

STATISTICA INDUSTRIALE

ora 1

Note Title

03/03/2015

Probabilità → prerequisite

Ing. gest 9/6 cfu → statistica multivariata (avanzata)

5/6/9 cfu → statistica

Contenuti : statistica di base + probs } statistica (lea)
regressione }
(ANOVA) analisi della varianza } metodi e modelli a
test di adattamento } supporto delle decisioni
(DOE) design of experiment } (lea)

+ argomenti a ricorrenza delle esigenze (text mining, big data)

... cluster analysis

... discriminant analysis

Bibliografia : Ross : "Prob & Stat per Ing e Scienze" Apogeo

Johnson Winchen : ...

↳ Sleeper : "Design for Six Sigma Statistics"
per "practitioner"

Corso molto pratico : poca teoria e poche dimostrazioni
poca astrazione e poca generalità

(Statistica per matematici : Mood : "..."
dispense Pratelli M. (Pisa))

Laboratorio : Excel / Minitab

Esame : scritto di laboratorio + orale (seminario?)

Storia della statistica

- Statistica descrittiva : informazioni complete sulla popolazione
1600 → (Istat & co)

- Statistica inferenziale : dedurre info parziali da un campione
~1930 lab scientifici

40 persone → 20 trattamento
 ↳ 20 controllo (placebo)

	g.	n.g.	
t.	16	4	20
c.	7	13	20

	g.	n.g.	
t.	14	6	20
c.	9	11	20

~ 1940/50 in azienda → NO

Deming → 1960 in Giappone (Toyota)
 "qualità" (aziendale) = statistica

~ 1980 Motorola "six sigma"

green belt, black belt, master black belt

~ oggi dipende dalle persone

STATISTICA DI BASE

(stimatori)

→ Intervalli di confidenza

→ Test statistici

STIMATORI

- Dati normali

X_1, X_2, \dots, X_n

~ $N(\mu, \sigma^2)$

legge normale

media

varianza

indipendenti

"campiono" : x_1, \dots, x_n

"popolazione" : $\mathcal{N}(\mu, \sigma^2)$

μ, σ^2 incognite (una o entrambe)

x_1, x_2, \dots, x_n realizzazione numerica del campione nota

$$\bar{X} := \frac{1}{n} \sum_{i=1}^n x_i \quad \text{media campionaria}$$

$\bar{X} \approx \mu$ è uno stimatore di μ

- Def una statistica è una funzione qualsiasi del campione (ad esempio \bar{X})
- Def uno stimatore corretto di un parametro θ è una statistica Θ tale che

$$E(\Theta) = \theta$$

altrimenti si chiama distorto e $E(\Theta) - \theta$ si chiama bias

★ \bar{X} è uno stimatore corretto di μ
 $E(\bar{X}) = E\left(\frac{1}{n} \sum x_i\right) = \frac{1}{n} \sum E(x_i) = \mu$

- Def uno stimatore consistente di un parametro θ è una statistica $\Theta_n = f_n(x_1, \dots, x_n)$ tale che

$$\Theta_n \xrightarrow[n \rightarrow \infty]{} \theta \quad \text{in probabilità o q.c.}$$

★ \bar{X} è uno stimatore consistente di μ
la LGN (forte) dice che x_1, x_2, \dots i.i.d. $x_i \in L^1$

allora $\frac{1}{n} \sum_{i=1}^n x_i \xrightarrow[n \rightarrow \infty]{} E(x_1)$ q.c.

Per una verifica diretta della convergenza in probabilità

- i. $E(\Theta_n) \xrightarrow{n \rightarrow \infty} \theta$
 ii. $\text{Var}(\Theta_n) \xrightarrow{n \rightarrow \infty} 0$
- + disug. di Chebyshev \Rightarrow conv in prob (HW)
 independent!

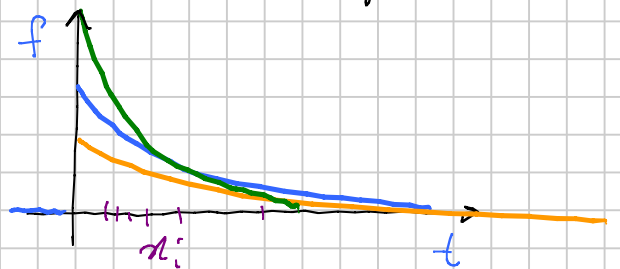
$$\star \text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n x_i\right) = \frac{1}{n^2} \text{Var}\left(\sum x_i\right) = \frac{1}{n^2} \sum \text{Var}(x_i) = \frac{\sigma^2}{n}$$

• Stimatori di massima likelihood (=verosimiglianza)

esempio $x_1, x_2, \dots, x_n \sim \text{expo}(\lambda)$ λ incognito

$$f(t) = \lambda e^{-\lambda t}, t \geq 0$$

$$X := (x_1, x_2, \dots, x_n)$$



$f_X: \mathbb{R}^n \rightarrow \mathbb{R}_+$ densità congiunta

$$f_X(t_1, t_2, \dots, t_n) = \prod_{i=1}^n (\lambda e^{-\lambda t_i}) = \lambda^n e^{-\lambda \sum t_i}$$

↑
indipendenza

$$L(\lambda) := f_X(x_1, x_2, \dots, x_n; \lambda) = \lambda^n e^{-\lambda \sum_1^n x_i}$$

↑
funzione di likelihood

↑
parametri

↑
variabile

cerco $\hat{\lambda}$ che massimizza L

$$l(\lambda) = \log L(\lambda) = n \log \lambda - \lambda \sum_1^n x_i$$

log-likelihood

$$\text{deriva } l'(\lambda) = \frac{n}{\lambda} - \sum_1^n x_i \quad \text{e pongo } l'(\hat{\lambda}) = 0$$

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

lo stimatore ML $\hat{\lambda} = \frac{1}{\bar{x}}$ non credo sia corretto (HW?)

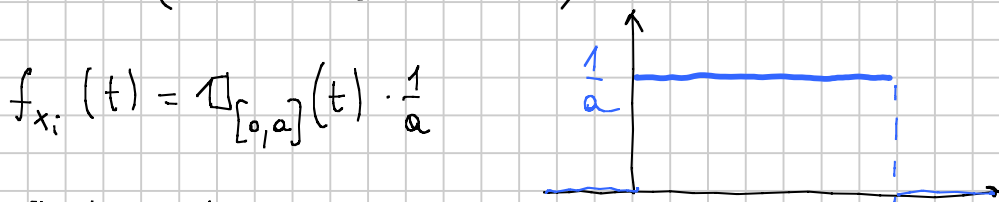
★ Ogni stimatore ML $\hat{\lambda}$ per forza consistente

HW: $X_1, X_2, \dots, X_n \sim \text{unif}[0, a]$ a incognito

- i. trovare ML est di a
- ii. verificare che non è corretto
- iii. proporre una modifica affinché diventi corretto (restando consistente)

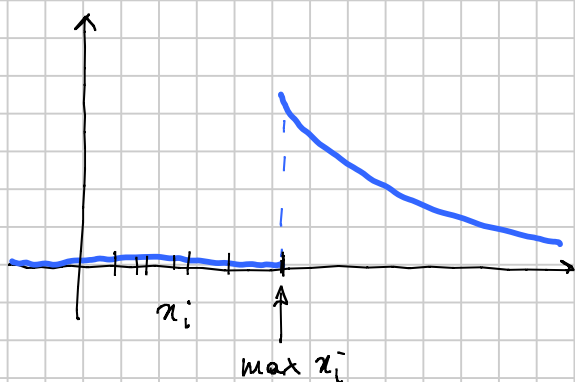
HW: $X_1, X_2, \dots, X_n \sim \text{unif}[0, a]$ a incognito

- i. trovare ML est di a
- ii. verificare che non è corretto
- iii. proporre una modifica affinché diventi corretto (restando consistente)



x_1, x_2, \dots, x_n

$$L(a) = \prod_{i=1}^n f_{X_i}(x_i; a) = \prod_{i=1}^n \left[\mathbb{1}_{[0, a]}(x_i) \cdot \frac{1}{a} \right] = a^{-n} \mathbb{1}_{[0, a]}(\max_i x_i)$$



$L(\hat{a})$ è massima per $\hat{a} = \max x_i$

MLE $\hat{a} := \max_i X_i$ è consistente

$$\hat{a}_n \xrightarrow[n \rightarrow \infty]{P} a$$

Non è corretto: $E(\hat{a}) \neq a$

$$E(\max_i X_i) = ?$$

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

$$F_{\hat{a}}(t) := P(\max_i X_i \leq t) = P(X_i \leq t \ \forall i=1, 2, \dots, n)$$

$$= P(\bigcap_i \{X_i \leq t\}) \stackrel{\text{indipendenza}}{=} \prod_i P(X_i \leq t) = [F_{X_i}(t)]^n = \begin{cases} \frac{t^n}{a^n}, & t \in [0, a] \\ 1 & t > a \\ 0 & t < 0 \end{cases}$$

$$F_{\hat{a}} \longrightarrow f_{\hat{a}} \quad f_{\hat{a}}(t) = \frac{n}{a^n} t^{n-1}$$

$$E(\hat{a}) = \int_0^a t \cdot \frac{n}{a^n} t^{n-1} dt = \frac{a^{n+1}}{n+1} \cdot \frac{n}{a^n} = a \frac{n}{n+1}$$

$$A := \frac{n+1}{n} \max_i X_i \quad \text{corretto e consistente}$$

$$A_n - \hat{a}_n = \frac{1}{n} \hat{a}_n \quad \hat{a}_n \xrightarrow{P} a$$

↘ 0

$$A_n \xrightarrow[n \rightarrow \infty]{P} a$$

\hat{a} stimatezza sistematicamente a

A e' corretto

HW : 1942 carri armati tedeschi catturati

77 carri tra marzo e ottobre 1942

1042, 1717, 348, ..., 4434, ..., 2345

↑
max 71° dato

numeri di serie

inizio marzo

fine ottobre

→ stimare la capacità produttiva mensile

INTERVALLI DI CONFIDENZA

- Caso campione normale σ nota, stimare μ
 $X_1, X_2, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$ indep

\bar{X} è stimatore c&c di μ

$\mu \approx \bar{x}$ stima puntuale

$\mu = \bar{x} \pm r$ stima intervallo / fornice ←

l'aspetto che $\mu \in [\bar{x} - r; \bar{x} + r]$ con livello di confidenza $1 - \alpha$

95% 90%
↙ ↘

Il significato è che " $1 - \alpha = P(\mu \in [\bar{x} - r; \bar{x} + r])$ "

Non è una vera probabilità $\bar{n} = 10,1$ $r = 0,3$

$\mu \in [9,8; 10,4]$ al 95% di confidenza

tutto deterministico: o è vera o no

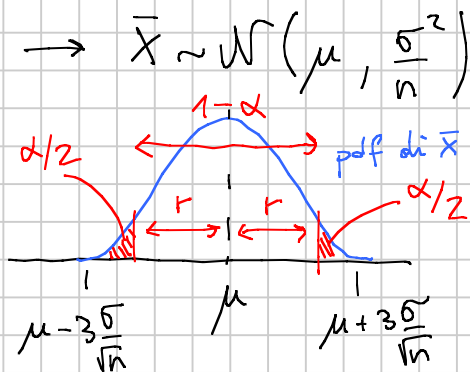
= fiducia nella affermazione

- Serve la legge di \bar{X}

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n x_i \quad x_i \sim \mathcal{N}(\mu, \sigma^2) \text{ indep.}$$

* la classe \mathcal{N} è riproducibile ovvero se sommo v.a.a. normali indipendenti (non occorre isonomie) ottengo una v.a. \mathcal{N}

(more on that)



99,73%

Esiste $r > 0$ t.c.

$$\begin{aligned} 1-\alpha &= P(\bar{X} \in [\mu-r; \mu+r]) \\ &= P(|\bar{X} - \mu| \leq r) \\ &= P(\mu \in [\bar{X}-r; \bar{X}+r]) \end{aligned}$$

$$1-\alpha = P(\mu \in [\bar{X}-r; \bar{X}+r]) \leftarrow \text{intervallo random}$$

$$\mu \in [\bar{x}-r; \bar{x}+r] \text{ con livello di conf. } 1-\alpha$$

- Come trovo r t.c. $P(|\bar{X} - \mu| \leq r) = 1-\alpha$?

$$\bar{X} \sim \mathcal{N}(\mu, \frac{\sigma^2}{n})$$

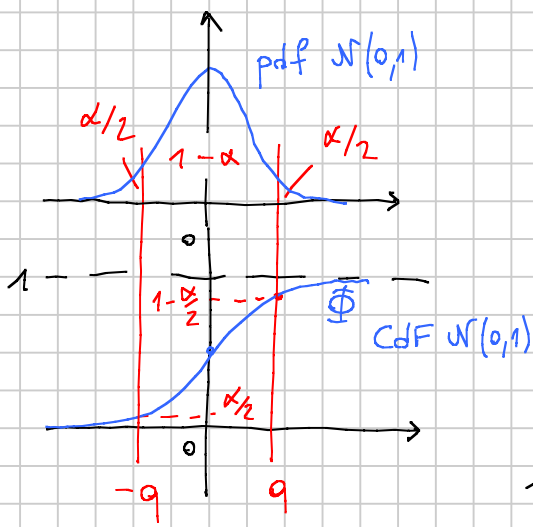
$$\bar{X} - \mu \sim \mathcal{N}(0, \frac{\sigma^2}{n})$$

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

$$1-\alpha = P\left(\left|\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}\right| \leq \frac{r}{\sigma/\sqrt{n}}\right)$$

9

$$\begin{aligned} E(a+bX) &= a+bE(X) \\ \text{Var}(a+bX) &= b^2 \text{Var}(X) \end{aligned}$$



$$q = F_{N(0,1)}^{-1} \left(1 - \frac{\alpha}{2} \right)$$

$$= \Phi^{-1} \left(1 - \frac{\alpha}{2} \right)$$

$$\Phi = F_{N(0,1)}$$

$$r = \frac{\sigma}{\sqrt{n}} q$$

$$1 - \alpha = 95\%$$

$$q = 1,96 \approx 2$$

$$= 90\%$$

$$q = 1,645$$

$$= 99,73\%$$

$$q = 3$$

* Valori di **quantili** sensati sono intorno a 2, o fra 1 e 3

$$\mu \in \bar{x} \pm 1,96 \frac{\sigma}{\sqrt{n}} \text{ al } 95\% \text{ di conf.}$$

ora 4

$$\Phi(t) := \int_{-\infty}^t \frac{1}{\sqrt{2\pi}} e^{-s^2/2} ds$$

non ha scrittura analitica

Funzioni Excel :

$$\Phi = \text{NORMSDIST}$$

DISTRIB. NORM. ST

$$\Phi^{-1} = \text{NORMSINV}$$

....

NORMDIST

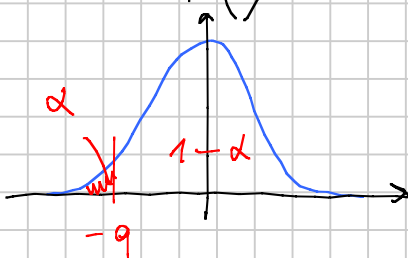
HW: prendere confidenza

HW: Ross Cap 7 pbm 8 e 9

• int. confid. unilaterale $X_i \sim N(\mu, \sigma^2)$ σ nota

voglio $\mu \leq UCL$ con lvl di conf $1 - \alpha$

$$1 - \alpha = P(\mu \leq \bar{X} + r) = P(\bar{X} - \mu \geq -r) = P\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \geq -r \frac{\sqrt{n}}{\sigma}\right)$$



$$q = \Phi^{-1}(1 - \alpha)$$

$$-q = \Phi^{-1}(\alpha)$$

$$\text{HW: } \Phi(-\alpha) = 1 - \Phi(\alpha)$$

$$\Phi^{-1}(1 - \alpha) = -\Phi^{-1}(\alpha)$$

• Generalizzazione : funzione ausiliaria

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

in generale \bar{X} è una funzione che dipende da :

- i. i dati
- ii. il parametro da stimare
- iii. altre grandezze note

e la cui distribuzione sia nota

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

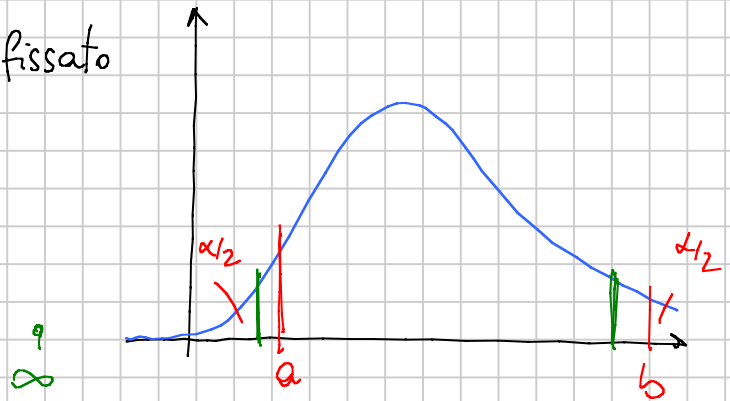
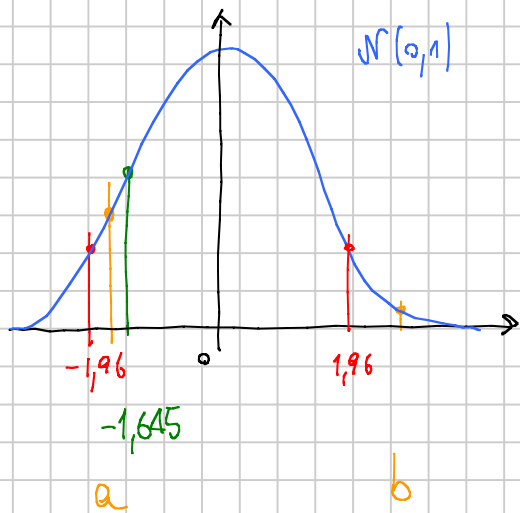
↓ ↓
quantile / i

$$1-\alpha = P\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \in I_q\right) = \dots = P(\mu \in I_{\bar{X}, q})$$

↑
ricavo μ

$\mu \in I_{\bar{X}, q}$ con lvl di conf...

• Generalizzazione : altri intervalli



$$1-\alpha = P(a \leq \mathcal{N}(0,1) \leq b)$$

★ L'intervallo di ampiezza minima e prob $1-\alpha$ è quello in cui il valore di f nei due estremi è uguale

• Caso Bernoulli/binomiale

faccio un sondaggio
n persone, chiedo
se seguono il calcio

Bernoulliana:

X_1, X_2, \dots, X_n

↓ ↓
1 0 0 ... 1 ... 0

1: segue il calcio

0: no

binomiale: su n persone X sono quelle che seguono il calcio

→ Legame: $X = \sum_{i=1}^n X_i$

$X_i \sim \text{bin}(1, p)$ indipendenti

$X \sim \text{bin}(n, p)$

dove p: probabilità che una persona a caso segua il calcio
p incognita, da stimare

★ Esempio industriale: arriva fornitura 10000 pezzi ^{popolazione} con
frazione di difetti $p \in [0, 1]$ incognita

Sceglie un campione di n (say 20) li testa e
conta i difetti → X

★ Stimo p puntualmente e con intervallo di conf.

→ stimatore di p: $\hat{p} := \frac{X}{n}$ (HW: MLE?)

$\hat{p} := \bar{X}$ a livello delle X_i

→ legge di \hat{p} ? $X \sim \text{bin}(n, p)$ $\hat{p} = \frac{1}{n} X$

$X \sim \mathcal{N}(np; np(1-p))$ $\hat{p} \sim \mathcal{N}(p; \frac{p(1-p)}{n})$ $\frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} \sim \mathcal{N}(0, 1)$

$X \sim \text{bin}(n, p)$

$\varphi_x(k) = P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}$ $k=0, \dots, n$

n=1 Bernoulli

$P(X=1) = p$ $P(X=0) = 1-p$

n ≥ 1 binomiale

$E(X) = np$

$\text{Var}(X) = np(1-p)$

Classe $\text{bin}(\cdot, p)$ è riproducibile

① Approssimazione normale

$$\frac{\sum_1^n x_i - nE(x_1)}{\sqrt{n \text{Var}(x_1)}} \underset{\sim}{\sim} \mathcal{N}(0,1)$$

approssimativamente distribuito come

$$\sum_1^n x_i \underset{\sim}{\sim} \mathcal{N}(nE(x_1); n \text{Var}(x_1))$$

→ Se x_1, \dots, x_n sono iid di media μ e var σ^2 , se n è abbastanza grande,

$$\sum_1^n x_i \underset{\sim}{\sim} \mathcal{N}(n\mu; n\sigma^2)$$

→ quanto grande n dipende fortemente dalla legge di X_i

$$X_i \sim \mathcal{N} \quad n=1$$

$$X_i \sim \text{bin}(1, 10^{-10}) \quad n \geq 10^{10}$$

→ l'approssimazione $\sum x_i \underset{\sim}{\sim} \mathcal{N}$ è migliore al centro della campana che sulle code

TLC: $x_1, x_2, \dots \in L^2$ iid

$$\frac{\sum_1^n x_i - nE(x_1)}{\sqrt{n \text{Var}(x_1)}} \xrightarrow[n \rightarrow \infty]{L} \mathcal{N}(0,1)$$

● Campione bernoulliano X : numero totale di "positivi" su n prove
 p la prob. incognita di una prova di essere "positiva"

$$X \sim \text{bin}(n, p) \quad E(X) = np \quad \text{Var}(X) = np(1-p)$$

$$X \approx \mathcal{N}(np; np(1-p)) \quad \text{approssimazione per il TLC}$$

→ intervallo di conf. per p

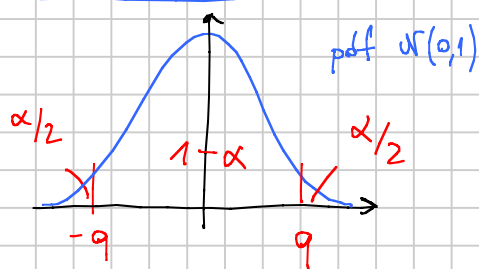
$$p \approx \hat{p} := \frac{X}{n} \quad \text{stimatore puntuale}$$

$$\hat{p} \approx \mathcal{N}\left(p; \frac{p(1-p)}{n}\right) \quad \text{distrib. stimatore (approx)}$$

$$\frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} \approx \mathcal{N}(0, 1) \quad \text{funz. ancillare (approx)}$$

$$\frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \approx \mathcal{N}(0, 1)$$

funz. ancillare (più approx)



$1-\alpha$ livello di conf assegnato
 tipicamente 95%

$$q = \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)$$

$$1-\alpha = P(-q \leq \mathcal{N}(0,1) \leq q) \approx P\left(-q \leq \frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \leq q\right) = P\left(p \in \hat{p} \pm q \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right)$$

$$\text{con lvl di conf } 1-\alpha \quad p \in \hat{p} \pm q \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

★ La II approssimazione non è necessaria, ma semplifica le formule

★ La I approx $\text{bin}(n, p) \approx \mathcal{N}(np; np(1-p))$ $np(1-p) \geq 5$

● Stimatori per la varianza

$$X \in L^2 \quad \text{Var}(X) = E[(X - E(X))^2] \stackrel{\text{(check)}}{=} E(X^2) - E(X)^2$$

X_1, X_2, \dots, X_n campione iid. $\mu = E(X_i)$ $\sigma^2 = \text{Var}(X_i)$

$$\sigma^2 \approx S^2 := \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \left[\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right] \quad \text{var. campionaria}$$

(check)

17, 14, 14, 18, 15, ...

$$S \approx 1,4 \quad \sum x_i^2 \approx 3000 \quad n\bar{x}^2 \approx 3000$$

$$10017, 10014, 10014, \dots \quad S \approx 0,0014$$

★ La II è mal condizionata

$$\star \bar{X}(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i \quad \bar{X}(a + bX_1, a + bX_2, \dots) = a + b\bar{X}$$

$$S^2(x_1, \dots, x_n) = \dots \quad S^2(a + bX_1, a + bX_2, \dots) = b^2 S^2$$

★ S^2 è uno stimatore corretto di σ^2

$$E(S^2) = \frac{1}{n-1} \sum_i E(x_i^2) - \frac{n}{n-1} E(\bar{X}^2) = \frac{n}{n-1} \left\{ \sigma^2 + \mu^2 - \frac{\sigma^2}{n} - \mu^2 \right\} = \sigma^2$$

$\uparrow \sigma^2 + \mu^2$ $\uparrow \text{Var}(\bar{X}) + E(\bar{X})^2 = \frac{\sigma^2}{n} + \mu^2$

→ Ho usato che $E(Y^2) = \text{Var}(Y) + E(Y)^2$

★ Nel caso (piuttosto raro) che si conosca μ

$$S_*^2 := \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

→ (stimatore corretto)

★ La media aritmetica $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ è il punto y che

$$\text{minimizza} \quad \sum_{i=1}^n (x_i - y)^2$$

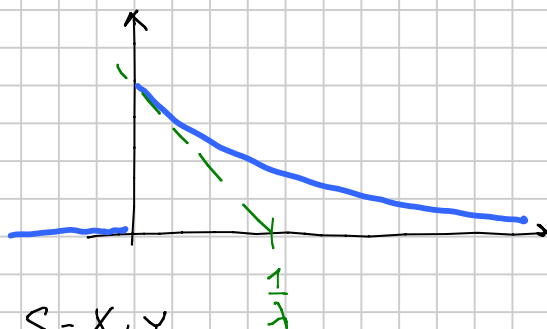
HW: Cosa si deve scegliere per minimizzare $\sum_{i=1}^n |x_i - y_i|$?

HW: S è uno stimatore distorto di σ (sottostima o sovrastima?)

