

METODI E MODELLI A SUPPORTO DELLE DECISIONI (II) ora 1

Titolo nota

29/09/2009

2 moduli : I Nicolodi → programmazione lineare 6 CFU
II Morandin → statistica multivariata 6 CFU

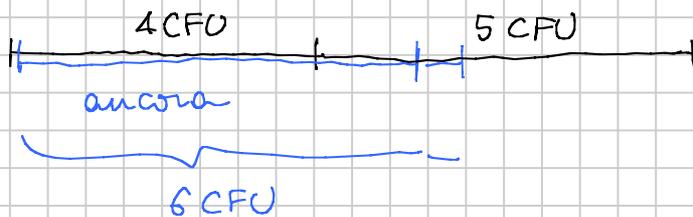
$$5 + 9 = 14 \quad \Rightarrow \quad 9 + (3) + 6 = 15 \quad (18)$$

II modulo :

lun 14.30 - 16.30 LAB
mar 13.30 - 16.30 teoria

} 60 ore = 24 + 36

rispetto all'anno scorso 24 lab + 24 teoria in comune
+12 teoria : *ripasso + approfondimento di statistica + altro*



PROGRAMMA :

Ross : 1 - 11 (1-8 già fatti quasi del tutto)
completare 8 : curva OC / potenza dei test / er I specie
9 : regressione

10 : ANOVA

11 : test di adattamento

approfondimenti :

→ uso di Excel

→ regressione avanzata : Afifi "Computer-aided multivariate analysis"

→ D.O.E. : Sleeper "Design for 6 sigma statistics"

BIBLIOGRAFIA :

Ross, Afifi, Sleeper

www.statsoft.com → stat. textbook

+ Middleton per Excel
Apogeo

dispense di Soliani → 2800 pagine

- STUDENTI → top → capire perfettamente → innovazione
- ↳ middle → imparare a passare l'esame → fatica per niente
- ↳ bottom → cercare di non bloccarsi

• Riciclaggio: martedì 16.30-18.30 dal 5 in poi
(questa settimana mer 9.00-10.00)

ESAME

- * test a crocette → idoneità per i due scritti (entrambi, ma soprattutto Morandini)
- * prova pratica (Excel) (Morandini)
- * prova scritta (Nicolodi)
- * orale (entrambi, ma soprattutto Nicolodi)

COSE PRATICHE

corsi.unipr.it → iscriversi a Metodi Statistici a Supporto delle Decisioni

- foto
- avvisi
- forum
- documenti
 - pdf (anche vecchi)
 - avi (scaricarli subito)
 - .xlsx (Office 2007)
 - scritti vecchi

ora 2

COME FUNZIONA EXCEL

	A	B	...	Z	AA	AB...	K
1	A1						
2			C2				
3							
⋮							
N							

$$N = 65536 = 2^{16} \text{ Office 2003 e precedenti}$$
$$= 1048576 = 2^{20} \text{ Office 2007}$$

$$k = 256 = 2^8 \text{ Office 2003 e precedenti}$$
$$= 16384 = 2^{14} \text{ Office 2007}$$

- riferimenti : a una cella $B4$ = $B4^{\wedge}2$
- a un rettangolo $A1:C2$ = $SOMMA(A1:C2)$

rettangoli "impropri" : AA:AB Z:Z

• Una cella Excel può contenere :

a) Numeri

b) Stringhe

c) Formule (che possono restituire numeri o stringhe)

d) (un pezzo di matrice)

• Formule : iniziano con = oppure -

si possono usare + - * / ^ ()

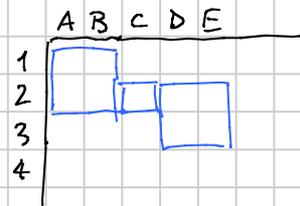
e le relazioni < > <= >= = <>

poi ci sono le funzioni interne

=SOMMA(A1:B2; C2; D2:E3)

il ; separa gli argomenti

=SUM(...)



in generale i comandi sono tradotti automaticamente se si carica il file con una versione di lingua diversa, ma:

CASUALE.TRA

RANDBETWEEN

FRAZ.ANNO

YEARFRAC

invece danno err

→ elenco comandi più usati :

SOMMA, MEDIA, VAR, DEV.ST, INV.NORM.ST,

DISTRIB.NORM.ST, INV.NORM, DISTRIB.NORM,

DISTRIB.T, DISTRIB.CHI, INV.T, INV.CHI, SE

CASUALE

★ CONSIGLIO : Fate molti esercizi dei Capitoli 7 e 8 del Ross usando Excel.

• Numeri : contenuto delle celle e visualizzazione sono due concetti distinti e complementari

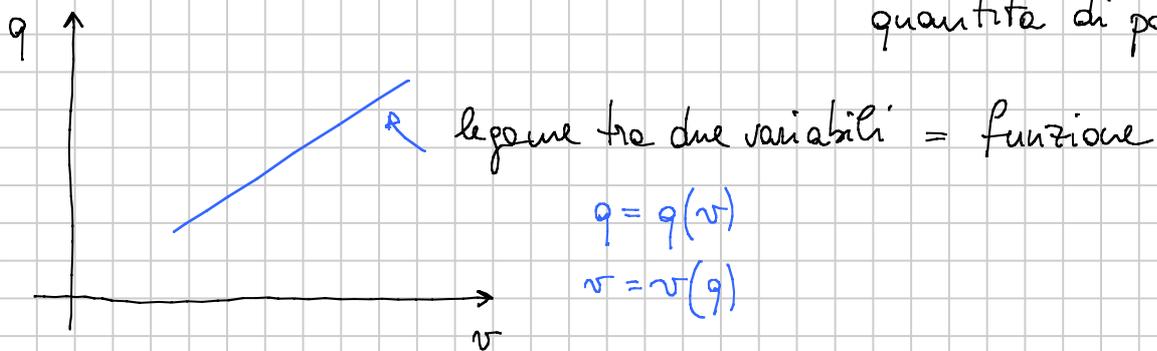
Cambio orario : lun laboratorio 2 e 3 dalle 16:30 alle 18:30

REGRESSIONE

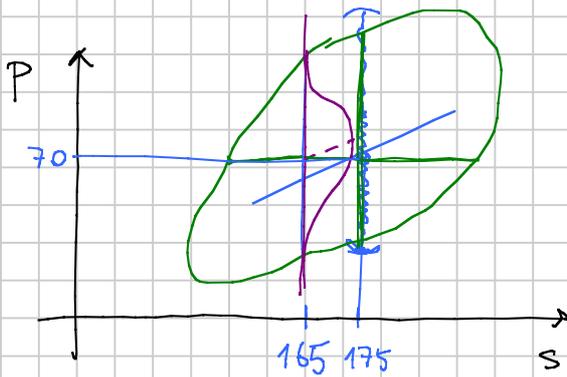
- ★ r. lineare semplice
- ★ r. n multiple
- ★ r. nonlineare
- ...

La regressione studia il rapporto tra due variabili numeriche

es 1 : dosatore automatico di polvere → vel di rotazione
 quantità di polvere



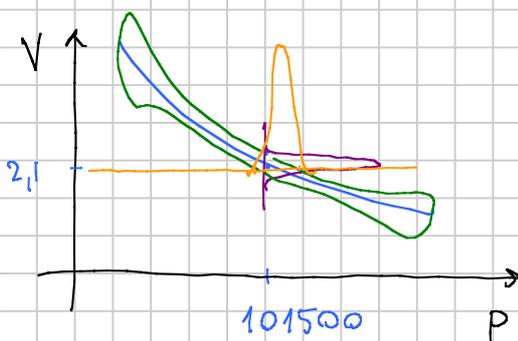
es 2 : misuro esseri umani → peso
 statura



il legame tra queste variabili è "casuale": sapere s permette di conoscere la distribuzione del peso; ad esempio

$$s = 165 \quad p \sim \mathcal{N}(60, 12^2)$$

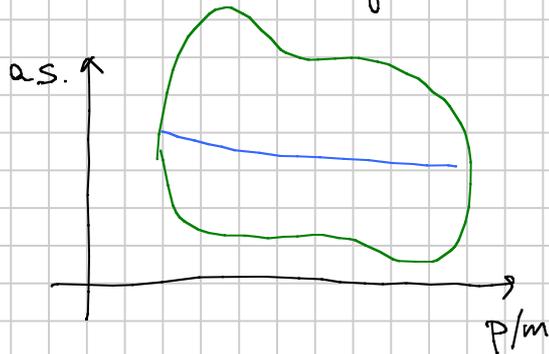
es 3 : barometro elettronico : il sensore trasforma → pressione
 diff di potenziale



(situazione comune a tutti gli strumenti di misurazione)

- La regressione generalizzata lo studio sul rapporto tra due variabili ai casi in cui c'è un margine di incertezza
- La prima risposta della regressione è se ci sia una qualche relazione tra le variabili, o se invece siano indipendenti

es 4: c'è un legame tra → pulsazioni a riposo e livello culturale (= anni di studi) ?



