (trattamenti)	$k-1$	$SS_B = \frac{1}{n} \sum_{j=1}^k y_j^2 - \frac{y_{..}^2}{nt}$	$MS_B = \frac{SS_B}{k-1}$	$\frac{MS_B}{MS_W}$
Errore sperimentale	$k(n-1)$	$SS_W = \sum_{i=1}^n \sum_{j=1}^k y_{ij}^2 - \frac{1}{n} \sum_{j=1}^k y_{.j}^2$	$MS_W = \frac{SS_W}{k(n-1)}$	
Totali	$nk-1$	$SS_{TC} = \sum_{i=1}^n \sum_{j=1}^k y_{ij}^2 - nk\bar{y}^2$		

rapporto tra le stime delle varianze

09/11/15 Analisi della varianza 14

↓
 sarà in modo diverso da precedenti ma uguale

Dubbio!!!!

L'errore sperimentale (variazione dovuta al caso) è grande!
 Non è che è stata attribuita al caso una variazione che invece può essere stata originata (se non tutta, almeno in parte) a qualche altra causa?

09/11/15

Analisi della varianza

19

	Fertiliz. 1	Fertiliz. 2	Fertiliz. 3	Fertiliz. 4	Tot. riga	Media riga
Varietà 1	12.6	13.1	12.7	12.9	51.3	12.825
Varietà 2	6.2	7.0	6.7	7.0	26.9	6.725
Varietà 3	4.8	5.1	4.8	5.2	19.9	4.975
Varietà 4	3.5	3.8	3.6	3.7	14.6	3.65
Varietà 5	2.9	3.0	2.7	2.9	11.5	2.875
Tot. colon.	30	32	30.5	31.7	124.2	
Media colonna	6	6.4	6.1	6.34		6.21

se sono
 sbagliati
 posso
 ripetere
 il test

Obiettivo: Cercare quale parte di variazione totale è dovuta alla coltivazione di varietà diverse: presenza di un secondo fattore

09/11/15

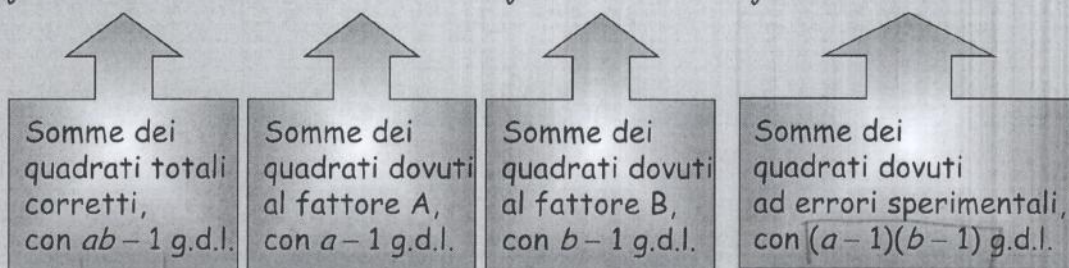
Analisi della varianza

20

Scomposizione della variabilità:

lo divido per...
 + ho la varianza
 totale?

$$\sum_{i=1}^a \sum_{j=1}^b (y_{ij} - \bar{y}_{..})^2 = b \sum_{i=1}^a (\bar{y}_{i.} - \bar{y}_{..})^2 + a \sum_{j=1}^b (\bar{y}_{.j} - \bar{y}_{..})^2 + \sum_{j=1}^b \sum_{i=1}^a (y_{ij} - \bar{y}_{.j} - \bar{y}_{i.} + \bar{y}_{..})^2$$



1° fattore
 2° fattore

colori nelle
 stime, quindi
 a-b

2° fattore
 stima
 2-1 = media
 complessiva
 che va a
 stimare la
 media

09/11/15

Analisi della varianza

23

1 parametro
 è stato stimato
 quindi: 2b-1

Tabella dell'Analisi della Varianza con due fattori sotto controllo

Origine della variazione	Gradi di libertà	Somme dei quadrati	Quadrati medi	F _{calc}
Tra colonne (fattore A)	a-1	$SS_A = \frac{1}{b} \sum_{i=1}^a y_{i.}^2 - \frac{Y_{..}^2}{ab}$	$MS_A = \frac{SS_A}{a-1}$	$\frac{MS_A}{MS_{exp}}$
Tra righe (fattore B)	b-1	$SS_B = \frac{1}{a} \sum_{j=1}^b y_{.j}^2 - \frac{Y_{..}^2}{ab}$	$MS_B = \frac{SS_B}{b-1}$	$\frac{MS_B}{MS_{exp}}$
Errore sperimentale	(a-1)(b-1)	$SS_{exp} = SS_{TC} - SS_A - SS_B$	$MS_{exp} = \frac{SS_{exp}}{(a-1)(b-1)}$	
Totale	ab-1	$SS_{TC} = \sum_{i=1}^a \sum_{j=1}^b y_{ij}^2 - ab\bar{y}_{..}^2$		

Se $F_{calc} = \frac{MS_A}{MS_{exp}} > F_{a-1, (a-1)(b-1), 1-\alpha}$ si rifiuta l'ipotesi nulla che l'effetto del fattore A sia non significativo, al livello di significatività α

Se $F_{calc} = \frac{MS_B}{MS_{exp}} > F_{b-1, (a-1)(b-1), 1-\alpha}$ si rifiuta l'ipotesi nulla che l'effetto del fattore B sia non significativo, al livello di significatività α .

09/11/15

Analisi della varianza

24