■ DISTRIBUZIONE GAMMA (CHI-QUADRO, ERLANG, ESPONENZIALE)

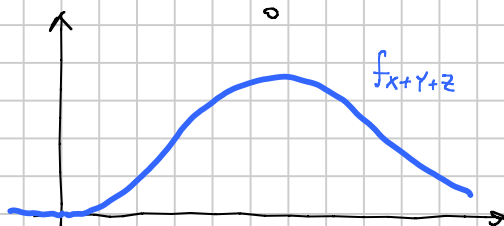
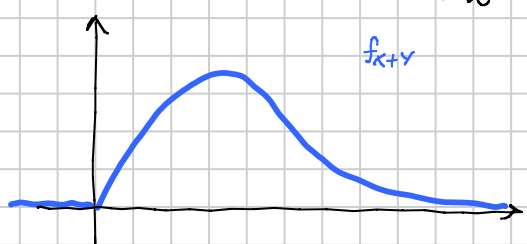
$X \sim \text{expo}(\lambda)$ λ rate / intensità $\lambda = \frac{1}{E(X)}$

$$f_x(t) = \begin{cases} \lambda e^{-\lambda t} & t \geq 0 \\ 0 & \text{altrove} \end{cases}$$



$X, Y \sim \text{expo}(\lambda)$ indip. $S = X + Y$

$$f_s(t) = f_x * f_y(t) = \int_{-\infty}^{\infty} f_x(s) f_y(t-s) ds = \int_0^t \lambda^2 e^{-\lambda t} ds = \lambda^2 t e^{-\lambda t} \quad t \geq 0$$



$$f(t) = \frac{\lambda^3}{2} t^2 e^{-\lambda t}$$

$X_1, X_2, \dots, X_n \sim \text{expo}(\lambda)$ iid. $S = X_1 + \dots + X_n$

$$f_s(t) = c(n) \lambda^n t^{n-1} e^{-\lambda t}, \quad t \geq 0 \quad \text{Erlang}(n, \lambda)$$

$$\int_0^{\infty} t^k e^{-t} dt = k! \quad \Rightarrow \quad c(n) = \frac{1}{(n-1)!}$$

• Legge gamma

$X \sim \text{gamma}(\alpha; \lambda)$

$X \sim \text{gamma}(\alpha; \beta)$

$$\beta = \frac{1}{\lambda}$$

$$\alpha > 0 \quad \beta > 0$$

$$f_x(t) = \frac{1}{\Gamma(\alpha)} \lambda^\alpha t^{\alpha-1} e^{-\lambda t}, \quad t \geq 0$$

* $E(X) = \alpha\beta$ $\text{Var}(X) = \alpha\beta^2$ $aX \sim \text{gamma}(\alpha; a\beta)$

* $\text{gamma}(\cdot, \beta)$ è riproducibile

- Legge chi-quadro $k = 1, 2, \dots$ numero di gradi di libertà

$$X \sim \chi^2(k) \sim \text{gamma}\left(\frac{k}{2}; 2\right)$$

$$E(X) = k \quad \text{Var}(X) = 2k$$

$$\chi^2(2) \sim \text{expo}\left(\frac{1}{2}\right)$$

- ★ $X_1, X_2, \dots, X_n \sim \mathcal{N}(0, 1)$ indipendenti

$$\sum_{i=1}^n x_i^2 \sim \chi^2(n)$$

- Distribuzione della varianza campionaria per dati normali

$$X_1, X_2, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2) \quad \sigma \text{ incognita}$$

$$\rightarrow \mu \text{ nota} \quad S_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

$$x_i \sim \mathcal{N}(\mu, \sigma^2)$$

$$x_i - \mu \sim \mathcal{N}(0, \sigma^2)$$

$$\frac{x_i - \mu}{\sigma} \sim \mathcal{N}(0, 1)$$

$$\star \quad \frac{S_x^2}{\sigma^2} n \sim \chi^2(n)$$

$$\rightarrow \mu \text{ incognita} \quad S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\star \quad \frac{S^2}{\sigma^2} (n-1) \sim \chi^2(n-1)$$

funzione ausiliare per σ

$$\rightarrow \frac{S^2}{\sigma^2} (n-1) \sim \text{gamma}\left(\frac{n-1}{2}; 2\right)$$

$$S^2 \sim \text{gamma}\left(\frac{n-1}{2}; \frac{2}{n-1} \sigma^2\right)$$

$$E(S^2) = \alpha\beta = \sigma^2 \quad \text{Var}(S^2) = \alpha\beta^2 = \frac{2}{n-1} \sigma^4 \quad \text{sim. c\&c}$$

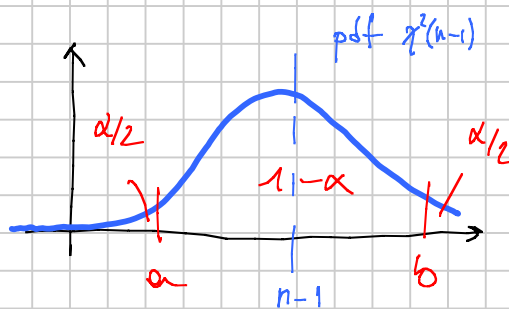
● Intervallo di conf. per la varianza di una popolaz. normale

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

$$\sigma^2 \approx S^2$$

$$\frac{S^2}{\sigma^2} (n-1) \sim \chi^2(n-1)$$

$1-\alpha$ livello di confidenza



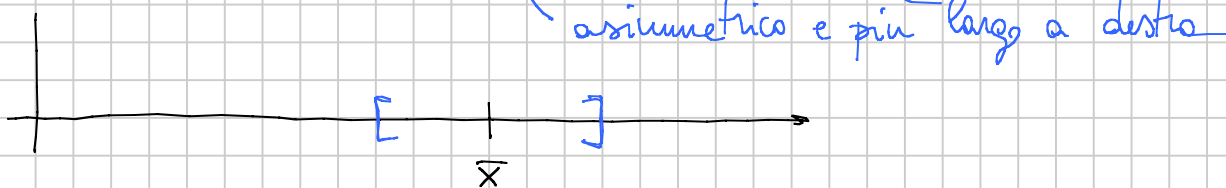
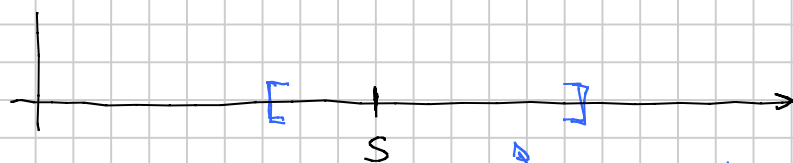
$$a < n-1 \quad \left\| \quad a = F_{\chi^2(n-1)}^{-1}\left(\frac{\alpha}{2}\right) = \text{INV.CH1}\left(1-\frac{\alpha}{2}; n-1\right)$$

$$b > n-1 \quad \left\| \quad b = F_{\chi^2(n-1)}^{-1}\left(1-\frac{\alpha}{2}\right) = \text{INV.CH1}\left(\frac{\alpha}{2}; n-1\right)$$

$$1-\alpha = P(a \leq \chi^2(n-1) \leq b) = P\left(a \leq \frac{S^2}{\sigma^2} (n-1) \leq b\right) = P\left(\frac{n-1}{b} S^2 \leq \sigma^2 \leq \frac{n-1}{a} S^2\right)$$

$$\sigma^2 \in \left[\frac{n-1}{b} S^2 ; \frac{n-1}{a} S^2 \right] \quad \text{con lvl di conf } 1-\alpha$$

$$\sigma \in \left[\underbrace{\sqrt{\frac{n-1}{b}} S}_{< S} ; \underbrace{\sqrt{\frac{n-1}{a}} S}_{> S} \right] \quad \text{con lvl di conf } 1-\alpha$$



HW: $Z \sim \mathcal{N}(0,1)$ $W = Z^2$ $f_W = ?$

$$F_W(t) = P(W \leq t) = P(Z^2 \leq t) = \dots = \Phi\left(\sqrt{t}\right) - \Phi\left(-\sqrt{t}\right)$$

1) $f_W = F_W'$ ci riconosco la gamma

2) W_1, W_2 iid $f_{W_1+W_2} = ?$ (expo $(\frac{1}{2})$)

3) $Z_1, Z_2 \sim \mathcal{N}(0,1)$ $f_{Z_1^2+Z_2^2} = ?$ (expo $(\frac{1}{2})$)

HW: Intervallo di confidenza per la media di un campione
esponenziale

$$X_1, X_2, \dots, X_n \sim \text{expo}\left(\frac{1}{\beta}\right)$$

$$\beta \approx \bar{X}$$

eccetera ...

TEST STATISTICI

● Caso normale σ nota μ incognita
 $X_1, X_2, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$ indep.

→ il test è sul parametro incognito (μ)

→ si fanno due ipotesi in opposizione e ci si domanda "quale sia vera"?

$$H_0: \mu = 10 \quad H_1: \mu \neq 10$$

μ_0
 valore target

AQL (acceptable quality level)

10,04 9,98 10,01 10,03 9,95 ...

Q: μ non può essere esattamente uguale a 10 **NONSENSE**

Non cerco quale è vera: so che è H_1

Cerco in effetti di capire:

se H_0 è plausibile, visti i dati
 è compatibile con i dati } H_0 può essere accettata (test negativo)

se c'è evidenza statistica che H_1 sia vera } H_0 va rifiutata (test positivo)

★ Come si effettua il test?

dati → black box → decisione

	test	
	0	1
0	✓	FP
1	fn	✓

fp: errori di I specie

fn: errori di II specie

α : probs. di un fp

β : probs. di un fn

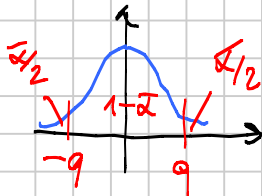
$\boxed{bb} \rightarrow \alpha \leq \bar{\alpha}$
tipicamente 5%
livello di significatività
(assegnato)

① Primo modo: calcolo la statistica del test

$$Z := \frac{\bar{X} - \mu_0}{\sigma / \sqrt{n}}$$

la funz. ancillare era $\frac{\bar{X} - \mu}{\sigma / \sqrt{n}} \sim \mathcal{N}(0,1)$

$$q = \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)$$



$\boxed{bb} \rightarrow \begin{cases} \text{se } |Z| \leq q & \text{accetto } H_0 \\ \text{se } |Z| > q & \text{rifiuto } H_0 \end{cases}$ vera H_0

$$\alpha = P(fp) = P(\text{dire } H_1 \mid \text{vera } H_0) = P(|Z| > q) = P\left(\left|\frac{\bar{X} - \mu}{\sigma / \sqrt{n}}\right| > q\right) = \alpha$$

$\mathcal{N}(0,1)$

→ Notazione: RA_T la regione di accettazione (di H_0)
per una statistica T

$$\textcircled{1}: Z := \frac{\bar{X} - \mu_0}{\sigma / \sqrt{n}} \quad RA_Z = [-q; q]$$

② secondo modo: uso lo stimatore \bar{X}

$$RA_{\bar{X}} = \mu_0 \pm q \frac{\sigma}{\sqrt{n}}$$

→ Attenzione: assomiglia all'intervallo di conf

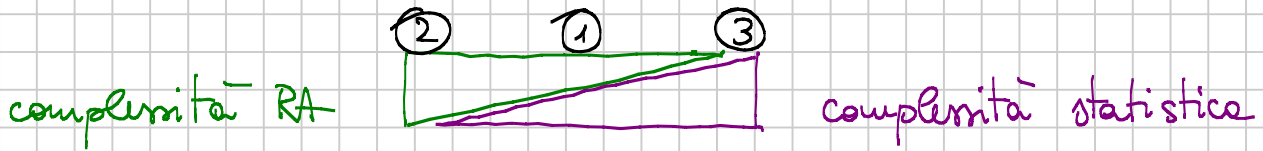
$$\mu \in \bar{x} \pm q \frac{\sigma}{\sqrt{n}}$$

③ terzo modo: calcolo il p dei dati (p-value)

$$\alpha^* := 2 - 2\Phi(|Z|) \in [0, 1]$$

$$RA_{\alpha^*} = [\bar{x}; 1]$$

HW: verificare che sia equivalente



★ Con il p dei dati non occorre decidere prima \bar{x}

→ se α^* è molto piccolo $\alpha^* = 10^{-4}$ → dico H_1

→ se α^* è grande $\alpha^* = 0,3$ → dico H_0

se α^* è vicino al 5% ovviamente occorre scegliere \bar{x} .

● Curva operativa caratteristica

permette di valutare le performance del test (\square) prima di utilizzarlo

$$h: \mathbb{R} \rightarrow [0, 1]$$

↑ spazio dove vive il parametro (μ)

$$h: y \mapsto P(\text{accettare } H_0) \text{ quando } \mu = y$$

↓ $N(y; \frac{\sigma^2}{n})$

$$h(y) = P_{\mu=y}(\text{accetto } H_0) = P_{\mu=y}\left(\left|\frac{\bar{X}-\mu_0}{\sigma/\sqrt{n}}\right| \leq q\right)$$

$$= P_{\mu=y}\left(\underbrace{-q + \frac{\mu_0 - \mu}{\sigma/\sqrt{n}}}_a \leq \underbrace{\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}}_{N(0,1)} \leq \underbrace{q + \frac{\mu_0 - \mu}{\sigma/\sqrt{n}}}_b\right) = \int_a^b f_{N(0,1)} dx$$

$$= \Phi(b) - \Phi(a)$$

$$b = b(y) = q + \frac{\mu_0 - y}{\sigma/\sqrt{n}}$$

$$a = a(y) = -q + \frac{\mu_0 - y}{\sigma/\sqrt{n}}$$



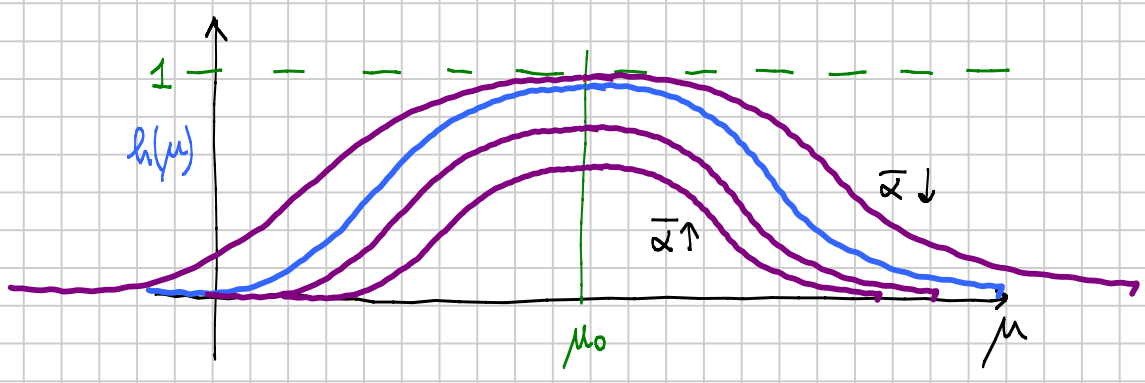
$\mu = \mu_0$ dico H_0 con prob $1-\alpha$

$\mu \approx \mu_0$ dico H_0 con prob poco minore

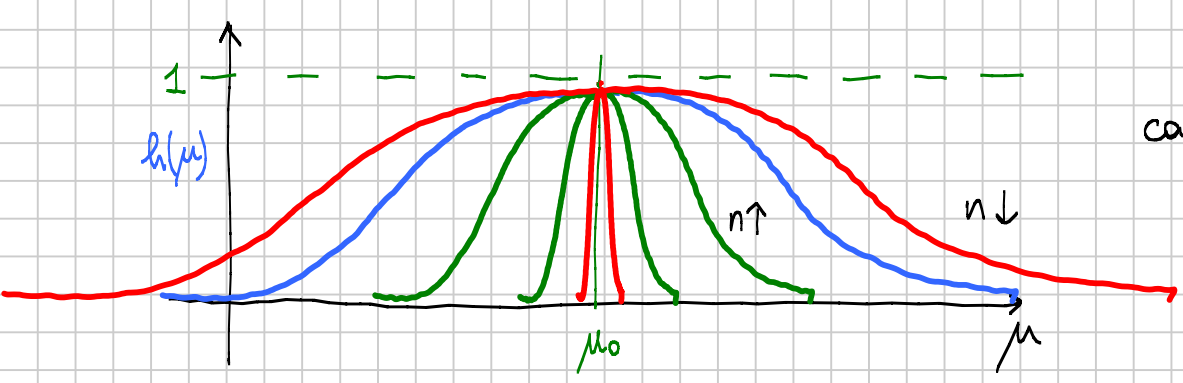
$\mu \neq \mu_0$ $h(\mu) = 1 - \beta(\mu)$ prob di errore di II specie (fn)

μ molto lontano da μ_0 dico H_1 con elevata probabilita

★ La potenza del test e $1 - \beta(\mu)$ la "capacita" di dire H_1



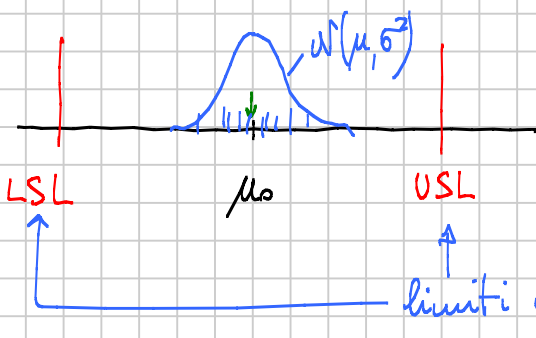
cambio α



cambio n

Se n aumenta il test e piu potente, ma se e troppo potente ricado nel **NONSENSE** dell'ora 7

PROCESS CAPABILITY



$$C_p = \frac{USL - LSL}{6\sigma}$$

$C_p = 1$ se i SL stanno a $\pm 3\sigma$
 $\approx 2,7\%$ di difetti

$$C_{pk} = \min\left(\frac{USL - \mu}{3\sigma}; \frac{\mu - LSL}{3\sigma}\right)$$

$C_p = \frac{4}{3}, \frac{3}{2}, \frac{5}{3}$ buoni

$C_p = 2$ SIX-SIGMA

$$C_{pk} \leq C_p$$

$$\text{Spesso } C_{pk}(\mu) = C_{pk}(\hat{\mu})$$

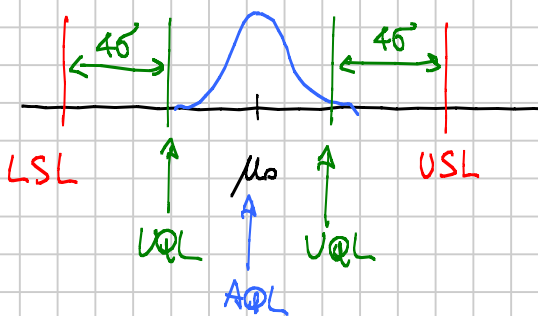
$\hat{\mu}$ stima di μ

★ Devo accorgermi che $\mu \neq \mu_0$ e dire H_1 tutte le volte che la gaussiana è così spostata che sto producendo troppi difetti

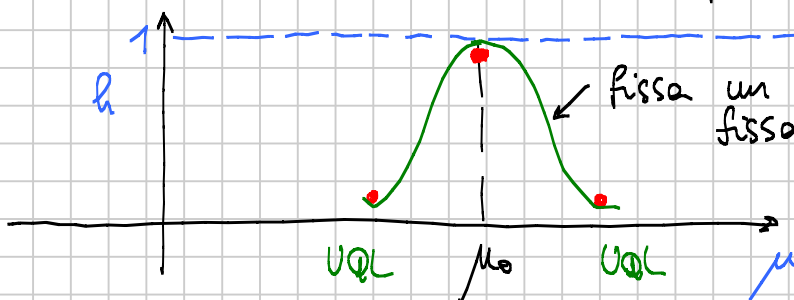
$C_{pk} < 1,33$ inaccettabile (mettiamo)

$$\frac{USL - \mu}{3\sigma} < 1,33 \Leftrightarrow \mu > USL - 4\sigma \quad UQL$$

$$\frac{\mu - LSL}{3\sigma} < 1,33 \Leftrightarrow \mu < LSL + 4\sigma \quad UQL$$



unacceptable quality level



fissa un n minimo
 fissa α
 fissa $\beta(UQL)$

★ Attenzione che nella realtà industriale non siano confusi i
concetti : SL / UQL / $RA_{\bar{x}}$ tutti intervalli $\mu \pm$ qualcosa

- Test per σ^2 in caso di campione normale

$$X_1, X_2, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$$

→ μ nota $S_x^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$

$$\frac{S_x^2}{\sigma^2} n \sim \chi^2(n)$$

→ μ incognita $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$

$$\frac{S^2}{\sigma^2} (n-1) \sim \chi^2(n-1)$$

★ Test bilaterale

$H_0: \sigma = \sigma_0$ dato

$H_1: \sigma \neq \sigma_0$

α livello di significatività

ipotesi nulla

ipotesi alternativa

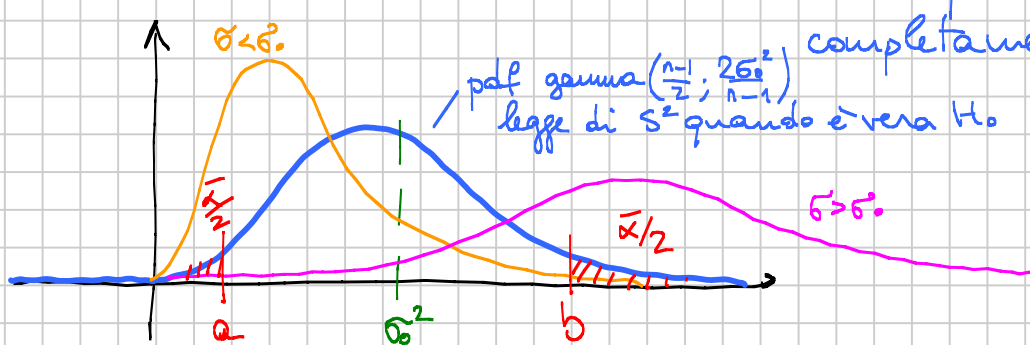
per costruire il test, si suppone vera H_0 e si studia la distribuzione di uno stimatore del parametro incognito determinandone i valori "verosimili"

$S^2 \approx \sigma^2$ suppongo vera $H_0: \sigma = \sigma_0$

$$\frac{S^2}{\sigma^2} (n-1) \sim \chi^2(n-1) \sim \text{gamma}\left(\frac{n-1}{2}; 2\right)$$

$$S^2 \sim \text{gamma}\left(\frac{n-1}{2}; \frac{2\sigma_0^2}{n-1}\right)$$

$$S^2 \sim \text{gamma}\left(\frac{n-1}{2}; \frac{2\sigma^2}{n-1}\right)$$



$$a = F_{\text{gamma}}^{-1}\left(\frac{\alpha}{2}\right) \quad RA_{S^2} = [a; b]$$

$$b = F_{\text{gamma}}^{-1}\left(1 - \frac{\alpha}{2}\right)$$

Se è vera H_0 , $S^2 \in RA_{S^2}$ con prob $1 - \alpha$

② Se $S^2 \in RA_{S^2}$ accetto H_0 , altrimenti la rifiuto perché è un forte indizio che sia vera H_1

① $\frac{S^2}{\sigma_0^2}(n-1) \sim \chi^2(n-1)$ $W := \frac{S^2}{\sigma_0^2}(n-1) \stackrel{H_0}{\sim} \chi^2(n-1)$

$a' = F_{\chi^2(n-1)}^{-1}(\bar{\alpha}/2)$

$b' = F_{\chi^2(n-1)}^{-1}(1 - \bar{\alpha}/2)$

$RA_W = [a' ; b']$

verifico equivalenza con ②

$a' \leq \frac{S^2}{\sigma_0^2}(n-1) \leq b' \iff \frac{a'}{n-1} \sigma_0^2 \leq S^2 \leq \frac{b'}{n-1} \sigma_0^2$

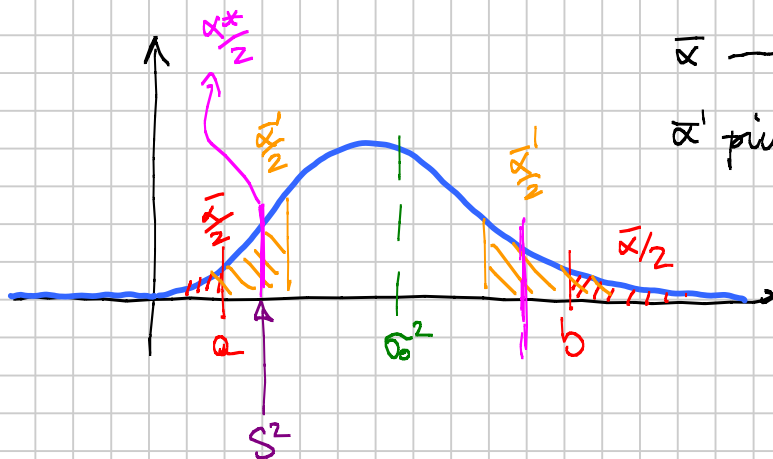
$RA_{S^2} = \left[\frac{a'}{n-1} \sigma_0^2 ; \frac{b'}{n-1} \sigma_0^2 \right]$

devo verificare che $a = \frac{a'}{n-1} \sigma_0^2$ e analogo per b

$a = F_{\text{gamma}(\frac{n-1}{2}; \frac{2\sigma_0^2}{n-1})}^{-1}(\frac{\bar{\alpha}}{2}) = \frac{\sigma_0^2}{n-1} F_{\text{gamma}(\frac{n-1}{2}; 2)}^{-1}(\frac{\bar{\alpha}}{2}) = \frac{\sigma_0^2}{n-1} a'$
 $\chi^2(n-1)$

● Il p-dei-dati

③

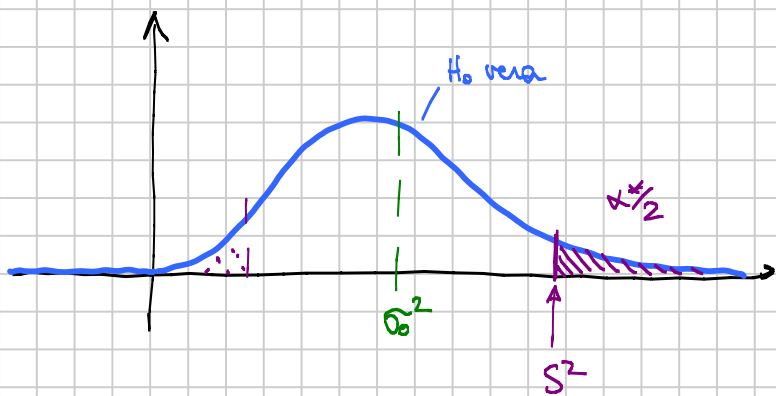


$\bar{x} \rightarrow$ dico H_0

\bar{x} più grande \rightarrow dico H_1

Fissato il campione, dico H_0 se il hl di sign. è molto piccolo e dico H_1 se è molto grande.