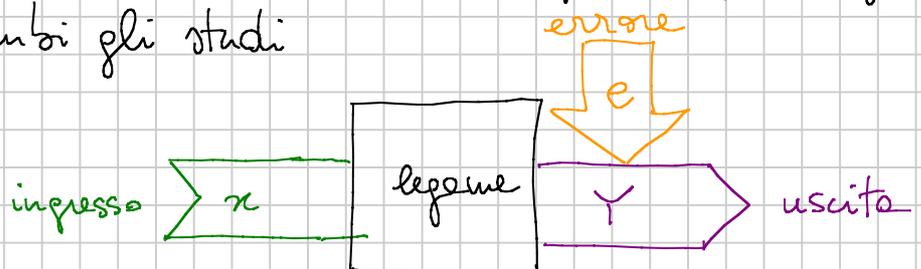
NB "correlazione non significa causalità"

- La regressione richiede di stabilire :
 - * una variabile "causa" che si chiama **variabile di ingresso**
 - o **indipendente** (vel ; stat ; p ; ?) 'X'
 - * una variabile "effetto" che si chiama **variabile di risposta**
 - o **dipendente** (qta ; peso ; V ; ?) 'Y'

→ nell'es. 3 si potrebbero anche scambiare le due var tenendo conto che l'elettronica "vede" V e deve indovinare p

→ nell'es. 4 non è chiaro.

- Non è una scelta innocente; nonostante questo a volte deve essere fatta con un certo arbitrio. Si può comunque sempre scegliere di condurre entrambi gli studi

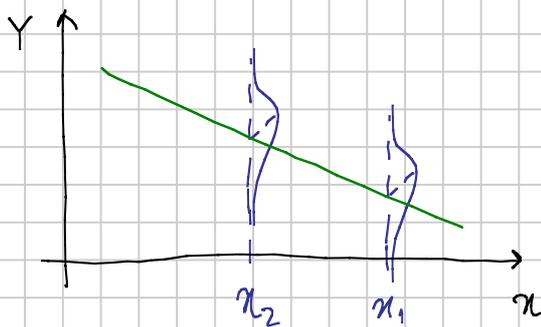


★ Si sceglie quindi quale variabile è "deterministica" e quale è "aleatoria" perché contiene l'incertezza o l'errore

→ Spesso questo punto di vista è irrealistico, ma si tratta di una semplificazione matematica di compromesso

ora 4

REGRESSIONE LINEARE SEMPLICE



legge lineare: " $Y = \alpha + \beta x$ "

considerando anche l'errore:

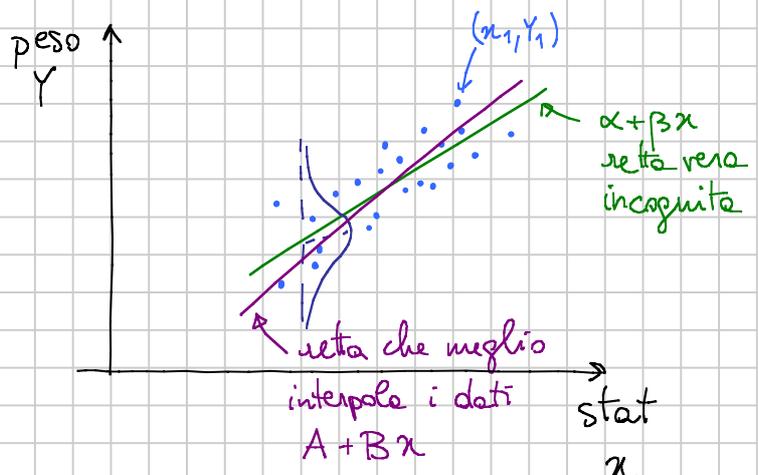
$$Y = \alpha + \beta x + e$$

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

★ σ la dev. standard di e si suppone incognita ma non dipendente da x (= modello omoschedastico)

● Lo scopo della regr in questo caso è trovare α , β e σ !
Questo avviene a partire da un po' di dati sperimentali (2 var per n esperimenti)

id	stat	peso
1	170	82
2	173	64
3	168	71
⋮	181	75
⋮	⋮	⋮
⋮	⋮	⋮
n	154	55



Dati i punti $(x_1, Y_1) \dots (x_n, Y_n)$ cerchiamo la retta $A + Bx$ che meglio li interpola, ovvero più "vicina" ai punti

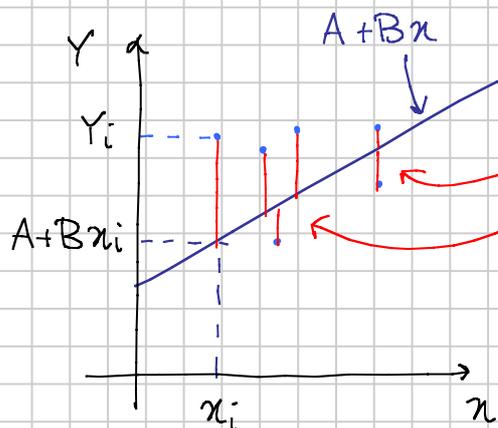
- Date una retta si misura il suo "scarto" dai punti e poi si cerca tra tutte le rette quella con "scarto" minimo.

→ Come misuro lo "scarto" ?



a) potrei sommare tutte le distanze

troppo complicato



$$R_i = Y_i - (A + Bx_i)$$

b) potrei provare a sommare i residui

$$\text{scarto (retta)} = \sum_i R_i$$

sbagliato

c) metto il valore assoluto

$$\text{scarto (retta)} = \sum_i |R_i|$$

troppo complicato

d) faccio i quadrati

$$SS_R = SS = \sum_{i=1}^n R_i^2 = \sum_{i=1}^n (Y_i - A - Bx_i)^2$$

- Cerchiamo i coefficienti A, B a cui corrisponde il minimo di $SS = SS(A, B)$. Come si trovano A e B che rendono minima SS ?

↳ con il "risolutore" di Excel (solver)

↳ facendo le derivate

$$\left\{ \begin{aligned} 0 &= \frac{\partial SS}{\partial A} = \frac{\partial}{\partial A} \sum R_i^2 = \sum \frac{\partial}{\partial A} R_i^2 = 2 \sum R_i \frac{\partial}{\partial A} R_i = -2 \sum R_i \end{aligned} \right.$$

$$\left\{ \begin{aligned} 0 &= \frac{\partial SS}{\partial B} = \frac{\partial}{\partial B} \sum R_i^2 = \sum \frac{\partial}{\partial B} R_i^2 = 2 \sum R_i \frac{\partial}{\partial B} R_i = -2 \sum R_i x_i \end{aligned} \right.$$

equazioni normali $\begin{cases} \sum R_i = 0 \\ \sum x_i R_i = 0 \end{cases} \quad \begin{cases} \sum_{i=1}^n (Y_i - A - Bx_i) = 0 \\ \sum_{i=1}^n x_i (Y_i - A - Bx_i) = 0 \end{cases}$ sist. lineare in A, B

$$\begin{cases} \sum Y_i = nA + B \sum x_i \\ \sum x_i Y_i = A \sum x_i + B \sum x_i^2 \end{cases}$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i \quad \overline{xY} = \frac{1}{n} \sum_{i=1}^n x_i Y_i \quad \overline{x^2} = \frac{1}{n} \sum_{i=1}^n x_i^2$$

$$\begin{cases} \bar{Y} = A + B\bar{x} \\ \overline{xY} = A\bar{x} + B\overline{x^2} \end{cases} \quad \begin{cases} A = \bar{Y} - B\bar{x} \\ \overline{xY} = (\bar{Y} - B\bar{x})\bar{x} + B\overline{x^2} \end{cases}$$

$$\begin{cases} A = \bar{Y} - B\bar{x} \\ B = \frac{\overline{xY} - \bar{x}\bar{Y}}{\overline{x^2} - \bar{x}^2} \end{cases}$$

coeff della retta che interpola meglio

★ $\bar{Y} = A + B\bar{x}$ significa che il punto (\bar{x}, \bar{Y}) , che è il baricentro di tutti i punti, sta sulla retta che interpola meglio. Ovvero $A + Bx$ passa per il baricentro