Def: il valore critico del hl di sign. per cui cambio idea è il p dei dati e si denota con x^*



Altrimenti detto x^* e' la probabilita' che, sotto ipotesi H_0 , si osservino dei dati "strani" come quelli del campione o piu'

Qui "strani" vuol dire che sembrano indicare H_1

$$x^* = 2 \min \left(F_{\text{gamma}}(S^2); 1 - F_{\text{gamma}}(S^2) \right)$$

↑ coda sx ↑ coda dx

$\rightarrow \bar{x} \leq x^*$ accetto H_0
 $\bar{x} > x^*$ rifiuto H_0

fissato \bar{x} questo vuol dire

$$RA_{x^*} = [\bar{x}; 1]$$

★ Fatto importante: se e' vera H_0 , $x^* \sim \text{unif}(0; 1)$

$$\alpha = P(F_p) = P(x^* \notin RA_{x^*}) = P(x^* < \bar{x}) = \bar{x}$$

se e' vera H_1 , x^* tende ad assumere valori piccoli con probabilita' piu' elevata

HW: 1) Se F e' la Cdf di una v.a. X , allora $F(X) \sim \text{unif}(0; 1)$

2) Se F e' una Cdf assegnata e $U \sim \text{unif}(0; 1)$
 allora $X := F^{-1}(U)$ e' una v.a. con Cdf F

Q: cosa si fa se F non e' invertibile?

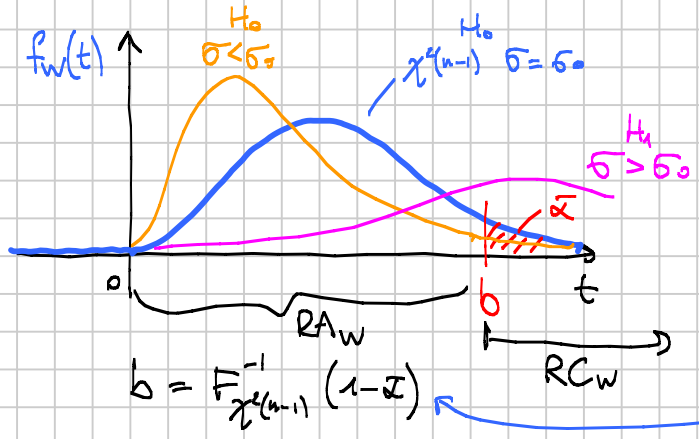
A: cercare Skorokhod sul Williams

3) $x^* \stackrel{H_0}{\sim} \text{unif}(0; 1)$ nel caso dell'esempio

TEST UNILATERALI

$H_0: \sigma \leq \sigma_0$ $H_1: \sigma > \sigma_0$ (caso gaussiano, test su σ , μ incogn.)

① funz. ausiliare $\frac{S^2}{\sigma^2}(n-1) \sim \chi^2(n-1)$
 → statistica $W := \frac{S^2}{\sigma_0^2}(n-1) \sim \chi^2(n-1)$ *se H_0 vera al bordo*



diverse possibili deviazioni di W a seconda di σ

$b = F_{\chi^2(n-1)}^{-1}(1-\bar{\alpha})$ $P(\chi^2(n-1) > b) = \bar{\alpha}$

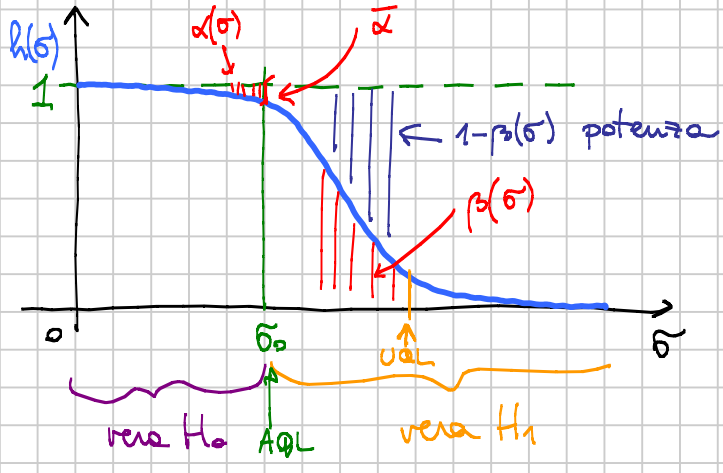
$\alpha = P(f_p) = P(W > b) = P\left(\underbrace{W \frac{\sigma_0^2}{\sigma^2}}_{f.anc.} > b \frac{\sigma_0^2}{\sigma^2}\right) = P\left(\chi^2(n-1) > b \frac{\sigma_0^2}{\sigma^2}\right)$

$H_0: \frac{\sigma_0^2}{\sigma^2} \geq 1 \Rightarrow \alpha \leq \bar{\alpha}$ uguaglianza al bordo

$H_1: \frac{\sigma_0^2}{\sigma^2} < 1 \Rightarrow P(\text{dire } H_1) > \bar{\alpha}$
↑ potenza

* Curva O.C.

$h(\sigma) = P(\text{dire } H_0) = P(W \leq b) = P\left(\chi^2(n-1) \leq b \frac{\sigma_0^2}{\sigma^2}\right)$



Un test ben realizzato dice raramente H_1 se $\sigma \leq AQL$ e dice raramente H_0 se $\sigma \geq UQL$, mentre nella regione intermedia è indeciso

★ Se scambiamo le ipotesi, $AQL < UQL$ posso ottenere un test equivalente

$$H_0: \bar{5} \leq AQL \quad H_1: \bar{5} > AQL \quad \text{lvl di sign. } \alpha \text{ assegnato}$$

... $UQL \quad \beta(UQL) =: \delta$

$$H_0: \bar{5} \geq UQL \quad H_1: \bar{5} < UQL \quad \text{lvl di sign. } \delta \text{ assegnato}$$

→ Il primo test dice H_0 se il secondo dice H_1

★ Invece il test $H_0: \bar{5} \geq AQL \quad H_1: \bar{5} < AQL$ è completamente diverso,

HW: Curva OC per il test bilaterale in $\bar{5}$

■ TEST PER BERNOULLI/BINOMIALE (CONTROLLO QUALITÀ)

Ho un lotto di N pezzi e i difetti dovrebbero essere una frazione inferiore a $3\% = AQL = p_0$. Per verificarlo, testo un campione di n pezzi e verifico che i difetti non siano troppi.

X : # di difetti nel campione

p : frazione difetti nel lotto

assumo che $N \gg 1 \Rightarrow X \sim \text{bin}(n, p)$

$$H_0: p \leq p_0 \quad H_1: p > p_0$$

(HW: perché?
invece cos'è?)

TEST PER BERNOULLI / BINOMIALE (CONTROLLO QUALITÀ)

Ho un lotto di N pezzi e i difetti dovrebbero essere una frazione inferiore a $3\% = AQL = p_0$. Per verificarlo, testo un campione di n pezzi e verifico che i difetti non siano troppi.

X : # di difetti nel campione

p : frazione difetti nel lotto

assumo che $N \gg n \Rightarrow X \sim \text{bin}(n, p)$

$H_0: p \leq p_0$ $H_1: p > p_0$ $\bar{\alpha}$ lvl di significatività

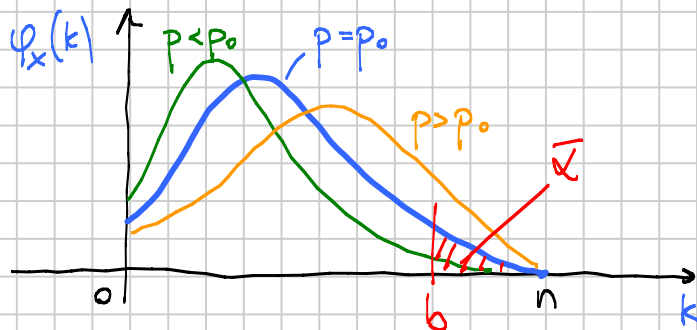
→ Approccio canonico tipo ① approssimato

$$\frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} \sim \mathcal{N}(0, 1)$$

$$Z := \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \sim \mathcal{N}(0, 1)$$

Ho al bordo
non occorre mettere \hat{p}

● Approccio moderno: X è la statistica del test
se H_0 è vera (al bordo) $p = p_0$, allora $X \sim \text{bin}(n, p_0)$



Scelgo il più piccolo b tale che
 $P(\text{bin}(n, p_0) \geq b) < \bar{\alpha}$

Se $X < b$ accetto H_0

Se $X \geq b$ rifiuto H_0

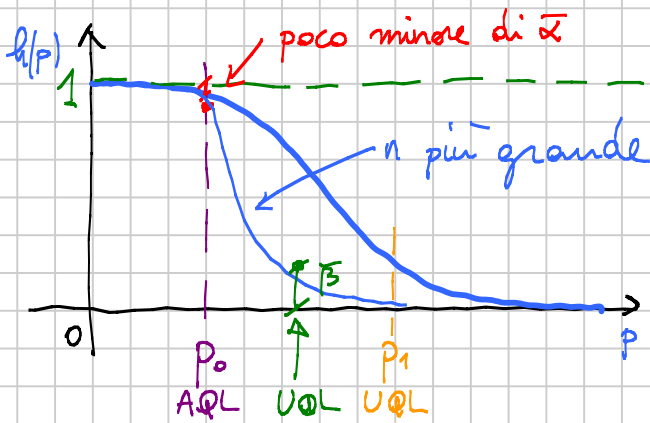
$f_p \Rightarrow p \leq p_0$

★ check: $P(f_p) = P(X \geq b) = P(\text{bin}(X, p) \geq b) \leq P(\text{bin}(X, p_0) \geq b) \leq \bar{\alpha}$

$RA_X = \{0, 1, 2, \dots, b-1\}$

● Curva OC

$$h = h(p) = P(\text{dice } H_0) = P(\text{bin}(n, p) \leq b-1) = F_{\text{bin}(n, p)}(b-1)$$



★ Diverse curve OC al variare di n e b hanno diverse performance che possono essere confrontate, ricavando a posteriori \bar{x} e $\bar{\beta}$ (e/o AQL e UQL che spesso non sono fissati in modo assoluto)

HW: formula del p -dei-dati per questo test

■ COLLEGAMENTO TRA INT. DI CONF. E TEST

Esempio int. conf. p Bernoulliana:

campione di n X : quelli che hanno la caratteristica

$$X \sim \text{bin}(n, p) \quad \hat{p} = \frac{X}{n}$$

$$p \in \hat{p} \pm q \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

↑
quantile gaussiano

Quanti sono i mancini? $n=5$ $X=0$ non funziona

→ Si può dare almeno un upper bound? Sì:

$$90\%? \quad p_0 = 90\%$$

$$H_0: p \geq p_0 \quad H_1: p < p_0 \quad \bar{x} = 5\%$$

$$\alpha^* = P(\text{bin}(n, p_0) = 0) = (1-p_0)^n = (0,1)^5 = 10^{-5} \ll \bar{x} \quad H_1!$$

Modifico p_0 fino a che il test non dice H_0

$$(1-p_0)^5 = 5\% \quad p_0 \approx 45\%$$

→ Questa procedura equivale a fare l'intervallo di confidenza $p \in [0; 45\%]$ con lvl di conf $1-\alpha = 95\%$

★ Fatto generale: l'intervallo di confidenza di un parametro θ a lvl di conf $1-\alpha$ è l'insieme dei valori θ_0 tali che il test $H_0: \theta = \theta_0$ (o unilaterale) accetta H_0 .

★ Fatto generale: se il test dice H_1 vuol dire che "è sicuro", indipendentemente da quanto piccolo sia n (tiene conto di quanto vale n , per dire H_1)

ora 12

● Funzioni ausiliarie

★ Campione $\mathcal{N}(\mu, \sigma^2)$

1) per μ, σ note $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim \mathcal{N}(0,1)$

2) per μ, σ incognita $\frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t(n-1)$

3) per σ, μ nota $\frac{S^2}{\sigma^2} n \sim \chi^2(n)$

4) per σ, μ incognita $\frac{S^2}{\sigma^2} (n-1) \sim \chi^2(n-1)$

Test di Cochran

★ Campione Bernoulli p

5) approssimazione \mathcal{N} $\frac{\hat{p} - p}{\sqrt{p(1-p)/n}} \sim \mathcal{N}(0,1)$

★ Campione esponenziale rate λ

6) $2n\lambda\bar{X} \sim \chi^2(2n)$

▣ THM DI COCHRAN

• Sulla distribuzione normale multivariata

- Def X v.a. a valori in \mathbb{R}^n si dice avere legge normale se $\forall a \in \mathbb{R}^n$ $a \cdot X$ ha legge normale univariata

- Thm se X ha legge normale multivariata $X: \Omega \rightarrow \mathbb{R}^n$ allora la sua densità congiunta è data da:

$$f: \mathbb{R}^n \rightarrow \mathbb{R}_+$$
$$f(x) = (2\pi)^{-n/2} \cdot \frac{1}{\sqrt{\det Q}} \exp \left\{ \frac{1}{2} \langle x - \mu, Q^{-1}(x - \mu) \rangle \right\}$$

dove $\mu \in \mathbb{R}^n$ $\mu = E(X)$ $\mu = (\mu_1, \mu_2, \dots, \mu_n)$ $\mu_i = E(x_i)$

e dove $Q \in M_n$ simmetrica def positiva

$$Q = C(X) \quad Q_{ij} := \text{Cov}(x_i; x_j)$$

matrice di covarianza

* Oppure X è concentrata su un sottospazio affine se $\det Q = 0$ e si può scrivere comunque la legge

- Cor Se Q è diagonale, le componenti x_i sono indipendenti
"per la legge di $\text{Cov} = 0 \Rightarrow$ indipendenza"

- Prop Se $X \sim \mathcal{N}(\mu; Q)$ e $N \in M_{k,n}$ allora

1) NX è normale

2) $NX \sim \mathcal{N}(N\mu; NQN^T)$

* Thm e Prop 1) discendono dal Thm di Lévy sulla funz. caratteristica (trasf. di Fourier)

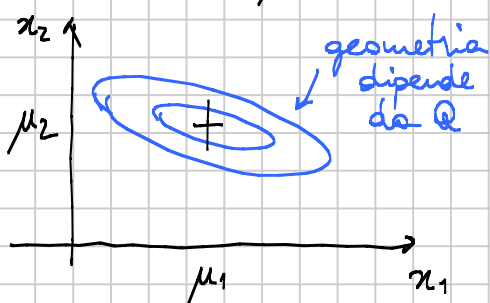
Dim 2) $Y = NX$ $Y_i = \sum_{j=1}^n N_{ij} X_j$

$E(Y_i) = \sum_{j=1}^n N_{ij} E(X_j) = \sum_{j=1}^n N_{ij} \mu_j = [N\mu]_i$

$[C(Y)]_{ij} := \text{Cov}(Y_i; Y_j) = \text{Cov}\left(\sum_{h=1}^n N_{ih} X_h; \sum_{l=1}^n N_{jl} X_l\right)$
 $= \sum_{h,l} N_{ih} N_{jl} \text{Cov}(X_h; X_l) = \sum_{h,l} N_{ih} N_{jl} [C(X)]_{h,l} = \sum_{h,l} N_{ih} Q_{h,l} N_{jl} = [NQNT^T]_{ij}$

HW: $v \in \mathbb{R}^n$ $v + X \sim \mathcal{N}(v + \mu; Q)$

* $n=2$ $X \sim \mathcal{N}(\mu, Q)$ $f: \mathbb{R}^2 \rightarrow \mathbb{R}_+$ $Q = \begin{pmatrix} \sigma_1^2 & c \\ c & \sigma_2^2 \end{pmatrix}$



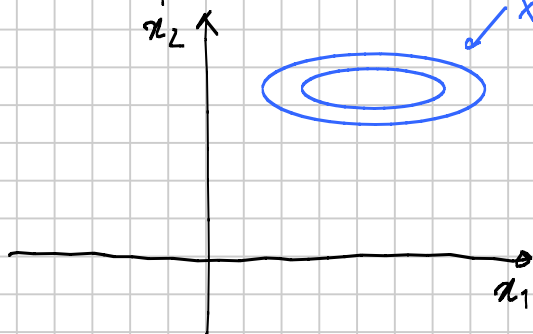
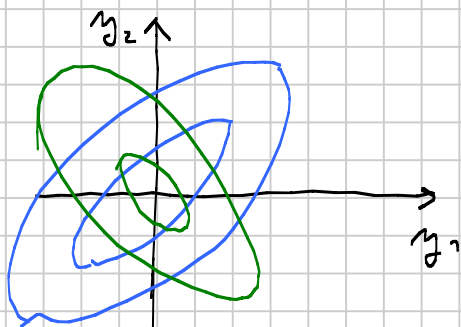
↑
curve di livello

standardizzato :

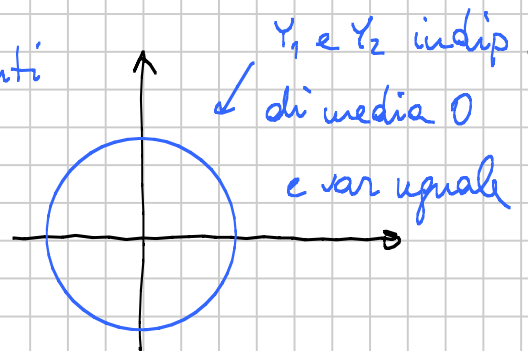
$Y_i := \frac{X_i - \mu_i}{\sigma_i} \quad i=1,2$

$E(Y_i) = 0 \quad \text{Var}(Y_i) = 1$

$Y \sim \mathcal{N}\left(0; \begin{pmatrix} 1 & \alpha \\ \alpha & 1 \end{pmatrix}\right)$



↑ X_1 e X_2 indipendenti



↑ Y_1 e Y_2 indep. di media 0 e var. uguale

* Dato $X \sim \mathcal{N}(\mu, Q)$ esiste una rotazione N di \mathbb{R}^n tale che NX ha componenti indipendenti (Thm spettrale)

* Se $X \sim \mathcal{N}(\mu, \sigma^2 I)$ allora qualunque rotazione non modifica l'indipendenza delle componenti

● Enunciato pratico del tlm di Cochran

i. $X \sim \mathcal{N}(\mu; \sigma^2 I)$ (μ e σ incognite)

ii. $\mu \in V$ sottospazio vettoriale di \mathbb{R}^n $k = \dim V$

Allora 1) la proiezione ortogonale $\pi_V(x)$ di x su V è lo stimatore di massima verosimiglianza di μ e minimizza la distanza $|x - \pi_V(x)|$ e' corretto e consistente

2) $SS := |x - \pi_V(x)|^2$ è una v.a. indipendente da $\pi_V(x)$ e inoltre $\frac{SS}{\sigma^2} \sim \chi^2(n-k)$

● Applicazione a campione normale

x_1, x_2, \dots, x_n $x_i \sim \mathcal{N}(\mu, \sigma^2)$ indep, μ, σ incognite

$X = (x_1, \dots, x_n) \sim \mathcal{N}(\underbrace{(\mu, \mu, \dots, \mu)}_{V \in \mathbb{R}^n}; \sigma^2 I)$

$V \in V := \text{Span}((1, 1, \dots, 1))$ $k = \dim(V) = 1$ $\pi_V(x) = (\gamma, \gamma, \dots, \gamma)$

1) $SS := |x - \pi_V(x)|^2 = \sum_{i=1}^n (x_i - [\pi_V(x)]_i)^2 = \sum_{i=1}^n (x_i - \gamma)^2$

$\frac{dSS}{d\gamma} = \sum_i -2(x_i - \gamma) = 0$ $\sum x_i = n\gamma \Rightarrow \gamma = \bar{x}$
 ↑
 punto (convergenza)

$\pi_V(x) = (\bar{x}, \bar{x}, \dots, \bar{x})$

2) $SS := \sum_{i=1}^n (x_i - \bar{x})^2$ *sum of squares*

$S^2 := \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{SS}{n-1}$

$\frac{S^2}{\sigma^2} (n-1) = \frac{SS}{\sigma^2} \sim \chi^2(n-1)$ indep da \bar{x}

Grazie all'indipendenza: $\frac{\bar{x} - \mu}{S/\sqrt{n}} = \frac{\frac{\bar{x} - \mu}{\sigma/\sqrt{n}}}{\frac{S}{\sigma}} =: \frac{Z}{\sqrt{\frac{W}{n-1}}}$

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

sono indipendenti

$$W := \frac{S^2}{\sigma^2}(n-1) \sim \chi^2(n-1)$$

Perciò la legge di $\frac{Z}{\sqrt{\frac{W}{n-1}}}$ è calcolabile e dipende solo da n

• Def Se $Z \sim \mathcal{N}(0,1)$ e $W \sim \chi^2(k)$ indipendenti

$\frac{Z}{\sqrt{\frac{W}{k}}} \sim t(k)$ ha legge t di Student con k gradi di libertà

★ Prossimamente facciamo qualche martedì 16:30-18:30

● Enunciato pratico del thm di Cochran

i. $X \sim \mathcal{N}(\mu; \sigma^2 I)$ (campioni omoschedastico) (μ e σ incognite)

ii. $\mu \in V$ sottospazio vettoriale di \mathbb{R}^n $k = \dim V$

Allora 1) la proiezione ortogonale $\pi_V(x)$ di x su V è lo stimatore di massima verosimiglianza di μ e minimizza la distanza $|x - \pi_V(x)|$ e corretto ~~e consistente~~

2) $SS := |x - \pi_V(x)|^2$ è una v.a. indipendente da $\pi_V(x)$ e inoltre $\frac{SS}{\sigma^2} \sim \chi^2(n-k)$

Dim v_1, v_2, \dots, v_k base ortonormale di V

① $\pi_V(x) := \sum_{i=1}^k \langle x; v_i \rangle v_i$ proiezione ortogonale su V

$$\langle x - \pi_V(x); v \rangle = 0 \quad \forall v \in V$$

$$\forall v \in V \quad |x - v|^2 = |(x - \pi_V(x)) - (v - \pi_V(x))|^2 = |x - \pi_V(x)|^2 + \underbrace{|v - \pi_V(x)|^2}_{\in V} \geq |x - \pi_V(x)|^2$$

$$\begin{aligned} \text{MLE} \quad \ell(\nu) &= \log f(x; \mu = \nu) = c - \frac{1}{2} \langle x - \nu; \sigma^2 I (x - \nu) \rangle \\ &= c - c_0 |x - \nu|^2 \quad c_0 > 0 \end{aligned}$$

$\pi_V(x)$ è MLE di μ

$$E(\pi_V(x)) = \pi_V(E(x)) = \pi_V(\mu) = \mu \quad \text{stimatore corretto}$$

$$\textcircled{2} \quad R = X - \pi_V(X) \quad R = (R_1, R_2, \dots, R_n) \quad \text{considerati residui}$$

v_{k+1}, \dots, v_n vettori che completano v_1, \dots, v_k a base ortonormale di \mathbb{R}^n

$$N = \begin{pmatrix} -v_1- \\ -v_2- \\ \dots \\ -v_n- \end{pmatrix} \quad [NX]_i = \langle v_i; X \rangle \quad \forall i$$

$$[N\pi_V(X)]_i = \left[N \sum_{j=1}^k \langle v_j; X \rangle v_j \right]_i = \begin{cases} \langle v_i; X \rangle & i=1, 2, \dots, k \\ 0 & i=k+1, \dots, n \end{cases}$$

$$[NR]_i = [NX - N\pi_V(X)]_i = \begin{cases} 0 & i=1, 2, \dots, k \\ \langle v_i; X \rangle & i=k+1, \dots, n \end{cases}$$

$$X \sim \mathcal{N}(\mu, \sigma^2 I) \Rightarrow NX \sim \mathcal{N}(N\mu; \sigma^2 I) \Rightarrow \text{componenti indipendenti}$$

\uparrow
 $N\sigma^2 I N^T$

$\Rightarrow N\pi_V(X)$ e NR indipendenti $\Rightarrow \pi_V(X) \in R$ indep.

$$\frac{SS}{\sigma^2} = \frac{|R|^2}{\sigma^2} = \frac{|NR|^2}{\sigma^2} = \sum_{i=k+1}^n \frac{\langle v_i; X \rangle^2}{\sigma^2} \sim \chi^2(n-k)$$

$$\langle v_i; X \rangle \sim \mathcal{N}(\langle v_i; \mu \rangle; \sigma^2) \quad \text{indipendenti per } i=k+1, \dots, n$$

\uparrow \uparrow
 I_V e_V

$$\frac{\langle v_i; X \rangle}{\sigma} \sim \mathcal{N}(0, 1) \quad \text{iid per } i=k+1, \dots, n$$

HW

$$X_1, X_2, \dots, X_n \quad X_i \sim \mathcal{N}(\mu_1, \sigma^2)$$

μ_1, μ_2, σ^2 incognite

$$Y_1, Y_2, \dots, Y_m \quad Y_i \sim \mathcal{N}(\mu_2, \sigma^2)$$

i. $\bar{X} - \bar{Y}$ stimatore MLE, corretto e consistente di $\mu_1 - \mu_2$

S_X^2, S_Y^2 le varianze campionarie

$$S_p^2 := \frac{n-1}{m+n-2} S_X^2 + \frac{m-1}{m+n-2} S_Y^2$$

stimatore pooled di σ^2

ii. S_p^2 stimatore corretto (e consistente) di σ^2

iii. $\frac{\bar{X} - \bar{Y} - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{m} + \frac{1}{n}}} \sim t(n+m-2)$ funz. ausiliaria per $\mu_1 - \mu_2$

iv. S_p^2 ha varianza minima tra tutte le combinazioni convexe di S_X^2 e S_Y^2

Cosa viene adesso:

① Regressione

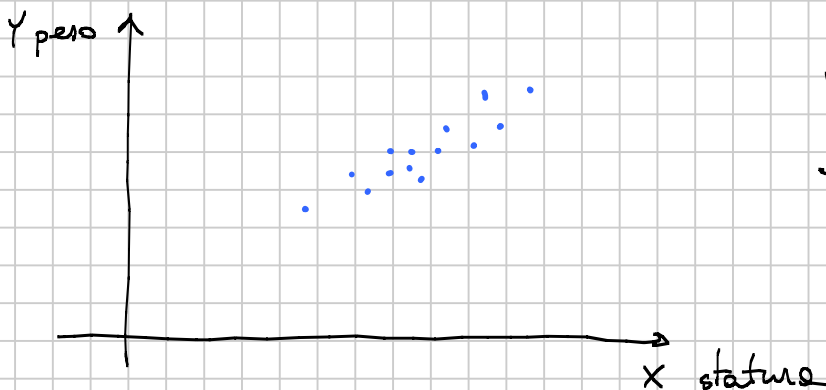
② Analisi della varianza

???

ora 14

REGRESSIONE

Campione bivariato $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$



Vogliamo studiare la relazione tra due variabili

REGRESSIONE LINEARE SEMPLICE

$$Y = \beta_0 + \beta_1 x + e \quad \leftarrow \text{errore additivo } e \sim \mathcal{N}(0, \sigma^2)$$

relazione
lineare

la variabile "causa"
(aka "indipendente"
o "di ingresso")

è immaginata deterministica
(a volte è impostata dallo
sperimentatore)



ipotesi di
omoschedasticità

Notazione

$$Y = \beta_0 + \beta_1 x + e$$

va intesa: $(x_1, Y_1), \dots, (x_n, Y_n)$ $Y_i \sim \mathcal{N}(\beta_0 + \beta_1 x_i; \sigma^2)$

indipendenti

oppure: $Y_i = \beta_0 + \beta_1 x_i + e_i$ $e_i \sim \mathcal{N}(0, \sigma^2)$ i.i.d. (equivalent)

→ ci sono 3 incognite β_0, β_1 e σ^2

→ possono essere stimati B_0, B_1 e S_e^2 saranno gli stimatori

HW: trovare B_0 e B_1 MLE di β_0 e β_1 e usare Cochran:

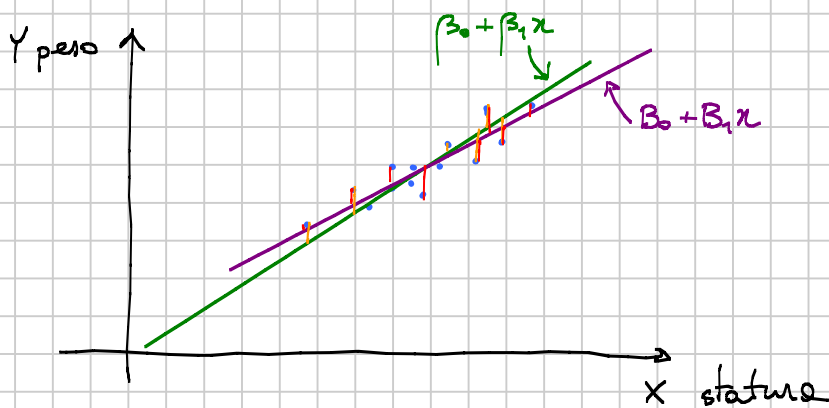
$$S_e^2 := \frac{\sum_{i=1}^n (Y_i - B_0 - B_1 x_i)^2}{n-2}$$

$$\frac{S_e^2}{\sigma^2} (n-2) \sim \chi^2(n-2) \text{ indep. da } B_0 \text{ e } B_1$$

Si trova: $B_1 = \frac{\overline{xY} - \bar{x}\bar{Y}}{\overline{x^2} - \bar{x}^2}$ $B_0 = \bar{Y} - B_1 \bar{x}$ $S_e^2 = \dots$

B_0 e B_1 sono funzioni lineari delle Y_i e quindi hanno legge normale

★ B_0 e B_1 non sono indipendenti: hanno covarianza non nulla



$$e_i = Y_i - (\beta_0 + \beta_1 x_i) \sim \mathcal{N}(0, \sigma^2)$$

$$R_i = Y_i - (B_0 + B_1 x_i) \sim \mathcal{N}(0, ?)$$

complicate e non sono indipendenti

La retta $B_0 + B_1 x$ si chiama *retta di regressione* oppure *best fit* (rispetto al campione)

$$SS = \sum_{i=1}^n R_i^2$$

$$S_e^2 := \frac{SS}{n-2}$$

$$\sigma^2 \approx \frac{\sum e_i^2}{n}$$

se si divide per n si sottostimerebbe σ^2

● Quanto sono precisi gli stimatori? (intervalli di confidenza e test statistici, quindi: funzioni ancillari)

$$B_1 \sim \mathcal{N}(\beta_1; \sigma^2 k_1) \quad B_0 \sim \mathcal{N}(\beta_0; \sigma^2 k_0)$$

dipende solo dalle x_i

faremo tutto per bene nel caso più generale

$$\text{Cov}(B_0; B_1) = \sigma^2 k_{01}$$

$$\frac{B_1 - \beta_1}{\sigma \sqrt{k_1}} \sim \mathcal{N}(0, 1)$$

Labels: $B_1 - \beta_1$ (par. incognito), σ (incognita), $\sqrt{k_1}$ (dati), $\mathcal{N}(0, 1)$ (distrib. nota)

non è f. anc.

$$\frac{B_1 - \beta_1}{S_e \sqrt{k_1}} \sim t(n-2)$$

funz. ancillare

$$Z := \frac{B_1 - \beta_1}{\sigma \sqrt{k_1}} \sim \mathcal{N}(0, 1)$$

$$W := \frac{S_e^2}{\sigma^2} (n-2) \sim \chi^2(n-2)$$

indipendenti per Cochran

$$\frac{B_1 - \beta_1}{S_e \sqrt{k_1}} = \frac{Z}{\sqrt{W/n-2}} \sim t(n-2) \quad \text{def di } t \text{ di Student}$$

$$\frac{B_0 - \beta_0}{\sigma \sqrt{k_0}} \sim \mathcal{N}(0,1)$$

$$\frac{B_0 - \beta_0}{S_e \sqrt{k_0}} \sim t(n-2)$$

- Test per vedere se "c'è regressione" (= Y dipende da x)
= le due variabili sono correlate / legate / ... ecc.

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

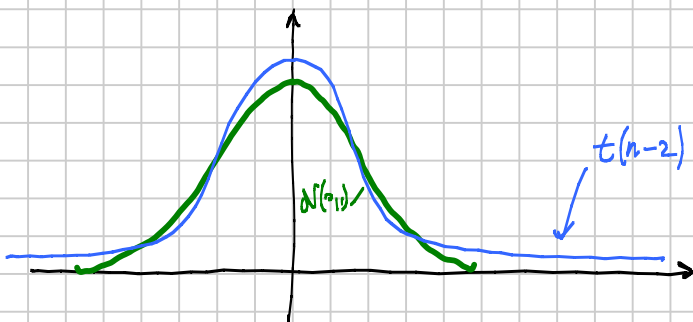
α lvl di significatività

↑
retta orizzontale $\rightarrow Y_i \sim \mathcal{N}(\beta_0, \sigma^2)$ non dipende da x

$$\textcircled{2} T := \frac{B_1}{S_e \sqrt{k_1}} \stackrel{H_0}{\sim} t(n-2)$$

$$RA_T = [-q; q]$$

$$q = F_{t(n-2)}^{-1} \left(1 - \frac{\alpha}{2} \right)$$



- i. è simile a $\mathcal{N}(0,1)$
- ii. è simmetrica rispetto 0
- iii. ha le code pesanti (va a zero in modo polinomiale)
- iv. per $n \rightarrow \infty$ $t(n) \rightarrow \mathcal{N}(0,1)$

★ I quantili della t di Student si trovano come quelli della $\mathcal{N}(0,1)$ usando la CdF opportuna.

Vengono sempre più larghi e la differenza è maggiore quando n è piccolo

● Intervallo di confidenza per la risposta media

x generico, $Y \sim \mathcal{N}(\beta_0 + \beta_1 x; \sigma^2)$ 1- α lvl di confidenza
↑ livello di ingresso ↑ risposta associata a x ↑ risposta media associata a x

→ $\beta_0 + \beta_1 x \in$ intervallo di confidenza

→ stimatore puntuale $\beta_0 + \beta_1 x \approx B_0 + B_1 x$ (corretto e consistente perché lo sono B_0 e B_1)

→ distribuzione stimatore $B_0 + B_1 x \sim \mathcal{N}(\beta_0 + \beta_1 x; \sigma^2 p(x))$

$$\begin{aligned} \text{Var}(B_0 + B_1 x) &= \text{Cov}(B_0 + B_1 x; B_0 + B_1 x) = \text{Cov}(B_0; B_0) + 2x \text{Cov}(B_0; B_1) + x^2 \text{Cov}(B_1; B_1) \\ &= \sigma^2 [k_0 + 2xk_1 + x^2k_1] \end{aligned}$$

polinomio di II grado
i coefficienti dipendono dalle x_i

→ funz. anc. $\frac{B_0 + B_1 x - (\beta_0 + \beta_1 x)}{\sigma \sqrt{p(x)}} \sim \mathcal{N}(0; 1)$

⇓ Cochran

$$\frac{B_0 + B_1 x - (\beta_0 + \beta_1 x)}{S_e \sqrt{p(x)}} \sim t(n-2)$$

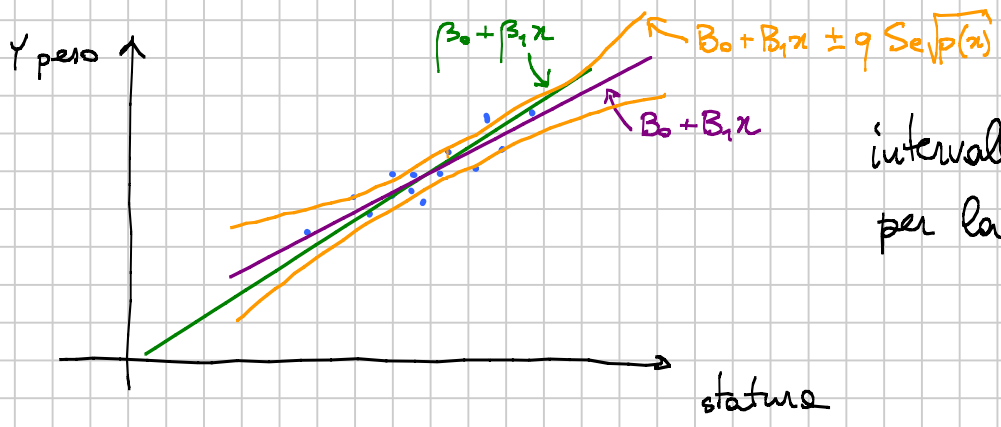
→ quantili $\pm q$ $q = F_{t(n-2)}^{-1}\left(1 - \frac{\alpha}{2}\right)$

→ intervallo di conf

$$\beta_0 + \beta_1 x \in B_0 + B_1 x \pm q S_e \sqrt{p(x)}$$

$$p(x) = \frac{1}{n} + \frac{(x - \bar{x})^2}{n(\bar{x}^2 - \bar{x}^2)}$$

(o qualcosa di simile)



intervallo di confidenza
per la retta verde

● Intervallo di predizione per risposte future

x generico $Y \sim \mathcal{N}(\beta_0 + \beta_1 x; \sigma^2)$

→ $Y \in$ intervallo di predizione

→ Y "viene previsto" da $\beta_0 + \beta_1 x$ tenendo conto dell'incertezza σ

→ $\beta_0 + \beta_1 x \approx B_0 + B_1 x \sim \mathcal{N}(\beta_0 + \beta_1 x; \sigma^2 p(x))$

→ $Y - (B_0 + B_1 x) \sim \mathcal{N}(0; \sigma^2(1 + p(x)))$

risposta di un esperimento futuro
dipendono da Y_1, \dots, Y_n

indipendenti

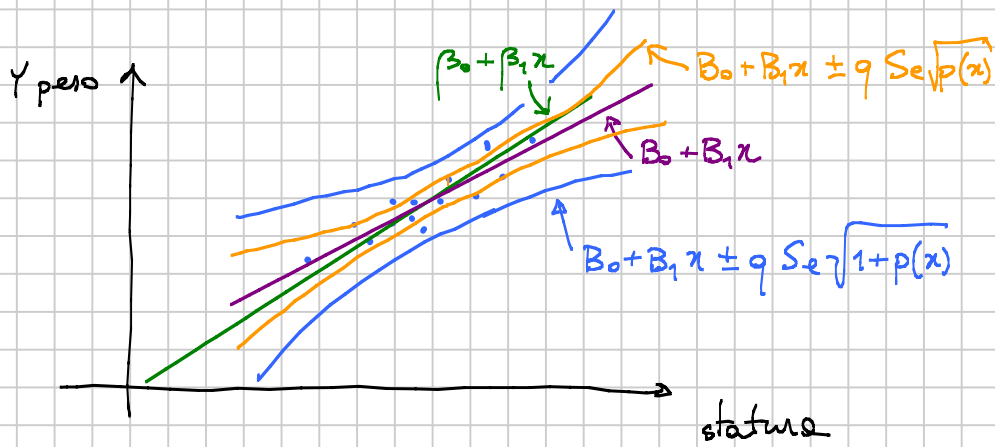
$$\text{Var}(X - Y) = \text{Var}(X) + \text{Var}(Y) - 2 \text{Cov}(X; Y)$$

→ f. ancillare (check)

→ intervallo di predizione:

$$Y \in B_0 + B_1 x \pm q \text{Se} \sqrt{1 + p(x)}$$

unica differenza con l'intervallo di confidenza per la risposta media



intervallo che
contiene una
frazione $1-\alpha$ dei
punti

★ Siccome l'ipotesi di linearità di $Y = \beta_0 + \beta_1 x + e$ è spesso una approssimazione valida localmente, è molto pericoloso usare questi intervalli al di fuori del range dei valori di x osservati nel campione.

HW: x_1, x_2, \dots, x_n campione $\mathcal{N}(\mu, \sigma^2)$ (dati osservati)
 x_{n+1}, \dots, x_{n+m} esperimenti futuri (incognite)

i. intervallo di predizione per $\sum_{i=1}^m x_{n+i}$

ii. intervallo di predizione per $\frac{1}{m} \sum_{i=1}^m x_{n+i}$

ora 16

mar 31 marzo 8:30-10:30
 mer 1 aprile 11:30-13:30] LABORATORIO

sede scientifiche di ingegneria

REGRESSIONE LINEARE MULTIPLA

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + e$$

$$e \sim \mathcal{N}(0, \sigma^2)$$

↳ la variabile di ingresso diventa un vettore

$$(x_{1,1}, x_{1,2}, x_{1,3}, \dots, x_{1,p}, Y_1) \quad \text{primo esperimento}$$

$$(x_{n,1}, x_{n,2}, \dots, x_{n,p}, Y_n) \quad \text{ultimo esperimento}$$

⚡ $p < n$
necessario

⚡ $p \ll n$
preferibile

↳ si aggiunge di solito una x_0 "virtuale" (dummy variable) che vale sempre 1

$$x_{i,0} \equiv 1 \quad i = 1, 2, \dots, n$$

→ scrittura matriciale : $Y = X\beta + E$

$$Y \sim \mathcal{N}(X\beta; \sigma^2 I)$$

dove : $Y = (Y_1, Y_2, \dots, Y_n)$ vettore colonna delle risposte

$X = (x_{i,j})_{i=1, \dots, n; j=0, \dots, p}$ matrice $n \times (p+1)$ dei valori di ingresso

$\beta = (\beta_0, \beta_1, \dots, \beta_p)$ vettore colonna dei coefficienti

$E \sim \mathcal{N}(0, \sigma^2 I)$ vettore gaussiano di \mathbb{R}^n

→ incognite : $\beta \in \mathbb{R}^{p+1}$, $\sigma > 0$

→ tesi di Cochran i. $Y \sim \mathcal{N}(X\beta; \sigma^2 I)$

ii. $\mu \in \text{Im}(X) = X(\mathbb{R}^{p+1}) = V$ ssv. di \mathbb{R}^n di dimensione $p+1$

$\pi_V(Y) =: \hat{Y}$ trovo \hat{Y} minimizzando $|R|^2$ dove $R = Y - \hat{Y}$

Sia $v \in V$ generico $v = Xb$ $b \in \mathbb{R}^{p+1}$ qualsiasi

$$SS = SS(b) = |R|^2 = |Y - Xb|^2 = \langle Y - Xb; Y - Xb \rangle = (Y^T - b^T X^T)(Y - Xb)$$

$$\nabla_b SS = -2(Y^T - b^T X^T)X = 0$$

$B \in \mathbb{R}^{p+1}$ quello che minimizza

$$Y^T X = B^T X^T X \quad \Leftrightarrow \quad X^T Y = X^T X B$$



$$B = (X^T X)^{-1} X^T Y$$

stimatore per β

★ X di rango $p+1 \Rightarrow X^T X$ invertibile

$$\hat{Y} = X B$$

valori previsti per Y

$\rightarrow Y = \alpha B$ equazione di regressione (analogo della retta viola)
 $\alpha \in \mathbb{R}^{p+1}$ vettore riga è un generico valore di ingresso
e Y la risposta associata

$$\rightarrow SS = |R|^2 = \sum_{i=1}^n R_i^2 = |Y - \hat{Y}|^2 = \sum_{i=1}^n \left(Y_i - \underbrace{\sum_{j=0}^p B_j x_{ij}}_{\hat{Y}_i \text{ previsto}} \right)^2$$

$$\frac{SS}{\sigma^2} \sim \chi^2(n-p-1)$$

indipendente da \hat{Y} e quindi da B

\rightarrow Stimatore per σ : errore standard S_e

$$S_e := \sqrt{\frac{SS}{n-p-1}}$$

$$\frac{S_e^2}{\sigma^2} (n-p-1) \sim \chi^2(n-p-1)$$

funz. anc.

HW: $E(S_e^2) = \sigma^2$; S_e^2 è consistente

\rightarrow Distribuzione di B

$$B = (X^T X)^{-1} X^T Y$$

$\therefore N$ matrice $p+1 \times n$

$$Y \sim \mathcal{N}(\mu; \Omega)$$

$$NY \sim \mathcal{N}(N\mu; N\Omega N^T)$$

$$B \sim \mathcal{N}(\beta; \sigma^2 (X^T X)^{-1})$$

$$\text{verifico: } N\mu = (X^T X)^{-1} X^T X \beta = \beta$$

$$NQNT = (X^T X)^{-1} X^T \sigma^2 I X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1}$$

• Esempio di inferenza

Test per vedere se la variabile x_j è significativa
(= se Y dipende da x_j)

$$H_0: \beta_j = 0 \quad H_1: \beta_j \neq 0 \quad \alpha \text{ lvl di significatività}$$

$$\rightarrow B_j \approx \beta_j \quad B_j \sim \mathcal{N}(\beta_j; \sigma^2 [(X^T X)^{-1}]_{jj})$$

$$\rightarrow \text{funz. ausiliare} \quad \frac{B_j - \beta_j}{\text{Se} \sqrt{[(X^T X)^{-1}]_{jj}}} \sim t(n-p-1)$$

$$\rightarrow \text{statistico} \quad T := \frac{B_j}{\text{Se} \sqrt{[(X^T X)^{-1}]_{jj}}} \quad RA_T = [-q; q] \quad q = F_{t(n-p-1)}^{-1}(1 - \frac{\alpha}{2})$$

SELEZIONE DELLE VARIABILI

$$Y = \sum_{j=0}^p \beta_j x_j + e \quad e \sim \mathcal{N}(0, \sigma^2)$$



candidati

x_j non impattante $\Rightarrow \beta_j = 0 \not\Rightarrow B_j = 0$

$B_j \neq 0$ sempre non nullo

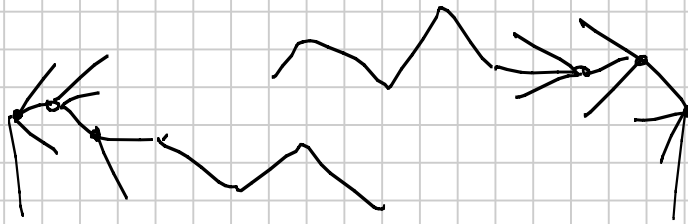
$$Y = B_0 + B_1 x_1 + \dots + B_j x_j + \dots + B_p x_p$$

- ★ Le variabili non impattanti andrebbero tutte tolte dal modello per evitare "rumore" se si usa per prevedere Y su dati futuri
- inoltre tipicamente la precisione sugli altri stimatori migliora, togliendo le variabili irrilevanti

Lo strumento base per la selezione è il test $H_0: \beta_j = 0$ $H_1: \beta_j \neq 0$ che va fatto per tutte le variabili, tranne il termine noto.

- ★ Se tolgo una variabile il risultato dei test per le altre può cambiare
- si toglie sempre una sola variabile alla volta
- l'ordine in cui si tolgono è influente: non per forza se x_j ha $\alpha^* = 80\%$ e x_k ha $\alpha^* = 70\%$ è giusto togliere prima x_j
- il metodo di selezione stepwise backward però funziona proprio così: toglie iterativamente la variabile con α^*_{max}

→ il metodo di selezione *stepwise forward* inserisce una var alla volta, scegliendo quella con impatto maggiore (confronta x^* globale delle regressioni su modelli diversi)



non sempre si incontrano

→ la forward si fa male con Excel: nel caso si sceglie, occorre uno strumento più sofisticato

★ Analisi della varianza di regressione

$$Y = X\beta + e$$

notazione matriciale

$$Y = XB + R$$

↓ toglie alcune variabili : $\tilde{B} = (\tilde{\beta}_0, \tilde{\beta}_1, \dots, \tilde{\beta}_k, 0, \dots, 0)$

$\tilde{\beta}_p$

$$Y = X\tilde{B} + \tilde{R}$$

$$|R| \leq |\tilde{R}|$$

perché B è ottimo

$$|R|^2 = SS = SS_R = \sum_{i=1}^n R_i^2$$

$$SS_{\tilde{R}} = \sum_{i=1}^n \tilde{R}_i^2$$

$$SS_D = SS_{\tilde{R}} - SS_R \geq 0$$

ipotesi : $\beta = (\beta_0, \beta_1, \dots, \beta_k, 0, \dots, 0)$

$$S_e^2 = \frac{SS_R}{n-p-1} \approx \sigma^2 \quad \text{stimatore corretto}$$

$$\tilde{S}_e^2 = \frac{SS_{\tilde{R}}}{n-k-1} \approx \sigma^2 \quad \text{stimatore corretto}$$

$$SS_D = (n-k-1)\tilde{S}_e^2 - (n-p-1)S_e^2 \approx (p-k)\sigma^2$$

$$S_D^2 := \frac{SS_D}{p-k} \approx \sigma^2 \quad \text{stimatore corretto}$$

$$H_0: \beta_{k+1} = \beta_{k+2} = \dots = \beta_p = 0$$

H_1 : non tutti nulli

$$F := \frac{S_D^2}{S_e^2} \quad \text{statistica del test}$$

→ in ipotesi H_0 fa circa 1

→ in ipotesi H_1 viene più grande

★ È sempre possibile usare questo test per confrontare due modelli di regressione con set di variabili contenute uno nell'altro.

ord 18

★ Gli strumenti di calcolo che fanno la stepwise backward e forward usano tipicamente test F.

Ad esempio: forward, passo $k+1$

ho già inserito x_i $i \in I_k = \{1, 2, \dots, n\}$ $\#I_k = k$

$\forall j \in \{1, \dots, n\} \setminus I_k$ confronto la regressione con variabili I_k con quella con var $I_k \cup \{j\}$ e ottengo un valore di F_j per il test (e un valore di α_j^*)

Scego j che massimizza F_j (o minimizza α_j^*)
 $I_{k+1} := I_k \cup \{j\}$

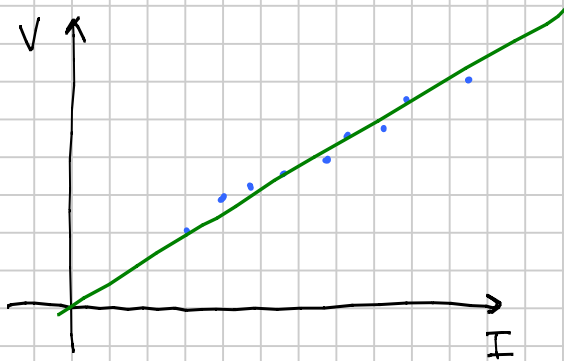
Mi fermo quando $\max_j F_j \leq \text{soglia}$ o $\min_j \alpha_j^* \geq \text{soglia}$

★ Stepwise generale: parte come forward, ma ad ogni step verifica anche se ci sono variabili da togliere (backward step)

★ Regola gerarchica: se nel modello includo un termine, devo includere anche tutti quelli "che lo dividono"

→ in particolare tengo il termine noto sempre

→ eccezione: se ho motivi teorici astratti per supporre che alcuni coefficienti siano nulli



legge di Ohm : $V = I \cdot R$

→ Esempio : $p = 3$

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

dummy : x_1^2 , $x_1 x_2$, $x_1^2 x_2$, x_2^2 , $x_1^2 x_2^2$

(Note: A red arrow points from x_1^2 to $x_1 x_2$. A red arrow points from $x_1 x_2$ to $x_1^2 x_2$. A blue bracket underlines $x_1^2 x_2$, x_2^2 , and $x_1^2 x_2^2$ with the label "non significativi".)