→ Inoltre $\sum_i R_i = 0$ quindi residui positivi e negativi si bilanciano

HW: Riscrivere la formula di B in funzione della varianza campionaria della x_i , S_x^2 e della covarianza campionaria di x e Y (Capitolo 2 del Ross)

ora 5

- Si trova che A, B sono stimatori corretti e consistenti di α e β , con distribuzione normale di parametri "accessibili"
- Hanno distribuzione normale perché sono combinazioni lineari di vva normali indipendenti:

$$A = \sum_{i=1}^n d_i Y_i \quad B = \sum_{i=1}^n c_i Y_i \quad c_i, d_i \in \mathbb{R}$$

- Nota bene : $Y_i = \alpha + \beta x_i + e_i$ con $e_i \sim \mathcal{N}(0, \sigma^2)$
 allora : $Y_i \sim \mathcal{N}(\alpha + \beta x_i, \sigma^2)$

Siccome gli errori e_i si possono supporre indipendenti, anche le Y_i si possono supporre indipendenti

- Torniamo a $B = \sum_{i=1}^n c_i Y_i$ e giustificiamolo

$$B = \frac{\bar{xY} - \bar{x}\bar{Y}}{\bar{x^2} - \bar{x}^2} = \frac{\sum x_i Y_i - \bar{x} \sum Y_i}{n(\bar{x^2} - \bar{x}^2)} = \sum_{i=1}^n \underbrace{\frac{x_i - \bar{x}}{n(\bar{x^2} - \bar{x}^2)}}_{c_i} Y_i =: \sum_{i=1}^n c_i Y_i$$

Vale anche per A :

$$A = \bar{Y} - B\bar{x} = \frac{1}{n} \sum_{i=1}^n Y_i - \bar{x} \sum_{i=1}^n c_i Y_i = \sum_{i=1}^n \underbrace{\left(\frac{1}{n} - c_i \bar{x} \right)}_{d_i} Y_i =: \sum_{i=1}^n d_i Y_i$$

★ I coefficienti c_i e d_i sono numeri deterministici che dipendono solo dalle x_i

- Sono stimatori corretti : $E(A) = \alpha$ $E(B) = \beta$

$$E(B) = E\left(\sum_i c_i Y_i\right) \stackrel{\substack{\uparrow \\ \text{linearità} \\ \text{di } E}}{=} \sum_i c_i E(Y_i) = \sum_i c_i (\alpha + \beta x_i) = \alpha \sum_i c_i + \beta \sum_i c_i x_i$$

$\sum_i c_i = 0$ $\sum_i c_i x_i = 1$

$$= \alpha \cdot 0 + \beta \cdot 1 = \beta$$

infatti $\sum_i c_i = \frac{\sum (x_i - \bar{x})}{n(\bar{x^2} - \bar{x}^2)} = \frac{n\bar{x} - n\bar{x}}{n(\bar{x^2} - \bar{x}^2)} = 0$

$$\sum_i c_i x_i = \frac{\sum (x_i - \bar{x}) x_i}{n(\bar{x^2} - \bar{x}^2)} = \frac{\sum x_i^2 - \bar{x} \sum x_i}{n(\bar{x^2} - \bar{x}^2)} = \frac{\bar{x^2} - \bar{x}^2}{\bar{x^2} - \bar{x}^2} = 1$$

$$E(A) = E(\bar{Y} - B\bar{x}) = E(\bar{Y}) - \bar{x} E(B) = \frac{1}{n} \sum E(Y_i) - \beta \bar{x}$$

$$= \frac{1}{n} \sum (\alpha + \beta x_i) - \beta \bar{x} = \frac{1}{n} \sum \alpha + \beta \frac{1}{n} \sum x_i - \beta \bar{x} = \alpha$$

★ Punto della situazione : $A \sim \mathcal{N}(\alpha; ?)$ $B \sim \mathcal{N}(\beta; ?)$

● Troviamo le varianze

$$\text{Var}(B) = \text{Var}\left(\sum_i c_i Y_i\right) = \text{Cov}\left(\sum_i c_i Y_i; \sum_j c_j Y_j\right) = \sum_i \sum_j c_i c_j \text{Cov}(Y_i; Y_j)$$

$$= \sum_{i=1}^n c_i^2 \text{Var}(Y_i) = \sigma^2 \sum_{i=1}^n c_i^2$$

↑ $\text{Cov}(Y_i; Y_j) = 0$ se $i \neq j$

$$\begin{aligned} \sum_i c_i^2 &= \sum_i \frac{(x_i - \bar{x})^2}{n(\bar{x}^2 - \bar{x}^2)^2} = \sum_i \frac{x_i^2 - 2\bar{x}x_i + \bar{x}^2}{n(\bar{x}^2 - \bar{x}^2)^2} = \frac{\sum x_i^2 - 2\bar{x}\sum x_i + n\bar{x}^2}{n(\bar{x}^2 - \bar{x}^2)^2} \\ &= \frac{\bar{x}^2 - 2\bar{x}^2 + \bar{x}^2}{n(\bar{x}^2 - \bar{x}^2)^2} = \frac{\bar{x}^2 - \bar{x}^2}{n(\bar{x}^2 - \bar{x}^2)^2} = \frac{1}{n(\bar{x}^2 - \bar{x}^2)} =: k_B \end{aligned}$$

$$\text{Var}(B) = \sigma^2 k_B \quad k_B = \sum_i c_i^2 = \frac{1}{n(\bar{x}^2 - \bar{x}^2)}$$

analogamente

$$\text{Var}(A) = \sigma^2 k_A \quad k_A = \sum_i d_i^2 = \frac{\bar{x}^2}{n(\bar{x}^2 - \bar{x}^2)}$$

● H_0 la distribuzione di A e B :

$$A \sim \mathcal{N}(\alpha, \sigma^2 k_A) \quad B \sim \mathcal{N}(\beta, \sigma^2 k_B)$$

★ NON sono indipendenti : $\text{Cov}(A; B)$ si può calcolare

● Mancano solo uno stimatore per σ

$$SS = \sum_{i=1}^n R_i^2 = \sum_{i=1}^n (Y_i - A - Bx_i)^2$$

... se invece ci fossero α e β :

$$Y_i \sim \mathcal{N}(\alpha + \beta x_i, \sigma^2)$$

$$\sum_{i=1}^n \left(\frac{Y_i - (\alpha + \beta x_i)}{\sigma} \right)^2 = \sum_{i=1}^n \left(\mathcal{N}(0, 1) \right)^2 \sim \chi^2(n) \quad \text{per definizione}$$

$$\frac{SS}{\sigma^2} = \sum_{i=1}^n \left(\frac{Y_i - A - Bx_i}{\sigma} \right)^2 \sim \chi^2(n-2)$$

↑
 then
 che faccio
 la prox volta

↑
 ho sostituito
 2 parametri con
 i relativi stimatori

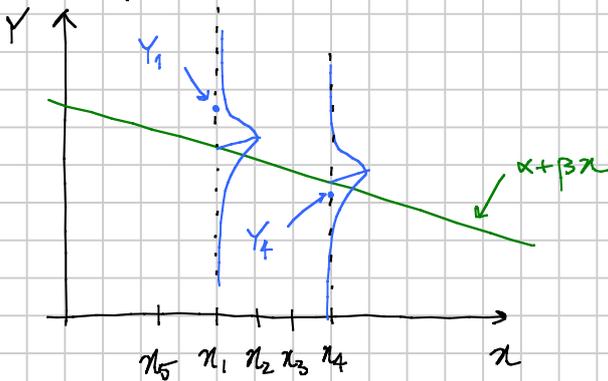
★ inoltre SS indep da A e SS indep da B

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

■ $\frac{SS}{n-2}$ è uno stimatore corretto di σ^2 , indipendente da A e B.

$$\begin{aligned} \text{Infatti } E\left(\frac{SS}{n-2}\right) &= \frac{1}{n-2} E(SS) = \frac{\sigma^2}{n-2} E\left(\frac{SS}{\sigma^2}\right) = \frac{\sigma^2}{n-2} E(\chi^2(n-2)) \\ &= \frac{\sigma^2}{n-2} \cdot (n-2) = \sigma^2 \end{aligned}$$

Regressione lineare semplice



$$Y_i \sim \mathcal{N}(\alpha + \beta x_i, \sigma^2)$$

Scopo: inferenza su α, β, σ

1) A, B stimatori per α, β

$$A \sim \mathcal{N}(\alpha, \sigma^2 k_A) \quad B \sim \mathcal{N}(\beta, \sigma^2 k_B)$$

2) $\frac{SS}{n-2}$ stimatore per σ^2

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

Vediamo ora le condizioni generali per dedurre che $\frac{SS}{\sigma^2} \sim \chi^2(n-2)$ e enunciati analoghi

→ Prima altri 2 esempi (e ripasso)

a) $X_1, X_2, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$ μ e σ incognite
 $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \approx \mu$ $\bar{X} \sim \mathcal{N}(\mu, \frac{\sigma^2}{n})$

$$\sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma} \right)^2 \sim \chi^2(n)$$

per definizione

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

$$\frac{S^2}{\sigma^2} (n-1) = \sum_{i=1}^n \left(\frac{X_i - \bar{X}}{\sigma} \right)^2 \sim \chi^2(n-1)$$

b) due campioni X_i e Y_j $X_i \sim \mathcal{N}(\mu_X, \sigma^2)$ $Y_j \sim \mathcal{N}(\mu_Y, \sigma^2)$

$$\bar{X} - \bar{Y} \sim \mathcal{N}(\mu_X - \mu_Y; \frac{\sigma^2}{m} + \frac{\sigma^2}{n})$$

m, n numerosità campioni

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

$$\sum_{i=1}^m \left(\frac{x_i - \mu_x}{\sigma} \right)^2 + \sum_{j=1}^n \left(\frac{Y_j - \mu_Y}{\sigma} \right)^2 \sim \chi^2(m+n) \quad \text{per def}$$



$$\sum_{i=1}^m \left(\frac{x_i - \bar{X}}{\sigma} \right)^2 + \sum_{j=1}^n \left(\frac{Y_j - \bar{Y}}{\sigma} \right)^2 \sim \chi^2(m+n-2) \quad \text{per thm}$$

$$\frac{S_P}{\sigma^2} \parallel \leftarrow \text{HW: verificare}$$

c) x_1, \dots, x_n numeri $Y_i \sim \mathcal{N}(\alpha + \beta x_i; \sigma^2)$ indipendenti

$$\sum_{i=1}^n \left(\frac{Y_i - (\alpha + \beta x_i)}{\sigma} \right)^2 \sim \chi^2(n) \quad \text{per def}$$



$$SS := \sum_{i=1}^n (Y_i - A - B x_i)^2$$

$$\frac{SS}{\sigma^2} = \sum_{i=1}^n \left(\frac{Y_i - (A + B x_i)}{\sigma} \right)^2 \sim \chi^2(n-2) \quad \text{per thm}$$

IN GENERALE

$H_P: Y_1, \dots, Y_n$ $Y_i \sim \mathcal{N}(\mu_i; \sigma^2)$ indipendenti

B_1, B_2, \dots, B_k ($k=2$ in b) e c)) combinazioni lineari a coeff noti delle Y_i :

$$B_1 = c_{11} Y_1 + c_{12} Y_2 + c_{13} Y_3 + \dots + c_{1n} Y_n$$

$$B_2 = c_{21} Y_1 + \dots$$

$$B_j = \sum_{i=1}^n c_{ji} Y_i$$

v^1, v^2, \dots, v^k vettori di \mathbb{R}^n linearmente indipendenti
 $v^j = (v_1^j, v_2^j, \dots, v_n^j) \in \mathbb{R}^n$

$$E \left[Y_i - (v_i^1 B_1 + v_i^2 B_2 + \dots + v_i^k B_k) \right] = 0$$

$$ts: W := \sum_{i=1}^n \left(\frac{Y_i - (v_i^1 B_1 + v_i^2 B_2 + \dots + v_i^k B_k)}{\sigma} \right)^2 \sim \chi^2(n-k)$$

inoltre W è indipendente da (B_1, B_2, \dots, B_k)

• c) usa il thm :

$$W = \sum_{i=1}^n \left(\frac{Y_i - (A + Bx_i)}{\sigma} \right)^2$$

$$k=2 \quad B_1=A, \quad B_2=B \quad \text{ok}$$

$$v_i^1 = 1 \quad i=1, 2, \dots, n$$

$$v^1 = (1, 1, \dots, 1)$$

$$v_i^2 = x_i \quad i=1, 2, \dots, n$$

$$v^2 = (x_1, x_2, \dots, x_n)$$

v^1 e v^2 sono lin. dipend. solo se $v^2 = \lambda v^1$ ovvero se

$$v^2 = (\lambda, \lambda, \dots, \lambda) \quad \text{ovvero se } x_i = \lambda \quad \forall i$$

→ Basta che gli x_i non siano tutti uguali

ok

$$B_1 = A = \sum_i d_i Y_i \quad B_2 = B = \sum_i c_i Y_i \quad \text{sono comb. lin.}$$

$$E(Y_i - A - Bx_i) = E(Y_i) - E(A) - x_i E(B)$$

$$= \alpha + \beta x_i - \alpha - x_i \beta = 0 \quad \text{ok}$$

Tesi : $\frac{SS}{\sigma^2} \sim \chi^2(n-2)$ inoltre SS è indipendente da (A, B)

HW: a) e b) usano il thm

ora 7

ERRORE STANDARD

$$\frac{SS}{\sigma^2} \sim \chi^2(n-2) \quad \Rightarrow \quad E\left(\frac{SS}{\sigma^2}\right) = n-2$$

$$\text{Var}\left(\frac{SS}{\sigma^2}\right) = 2n-4$$

$$\Rightarrow E(SS) = \sigma^2(n-2) \quad \Rightarrow \quad E\left(\frac{SS}{n-2}\right) = \sigma^2$$

$$\text{Var}\left(\frac{SS}{n-2}\right) = \text{Var}\left(\frac{SS}{\sigma^2} \cdot \frac{\sigma^2}{n-2}\right) = \left(\frac{\sigma^2}{n-2}\right)^2 \text{Var}\left(\frac{SS}{\sigma^2}\right) = \frac{2\sigma^4}{n-2} \quad \text{tende a 0}$$

$n \rightarrow \infty$

- $\frac{SS}{n-2}$ è uno stimatore corretto e consistente di σ^2

HW: verificare se A e B sono stimatori consistenti di α, β .

- $S_e := \sqrt{\frac{SS}{n-2}}$ è uno stimatore (non corretto) di σ .

Prende il nome di errore standard

HW: capire se $E(S_e) \geq \sigma$

★ S_e è indipendente da (A, B)

■ FUNZIONI ANCILLARI

→ funzione che dipende dai dati, dai parametri noti e da un solo parametro incognito (quello su cui si vuole fare inferenza) di cui sia nota la distribuzione

→ esempi: $\frac{\bar{X} - \mu}{\sigma / \sqrt{n}} \sim \mathcal{N}(0, 1)$ se σ è nota

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

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

$$\boxed{\frac{\bar{X} - \mu}{S / \sqrt{n}} \sim t(n-1)}$$

grazie al fatto che X e S sono indipendenti

• Veniamo ad A, B:

$$A \sim \mathcal{N}(\alpha; \sigma^2 k_A)$$

$$k_A = \frac{1}{n} \cdot \frac{\bar{x}^2}{\bar{x}^2 - \bar{\bar{x}}^2}$$

$$\frac{A - \alpha}{\sigma \sqrt{k_A}} \sim \mathcal{N}(0, 1)$$

non è ancillare

Tuttavia, visto che S_e è stimatore di σ indipendente da A e visto che $\frac{S_e^2}{\sigma^2} (n-2) = \frac{SS}{\sigma^2} \sim \chi^2(n-2)$, si ha che:

$$\boxed{\frac{A - \alpha}{S_e \sqrt{k_A}} \sim t(n-2)}$$

funz ancillare

Analogamente : $B \sim \mathcal{N}(\beta, \sigma^2 k_B)$ $k_B = \frac{1}{n} \cdot \frac{1}{\bar{x}^2 - \bar{x}^2}$

$$\frac{B - \beta}{\sigma \sqrt{k_B}} \sim \mathcal{N}(0, 1)$$

$$\boxed{\frac{B - \beta}{Se \sqrt{k_B}} \sim t(n-2)} \quad \text{funz. ausiliaria}$$

Infine

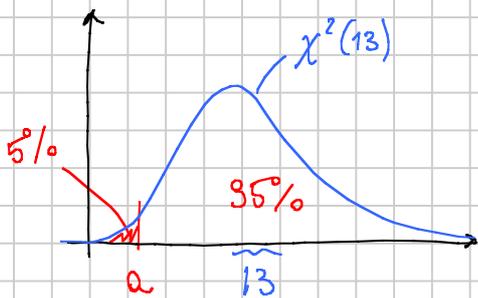
$$\boxed{\frac{SS}{\sigma^2} = \frac{Se^2}{\sigma^2} (n-2) \sim \chi^2(n-2)} \quad \text{funz. ausiliaria}$$

INFERENZA SULLA REGR. LIN. SEMPLICE

Siamo in grado di fare:

2) Intervalli di confidenza su α , β , σ

→ ad es.: int. conf. unil. sx su σ con lvl di conf 95%



$$n = 15, \quad Se = 0,7$$

$$a = \text{INV. CHI}(95\%; 13) \approx 5,89$$

$$95\% = P(\chi^2(13) > a) = P\left(\frac{SS}{\sigma^2} > a\right)$$

$$= P(SS > a \sigma^2) = P\left(\frac{SS}{a} > \sigma^2\right) = P\left(\sigma < \sqrt{\frac{SS}{a}}\right)$$

Con il 95% di conf $\sigma < Se \sqrt{\frac{n-2}{a}} \approx 1,04$

b) Test statistici su α , β , σ

→ ad es.: $H_0: \beta = 0$ vs $H_1: \beta \neq 0$ (p dei dati) $n=15$

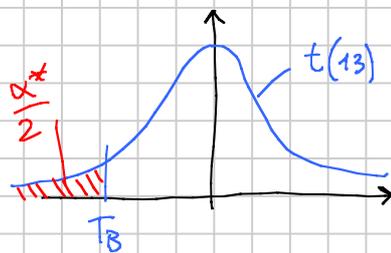
la statistica del test è

$$\frac{B - \beta}{Se \sqrt{k_B}} \sim t(n-2) \quad \Rightarrow \quad \frac{B}{Se \sqrt{k_B}} \stackrel{\beta_0=0}{=} T_B \stackrel{H_0}{\sim} t(n-2)$$

Calcolo $T_B = -2,314$

$$\alpha^* = 2 P(t(13) \leq -2,314)$$

$$\alpha^* = 2 - 2 P(t(n-2) \leq |T_B|)$$



$$\alpha^* = 2 * \text{DISTRIB.T}(-2,314; 13) = \text{NO}$$

$$= \text{DISTRIB.T}(2,314; 13; 2) \approx 3,8\%$$

ora 8

⊙ Il test più importante: $H_0: \beta = 0$ vs $H_1: \beta \neq 0$

Q: La var Y dipende dalla var x ? ($H_1 = \text{sì}$; $H_0 = \text{no}$)

$$\alpha^* = \text{DISTRIB.T}(\text{ABS}(T_B); n-2; 2)$$

★ Quando facciamo la regressione con gli strumenti automatici di Excel, questo valore è calcolato automaticamente

★ $H_1 \rightarrow$ tutto ok, la regressione è utile

$H_0 \rightarrow$ vuol dire che $Y_i \sim \mathcal{N}(\mu, \sigma^2)$ è un campione omogeneo che si può trattare con tecniche elementari

⊙ Analoghi su α $H_0: \alpha = 0$ vs $H_1: \alpha \neq 0$

Q: c'è termine noto nella relazione ($H_1 = \text{sì}$; $H_0 = \text{forse no}$)

a) Ci sono grandezze che hanno una rel lineare senza termine noto

\rightarrow testo un resistore $V = IR$

$$I = \frac{1}{R} V$$

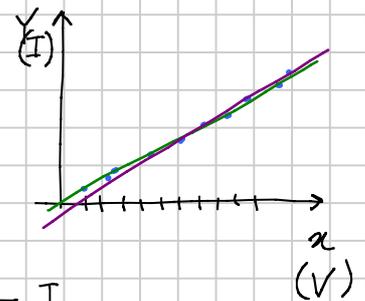
Posso stimare $\frac{1}{R}$ (e quindi R)

usando la regressione $x \leftarrow V$ $Y \leftarrow I$

$$Y = \alpha + \beta x$$

$$I = \alpha + \beta V = 0 + \frac{1}{R} V \quad \text{perci\u00f2} \quad \alpha = 0 \quad \beta = \frac{1}{R}$$

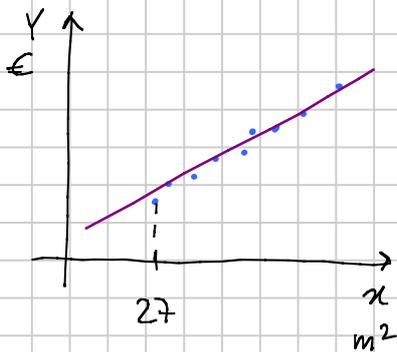
B stima $\frac{1}{R}$ B^{-1} stima R



★ In questo caso sarebbe più sensato fare la regressione con un modello ancora più semplice: $Y = \beta x$

b) Ci sono vari dubbi: non so se $\alpha = 0$ o $\alpha \neq 0$

ad esempio: appartamenti \rightarrow superficie [m^2]; prezzo [€]
 x Y



$$Y = \alpha + \beta x$$

in questo caso: può avere senso che $\alpha = 0$?
se no, non faccio il test e tengo il termine
noto anche se pdd è molto grande
se si, faccio il test:

se $H_1 \rightarrow$ tutto ok, continuo così

se $H_0 \rightarrow$ devo rifare la regressione con $Y = \beta x$

HW: Costruire da zero fino all'inferenza la regressione senza α .
 $Y = \beta x + e$ $e \sim \mathcal{N}(0, \sigma^2)$

• Inferenza sulla risposta media e sulle risposte future

\rightarrow esempio: regressione peso vs statura
 Y x

Q: Quanto dovrebbe pesare una persona alta 180 cm?

poniamo $x_0 = 180$; facciamo la regressione

$Y = \alpha + \beta x$ sui dati a disposizione; troviamo A, B, S_e

i. $\alpha + \beta x_0$ (incognito) è il peso medio che hanno
le persone alte x_0

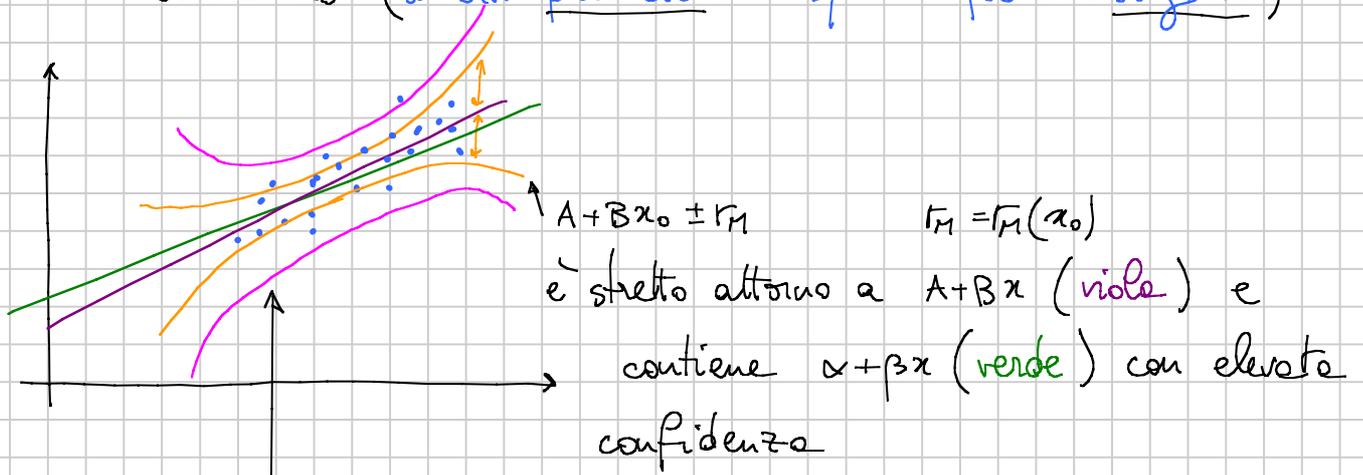
ii. $\mathcal{N}(\alpha + \beta x_0; \sigma^2)$ (incognita) è la distribuzione di
tutti i pesi delle persone alte x_0

iii. $A + B x_0$ stima puntuale del peso medio

\rightarrow iv. $A + B x_0 \pm t_m$ (t_m piccolo se n grande) intervallo di
confidenza per $\alpha + \beta x_0$ (interv. conf risposta media)