● Metodi globali di selezione delle variabili

Si definisce uno "score" globale di regressione e poi si cerca il modello tra tutti quelli possibili che lo rende massimo

→ Tutti i modelli sono tanti : 2^p come minimo

($2^{2^p + \binom{p}{2}}$ se includo un po' di dummy)

★ Esempio di "score" che non funziona : minimizzare $|R|$

→ con tutte le variabili il valore di $|R|$ è sempre minimo

★ Uno "score" che funziona è Se (minimizzato)

★ Un altro citato spesso si chiama AIC

● Coefficiente di determinazione (corretto o no)

$$(R^2) \quad R_d^2 := 1 - \frac{SSR}{SS_Y} \quad SS_Y = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

↑ = SS_Y con $k=0$

$0 \leq R_d^2 \leq 1$ è grande se l'ordine di grandezza dei residui è molto minore degli scarti di Y_i da \bar{Y} (stima senza variabili x_j di regressione)

★ massimizzare $R_d^2 \Leftrightarrow$ minimizzare $|R|$

non si può usare per la selezione delle variabili

$$\star \quad R_a^2 := 1 - \frac{Se^2}{S_Y^2} \quad S_Y^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

varianza campionaria

↑ = Se^2 con $k=0$

HW: $R_a^2 \leq R_d^2$

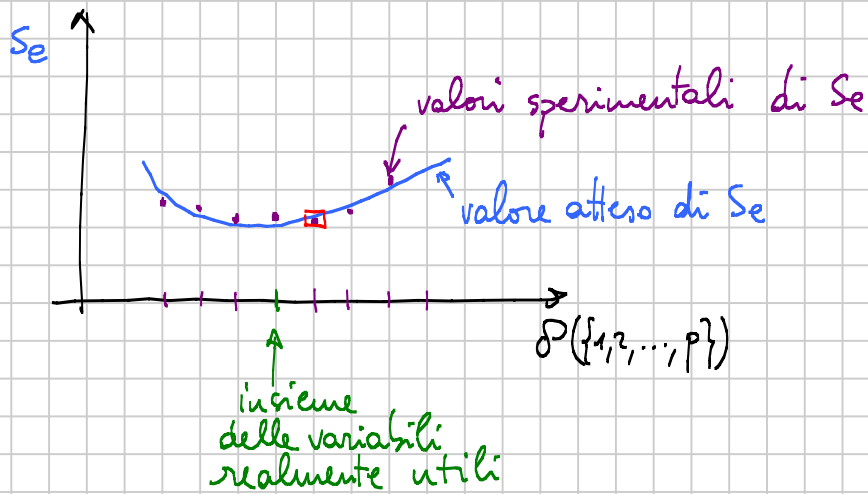
→ misura la qualità della regressione

★ massimizzare $R_a^2 \Leftrightarrow$ minimizzare Se

si può usare per la selezione delle variabili

!!! Sono canali!!!

● Selezione variabili con metodi globali



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

★ Non è opportuno cercare il minimo di una funzione usando una sola stima casuale

★ Molte volte tutti i metodi concordano.

■ GESTIONE DELLE VARIABILI

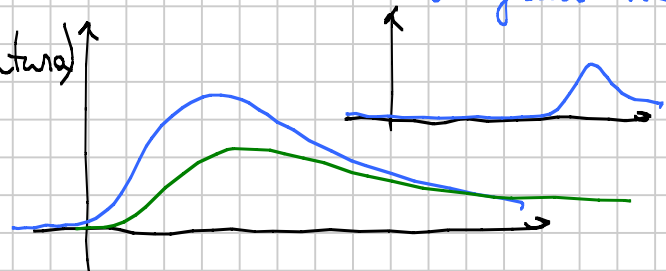
- categoriche nominali (regione, sesso)

- categoriche ordinali (titolo di studio)

↳ niente
↳ licenze media
↳ diploma
↳ laurea

- numeriche rapporto (reddito, popolazione città) diversi ordini di grandezza

- numeriche differenza (temperatura)



● Come si trattano le categoriche nominali

→ se le categorie sono solo 2, si codifica con 0/1 o -1/+1

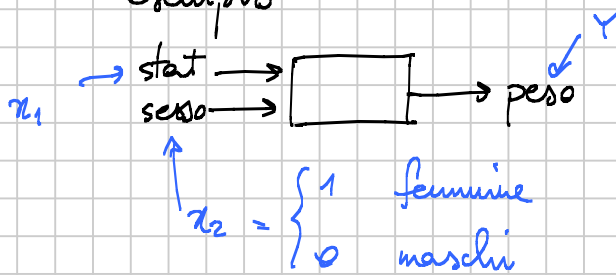
↳ diventa una variabile numerica

↳ i numeri scelti per la codifica non hanno nessuna importanza

* Consiglio comunque 0/1 usando lo 0 per la categoria più comune (di default) oppure per il "no" se le categorie sono sì/no

↳ Interpretazione risultati

esempio



$$Y = B_0 + B_1 x_1 + B_2 x_2$$

$$B_2 = -2,41$$



i maschi pesano in media e a parità di altezza 2,41 kg in più.

→ se le categorie sono $k > 2$

* Errore classico: codificare con i numeri da 1 a k

regione → $x_1 = \begin{cases} 1 & \text{valle d'Aosta} \\ 2 \\ \vdots \\ 20 & \text{Sardegna} \end{cases}$

$$Y = B_0 + B_1 x_1$$



l'ordine è del tutto arbitrario

★ Esploso in dicotomiche

↳ selgo una categoria di default

(a volte non è facile e occorre fare diverse prove)

↳ le altre $k-1$ diventano altrettante variabili 0/1

zone: nord/centro/sud

centro: default

id	zone	nord	sud
1	centro	0	0
2	sud	0	1
⋮	centro	0	0
⋮	nord	1	0
⋮	sud	0	1
⋮	⋮	⋮	⋮
n	centro	0	0

default → $\overbrace{0, 0, 0, \dots, 0}^{k-1}$

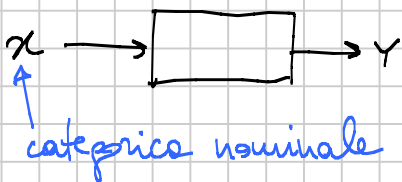
j -esima categ. → $\overbrace{0, 0, \dots, 1, 0, \dots, 0}^{k-1}$
 ↑
 posiz j

★ I coefficienti della regressione corrispondenti esprimono la differenza media attesa nella var di risposta tra la categoria in esame e quella di default

★ La scelta della codifica non è arbitraria: cambiare il default o cambiare altro (tipo imporre somma 0 su ogni riga: $-\epsilon, -\epsilon, \dots, 1, \epsilon, \epsilon, \dots, \epsilon$) fa cambiare i risultati,

★ Selezione variabili → si collapsano categorie assieme a quelle di default semplicemente togliendo variabili dicotomiche non significative.

Eventualmente cambiare il default per collapsare altri gruppi.



si studia (meglio) anche con la ANOVA ad una via.

* sono abbastanza equivalenti

● Categorie ordinali

1) come le nominali, ignorando l'ordine

2) codificare rispettando l'ordine

la scala non importa, ma i rapporti tra le differenze

si: la regressione può venire diversa

tit di studio	C1	C2	C3	C4	C5
niente	0	1	1	0	5
medie	1	2	3	1	8
diploma	2	3	5	3	13
laurea	3	4	7	4	17

equivalenti (under C1-C3)
 diversa (under C4-C5)

● Variabili numeriche

→ attenzione che non sia una categoria nascosta (tipo di pompa: 0, 1, 2)

Trasformazioni lineari di singole variabili

$$x_j \mapsto a_j + b_j x_j = x'_j \quad j=1, 2, \dots, p$$

$$Y \mapsto a_0 + b_0 Y = Y'$$

$$Y = B_0 + \sum_{j=1}^p B_j x_j + e$$

$$\begin{aligned}
 Y' &= a_0 + b_0 Y = a_0 + b_0 B_0 + \sum_{j=1}^p b_0 B_j \frac{x_j - a_j}{b_j} + b_0 e \\
 &= \underbrace{a_0 + b_0 B_0 - b_0 \sum_{j=1}^p \frac{B_j a_j}{b_j}}_{B_0'} + \sum_{j=1}^p \underbrace{\frac{b_0 B_j}{b_j}}_{B_j'} x_j + \underbrace{b_0 e}_{e'} = B_0' + \sum_{j=1}^p B_j' x_j + e' \\
 & \qquad \qquad \qquad e' \sim \mathcal{N}(0, b_0^2 \sigma^2)
 \end{aligned}$$

$B_j = 0 \Leftrightarrow B_j' = 0$ la relet delle var non cambia

HW: R_d^2 non cambia

★ trasformazioni lineari non hanno essenzialmente alcun effetto (Equivalenze viste sopra $d1 - -1/+1$ e $c_1 \sim c_2 - c_3$)

★ Standardizzazione (o normalizzazione) delle variabili

$$x_j : x_{1j}, x_{2j}, \dots, x_{nj} \qquad \bar{x}_j := \bar{x}_{*j} := \frac{1}{n} \sum_{i=1}^n x_{ij}$$

$$\text{devianze} \rightsquigarrow SS_j := \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2$$

$$\text{varianza} \rightsquigarrow S_j^2 := \frac{1}{n-1} SS_j$$

$$x_{ij} \mapsto x'_{ij} := \frac{x_{ij} - \bar{x}_j}{S_j} \qquad x'_{*j} = 0 \quad S_j' = 1$$

→ ottengo coefficienti di regressione confrontabili

$B_j = 2,4$ $B_k = 6,1$ significa che x_k ha impatto maggiore su Y di x_j

↳ in realtà per questo basta $\frac{x_{ij}}{S_j}$.

★ Forse da evitare sulle dicotomiche e sulle variabili rapporto.

Trasformazioni nonlineari

→ se ne applicate una, cambia tutto

→ i modelli possono essere confrontati solo qualitativamente

↳ grafici dei residui

↳ numero di variabili selezionate alla fine

→ se si trasformano (alcune delle) x_j ma non Y ,
i valori di R_d^2 , R_a^2 ed S_e sono confrontabili

* trasformazioni tipiche: $\log(x)$, \sqrt{x} , $\frac{1}{x}$, x^a

power transformation

→ In natura oltre a \mathcal{N} (contributi indipendenti additivi)
si trovano spesso anche leggi **lognormali** (contributi
indipendenti moltiplicativi): sono le classiche var. rapporto
→ conviene trasformarle

Def X ha legge lognormale di parametri μ, σ
se $\log X \sim \mathcal{N}(\mu, \sigma^2)$

$$\text{HW: } E(X) = e^{\mu + \frac{\sigma^2}{2}}$$

→ Anche in tutti i casi in cui un aspetto che gli errori
e siano proporzionali a Y , ha senso fare il $\log(Y)$

$$Y' = \log(Y)$$

$$Y' = \beta_0 + \beta_1 x_1 + \varepsilon$$

$$\varepsilon \sim \mathcal{N}(0, \sigma^2)$$

omoschedastico

$$Y = \exp Y' = e^{\beta_0} \cdot e^{\beta_1 x_1} \cdot \underbrace{e^{\varepsilon}}$$

ricome è moltiplicata
l'errore è proporzionale a Y

16.30 - 18.30 lab inf base 4 sede didattica ing.

HW dall'ora 7

$X_1, X_2, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$ σ nota μ_0 valore target \bar{x} val sign.

$H_0: \mu = \mu_0$

$H_1: \mu \neq \mu_0$

bilaterale

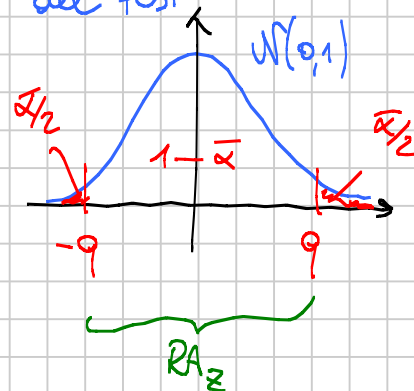
→ test con funzione ancillare : tre livelli / modi per farli

① $Z := \frac{\bar{X} - \mu_0}{\sigma / \sqrt{n}} \stackrel{H_0}{\sim} \mathcal{N}(0, 1)$

statistica del test

$RA_Z: [-q; +q]$

$q = \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)$



③ $\alpha^* = 2 - 2\Phi(|Z|)$

$RA_{\alpha^*}: [\bar{x}; 1]$

devo verificare che $Z \in RA_Z \Leftrightarrow \alpha^* \in RA_{\alpha^*}$

$Z \in RA_Z \Leftrightarrow |Z| \leq q \Leftrightarrow \Phi(|Z|) \leq \Phi(q) = 1 - \frac{\alpha}{2}$

$\Leftrightarrow \frac{\bar{x}}{2} \leq 1 - \Phi(|Z|)$

↑ Φ monotona

$\Leftrightarrow \bar{x} \leq 2 - 2\Phi(|Z|) =: \alpha^* \Leftrightarrow \alpha^* \in RA_{\alpha^*}$

* In generale si possono sempre fare i paraggi da $S \in RA_S$ generico per ricavare \bar{x} : la statistica confrontata con \bar{x} è il p-dei-dati

HW dall'ora 11

$X \sim \text{bin}(n, p)$ p_0 valore AQL $\bar{\alpha}$ lvl di sign.

$H_0: p \leq p_0$ $H_1: p > p_0$ test unilaterale no funz. ausiliare

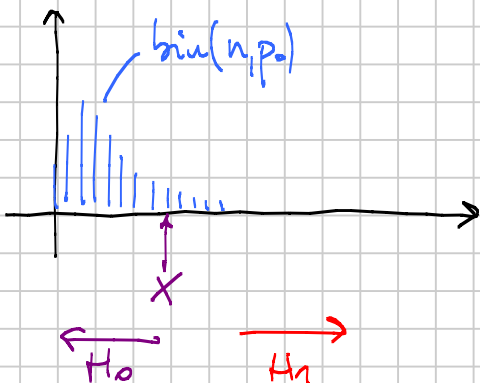
① Scelgo il più piccolo b tale che $P(\text{bin}(n, p_0) \geq b) < \bar{\alpha}$
 X statistica $RA_X = \{0, 1, \dots, b-1\}$ ↑ def di b

② Calcolo il p-dei-dati

$$X \leq b-1 \Leftrightarrow F_{\text{bin}(n, p_0)}(X-1) \leq F_{\text{bin}(n, p_0)}(b-2) \leq 1 - \bar{\alpha}$$

$$P(\text{bin}(n, p_0) \geq b-1) \geq \bar{\alpha} \Leftrightarrow F_{\text{bin}(n, p_0)}(b-2) \leq 1 - \bar{\alpha}$$

$$\bar{\alpha} \leq 1 - F_{\text{bin}(n, p_0)}(X-1) =: \alpha^*$$



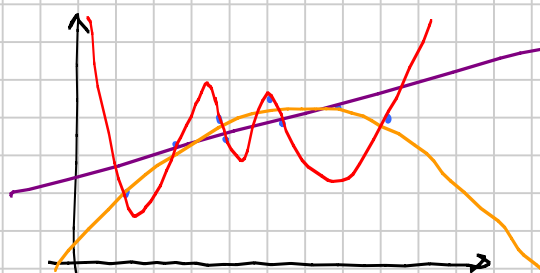
Se purtroppo vera H_0 , valori di X grandi sono "strani".
 α^* è la prob (sotto H_0) di osservare valori strani come X o più strani

$$\alpha^* = P(\text{bin}(n, p_0) \geq X)$$

OVERFITTING NELLA REGRESSIONE

Esempio: campione bivariato, fit polinomiale

id	x	x^2	$x^3 \dots$	Y
1				
⋮				
n				



$n=7$ se arrivo a x^6 trovo una curva che passa per tutti i punti

Overfitting: quando il modello interpola bene i punti del campione ma male quelli futuri

Diagnosi possibile: **cross-validation**

divido il campione in due sottocampioni A, B

fisso p (grado massimo)

regressione su A, calcolo i residui dai punti di B

regressione su B, calcolo i residui dai punti di A

sommo i quadrati di tutti i residui $\rightarrow SS_T$

$SS_T = SS_T(p)$ è uno **score**

compro p e rifaccio

alla fine scelgo p che massimizza lo score
regressione su tutto

* SS_T è comunque una v.a. quindi va massimizzata con buon senso, tenendo conto che non sappiamo la sua variabilità e che comunque a parità di tutto un modello con meno variabili è preferibile

ora 22

• Più in generale c'è overfitting nella regr. multiple tutte le volte che p è dell'ordine di grandezza di n

\rightarrow se $p+1 = n$ si dice che il modello è **saturo**

(residui nulli, impossibile fare inferenza: $n-p-1 = 0$)

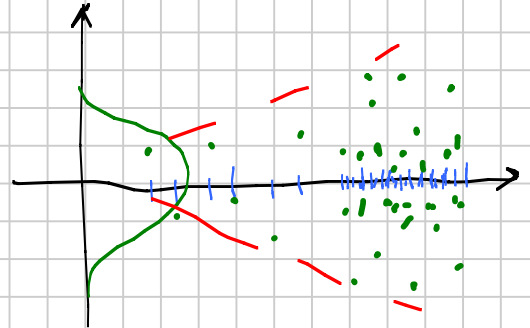
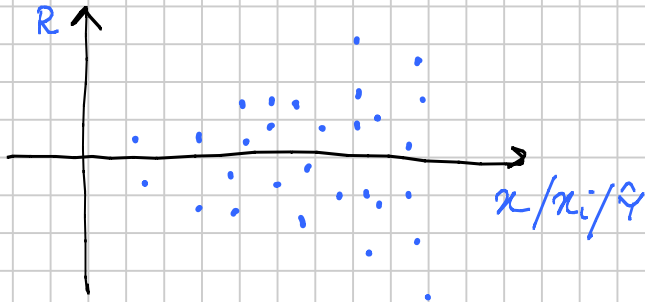
\rightarrow rule of thumb: $\frac{n}{p} \geq 5$ dovrebbe garantire che non ci sia overfitting

\rightarrow in caso di overfitting l'inferenza è inaffidabile
per questo la stepwise forward andrebbe preferita a quella backward.

REGRESSIONE PESATA

generalizzazione per quando il modello omoschedastico non è ritenuto adeguato

Diagnostici: analisi dei residui



- regr. lin. semplice $\rightarrow x$ in ascissa
- regr. lin. multiple \rightarrow vanno fatti molti grafici

$$x_1, x_2, \dots, x_p, \hat{Y}$$

\uparrow
previsi

- si cerca di "vedere" un andamento e proporre un modello:

$$\sigma \propto \sqrt{x}$$

$$\sigma \propto \sqrt{x_i} \text{ per un } i \text{ particolare}$$

$$\sigma \propto \hat{Y}$$

questi casi hanno a volte una motivazione teorica

x distanza da percorrere in città

Y tempo di percorrenza in auto con traffico

$$n: \# \text{ semafori} \propto x \quad Y = T_1 + T_2 + \dots + T_n \quad \text{Var}(Y) \propto n \propto x$$

$$\sigma \propto \sqrt{x}$$

nei casi di Y "log-normali" o "misti" può essere ragionevole ~~$\sigma \propto Y$~~ $\sigma \propto \hat{Y}$

★ Attenzione alle diverse densità dei punti in ascisse che può ingannare

★ In genere i risultati non dipendono molto dal modello scelto.

● Modello eteroschedastico

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + e \quad e \sim \mathcal{N}(0, \sigma^2(x_1, x_2, \dots, x_p))$$

σ ancora incognita, ma ipotizzo il tipo di dipendenza

$$v = (x_1, \dots, x_p) \quad \text{modello noto}$$

$\sigma \propto v$ ipotesi

→ definisco dei pesi : $i = 1, 2, \dots, n$

$$w_i := \frac{1}{v_i^2} = (v(x_{i,1}, \dots, x_{i,p}))^{-2} \quad \text{peso del punto } i$$

→ somma quadrati residui pesati

$$SS_w := \sum_{i=1}^n (Y_i - \sum_{j=0}^p \beta_j x_{ij})^2 w_i \quad \text{da minimizzare per trovare } \beta_j$$

→ alternativa equivalente più pratica

$$Y'_i = Y_i \sqrt{w_i} \quad x'_{i,j} = x_{ij} \sqrt{w_i} \quad i = 1, \dots, n \quad j = 0, 1, \dots, p$$

$$Y'_i = \beta_0 \sqrt{w_i} + \beta_1 x'_{i1} + \dots + \beta_p x'_{ip} + e'_i \quad e'_i \sim \mathcal{N}(0, \sigma^2)$$

$$Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \underbrace{e'_i / \sqrt{w_i}}_{=: e_i} \quad e_i \sim \mathcal{N}(0, \frac{\sigma^2}{w_i})$$

★ Studiare il modello "privato" come regressione omoschedastica
è equivalente a studiare il modello originale con quelle pesate

ANALISI DELLA VARIANZA (cap 10 Ross)

x_i numeriche \rightarrow $Y \sim \mathcal{N}(\mu; \sigma^2)$ *regressione*

1 var. categorica \rightarrow $Y \sim \mathcal{N}(\mu; \sigma^2)$ *ANOVA a 1 via*

2 var. categorica \rightarrow $Y \sim \mathcal{N}(\mu; \sigma^2)$ *ANOVA a 2 vie*

ANOVA a 1 VIA

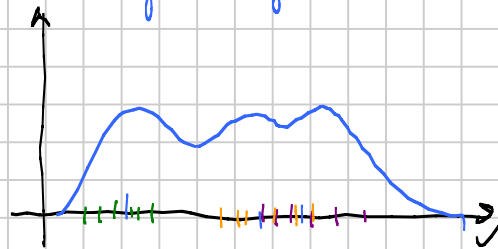
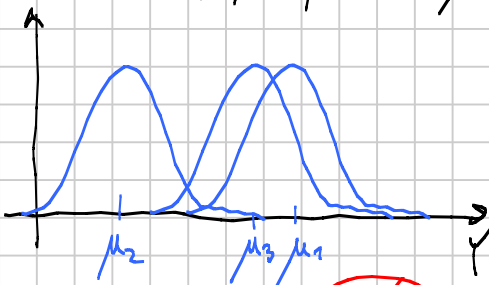
x categorica $x \in \{1, 2, \dots, m\}$ *codifica come etichette*

m il numero di categorie

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

$$x = 1, 2, \dots, m \quad \mu(x) = \mu_x \in \mathbb{R}$$

ogni categoria ha la sua media



Campione

~~$(x_i; y_i)$~~

$$Y_{1,1}, \dots, Y_{1,n_1} \sim \mathcal{N}(\mu_1; \sigma^2) \quad \text{iid}$$

$$Y_{2,1}, \dots, Y_{2,n_2} \sim \mathcal{N}(\mu_2; \sigma^2) \quad \text{iid}$$

...

$$Y_{m,1}, \dots, Y_{m,n_m} \sim \mathcal{N}(\mu_m; \sigma^2) \quad \text{iid}$$

* tipicamente Y nell'intero campione ha distribuzione multimodale

id	x_1	x_2	x_3	...	x_k
1					
...					
n					

categorica x

● Test fondamentale

H_0 : Y non dipende da x

H_1 : Y dipende da x

H_0 : $\mu_1 = \mu_2 = \dots = \mu_m$

H_1 : non tutte uguali

● Diverse "varianze campionarie"

1) Varianza *within* (aka "entro i campioni")

Si stima σ^2 dentro ogni campione e poi si fa la media

→ Not. standard

$$i = 1, 2, \dots, m \quad Y_{i,*} = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij} \quad \text{media campionaria}$$

$$S_i^2 = \frac{1}{n_i - 1} \sum_{j=1}^{n_i} (Y_{ij} - Y_{i,*})^2 \quad \text{var. campionaria}$$

$$S_W^2 = \sum_{i=1}^m \pi_i S_i^2 \quad \pi_i = \frac{n_i - 1}{N - m} \quad N = \sum_{i=1}^m n_i$$

HW: Per ogni scelta di π_i : $\sum \pi_i = 1$ la media pesata \bar{y} è uno stimatore corretto di σ^2 . La scelta

$\pi_i \propto$ gdl della categoria i

è quella che minimizza la varianza

HW: *Thm Cochran* \Rightarrow distribuzione di S_W^2

$$\frac{S_W^2}{\sigma^2} (N - m) \sim \chi^2(N - m)$$

S_W^2 indipendente da $Y_{1,*}, Y_{2,*}, \dots, Y_{m,*}$

2) Varianza between (aka "tra i campioni")

$\mu_i \approx Y_{i*}$ $\frac{1}{m-1} \sum (Y_{i*} - Y_{**})^2$ non proprio ma quasi

~~$Y_{**} = \frac{1}{m} \sum_{i=1}^m Y_{i*}$~~

$Y_{**} = \frac{1}{N} \sum_{i,j} Y_{ij} = \sum_{i=1}^m \frac{n_i}{N} Y_{i*}$

(check)

$Y_{i*} \sim \mathcal{N}(\mu_i; \frac{\sigma^2}{n_i})$ $(Y_{i*} - \mu_i)^2 \approx \text{Var}(Y_{i*}) = \frac{\sigma^2}{n_i}$

$S_B^2 = \sum_{i=1}^m \frac{n_i}{m-1} (Y_{i*} - Y_{**})^2$

★ Se H_0 è vera

S_B^2 è uno stimatore corretto di σ^2 indipendente da S_W^2

★ Se H_0 è falsa

S_B^2 è tendenzialmente più grande (e comunque indep. da S_W^2)

→ Verifichiamo la prima

$H_0: \mu_i \equiv \mu \quad i=1,2,\dots,m$ $Y_{ij} \sim \mathcal{N}(\mu, \sigma^2)$

$Y_{i*} = \frac{1}{n_i} \sum_j Y_{ij} \sim \mathcal{N}(\mu, \frac{\sigma^2}{n_i})$

$Y_{**} = \frac{1}{N} \sum_{i,j} Y_{ij} \sim \mathcal{N}(\mu, \frac{\sigma^2}{N})$

$X_i = \sqrt{n_i} Y_{i*} \sim \mathcal{N}(\mu\sqrt{n_i}; \sigma^2)$

$X = (X_1, \dots, X_m) \sim \mathcal{N}(v; \sigma^2 I)$

$v_i = \mu\sqrt{n_i}$

$v \in \text{Span}(\sqrt{n_1}, \sqrt{n_2}, \dots, \sqrt{n_m})$

SSV di \mathbb{R}^m di dim 1

$\sum_{i=1}^m (X_i - \mu\sqrt{n_i})^2 \rightsquigarrow \text{minimizzato} \quad \hat{\mu} = \frac{\sum \sqrt{n_i} X_i}{\sum n_i} = \frac{\sum n_i Y_{i*}}{\sum n_i} =: Y_{**}$

$0 = \sum_i \sqrt{n_i} (X_i - \hat{\mu}\sqrt{n_i})$

$$\frac{1}{\sigma^2} \sum_{i=1}^m (x_i - Y_{x,\mu} \sqrt{n_i})^2 \sim \chi^2(m-1) \quad \text{indipendente da } Y_{x,\mu}$$

$$\frac{1}{\sigma^2} \sum_{i=1}^m n_i (Y_{i,x} - Y_{x,\mu})^2 =: \frac{S_B^2}{\sigma^2} (m-1) \sim \chi^2(m-1)$$

quindi S_B^2 è
uno stimatore corretto
di σ^2

ANOVA A UNA VIA

- Se H_0 è falsa S_B^2 è tendenzialmente maggiore di σ^2

$$S_B^2 = \frac{1}{m-1} \sum_{i=1}^m n_i (Y_{i*} - Y_{***})^2 \quad SS_B = (m-1) S_B^2 \quad \text{devianza between}$$

$$SS_B = \sum_{i=1}^m n_i (Y_{i*} - \mu_i + \mu_i - \mu + \mu - Y_{***})^2 \quad \text{dove } \mu_i = \mu_{**} = \frac{1}{N} \sum_{j=1}^m \mu_i$$

$$= \sum_i n_i [(Y_{i*} - \mu_i)^2 + (\mu_i - \mu)^2 + (\mu - Y_{***})^2]$$

$$+ 2 \sum_i n_i (Y_{i*} - \mu_i)(\mu_i - \mu) + 2 \sum_i n_i (Y_{i*} - \mu_i)(\mu - Y_{***}) + 0$$

$$= \frac{1}{N} \sum_{i=1}^m n_i \mu_i$$

$$ESS_B = \sum_i n_i \underbrace{\text{Var}(Y_{i*})}_{\sigma^2} + \sum_i n_i (\mu_i - \mu)^2 + N \text{Var}(Y_{***}) + 0 - 2N \text{Var}(Y_{***})$$

$$(\mu - Y_{***}) \sum n_i (Y_{i*} - \mu_i) = (\mu - Y_{***}) (N Y_{***} - N \mu) = -N (Y_{***} - \mu)^2$$

$$Y_{i*} = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij} \quad Y_{ij} \sim \mathcal{N}(\mu_i, \sigma^2) \quad Y_{i*} \sim \mathcal{N}(\mu_i, \frac{\sigma^2}{n_i})$$

$$Y_{***} = \frac{1}{N} \sum_{i=1}^m n_i Y_{i*} \quad Y_{***} \sim \mathcal{N} \quad E(Y_{***}) = \frac{1}{N} \sum_{i=1}^m n_i \mu_i = \mu$$

$$\text{Var}(Y_{***}) = \sum_{i=1}^m \frac{n_i^2}{N^2} \frac{\sigma^2}{n_i} = \frac{\sigma^2}{N^2} \sum_i n_i = \frac{\sigma^2}{N}$$

$$ESS_B = m \sigma^2 + \sum_{i=1}^m n_i (\mu_i - \mu)^2 - \sigma^2$$

$$ES_B^2 = \frac{1}{m-1} ESS_B = \sigma^2 + \frac{1}{m-1} \sum_{i=1}^m n_i (\mu_i - \mu)^2 > \sigma^2 \quad \text{se è vera } H_1$$

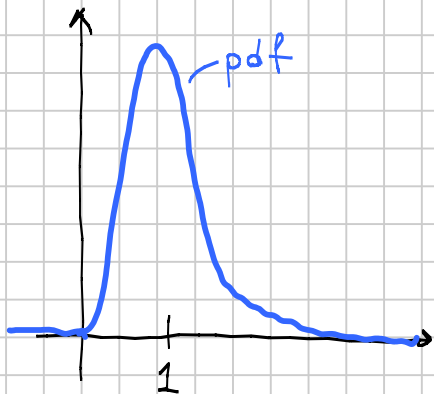
★ Prima di fare il test introduciamo la F di Fisher

● Distribuzione F di Fisher

Def: se $V \sim \chi^2(m)$ $W \sim \chi^2(n)$ indipendenti

$$\frac{V/m}{W/n} \sim F(m; n)$$

legge F di Fisher con m gdl al numeratore e n gdl al denominatore



★ Confronto delle varianze di due campioni normali (Cap 8)

$X_1, X_2, \dots, X_m \sim \mathcal{N}(\mu, \sigma^2)$
 $Y_1, Y_2, \dots, Y_n \sim \mathcal{N}(\nu, \tau^2)$ indipendenti

$H_0: \sigma = \tau$

$H_1: \sigma \neq \tau$

(altri esempi: $H_0: \sigma \geq \tau$, $H_0: \sigma \leq \tau$)

$\hat{\sigma}^2 \approx \frac{S_X^2}{m}$

$\frac{S_X^2}{\sigma^2} (m-1) \sim \chi^2(m-1)$

$\frac{S_Y^2}{\tau^2} (n-1) \sim \chi^2(n-1)$ indep

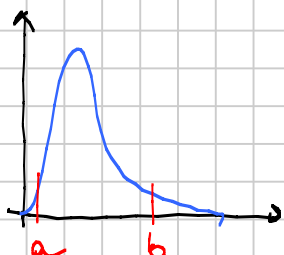
$\frac{\frac{S_X^2}{m}}{\frac{S_Y^2}{n}} \sim F(m-1; n-1)$

funz. ausiliare
 $\frac{S_X^2}{S_Y^2} \left(\frac{\sigma}{\tau}\right)^{-2} \sim F(m-1; n-1)$

$H_0: \hat{\sigma}^2 = 1$

statistica del test

$R := \frac{S_X^2}{S_Y^2} \stackrel{H_0}{\sim} F(m-1; n-1)$



α livello di sign. assegnato

$a = F_{F(m-1; n-1)}^{-1} \left(\frac{\alpha}{2} \right) = \text{INV.F} \left(1 - \frac{\alpha}{2}; m-1; n-1 \right)$

$b = F_{F(m-1; n-1)}^{-1} \left(1 - \frac{\alpha}{2} \right) = \text{INV.F} \left(\frac{\alpha}{2}; m-1; n-1 \right)$

(check)

$a^* = F^{-1} \left(\frac{\alpha^*}{2} \right)$

$b^* = F^{-1} \left(1 - \frac{\alpha^*}{2} \right)$

$R = a^* \circ R = b^*$

$F(R) = \frac{\alpha^*}{2} \circ F(R) = 1 - \frac{\alpha^*}{2}$

$$\alpha^* = 2F(R) \quad \text{o} \quad \alpha^* = 2 - 2F(R)$$

$$\alpha^* = 2 \min(F(R); 1 - F(R))$$

p dei dati

HW: curva OC

★ ANOVA a una via

$$R = \frac{S_B^2}{S_W^2} \quad \text{statistica del test}$$

- se è vera H_0 $\frac{S_B^2}{\sigma^2}(m-1) \sim \chi^2(m-1)$ $\frac{S_W^2}{\sigma^2}(N-m) \sim \chi^2(N-m)$

$$\Rightarrow R = \frac{\frac{S_B^2}{\sigma^2}}{\frac{S_W^2}{\sigma^2}} \sim F(m-1; N-m)$$

- se è vera H_1 R è tendenzialmente più grande

→ allora si imposta un test unilaterale



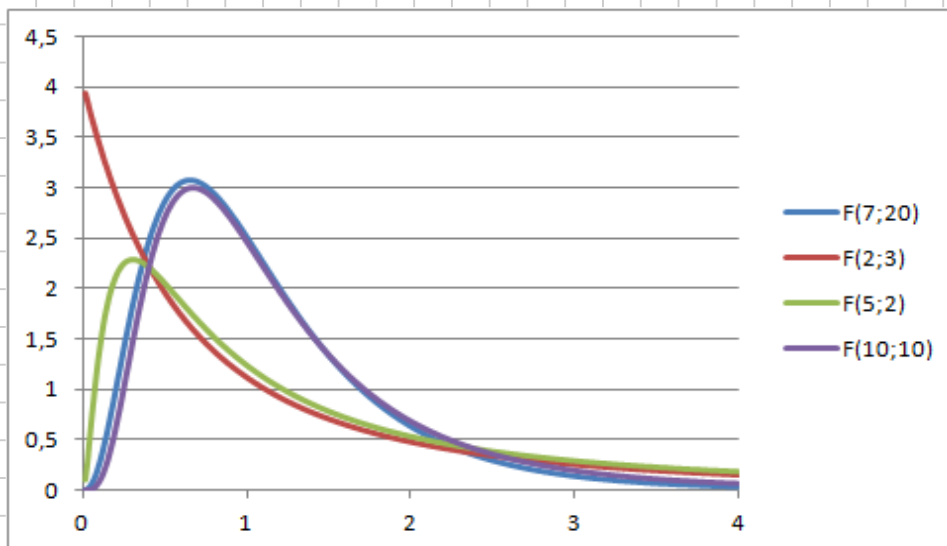
$$q = F_{F(m-1; N-m)}^{-1}(1 - \bar{\alpha}) = \text{INV.F}(\bar{\alpha}; m-1; N-m)$$

$$RA_R = [0; q]$$

$$\alpha^* = 1 - F(R)$$

(check)

ora 25



$$S_B^2 = \frac{1}{m-1} \sum_{i=1}^m n_i (Y_{i*} - Y_{***})^2$$

$$S_W^2 = \sum_{i=1}^m \frac{n_i - 1}{N - m} S_i^2$$

- Identità delle varianze (algebraica)

$$SS_Y = SS_W + SS_B$$

dove $S_Y^2 = \frac{1}{N-1} \sum_{i,j} (Y_{ij} - Y_{***})^2$

var. camp. complessiva

gdl_Y = N-1

dove $SS_i = gdl_i \cdot S_i^2$

* SS_Y e SS_W sono più veloci da ricavare di SS_B

Dim $SS_B + SS_W = \sum_{i=1}^m n_i (Y_{i*} - Y_{***})^2 + \sum_{i=1}^m (n_i - 1) S_i^2$

$$= \sum_{i,j} (Y_{i*} - Y_{***})^2 + \sum_{i,j} (Y_{ij} - Y_{i*})^2$$

$$\frac{1}{n_i - 1} \sum_{j=1}^{n_i} (Y_{ij} - Y_{i*})^2$$

$$\stackrel{?}{=} \sum_{i,j} (Y_{ij} - Y_{***})^2 = SS_Y$$

B.D.C. $\sum_{i,j} (Y_{i*} - Y_{***})(Y_{ij} - Y_{i*}) \stackrel{?}{=} 0$

$$\sum_{i,j} [Y_{i*} Y_{ij} - Y_{i*}^2 - Y_{***} Y_{ij} + Y_{***} Y_{i*}]$$

si semplificano sommando su j
si semplificano sommando su i

- Altra inferenza relativa a questo modello (in ipotesi H_1)

→ inferenze su μ_i

$$\mu_i \approx Y_{i*} \sim \mathcal{N}(\mu_i; \frac{\sigma^2}{n_i})$$

$$\frac{Y_{i*} - \mu_i}{\sigma / \sqrt{n_i}} \sim \mathcal{N}(0, 1)$$

$$\frac{Y_{i*} - \mu_i}{S_W / \sqrt{n_i}} \sim t(N - m)$$

→ int conf

→ test statistici

→ inferenza su $\mu_i - \mu_j$

$$\mu_i - \mu_j \approx Y_{i*} - Y_{j*} \sim \mathcal{N}(\mu_i - \mu_j; \frac{\sigma^2}{n_i} + \frac{\sigma^2}{n_j})$$

$$\frac{Y_{i*} - Y_{j*} - (\mu_i - \mu_j)}{S_w \cdot \sqrt{n_i^{-1} + n_j^{-1}}} \sim t(N-m)$$

→ int conf

→ test statistic

★ Se il test fondamentale dice H_0 (e se ci credo: dipende dalla potenza / numerosità) allora si fa inferenza più elementare:

$$Y_{i,j} \sim \mathcal{N}(\mu; \sigma^2) \quad \text{un campione normale}$$

$$\mu \approx Y_{***} = \bar{Y} \quad \frac{Y_{***} - \mu}{S / \sqrt{N}} \sim t(N-1)$$

$$\sigma^2 \approx S^2 \quad \frac{S^2}{\sigma^2} (N-1) \sim \chi^2(N-1)$$

● Altre considerazioni pratiche

1) Fare sempre se possibile l'analisi dei residui

→ va fatta a mano

→ controllare: outliers; omoschedasticità; normalità

$$R_{i,j} = Y_{i,j} - Y_{i*} \quad \text{e graficarli in qualche modo}$$

↑ miglior predittore per $Y_{i,j}$

2) In caso di eteroschedasticità: provare trasformazioni nonlineari

a) dati lognormali → $\log(Y_{i,j})$

b) dati conteggi ammissibili alla legge di Poisson → $\sqrt{Y_{i,j}}$

★ Se $X \sim \text{Pois}(\nu)$ $P(X=k) = \frac{\nu^k}{k!} e^{-\nu}$, $k=0,1,\dots$ $E(X)=\nu$, $\text{Var}(X)=\nu$

\sqrt{X} ha media che dipende da ν e varianza $\approx \frac{1}{4}$

$$Y = \sqrt{X} \quad X \sim \text{Pois}(\nu)$$

$$\text{Var}(Y) = E(Y^2) - E(Y)^2 = E(X) - E(Y)^2$$

$$E(Y) = E(\sqrt{X}) = \sum_{k=0}^{\infty} \sqrt{k} \frac{\nu^k}{k!} e^{-\nu}$$

$$E(Y)^2 = \sum_{i,j=0}^{\infty} \sqrt{i} \sqrt{j} \frac{\nu^{i+j}}{i!j!} e^{-2\nu} = \sum_{k=0}^{\infty} \sum_{i+j=k} \sqrt{i(k-i)} \frac{(2\nu)^k}{k!} \binom{k}{i} e^{-2\nu} \cdot 2^{-k}$$

$$= \sum_{k=0}^{\infty} e^{-2\nu} \frac{(2\nu)^k}{k!} \underbrace{\sum_{i=0}^k \binom{k}{i} i^{1/2} (k-i)^{1/2}}_{\frac{k}{2} + \frac{1}{4}} \cdot 2^{-k}$$

$$X \sim \text{Pois}(\nu) \quad Y = \sqrt{X}$$

$$(E(Y))^2 = \sum_{k=0}^{\infty} e^{-2\nu} \frac{(2\nu)^k}{k!} \sum_{i=0}^k \binom{k}{i} \cdot 2^{-k} i^{\frac{1}{2}} (k-i)^{\frac{1}{2}}$$

legge Poisson 2ν
legge bin $(k, \frac{1}{2})$

$$\sum_{i=0}^k \binom{k}{i} \cdot 2^{-k} i^{\frac{1}{2}} (k-i)^{\frac{1}{2}} = E\left[\sqrt{Z(k-Z)}\right] \quad \text{dove } Z \sim \text{bin}\left(k; \frac{1}{2}\right)$$

$$\approx \frac{k}{2} = \frac{k}{2} E\left[\sqrt{\frac{2}{k} Z \left(2 - \frac{2}{k} Z\right)}\right] \quad W := \frac{2}{k} Z - 1$$

$$= \frac{k}{2} E\left[\sqrt{(1+W)(1-W)}\right] = \frac{k}{2} E\left[\sqrt{1-W^2}\right]$$

$$k \gg 1 \quad Z \sim \mathcal{N}\left(\frac{k}{2}; \frac{k}{4}\right) \quad W \sim \mathcal{N}\left(0; \frac{1}{k}\right)$$

$$k \gg 1 \quad \sqrt{1-W^2} \approx 1 - \frac{1}{2} W^2 \Rightarrow$$

$$\sum_{i=0}^k \binom{k}{i} \cdot 2^{-k} i^{\frac{1}{2}} (k-i)^{\frac{1}{2}} \approx \frac{k}{2} E\left(1 - \frac{1}{2} W^2\right)$$

$$= \frac{k}{2} - \frac{k}{2} \cdot \frac{1}{2} \cdot \frac{1}{k} = \frac{k}{2} - \frac{1}{4}$$

HW: formalizzare con un limite

$$(E(Y))^2 = \sum_{k=0}^{\infty} e^{-2\nu} \frac{(2\nu)^k}{k!} \left(\frac{k}{2} - \frac{1}{4} + r(k)\right) = \frac{1}{2} \cdot 2\nu - \frac{1}{4} + \sum_{k=0}^{\infty} e^{-2\nu} \frac{(2\nu)^k}{k!} r(k)$$

↑
tende a zero

$$\text{Var}(Y) = \nu - E(Y)^2 = \frac{1}{4} - \sum_{k=0}^{\infty} e^{-2\nu} \frac{(2\nu)^k}{k!} r(k)$$

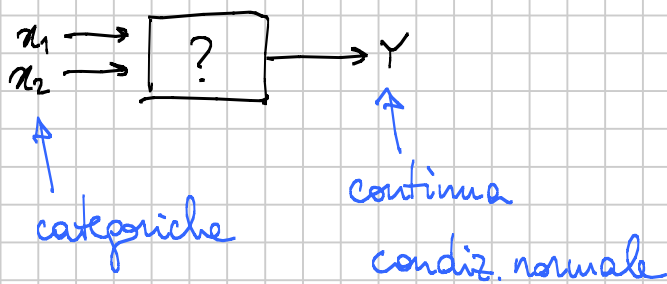
tende a 0 per $\nu \rightarrow \infty$

$$V_\nu \sim \text{Pois}(2\nu) \sim \mathcal{N}(2\nu; 2\nu)$$

$$\text{termine corr} = E(r(V_\nu)) \xrightarrow{\nu \rightarrow \infty} 0$$

HW: formalizzare

ANOVA A DUE VIE



x_1 m categorie $\{1, 2, \dots, m\}$
 x_2 n categorie $\{1, 2, \dots, n\}$

$m \times n$ combinazioni
 per ciascuna ho un certo numero l_{ij} di osservazioni

Struttura dei dati

id	x_1	x_2	Y
1			
⋮			
N			

$Y_{i,j,k}$ $i = 1, 2, \dots, m$
 $j = 1, 2, \dots, n$
 $k = 1, 2, \dots, l_{ij}$

★ Condizione di ortogonalità / buona decomposizione delle devianze

l_{ij} non dipende da $i, j \equiv l$

	x_1			
	1	2		m
1	7	4		
2		2	0	
n			5	

$l_{ij} \rightarrow$ devo buttare via dati e imporre
 $l := \min_{i,j} l_{ij}$

★ Se una delle caselle è vuota $l_{ij} = 0$ l'unica possibilità è buttare via la colonna o la riga.

★ Altri approcci generali (si trovano nei software) vanno considerati qualitativi.

$Y_{i,j,k} \sim \mathcal{N}(\mu_{ij}; \sigma^2)$ indipendenti

\rightarrow due casi : $l = 1$ (senza repliche) $l > 1$ (con repliche)

1WAY ANOVA \rightsquigarrow 2WAY ANOVA W/O R \rightsquigarrow 2WAY ANOVA W/R

SENZA REPLICHE ($l=1$)

• Ipotesi fondamentale di linearità

$$\mu_{ij} = \mu + \alpha_i + \beta_j$$

\uparrow \uparrow \uparrow
 media globale effetto riga effetto colonna

→ ipotesi tecnica (wlog) : $\sum_{i=1}^m \alpha_i = \sum_{j=1}^n \beta_j = 0$

HW: check μ, α_i, β_j qualsiasi → $\mu', \alpha'_i, \beta'_j := \dots$ t.c. $\mu'_{ij} = \mu'_{ij} + \text{ip. tec.}$

$$Y_{ij} = \mu + \alpha_i + \beta_j + e_{ij} \quad e_{ij} \sim \mathcal{N}(0, \sigma^2)$$

$\mu, \alpha_i, \beta_j, \sigma^2$ incognite

• Test principali

a) $H_0: \alpha_i \equiv 0$

Y non dipende da α_1

H_1 : non tutti zero

Y dipende da α_1

test per l'effetto riga

b) $H_0: \beta_j \equiv 0$

Y non dipende da α_2

H_1 : non tutti zero

Y dipende da α_2

test per l'effetto colonna

ord 27

• Stimatori

$\mu + \alpha_i \approx Y_{i*} \sim \mathcal{N}(\mu + \alpha_i; \frac{\sigma^2}{n})$ ✗

$\mu + \beta_j \approx Y_{*j} \sim \mathcal{N}(\mu + \beta_j; \frac{\sigma^2}{m})$ ✗

$\mu \approx Y_{**} \sim \mathcal{N}(\mu; \frac{\sigma^2}{mn})$ ✓

$\alpha_i \approx Y_{i*} - Y_{**} \sim \mathcal{N}(\alpha_i; \frac{n-1}{mn} \sigma^2)$ ✓

NON sono indipendenti

vedi prossima pagina

$$Y_{i*} - Y_{**} = Y_{i*} - \frac{1}{n} \sum_{a=1}^n Y_{a*} = Y_{i*} \left(1 - \frac{1}{n}\right) - \frac{1}{n} \sum_{a \neq i} Y_{a*}$$

↑
INDIPENDENTI

$$\left(\frac{n-1}{n}\right)^2 \frac{\sigma^2}{m} + \frac{1}{n^2} (n-1) \frac{\sigma^2}{m} = \frac{n-1}{mn} \sigma^2$$

$$\beta_j \approx Y_{*j} - Y_{**} \sim \mathcal{N}\left(\beta_j; \frac{n-1}{mn} \sigma^2\right) \quad \checkmark$$

* Per σ^2 , come nella regressione (e nell'ANOVA a 1 VIA) si fa una media quadratica dei residui

↳ previsti? $Y_{ij} \approx \mu + \alpha_i + \beta_j$ *miglior predittore*
 $\approx Y_{i*} + Y_{*j} - Y_{**}$ *HW: trovare la varianza*

$$R_{ij} = Y_{ij} - Y_{i*} - Y_{*j} + Y_{**}$$

$$R_{ij} = Y_{ij} \left(1 - \frac{1}{n} - \frac{1}{m} + \frac{1}{mn}\right) + \sum_{b \neq j} Y_{ib} \left(-\frac{1}{n} + \frac{1}{mn}\right) + \sum_{a \neq i} Y_{aj} \left(-\frac{1}{m} + \frac{1}{mn}\right) + \sum_{\substack{a \neq i \\ b \neq j}} Y_{ab} \cdot \frac{1}{mn}$$

$$\text{Var}(R_{ij}) = \frac{\sigma^2}{m^2 n^2} \left\{ (m-1)^2 (n-1)^2 + (m-1)^2 (n-1) + (m-1) (n-1)^2 + (m-1) (n-1) \right\}$$

$$= \frac{(m-1)(n-1)}{mn} \sigma^2$$

$$SS_e := \sum_{i,j} R_{ij}^2$$

$$E(SS_e) = mn \text{Var}(R_{ij}) = (m-1)(n-1) \sigma^2$$

$$S_e^2 := \frac{SS_e}{(m-1)(n-1)}$$

stimatore corretto di σ^2

HW: Cochran \Rightarrow

$$\frac{S_e^2}{\sigma^2} (m-1)(n-1) \sim \chi^2((m-1)(n-1))$$

indipendente da Y_{i*}, Y_{*j}

• Altri stimatori di σ^2

$$\star S_R^2 := \frac{1}{m-1} \sum_{i=1}^m n(Y_{i\cdot} - Y_{\cdot\cdot})^2$$

$$H_0 a) \Rightarrow S_R^2 \approx \sigma^2 \quad \frac{S_R^2}{\sigma^2} (m-1) \sim \chi^2(m-1)$$

$$H_1 a) \Rightarrow S_R^2 \geq \sigma^2$$

★ Analogo per le colonne

• Come si fanno i test

$$\star U_R := \frac{S_R^2}{S_e^2} \stackrel{H_0 a)}{\sim} F(m-1; (m-1)(n-1))$$

$H_1 a) \Rightarrow U_R$ tipicamente grande \rightarrow test unilaterale

$$\alpha_R^* = \text{DISTRIB. F}(U_R; m-1; (m-1)(n-1))$$

★ Analogo per le colonne

• Identità delle varianze (conseguenza ipotesi di ortogonalità)

$$SS_Y = SS_e + SS_c + SS_r$$

sono tutte indipendenti

$$\frac{SS_Y}{\sigma^2} \stackrel{H_0 b)}{\sim} \chi^2(mn-1)$$

$$\frac{SS_R}{\sigma^2} \stackrel{H_0 a)}{\sim} \chi^2(m-1)$$

$$\frac{SS_c}{\sigma^2} \stackrel{H_0 b)}{\sim} \chi^2(n-1)$$

$$\frac{SS_e}{\sigma^2} \sim \chi^2((m-1)(n-1))$$

$$mn-1 = (m-1) + (n-1) + (m-1)(n-1)$$

• Altre inferenze:

$$\star \alpha_i \approx Y_{i\cdot} - Y_{\cdot\cdot}$$

$$\frac{Y_{i\cdot} - Y_{\cdot\cdot} - \alpha_i}{S_e \sqrt{\frac{m-1}{mn}}} \sim t((m-1)(n-1))$$

★ $\mu, \beta_j, \mu + \alpha_i, \mu + \beta_j \dots$ analoghe

$$H_0: \alpha_1 = \alpha_2 \quad H_1: \alpha_1 \neq \alpha_2$$

$$\alpha_1 - \alpha_2 \approx \bar{Y}_{1\cdot} - \bar{Y}_{2\cdot} \sim \mathcal{N}(\alpha_1 - \alpha_2; ???) \rightarrow \frac{\bar{Y}_{1\cdot} - \bar{Y}_{2\cdot} - (\alpha_1 - \alpha_2)}{S_e \sqrt{???}} \sim t((m-1)(n-1))$$

→ e così via ...