\rightarrow v. $A + B x_0 \pm F_F$ (F_F sempre maggiore di S_e) intervallo che
con elevate probabilità conterrà il peso di una persona

alta x_0 (interv. predizione risposta futura singola)



$$A + Bx_0 \pm r_1 \quad r_1 = r_1(x_0)$$

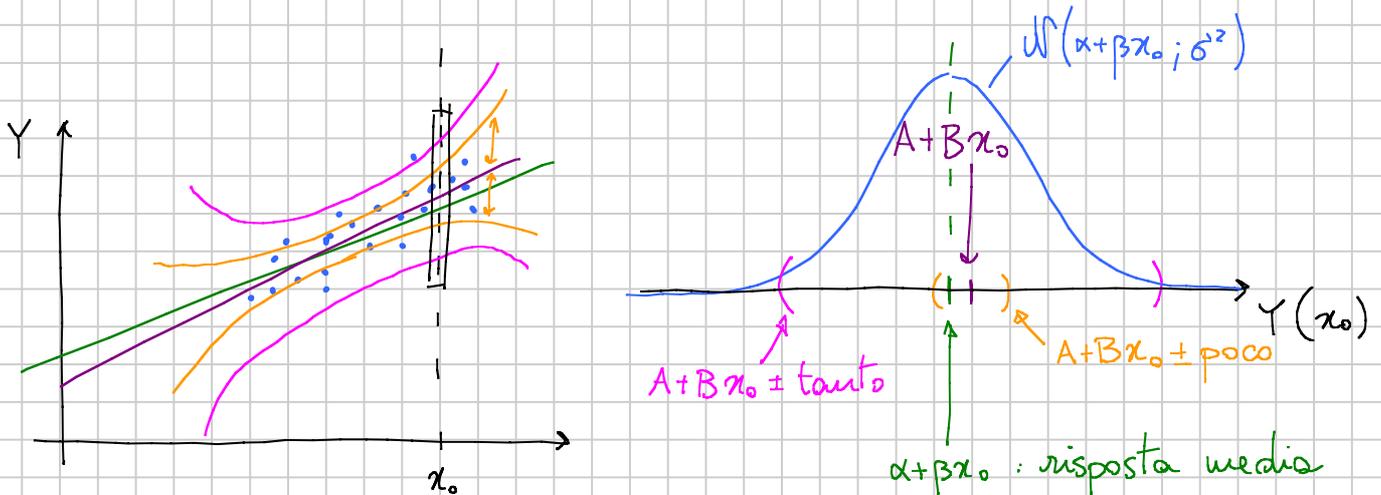
è stretto attorno a $A + Bx$ (viola) e

contiene $\alpha + \beta x$ (verde) con elevata
confidenza

$$A + Bx_0 \pm r_F \quad r_F = r_F(x_0)$$

è centrato attorno a $A + Bx$ (viola) e contiene
i singoli punti (vecchi e futuri) con elevata probabilità

INFERENZA PER RISPOSTA MEDIA E FUTURA



$$Y \sim \mathcal{N}(x + \beta x_0; \sigma^2)$$

$A + Bx_0 \pm a S_e \sqrt{k_m}$: (arancione) intervallo di confidenza
per la risposta media - può essere molto
stretto, soprattutto se n è grande

$A + Bx_0 \pm a S_e \sqrt{1 + k_m}$: (fucsia) intervallo di predizione :
una singola risposta futura Y relativa ad un livello
di ingresso x_0 apparterrà a questo intervallo con
probabilità assegnata

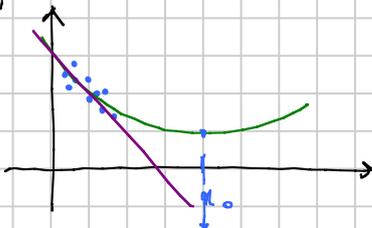
$$k_m = k_m(x_0) = \frac{1}{n} \left[1 + \frac{(x_0 - \bar{x})^2}{\bar{x}^2 - \bar{x}^2} \right]$$

→ è minima quando $x_0 = \bar{x}$ e si allarga verso i bordi

→ non è mai minore di $\frac{1}{n}$

→ tende a zero per $n \rightarrow \infty$

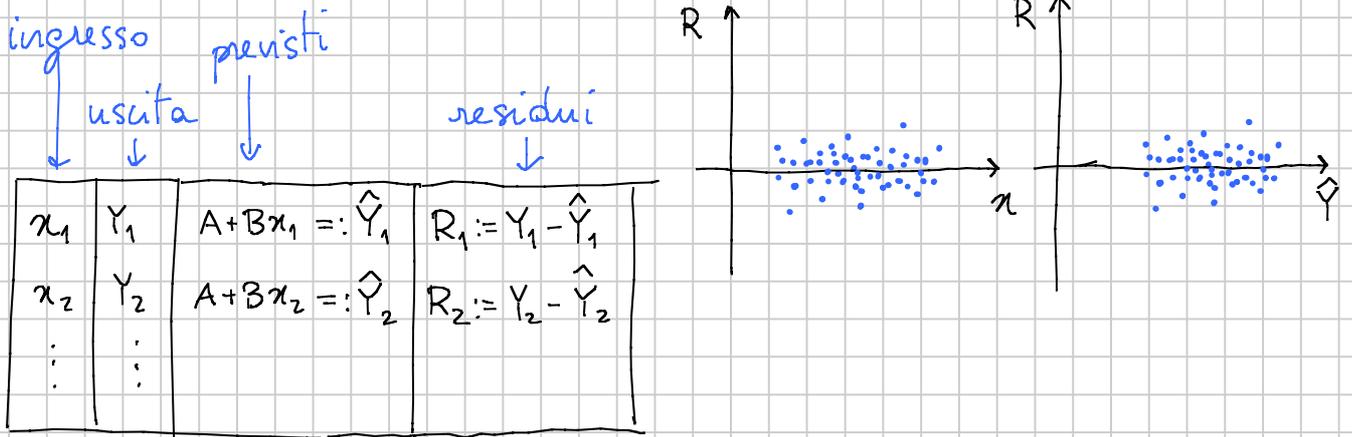
★ In generale l'affidabilità dell'inferenza a ingresso x_0 non è
garantita per valori esterni a quelli dei dati o vicini ai bordi.



- a deve essere un quantile t di Student con $n-2$ g.d.l.
(Si può prendere bilaterale o unilaterale...)

ANALISI DEI RESIDUI

Il modo migliore per capire se la regressione è affidabile è analizzare qualitativamente il diagramma di dispersione dei residui vs la var di ingresso o il valore previsto



A, B, Se

$$SS = \sum_i (y_i - (A+Bx_i))^2 = \sum_i (y_i - \hat{y}_i)^2 = \sum_i R_i^2$$

- Intendiamo che la regressione è affidabile se davvero:

$$Y \sim \mathcal{N}(\alpha + \beta x; \sigma^2)$$

la varianza non dipende da x
(modello omoschedastico)

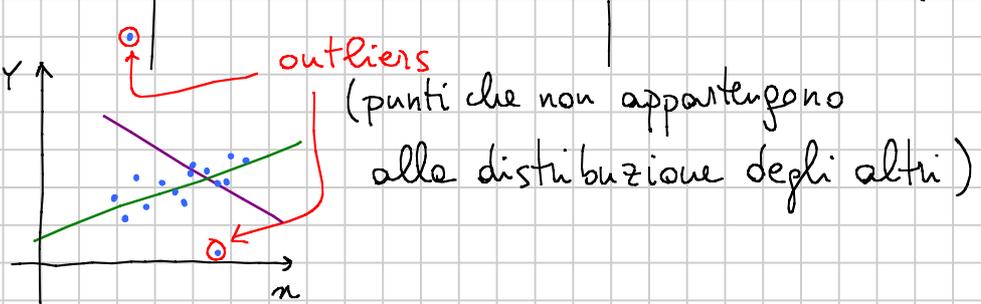
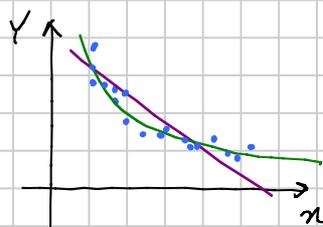
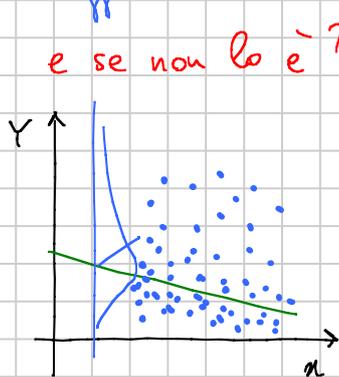
funzione lineare di x

e se invece $\sigma = \sigma(x)$?

e se non è lineare?

legge normale

e se non lo è?