★ Intervallo di predizione per osservazioni ulteriori in (i, j)

$$\tilde{Y}_{ij} \sim \mathcal{N}(\mu + \alpha_i + \beta_j; \sigma^2) \quad \text{indip da } Y_{a,b} \quad a=1, \dots, m, \quad b=1, \dots, n$$

$$\tilde{Y}_{ij} \approx \mu + \alpha_i + \beta_j \approx \bar{Y}_{i\cdot} + \bar{Y}_{\cdot j} - \bar{Y}_{\cdot\cdot}$$

$$\rightarrow \tilde{Y}_{ij} - \bar{Y}_{i\cdot} - \bar{Y}_{\cdot j} + \bar{Y}_{\cdot\cdot} \sim \mathcal{N}(0; ???)$$

\uparrow σ^2 $\underbrace{\qquad\qquad\qquad}_{\frac{m+n-1}{mn} \sigma^2}$ \uparrow $\frac{mn+m+n-1}{mn} \sigma^2$

$$\tilde{Y}_{ij} \in \bar{Y}_{i\cdot} + \bar{Y}_{\cdot j} - \bar{Y}_{\cdot\cdot} \pm q \cdot S_e \cdot \underbrace{\sqrt{\frac{mn+m+n-1}{mn}}}_{>1}$$

quantile $t((m-1)(n-1))$

ANOVA A DUE VIE CON REPLICHE

(aka "CON INTERAZIONI", "NONLINEARE")

$$Y_{ijk} \sim \mathcal{N}(\mu_{ij}; \sigma^2) \quad i=1,2,\dots,m, \quad j=1,2,\dots,n, \quad k=1,2,\dots,l$$

→ $l=1$ $\mu_{ij} = \mu + \alpha_i + \beta_j$ ipotesi di linearità

→ $l \geq 2$ μ_{ij} qualsiasi

$$\mu_{ij} = \mu + \alpha_i + \beta_j + \gamma_{ij}$$

↑ interazioni

$$0 = \sum_i \alpha_i = \sum_j \beta_j = \sum_i \gamma_{i,b} = \sum_j \gamma_{a,j} \quad \forall a, \forall b$$

HW: check che dati μ_{ij} , gli altri parametri sono univocamente determinati, grazie ai vincoli

$$Y_{ijk} \sim \mathcal{N}(\mu + \alpha_i + \beta_j + \gamma_{ij}; \sigma^2)$$

Test di linearità

$H_0: \gamma_{ij} = 0 \quad \forall i, j$ $H_1: \text{non tutti nulli}$

di nuovo si basa sul rapporto di certe "varianze" e sulla F di Fisher.

$$SS_E := \sum_{i,j,k} (Y_{ijk} - Y_{ij*})^2 \quad S_E^2 := \frac{SS_E}{(l-1)mn}$$

$$\frac{S_E^2 (l-1)mn}{\sigma^2} \sim \chi^2((l-1)mn)$$

$S_E^2 \approx \sigma^2$ sempre corretto

$$SS_R + SS_C + SS_{IN} + SS_E = SS_Y$$

$$(m-1) + (n-1) + (m-1)(n-1) + (l-1)mn = lmn - 1$$

identità delle devianze
gradi di libertà

→ Test linearità $\frac{S_{IN}^2}{S_E^2} \stackrel{H_0}{\sim} F((m-1)(n-1); (l-1)mn)$

→ Test effetto riga $\frac{S_R^2}{S_E^2} \stackrel{H_0}{\sim} F(m-1; (l-1)mn)$

→ Test effetto colonna $\frac{S_C^2}{S_E^2} \stackrel{H_0}{\sim} F(n-1; (l-1)mn)$

● Altra inferenza

$$\mu_{ij} \approx \bar{Y}_{ij*} \sim \mathcal{N}\left(\mu_{ij}, \frac{\sigma^2}{l}\right)$$

$$\frac{\bar{Y}_{ij*} - \mu_{ij}}{S_E / \sqrt{l}} \sim t((l-1)mn)$$

(esempio)

RELAZIONE TRA I TRE TIPI DI ANOVA

In genere si passa da ZWA con repliche a ZWA senza a 1w

$$Y_{ijk} \quad i=1, \dots, m, \quad j=1, \dots, n, \quad k=1, \dots, l$$

→ test di linearità

↳ se viene H_1 mi fermo

↳ se viene H_0 : $Y_{ij}^1 := \bar{Y}_{ij*} \sim \mathcal{N}\left(\mu + \alpha_i + \beta_j; \frac{\sigma^2}{l}\right)$

e rifaccio come ZWA senza repliche

* in questo modo posso verificare due volte se vi sia effetto riga o colonna e il secondo metodo in generale è più potente

↳ se effetto riga ed effetto colonna sono significativi mi fermo

↳ se almeno uno non lo è: $Y_i^1 := \bar{Y}_{i*} \sim \mathcal{N}\left(\mu + \alpha_i; \frac{\sigma^2}{n}\right)$

oppure $Y_j^1 := \bar{Y}_{*j} \sim \mathcal{N}\left(\mu + \beta_j; \frac{\sigma^2}{m}\right)$

e rifaccio come 1WA

* se nelle ZWA era significativa solo una variabile (say 'R') cancello l'altra e studio Y_i^1

NON per vedere se μ è dipendente da 'i' (lo so già)
 ma per rendere più accurate le stime e più precisi/potenti
 gli ALTRI tipi di inferenza

* se nella 2WA non era significativa nessuna delle
 due, ne cancello una (poi provo con l'altra)
 e provo 1WA (tipo stepwise backward nella
 regressione)

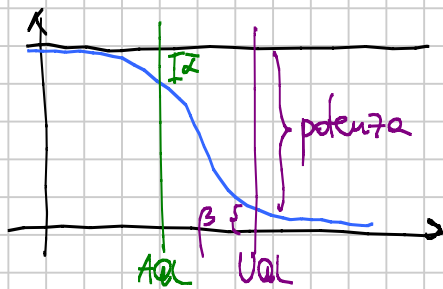
ora 29

☒ CORREZIONI PER TEST MULTIPLI

→ Se si fanno n test, ciascuno con livello di significatività α
 la probabilità che **almeno uno** dia H_1 , anche se sono
 vere tutte le ipotesi H_0 , è più alta di α

↳ esempio jelly beans - acne (xkcd/882)

↳ esempio Alzheimer - 80 geni



* Devo calare α ma questo
 diminuisce la potenza

↳ esempio ANOVA di tre tipi → 3 test effetto riga

↳ esempio: stepwise backward $\sim \frac{p^2}{2}$ test
 (anche nelle forward)

↳ esempio: $X_i \sim \mathcal{N}(\mu, \sigma^2)$ $Y_j \sim \mathcal{N}(\nu, \sigma^2)$
 $H_0: \mu = \nu$ $H_1: \mu \neq \nu$ (cap 8)

test $\frac{\bar{X} - \bar{Y}}{S_p \sqrt{\frac{1}{m} + \frac{1}{n}}} \stackrel{H_0}{\sim} t(m+n-2)$

→ $X_i \sim \mathcal{N}(\mu, \sigma^2)$ $Y_j \sim \mathcal{N}(\nu, \sigma^2)$ $Z_k \sim \mathcal{N}(\xi, \sigma^2)$

$$\mu = \nu \quad \text{poi} \quad \mu = \xi \quad (\text{poi} \quad \nu = \xi)$$

SBAGLIATO!!

↳ ANOVA 1 VIA è il test giusto

• Correzione di Bonferroni

Se devo fare n test e voglio un livello di significatività globale $\bar{\alpha}$ devo eseguire i singoli test con $\alpha' = \frac{\bar{\alpha}}{n}$

$i = 1, 2, \dots, n$ i diversi test $H_0^{(i)}$ $H_1^{(i)}$ e ipotesi
uso il p dei dati per confrontarli in modo semplice
 $\{U_i\}_{i=1, \dots, n}$ i diversi p dei dati

se $H_0^{(i)}$ è vera, $U_i \sim \text{unif}[0; 1]$

se $H_0^{(i)}$ è falsa, U_i è "piccolo"

sia t il lvl di significatività sui singoli test e α la probabilità di errore di I specie globale (almeno un test dà H_1 , anche se tutte le $H_0^{(i)}$ sono vere)

subadditività

$$\alpha = P(\exists i : U_i < t) = P(\min_i U_i < t) = P(\bigcup_i \{U_i < t\}) \leq \sum_i P(U_i < t) = nt$$

$$\alpha \leq nt$$

$$\text{Bonferroni: } t := \frac{\bar{\alpha}}{n} \Rightarrow \alpha \leq \bar{\alpha}$$

okay!

$H_0^{(i)}$ tutte vere

• Caso dei test indipendenti

indipendenza

$$1 - \alpha = 1 - P(\bigcup_i \{U_i < t\}) = P(\bigcap_i \{U_i \geq t\}) = \prod_i P(U_i \geq t) = (1 - t)^n$$

$$\text{Ipotesi di indipendenza: } t := 1 - (1 - \bar{\alpha})^{1/n} \Rightarrow \alpha = \bar{\alpha}$$

In genere non è ragionevole

$$\star \text{ Per } \bar{\alpha} \text{ piccolo } 1 - (1 - \bar{\alpha})^{1/n} \approx 1 - (1 - \frac{\bar{\alpha}}{n}) = \frac{\bar{\alpha}}{n}$$

● Correzione di Holm

Si ordinano gli U_i dal minore al maggiore

$$U_{(1)} \leq U_{(2)} \leq \dots \leq U_{(n)}$$

notazione standard statistiche di ordine

dico H_1 se $U_{(i)} < \frac{\alpha}{n-i+1}$ per qualche i

$$\alpha = P\left(\bigcup_i \left\{U_{(i)} < \frac{\alpha}{n-i+1}\right\}\right) = \dots \leq \dots = \alpha$$

HW: provate (domani lo faccio)

* Uniformemente migliore di Bonferroni (più potente)

$$B: H_1^{(B)} \text{ se } \min_i U_i < \frac{\alpha}{n} \Leftrightarrow U_{(1)} < \frac{\alpha}{n}$$

cfr ora 30

$$H: H_1^{(H)} \text{ se } U_{(1)} < \frac{\alpha}{n} \circ U_{(2)} < \frac{\alpha}{n-1} \circ \dots \circ U_{(n)} < \alpha$$

$$H_1^{(B)} \Rightarrow H_1^{(H)}$$

$$\alpha^{(B)} \leq \alpha^{(H)} \leq \alpha$$

▣ Prossimo argomento: Design of Experiment

→ pianificazione esperimenti

Sleeper (Cap 10)

Montgomery (tanti libri sul DoE)

→ esperimenti: soprattutto regressione, ma anche ANOVA

→ due "livelli" per ogni variabile (ipotesi di linearità)

→ si pianificano quali combinazioni fare (e quante, e come, e quando)

* è un elenco di "best practices" che escono dal mondo aziendale

● Correzione di Holm

Si ordinano gli U_i dal minore al maggiore

$$U_{(1)} \leq U_{(2)} \leq \dots \leq U_{(n)}$$

notazione standard statistiche di ordine

$$(1) \quad U_{(i)} \stackrel{H_0^{(i)}}{\leq} \frac{\bar{x}}{n-i+1}$$

★ Sia h il minimo i per cui dico $H_0^{(i)}$
 Il criterio di Holm prevede di dire

$$\underbrace{H_1^{(1)}, H_1^{(2)}, \dots, H_1^{(h-1)}, H_0^{(h)}}_{\text{secondo la disug. (1)}} , \underbrace{H_0^{(h+1)}, \dots, H_0^{(n)}}_{\text{forzati } H_0 \text{ indipendentemente da (1)}}$$

Dim (che $\alpha \leq \bar{\alpha}$)

Sia $I = \{1, 2, \dots, n\}$ ipotesi in cui è vera H_1

1) $\#I = 0 \quad I = \emptyset$ vera $H_0^{(i)}$ per ogni i

$$\alpha = P(\text{dire } H_1^{(i)} \text{ per qualche } i) = P(U_{(1)} < \frac{\bar{x}}{n}) \leq \bar{\alpha}$$

2) $\#I = 1 \quad I = \{j\}$ vera $H_0^{(i)}$ $i \neq j$ vera $H_1^{(j)}$

$$\alpha = P(\text{dire } H_1^{(i)} \text{ per qualche } i \neq j) \leq P(U_i \leq \frac{\bar{x}}{n-1} \text{ per qualche } i \neq j)$$

$$\alpha \quad (1) \neq j \quad f_p \Leftrightarrow \text{dire } H_1^{(2)} \Rightarrow \min_{i \neq j} U_i \leq \frac{\bar{x}}{n} < \frac{\bar{x}}{n-1}$$

$$\text{se } (1) = j \quad f_p \Leftrightarrow \text{dire } H_1^{(2)} \Leftrightarrow \min_{i \neq j} U_i \leq \frac{\bar{x}}{n-1}$$

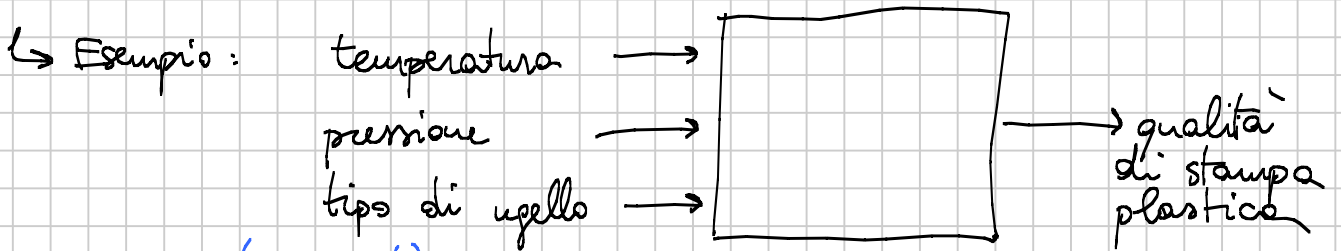
$$= P(\min_{i \neq j} U_i \leq \frac{\bar{x}}{n-1}) \leq \bar{\alpha}$$

Bonferroni

3) Analogamente (check) #I qualitativi

DESIGN OF EXPERIMENT

- pianificare esperimenti di regressione (decido X)
 - ↳ quanti esperimenti (\$)
 - ↳ quante variabili
 - ↳ quali variabili
 - ↳ quali esperimenti (imposto i valori delle x_i)
- nei casi più comuni tutte le var di ingresso vengono forzate ad assumere solo due valori ($T=180$, $T=200$)
(categoriche solo con 2 scelte: se ho più valori da testare ne scelgo solo 2 di rappresentative)
- fissati i valori delle x_i si fanno alcune prove (repliche), mai una sola e mai fatte consecutivamente
- idea base: provare tutte le combinazioni dei due livelli delle variabili di ingresso



(uncoded)
temp: 180 ~ 200

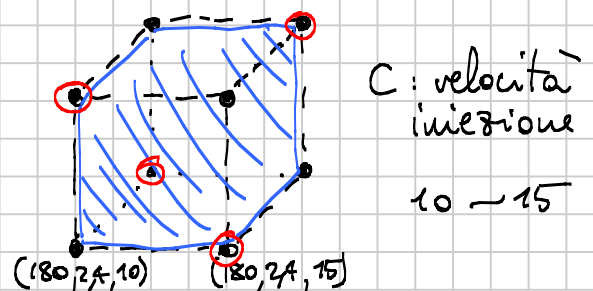
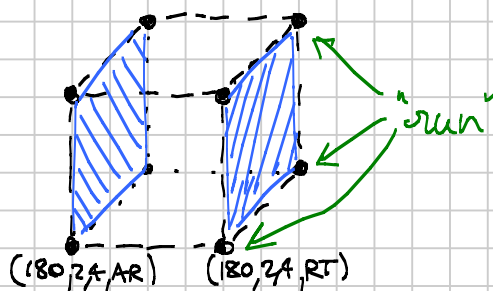
pressione: 2,4 ~ 2,8

ugello: AR, RT

fattore A: $\begin{matrix} 180 & 200 \\ -1 & +1 \end{matrix}$ (coded)

B: -1 +1

C: -1 +1 (no valori intermedi)

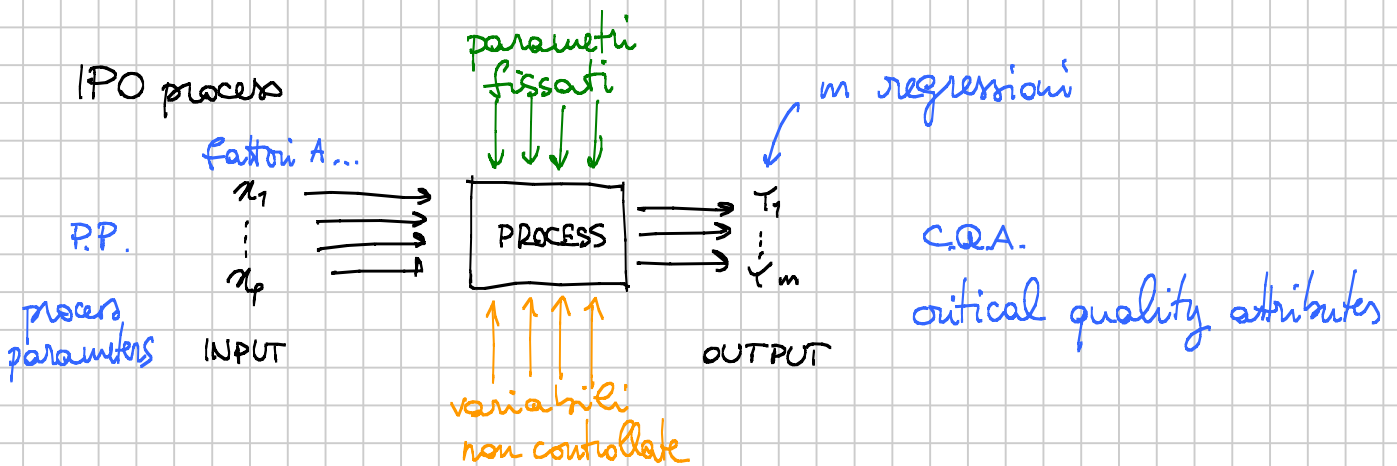


(due esempi di design space)

- Ogni configurazione che sperimento è una **run**.
Per ogni run faccio un po' (lo stesso numero!) di **repliche**
↳ es: 3 fattori $\rightarrow 8 = 2^3$ run \rightarrow (5 repliche) \rightarrow 40 esper.
- Siccome 2^p cresce troppo si introducono i design **frazionari** che riducono il numero delle run (di $\frac{1}{2}$, $\frac{1}{4}$, ... ecc)

ora 31

Scelta delle variabili



- CQA: devono essere numeriche, informative, precise (con errore gaussiano) NON dicotomiche o discrete

- PP: QUANTI? si distinguono

★ esperimenti di **screening** \rightarrow tanti fattori, poche run
 ↳ pesantemente frazionari
 ↳ meno informativi
 \rightarrow servono a selezionare le variabili più importanti

★ esperimenti di **modeling** (meno fattori, meno fraz, più inf)

\rightarrow servono per trovare un modello $Y = \sum_j \beta_j x_j$
 più preciso

● PP: livelli

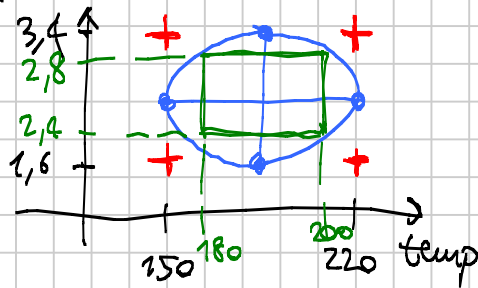
Criteri (ni usano con l'esperto / ingegnere per decidere)

a) linearità (se spazio è troppo grande fallisce)

b) SNR (signal-to-noise ratio)

da -1 a +1 → effetto in γ se è $\ll \sigma$ non lo vede
(se spazio è troppo piccolo fallisce)

c) funzionamento ai bordi (ost grande f)



d) utilità del design space (ost piccolo f)

■ Scelta del design

Vari tipi: Taguchi, Plackett-Burman, ...

↑
buoni solo per esperimenti di screening

Taguchi: L4, L8, L16, L32, ...

↑ ↑ ↑ ↑
di run

Design L8 di Taguchi

5-7 fattori	A	B	C	D	E	F	G
4 fattori	A	B	C				D
3 fattori	A	B	C	AB	AC	BC	ABC
1	-1	-1	-1	1	1	1	-1
2	1	-1	-1	-1	-1	1	1
3	-1	1	-1	-1	1	-1	1
4	1	1	-1	1	-1	-1	-1
5	-1	-1	1	1	-1	-1	1
6	1	-1	1	-1	1	-1	-1
7	-1	1	1	-1	-1	1	-1
8	1	1	1	1	1	1	1

si usano solo le colonne che servono

→ con la dummy $x_0 \equiv 1$ fattore I si ottiene la matrice X

Termini nonlineari e interazioni

L8 di Taguchi (ad esempio)

run	I	A	B	C	A²	AB	AC	BC	ABC
1	1	-1	-1	-1	1	1	1	1	-1
2	1	1	-1	-1		-1	-1	1	1
3	1	-1	1	-1		-1	1	-1	1
4	1	1	1	-1		1	-1	-1	-1
5	1	-1	-1	1		1	-1	-1	1
6	1	1	-1	1		-1	1	-1	-1
7	1	-1	1	1		-1	-1	1	-1
8	1	1	1	1	1	1	1	1	1

↑ tutte diverse e ortogonali

↑
 1 termine quadratico
 vanno bene solo con DoE
 a ≥ 3 livelli

5-7 fattori	A	B	C	D	E	F	G
4 fattori	A	B	C				D
3 fattori	A	B	C	AB	AC	BC	ABC
1	-1	-1	-1	1	1	1	-1
2	1	-1	-1	-1	-1	1	1
3	-1	1	-1	-1	1	-1	1
4	1	1	-1	1	-1	-1	-1
5	-1	-1	1	1	-1	-1	1
6	1	-1	1	-1	1	-1	-1
7	-1	1	1	-1	-1	1	-1
8	1	1	1	1	1	1	1

★ Se uso 3 fattori (full factorial) posso stimare i coefficienti di tutte le interazioni

★ Se uso 4 fattori (o più) alcune interazioni hanno la stessa colonna di un fattore

run	ABCD	BCD	ACD	ABD	ABC	AB CD	AC BD	BC AD
1	1	-1	-1	-1	-1	1	1	1
2	1	1	-1	-1	1	-1	-1	1
3	1	-1	1	-1	1	-1	1	-1
4	1	1	1	-1	-1	1	-1	-1
5	1	-1	-1	1	1	1	-1	-1
6	1	1	-1	1	-1	-1	1	-1
7	1	-1	1	1	-1	-1	-1	1
8	1	1	1	1	1	1	1	1

★ C'è una struttura di ALIAS

5f D=AB E=AC

I + ABCD

I + + +

A + BCD

A + BD + CE +

B + ACD

B + AD + +

C + ABD

C + AE + +

D + ABC

D + AB + +

AB + CD

E + AC + +

AC + BD

BC + DE + +

AD + BC

BE + CD + ABC +

→ # righe alias structure = run (talvolta < run)

→ in ogni riga lo stesso numero di termini:

1 se il design è full factorial (no alias)

2 se il design è $\frac{1}{2}$ - frazionario (ad es $L8 = \frac{1}{2} \cdot 16 = 2^4$)

4 se il design è $\frac{1}{4}$ - frazionario (ad es $L8 = \frac{1}{4} \cdot 32 = 2^5$)

....

Design Generators: $D = AB$; $E = AC$

Alias Structure

$I + ABD + ACE + BCDE$

A	+	BD	+	CE	+	ABCDE	3	} III risoluzione
B	+	AD	+	CDE	+	ABCE		
C	+	AE	+	BDE	+	ABCD		
D	+	AB	+	BCE	+	ACDE		
E	+	AC	+	BCD	+	ABDE		
BC	+	DE	+	ABE	+	ACD	4	
BE	+	CD	+	ABC	+	ADE		

● Risoluzione di un design frazionario

Def: è il minimo al variare della riga dell'alias structure della somma dei gradi dei due termini di grado più basso in alias tra loro

III: ci sono fattori singoli in alias con interazioni fra due

IV: ci sono interazioni tra due fattori in alias fra loro (o 1+3)

≥ V: non ci sono fattori singoli o interazioni tra due in alias fra loro

● Significato pratico della alias structure

★ Ipotesi: interazioni a tre o più non capitano nel mondo industriale

→ Se due termini sono in alias, il coefficiente stimato corrispondente sarà la somma delle stime dei due coefficienti → non è possibile separarli

→ In pratica di solito si assume che uno dei due sia zero e che la stima sia quella dell'altro.

↳ un po' arbitrario ma non assurdo: in molte situazioni pratiche è proprio così

★ Esperimenti di *modeling* vanno fatti di IV o V o più in generale non ci può essere aliasing tra due termini potenzialmente non nulli

★ Esperimenti di *screening* si possono fare anche di III (con qualche falso positivo da mettere in conto in caso di interazioni a due rilevanti e aliasing)

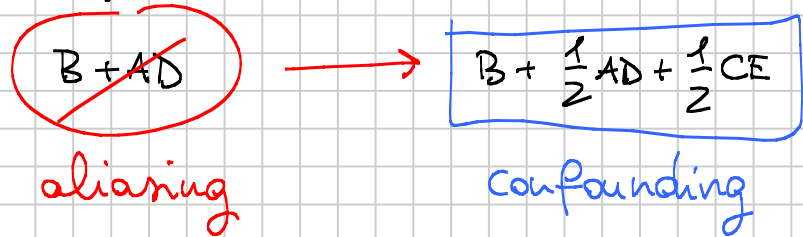
Available Factorial Designs (with Resolution)

Runs	2	3	4	5	6	7	8	9	10	11	12	13	14	15
4	Full	III												
8		Full	IV	III	III	III								
16			Full	V	IV	IV	IV	III	III	III	III	III	III	III
32				Full	VI	IV	IV	IV	IV	IV	IV	IV	IV	IV
64					Full	VII	V	IV	IV	IV	IV	IV	IV	IV
128						Full	VIII	VI	V	V	IV	IV	IV	IV

Available Resolution III Plackett-Burman Designs

Factors	Runs	Factors	Runs	Factors	Runs
2-7	12,20,24,28,...,48	20-23	24,28,32,36,...,48	36-39	40,44,48
8-11	12,20,24,28,...,48	24-27	28,32,36,40,44,48	40-43	44,48
12-15	20,24,28,36,...,48	28-31	32,36,40,44,48	44-47	48
16-19	20,24,28,32,...,48	32-35	36,40,44,48		

★ Design di Plackett-Burman



↳ ancora meglio per lo screening
(approfondimento per l'esame?)