Informazioni di servizio :

- gio 29 ~~Nicolodi~~ → Morandin
- gio 22 Nicolodi normalmente
- file .xlsx → è un file .zip ... click col dx → apri con... Excel Office XP - 2003 → google → open file office 2007

Bizzarrie da evitare

R_1, R_2, \dots, R_n sono un campione "quasi" gaussiano

$$\bar{R} = \frac{1}{n} \sum_{i=1}^n R_i \quad \text{media campionaria}$$

è sempre 0

$(x_i, R_i) \rightarrow$ faccio la regressione

viene sempre $A_R = 0 = B_R$

$$S_R^2 := \frac{1}{n-1} \sum_i (R_i - \bar{R})^2 = \frac{1}{n-1} \sum_i R_i^2 = \frac{SS}{n-1} \quad \text{è una stimate di } \sigma^2 \text{ che sottostima}$$

$$\sigma^2 \approx \frac{SS}{n-2}$$

COEFFICIENTE DI DETERMINAZIONE

$$R^2 = 1 - \frac{SS}{S_{YY}}$$

dove $SS = \sum_{i=1}^n (Y_i - (A + Bx_i))^2$

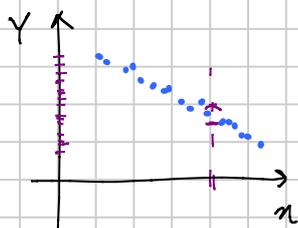
imprecisione sulla conoscenza di Y_i , conoscendo le x_i oppure no

$$S_{YY} = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

$\frac{SS}{S_{YY}}$ è la frazione di imprecisione che non viene rimossa anche se si conoscono le x_i

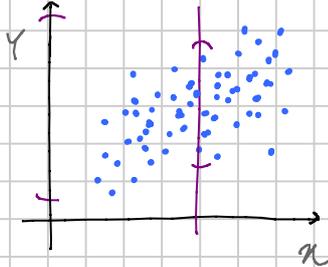
$R^2 = 1 - \frac{SS}{S_{YY}}$ è la fraz. di imprecisione rimossa dalla conoscenza delle x_i , o anche "la frazione di variante spiegata dalla regressione"

★ Un elevato valore di R^2 significa che la regressione è buona e utile ovvero che le nostre previsioni e stime su Y sono molto migliori quando conosciamo x



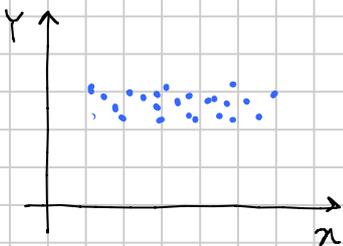
$R^2 \approx 97\%$ grande

p dei dati $H_0: \beta = 0$ 10^{-40} piccolo



$R^2 \approx 6\%$ piccolo

p dei dati $H_0: \beta = 0$ 10^{-20} piccolo



$R^2 \approx 3\%$ piccolo

p dei dati $H_0: \beta = 0$ 0,72 grande

★ $R^2 \in [0,1]$

★ Nel caso della regressione lineare semplice $R^2 = r^2$
dove r è il coefficiente di correlazione lineare campionaria (Cap 2)

$$\rho(x, Y) = \frac{\text{Cov}(x, Y)}{\sqrt{\text{Var}(x) \text{Var}(Y)}}$$



$$r(x_i, Y_i) = \dots$$

→ Se due variabili hanno una correlazione di 0,4 ... non è gran che perché $0,4^2 = 0,16 = 16\%$

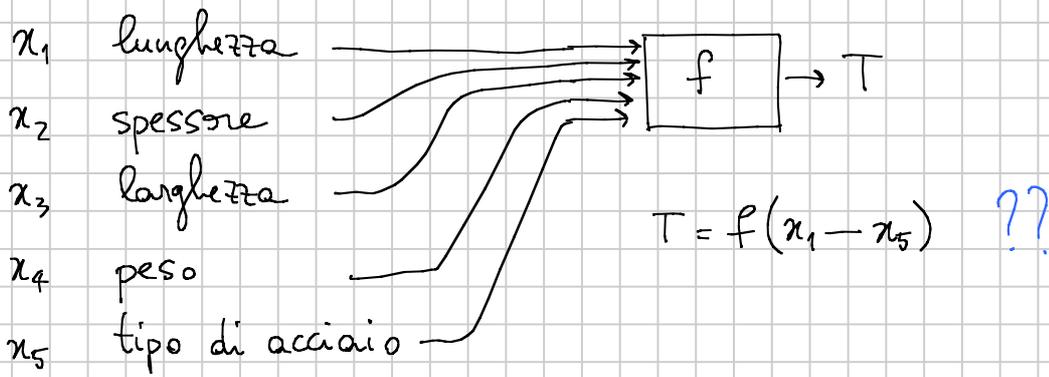
REGRESSIONE LINEARE MULTIPLA

→ Esempio: acciaieria → linea di lavorazione → programma



coil di acciaio

↓
l'ingresso viene occupato per tempo T non prevedibile facilmente



★ La repr. lineare multiple permette di studiare le relazioni tra $p \geq 1$ variabili di ingresso e una variabile di uscita

ora 11

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

Dati:

id	Y	x_1	x_2	...	x_p
1	$Y(1)$	$x_1(1)$	$x_2(1)$		$x_p(1)$
2	$Y(2)$	$x_1(2)$	\vdots		\vdots
\vdots	\vdots	\vdots	\vdots		\vdots
n	$Y(n)$	$x_1(n)$...		$x_p(n)$

$Y(i)$ risposta del punto i
 $x_j(i)$ var di ingresso j -esima del punto i

altra notazione:

- $Y(i) = Y_i$
- $x_j(i) = x_{ij}$ (come matrice)

$$Y(i) \sim \mathcal{N}\left(\beta_0 + \sum_{j=1}^p \beta_j x_j(i); \sigma^2\right)$$

modello lineare, gaussiano, omoschedastico

★ Molto spesso si inventa ed inserisce $x_0 \equiv 1$ in modo che sia:

$$Y \sim \mathcal{N}\left(\sum_{j=0}^p \beta_j x_j; \sigma^2\right)$$

Sia X la matrice x_{ij} $i = 1, 2, \dots, n$ $j = 0, 1, \dots, p$

$$\begin{pmatrix} 1 & x_1(1) & \dots & x_p(1) \\ 1 & x_1(2) & & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \dots & & x_p(n) \end{pmatrix} =: X \in M_{n, p+1}$$

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} =: Y \in \mathbb{R}^n$$

$$\begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} =: \beta \in \mathbb{R}^{p+1} \quad \parallel \quad \begin{pmatrix} B_0 \\ B_1 \\ \vdots \\ B_p \end{pmatrix} =: B \in \mathbb{R}^{p+1} \quad \parallel \quad B_i \approx \beta_i$$

stimatori

• Previsi : $\hat{Y}(i) = \sum_{j=0}^p B_j x_j(i) = \sum B_j x_{ij} = \sum_i x_{ij} B_j = [XB]_i$

• Residui : $R(i) = Y(i) - \hat{Y}(i) = Y_i - [XB]_i = [Y - XB]_i$

• SS : $SS = \sum_i R(i)^2 = \|Y - XB\|^2$

• Cerchiamo i coefficienti B_i che minimizzano SS

$$0 = \frac{\partial SS}{\partial B_j} = \frac{\partial}{\partial B_j} \sum_i \left(Y(i) - \sum_{k=0}^p B_k x_k(i) \right)^2 = 2 \sum_i \left(Y(i) - \sum_{k=0}^p B_k x_k(i) \right) (-x_j(i))$$

$$0 = \sum_i \left(Y(i) x_j(i) - \sum_{k=0}^p B_k x_k(i) x_j(i) \right) = \sum_i Y(i) x_j(i) - \sum_i x_j(i) [XB]_i$$

$$= \sum_i Y_i x_{ij} - \sum_i [XB]_i x_{ij} = [Y^T X]_j - [(XB)^T X]_j = [Y^T X - B^T X^T X]_j \quad \forall j$$

quindi $Y^T X - B^T X^T X = 0$ traspongo : $X^T Y - X^T X B = 0$

• Equazioni normali : $X^T Y = X^T X B$

• Ricavo B : $(X^T X)^{-1} X^T Y = B$

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

• Funzioni da usare su Excel

=MATR.PRODOTTO (matrice 1 ; matrice 2)

=MMULT(...)

=MATR.TRASPOSTA (matrice)

=TRANSPOSE(...)

=MATR.INVERSA (matrice)

=MINVERSE(...)

=MATR.SOMMAPRODOTTO (vett 1 ; vett 2)

prodotto scalare

=SUMPRODUCT(...)