● Efficienza di un design

- I + ABCDEF
- A + BCDEF
- B + ACDEF
- C + ABDEF
- D + APCDF
- E + ASCDF
- F + ABCDE
- AB + CDEF
- AC + BCDF
- AD + BCEF
- AE + BCDF
- AF + BCDE
- BC + ADEF
- BD + ACEF
- BE + ACDF
- BF + ACDE
- CD + ABDF
- CE + ABDF
- CF + ABDE
- DE + ABCF
- DF + ABCE
- EF + ABCD
- ABC + DEF
- ABD + CEF
- ABE + CDF
- ABF + CDE
- ACD + BEF
- ACE + BDF
- ACF + BDE
- ADE + BCF
- ADF + BCE
- AEF + BCD

ok

15 coefficienti quasi fatti nulli

10 coefficienti specifici

- I + ABCD + ADEF + BCDF
- A + BDE + DEF + ABCDF
- B + ACE + CDE + ABDEF
- C + ABE + BDF + ACDEF
- D + ADE + BCF + ABCDE
- E + APC + ADI + BCDEF
- F + ADE + ECD + ABCDE
- AB + CE + ACDF + BDEF
- AC + BE + ABCD + CDEF
- AD + EF + ABCF + BCDE
- AE + BC + DF + ABCDE
- AF + DE + ABCD + BCDF
- BD + CF + ABDE + ACDE
- BE + CD + ABDE + ACEF
- ABD + ACE + BDF + CDE
- ABE + ACD + EDE + CEF

ok

2 o 3 inter. in alias

2 coeff. specifici

L 16 6 fattori

L 32 6 fattori

$$\frac{22}{32} = \frac{11}{16} \approx 0,688$$

$$\frac{10}{16} \approx 0,625$$

★ Def: l'efficienza generale di un design è la media per righe dello "score" seguente

- 1 se c'è un solo EOI
- 1/2 se ci sono due EOI
- 0 se ci sono 0 o ≥ 3 EOI

↳ EOI = effect of interest

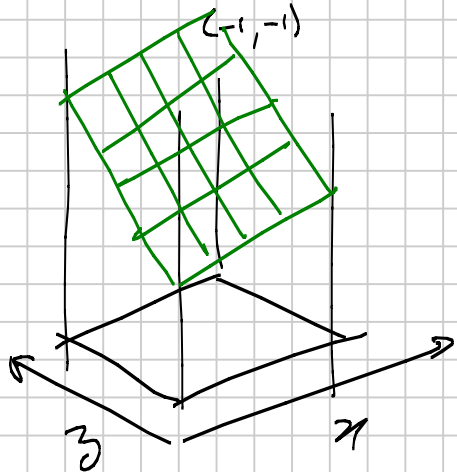
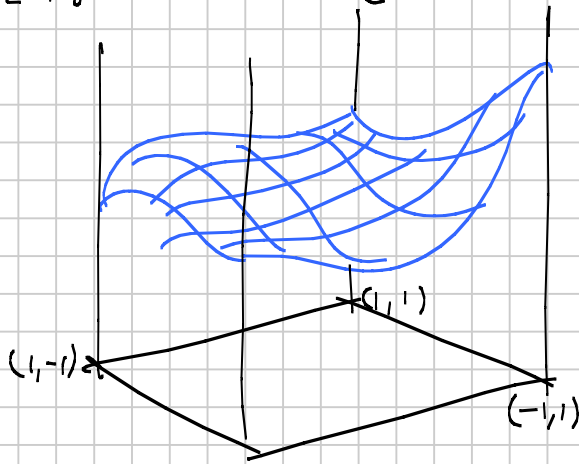
sono i termini che vogliamo inserire nel modello quindi in generale: fatti i fattori singoli e le interaz. a 2

★ L'efficienza specifico si calcola partendo da un set di EOI che sia un sottoinsieme di quelli sopra, dipendente dal particolare esperimento

Treatment Structure	Number of Factors						
	3	4	5	6	7	8	9
L4	0.375						
L8	0.875	0.813	0.500	0.125	0.125		
L16		0.688	1.000	0.625	0.500	0.563	0.313
L32			0.500	0.688	0.766	0.781	0.750
L64				0.344	0.453	0.578	0.641
L128					0.227	0.289	0.359

● Sulle interazioni

$f: [-1,1]^2 \rightarrow \mathbb{R}$ (due soli fattori)

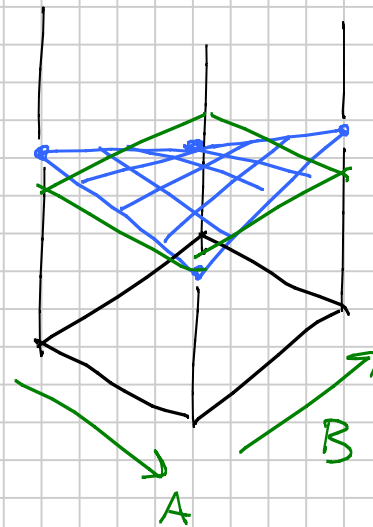


no interazioni:

→ effetto di A non dipende da B
e viceversa

$$f(x, y) = \alpha + \beta x + \gamma y$$

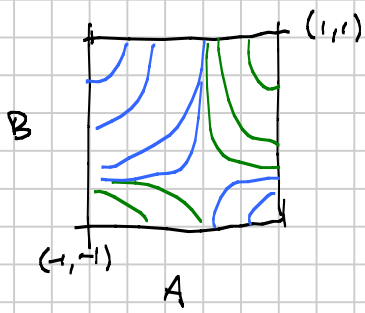
$$f(x, y) = \alpha + \beta x + \delta y + \delta xy$$



interazioni

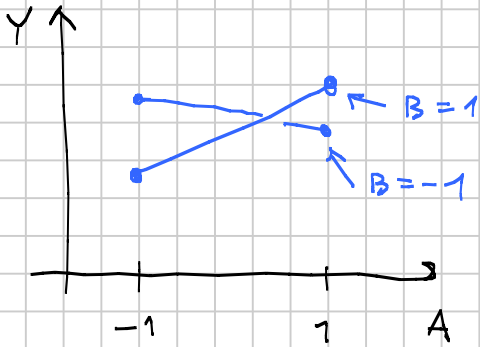
→ effetto di A dipende
dal valore di B e
viceversa

* Contour plot:



curve di livello di
 $c_1A + c_2B + c_3AB$

* Altra rappresentazione:

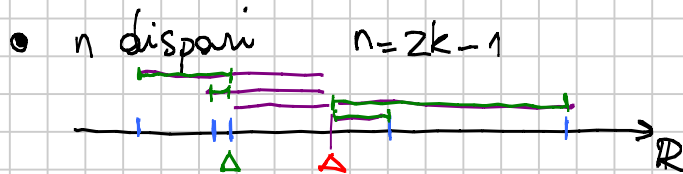


HW: Cosa si deve scegliere per minimizzare $\sum_{i=1}^n |x_i - y|$?

x_1, x_2, \dots, x_n reali

(dall'ora 5)

$y = \operatorname{argmin} f(y) \quad f(y) = \sum_{i=1}^n |x_i - y|$



$x_1 < x_2 < \dots < x_{2k-1}$ wlog

x_k mediana campionaria

$x_i < y < x_{i+1} \quad i \geq k$ wlog

$$f(y) - f(x_k) = \sum_{j=1}^k (|x_j - y| - |x_j - x_k|) + \sum_{k+1}^i (|x_j - y| - |x_j - x_k|) + \sum_{i+1}^{2k-1} (|x_j - y| - |x_j - x_k|)$$

$$= \sum (y - x_j - x_k + x_j) + \sum (y - x_j - x_j + x_k) + \sum (x_j - y - x_j + x_k)$$

$$\geq (y - x_k) \underbrace{(k - 2k + 1 + i)}_{i+1-k} - (y - x_k)(i - k) \geq y - x_k$$

$|x_j - y| - |x_j - x_k| \geq -|y - x_k|$

• n pari (analogo) $n = 2k$

$\forall x \in [x_k; x_{k+1}] \quad x = \operatorname{argmin}_y f(y)$

▣ Torniamo al DoE

• Numero di repliche

- più numerose sono, più precisa è l'analisi e più alti i costi

- sarebbe bene che fossero almeno 3 per ogni run

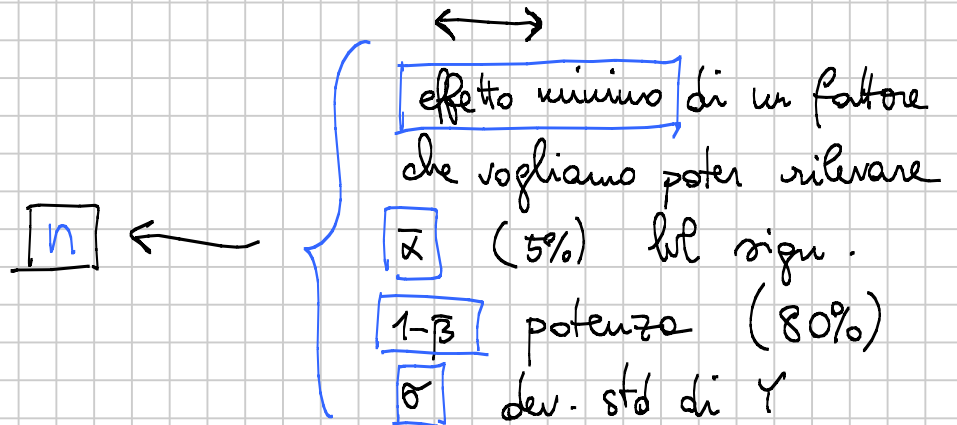
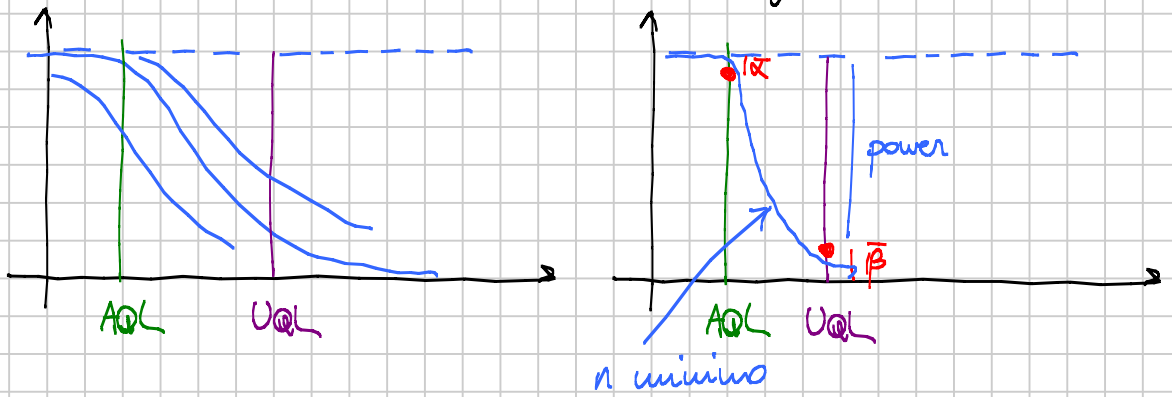
≥ 3 rilievo e identifico outliers

$= 2$ rilievo ma non identifico outliers

$= 1$ non rilievo gli outliers



- per determinare il numero minimo adeguato



* Complicato, richiede σ , effetto minimo (?)

↳ Scosciatoia: "rule of thumb"

$$\left\lceil \frac{\text{run} + 32}{\text{run}} \right\rceil = \text{repliche per ogni run}$$

L4 → 9 repliche

L8 → 5 repliche

L16 → 3 repliche

≥ L32 → 2 repliche

● Center point

I A B C D

1 -1 -1 -1 -1

.....

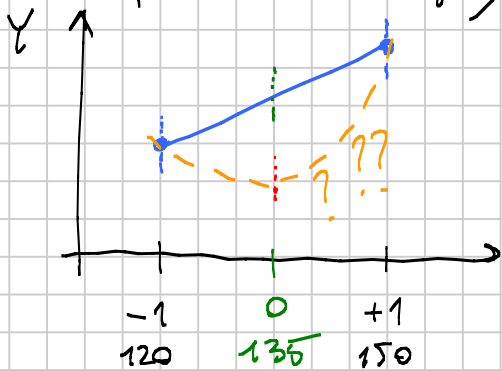
1 1 -1 -1 1

1 0 0 0 0 ← center point

* Punto aggiuntivo L4 → run = 5 → 8 repliche × 5 run = 40

↳ va replicato anche lui

★ Se vi sono variabili categoriche non si fa (o meno di complicare il design)



★ Serve a verificare (a posteriori) che il modello lineare sia adeguato

● Randomizzare le repliche

→ es: L4 3 fattori 4 repliche

run	I	A	B	C	r1	r2	r3	r4	— order —				— stesso case —			
1	1	1	1	1	~	~	~	~	16	1	13	4	1	4	13	16
2	1	1	-1	-1	~	~	~	~	14	10	9	6	6	9	10	14
3	1	-1	-1	1	~	~	~	~	2	7	3	15	2	3	7	15
4	1	-1	1	-1	~	~	~	~	8	12	5	11	5	8	11	12

★ si spreca tempo e denaro per randomizzare

★ Serve a "schermare" l'effetto del tempo

tempo = ordine degli esperimenti = x_1 variabile nascosta

↳ spesso viene significativa

i. usura → curve d'apprendimento

ii. usura, logorio, spreco

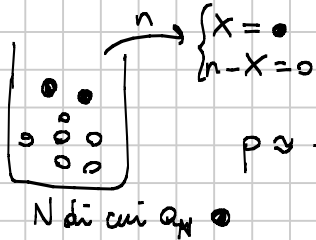
iii. variazioni macroscopiche (petto sostituito)



★ Se non randomizzato, qualunque effetto di questo tipo si può ripercuotere su uno o più fattori (nell'esempio, sicuramente A e forse C)

Parenteri: HW alla fine dell'ora 10

→ se estraggo n pezzi da una popolazione di N (senza rimessa) in cui vi è una frazione p di difetti, il numero di difetti estratti X ha legge ipergeometrica che può essere approssimata con quella binomiale



$$p \approx \frac{a_N}{N} = p_N \quad p_N \xrightarrow{N \rightarrow \infty} p$$

$$P(X=k) = \frac{\binom{a_N}{k} \binom{N-a_N}{n-k}}{\binom{N}{n}} = \frac{n!(N-n)!}{N!} \cdot \frac{a_N!}{k!(a_N-k)!} \cdot \frac{(N-a_N)!}{(n-k)!(N-a_N-n+k)!}$$

$$= \binom{n}{k} \underbrace{\frac{a_N(a_N-1) \dots (a_N-k+1)}{N(N-1) \dots (N-k+1)}}_{k \text{ fattori}} \cdot \underbrace{\frac{(N-a_N) \dots (N-a_N-n+k+1)}{(N-k) \dots (N-n+1)}}_{n-k \text{ fattori}}$$

$$\frac{a_N - c}{N - c} \xrightarrow{N \rightarrow \infty} p$$

$$\frac{N - a_N - c}{N - k - c} \xrightarrow{N \rightarrow \infty} 1 - p$$

$$\approx \binom{n}{k} p^k (1-p)^{n-k}$$

DE: si fanno gli esperimenti

- presenza
- documentazione (diario)
- esperimenti, misurazioni, dati

● Analisi dei dati

- regressione con alcune differenze

$$X^T X = cI$$

B_i hanno covarianza nulla

se cambiano le variabili non variano i B_i calcolati

i B_i sono tutti confrontabili

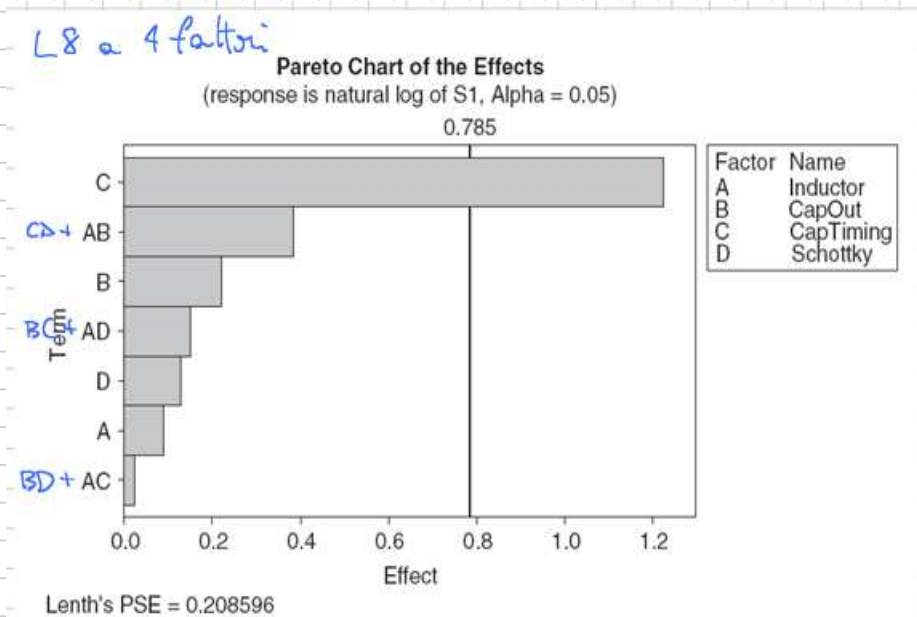
$$Y = 12,4 + 1,2A - 0,7B + 2,5C - 2,1AB$$

coefficienti: ± 1

effetti: C 5 ←
 AB 4,2
 A 2,4
 B 1,4

valore assoluto del doppio del coeff.

★ Selezione delle variabili



→ Si fa qualitativo, tenendo conto di cosa misura Y e di qual è un effetto rilevante

→ Sclte le variabili (regole gerarchica + aliasing)
 si riduce il modello e si ricalciano i residui e si controllano

TEST DEL CHI-QUADRO

φ legge vera di X_1, \dots, X_n iid (legge discreta con pochi valori k)

φ_0 legge ipotetica

$$P(X_i = j) = \varphi(j)$$

Test: $H_0: \varphi = \varphi_0$ $H_2: \varphi \neq \varphi_0$

$$O_j = \#\{i \leq n : X_i = j\} \quad j = 1, 2, \dots, k \text{ wlog}$$

→ Che legge ha O_j ? $O_j \sim \text{bin}(n; \varphi(j))$

$$E(O_j) = n\varphi(j) \text{ vero}$$

$$n\varphi_0(j) = A_j = E(O_j)$$

↑ sotto ipotesi H_0

★ Che legge ha (O_1, O_2, \dots, O_k) ?
 le componenti sono binomiali ma non indipendenti
 $\sum_j O_j = n$ ad esempio

→ Si definisce la legge congiunta **multinomiale**

$$P((O_1, \dots, O_k) = (h_1, \dots, h_k)) = \binom{n}{h_1, h_2, \dots, h_k} \varphi(1)^{h_1} \varphi(2)^{h_2} \dots \varphi(k)^{h_k}$$

$\sum_j h_j = n$
 $n!$
 $h_1! h_2! \dots h_k!$

coefficiente multinomiale

(conta gli anagrammi)

"ABRACADABRA" ha $\binom{11}{5, 2, 2, 1, 1}$ anagrammi

★ Se prendo O_1, \dots, O_k indipendenti con $O_j \sim \text{Pois}(\nu_j)$

la legge di (O_1, \dots, O_k) condizionata all'evento $\{\sum_j O_j = n\}$

è quella multinomiale

• Multinomiale e Poisson

$$S \sim \text{Pois}(\sum_j \nu_j)$$

$O = (O_1, \dots, O_k)$ $O_j \sim \text{Pois}(\nu_j)$ indipendenti

$$S = O_1 + \dots + O_k$$

$$P(O = (a_1, \dots, a_k)) = \prod_{j=1}^k P(O_j = a_j) = \prod_{j=1}^k \frac{\nu_j^{a_j}}{a_j!} e^{-\nu_j}$$

$$\nu := \sum_{j=1}^k \nu_j$$

$$= \frac{(\sum a_j)!}{a_1! a_2! \dots a_k!} \frac{(\sum \nu_j)^{\sum a_j}}{(\sum a_j)!} \cdot \frac{\nu_1^{a_1} \dots \nu_k^{a_k}}{(\sum \nu_j)^{a_1} \dots (\sum \nu_j)^{a_k}} e^{-\sum \nu_j}$$

$$n = \sum_{j=1}^k a_j$$

$$= \underbrace{\binom{n}{a_1, \dots, a_k}}_{\text{legge multinomiale}} \prod_{i=1}^k \left(\frac{\nu_i}{\nu}\right)^{a_i} \underbrace{\frac{\nu^n}{n!} e^{-\nu}}_{\text{legge di Poisson}}$$

$$= P(S = n) P(\text{multin}(n; \frac{\nu_1}{\nu}, \frac{\nu_2}{\nu}, \dots, \frac{\nu_k}{\nu}) = (a_1, \dots, a_k))$$

$$= P(S = n) P(O = (a_1, \dots, a_k) | S = n)$$

$$P(O = (a_1, \dots, a_k)) = P(\{O = (a_1, \dots, a_k)\} \cap \{S = n\}) = P(O = (a_1, \dots, a_k) | S = n) P(S = n)$$

$$= f(a_1, \dots, a_k; n) g(n)$$

★ Se $\sum_n g(n) = 1$ e $\sum_a f(a; n) \neq n$ allora $g(n) = P(S = n)$
 e $f(a; n) = P(O = a | S = n)$
 (HW: check)

★ $O = (O_1, \dots, O_k)$ $O_j \sim \text{Pois}(\nu_j)$ indipendenti $S = \sum_{j=1}^k O_j$ $\nu = \sum_{j=1}^k \nu_j$

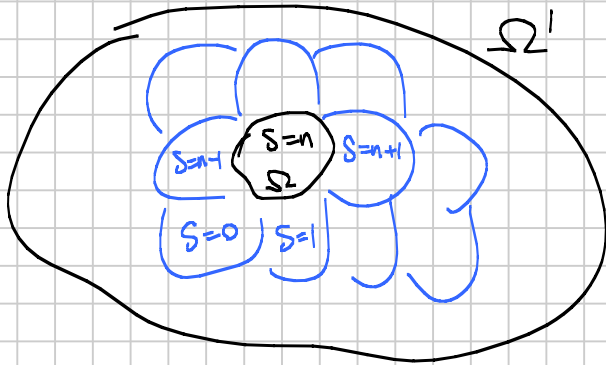
La legge di O condizionata a $\{S = n\}$ è multinomiale $(n; \frac{\nu_1}{\nu}, \dots, \frac{\nu_k}{\nu})$

• Chi-quadro.

X_1, \dots, X_n con legge φ indep. φ ha k valori possibili

O_1, \dots, O_k $O_j = \#\{i \leq n : X_i = k\}$ $S = \sum_{j=1}^k O_j$ $S = n$ q.c.

$(O_1, \dots, O_k) \sim \text{multin}(n; \varphi(1), \dots, \varphi(k))$



Su Ω' O_1, \dots, O_k sono Poisson indipendenti di medie

ν_1, \dots, ν_k proporzionali a $\varphi(1) \dots \varphi(k)$

Per comodità: $\nu_j = \varphi(j)n$

in modo che $E(S) = n$, $S = n$ credibile

★ Statistica: $W = \sum_{j=1}^k \frac{(O_j - A_j)^2}{A_j}$

dove $A_j = n\varphi_0(j)$

sotto ipotesi $H_0: \varphi = \varphi_0$

$W = \sum_{j=1}^k \frac{(O_j - \nu_j)^2}{\nu_j} = \sum_{j=1}^k Z_j^2$

$Z_j := \frac{O_j - \nu_j}{\sqrt{\nu_j}}$ $E(Z_j) = 0$
 $\text{Var}(Z_j) = 1$

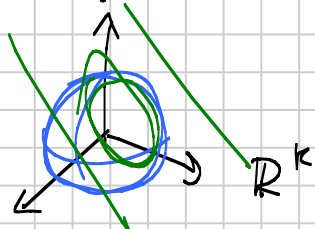
TLC: $\nu_j \gg 1 \Rightarrow Z_j \sim \mathcal{N}(0; 1)$

Nello spazio Ω' , Z_j sono indipendenti: se $\nu_j \gg 1 \forall j$

$W \sim \chi^2(k)$

La legge di W nello spazio Ω è la stessa che in Ω' ma condizionata a $\{S = n\}$

$S = n \Leftrightarrow \sum_{j=1}^k O_j = n \Leftrightarrow \sum_{j=1}^k (\nu_j + \sqrt{\nu_j} Z_j) = n \Leftrightarrow \sum_j \alpha_j Z_j = c$

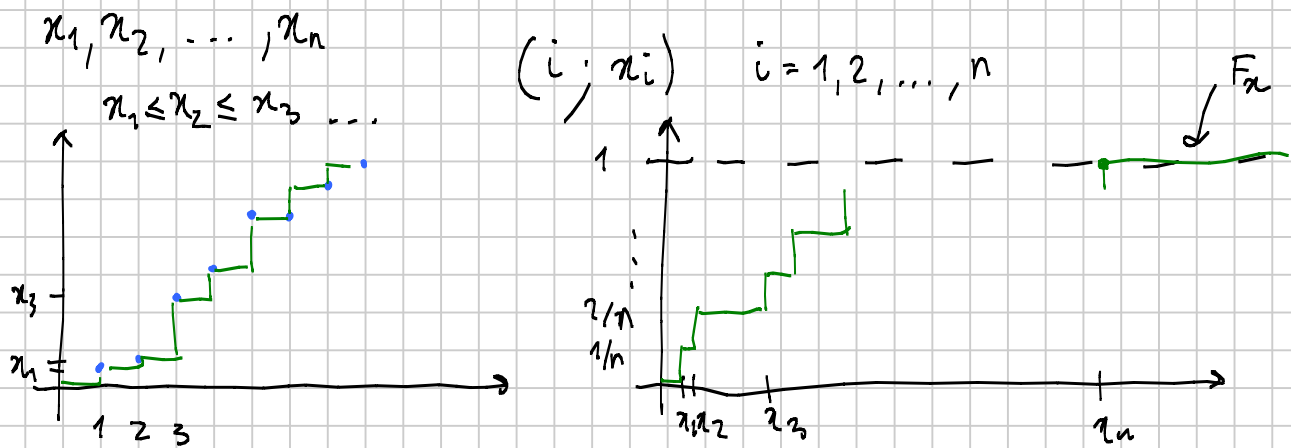


$W \sim \chi^2(k-1)$ su Ω

relazione lineare

Vedi pdf 2011

FINE ?



$$F_n(t) = \frac{1}{n} \#\{i \leq n : x_i \leq t\} \approx P(x_1 \leq t) = F_{X_1}(t)$$

Funzione di ripartizione empirica