• Distribuzione coefficienti

$$1) B_j = [(X^T X)^{-1} X^T Y]_j = \sum_i [(X^T X)^{-1} X^T]_{ji} Y_i \quad \text{è comb. lineare delle } Y_i$$

Quindi B_j ha distribuzione ^{$p+1 \times n$} normale

$$2) E(B_j) = \dots = \beta_j \quad \text{stimatore corretto}$$

$$3) \text{Var}(B_j) = \dots$$

è più facile e più utile : $\text{Cov}(B_j; B_k) =: \Sigma_{jk}$

Σ 'Sigma' matrice di covarianza $\in M_{p+1, p+1}$

$$\Sigma = \dots = \sigma^2 (X^T X)^{-1}$$

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

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

$$\text{Cov}(B; B) =$$

Regr. multilineare \rightarrow matrici

var di ingresso

$$X = \begin{pmatrix} 1 & x_1(1) & x_2(1) & \dots & x_p(1) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1(n) & x_2(n) & \dots & x_p(n) \end{pmatrix} \in M_{n,p+1}$$

$$x_0(i) = 1 \quad \forall i = 1, 2, \dots, n$$

$$\beta = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_p \end{pmatrix} \in \mathbb{R}^{p+1}$$

$$B = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_p \end{pmatrix} \in \mathbb{R}^{p+1}$$

var di risposta

$$Y = \begin{pmatrix} Y(1) \\ \vdots \\ Y(n) \end{pmatrix} \in \mathbb{R}^n$$

$$(X^T X)^{-1} \in M_{p+1,p+1}$$

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

• Distribuzione di $B = (\beta_0, \beta_1, \dots, \beta_p)$ è un vettore \rightarrow distrib. congiunta

Per fortuna non occorre scrivere $f_B(t_0, t_1, \dots, t_p) = ???$

Infatti B ha distribuzione gaussiana multidimensionale ...

... basta conoscere quali sono i parametri che la identificano :

\rightarrow la media di ciascuna componente $E(\beta_i), i = 0, 1, \dots, p$

\rightarrow la varianza di ciascuna componente $\text{Var}(\beta_i), i = 0, \dots, p$

\rightarrow la covarianza di ogni coppia di comp. $\text{Cov}(\beta_i, \beta_j), i, j = 0, \dots, p$

★ Ovviamente $\text{Var}(\beta_i) = \text{Cov}(\beta_i, \beta_i)$ quindi la \mathbb{II} è inclusa nella \mathbb{III}

a) $B = NY \quad N := (X^T X)^{-1} X^T \in M_{p+1,n}$

B è una trasf. lineare di Y , quindi è normale.

b) $E(B) = E(NY) = NE(Y) = NX\beta = (X^T X)^{-1} X^T X \beta = \beta$

infatti :

stimatore corretto

$$Y(i) \sim \mathcal{N}(\beta_0 + \beta_1 x_1(i) + \beta_2 x_2(i) + \dots + \beta_p x_p(i); \sigma^2) \sim \mathcal{N}(X\beta; \sigma^2)$$

$$E(Y(i)) = [X\beta]_i \Rightarrow E(Y) = X\beta$$

c) $\Sigma := \text{Cov}(B) \in M_{p+1, p+1}$

matrice di covarianza di B

$$\Sigma_{ij} := \text{Cov}(B_i; B_j)$$

$$\text{Cov}(B; B) = \text{Cov}(NY; NY) = ?$$

Più in generale: X, Y vettori aleatori con $S = \text{Cov}(X; Y)$
matrice di covarianza; M, N matrici:

$$[\text{Cov}(MX; NY)]_{ij} := \text{Cov}([MX]_i; [NY]_j) = \text{Cov}\left(\sum_k M_{ik} X_k; \sum_l N_{jl} Y_l\right)$$

$$= \sum_k \sum_l M_{ik} N_{jl} \text{Cov}(X_k; Y_l) =: \sum_k \sum_l M_{ik} N_{jl} S_{kl} = \sum_{kl} M_{ik} \cdot S_{kl} \cdot N_{jl} =$$

$$=: [MSN^T]_{ij}$$

$$\boxed{\text{Cov}(MX; NY) = M \text{Cov}(X; Y) N^T}$$

$$\text{Cov}(B; B) = \text{Cov}(NY; NY) = N \text{Cov}(Y; Y) N^T = N \sigma^2 I N^T = \sigma^2 N N^T$$

$$\text{Cov}(Y_i; Y_j) = \begin{cases} 0 & i \neq j \\ \sigma^2 & i = j \end{cases}$$

$$N N^T = (X^T X)^{-1} X^T X (X^T X)^{-1} = (X^T X)^{-1}$$

$$\boxed{\text{Cov}(B; B) = \sigma^2 (X^T X)^{-1}}$$

• Serve uno stimatore di σ^2

$$S_e^2 := \frac{SS}{n - (p+1)} \approx \sigma^2$$

stimatore consistente e corretto di σ^2

$$\rightarrow \frac{SS}{\sigma^2} \sim \chi^2(n - p - 1) \quad \text{indipendente da } B$$

HW: Verificate l'applicabilità del thm del chi-quadro

INFERENZA SUI β_k

$$B_k \sim \mathcal{N}\left(\beta_k; \sigma^2 \cdot (X^T X)^{-1}_{k,k}\right)$$

$$\Rightarrow \frac{B_k - \beta_k}{\sigma \sqrt{(X^T X)^{-1}_{k,k}}} \sim \mathcal{N}(0, 1) \quad \text{se } \sigma \text{ incognita...}$$

$$\Rightarrow \frac{B_k - \beta_k}{S_e \sqrt{(X^T X)^{-1}_{k,k}}} \sim t(n-p-1) \quad \text{funzione ausiliare}$$

$\rightarrow H_0: \beta_k = 0 \quad H_1: \beta_k \neq 0$ test per vedere se $\beta_k = 0$

statistica $T_k := \frac{B_k}{S_e \sqrt{(X^T X)^{-1}_{k,k}}} \stackrel{H_0}{\sim} t(n-p-1)$

• Dopo la regressione preliminare, andrebbero eliminate le variabili per cui risulta plausibile $H_0: \beta_k = 0$.

In effetti tali variabili non vanno mai rimosse tutte assieme: se ne toglie una e si rifà la regressione: può essere che a questo punto una o più delle altre risultino in $H_1: \beta_k \neq 0$.

ora 13

INFERENZA risposta media e futura

$$p=1 \quad x_0 \rightarrow Y(x_0) \sim \mathcal{N}(\alpha + \beta x_0, \sigma^2)$$

$\alpha + \beta x_0$??

$$p \geq 1 \quad \tilde{x} \in \mathbb{R}^{p+1}, \tilde{x} = (1, \tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_p) \rightarrow Y(\tilde{x}) \sim \mathcal{N}(\tilde{x} \cdot \beta; \sigma^2)$$

$$\tilde{x} \cdot \beta = \beta_0 + \beta_1 \tilde{x}_1 + \dots + \beta_p \tilde{x}_p \quad ??$$

\rightarrow Stimatore puntuale $\tilde{x} \cdot B$

\rightarrow Che distribuzione ha? Certamente ha legge normale (scalare)

\rightarrow Parametri: $E(\tilde{x} \cdot B) = \tilde{x} \cdot E(B) = \tilde{x} \cdot \beta$

$\tilde{x} \cdot B$ è uno stimatore corretto

$$\text{Var}(\tilde{x} \cdot B) = \text{Cov}(\tilde{x} \cdot B; \tilde{x} \cdot B) = \text{Cov}(MB; MB) = \textcircled{*}$$

$$\text{dove } M = \tilde{x}^T = (1, \tilde{x}_1, \dots, \tilde{x}_p) \in M_{1, p+1}$$

$$\textcircled{*} = M \text{Cov}(B; B) M^T = \tilde{x}^T \Sigma \tilde{x} = \tilde{x} \cdot \Sigma \tilde{x} = \langle \tilde{x}, \Sigma \tilde{x} \rangle$$

$$\text{dove } \Sigma = \sigma^2 (X^T X)^{-1}$$

$$\text{Var}(\tilde{x} \cdot B) = \sigma^2 \tilde{x}^T (X^T X)^{-1} \tilde{x}$$

→ funz. ausiliaria:

$$\frac{\tilde{x} \cdot B - \tilde{x} \cdot \beta}{\sigma \sqrt{\tilde{x}^T (X^T X)^{-1} \tilde{x}}} \sim \mathcal{N}(0, 1) \Rightarrow \frac{\tilde{x} \cdot B - \tilde{x} \cdot \beta}{S_e \sqrt{\tilde{x}^T (X^T X)^{-1} \tilde{x}}} \sim t(n-p-1)$$

→ int. di conf.

$$\tilde{x} \cdot \beta \in \tilde{x} \cdot B \pm q S_e \sqrt{\tilde{x}^T (X^T X)^{-1} \tilde{x}}$$

↑ quantile $t(n-p-1)$ con lvl di conf prescelto

• Risposta futura:

\tilde{x} livello di ingresso $\tilde{x} = (1, \tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_p)$

$Y(1), Y(2), \dots, Y(n)$ con ingressi $x(i)$ diversi

$\tilde{x} \rightarrow Y(n+1)$ risposta per \tilde{x}

$Y(n+1) \sim \mathcal{N}(\tilde{x} \cdot \beta; \sigma^2)$ indipendente da Y e quindi da B

$$Y(n+1) - \tilde{x} \cdot B \sim \mathcal{N}(0; \sigma^2 + \sigma^2 \tilde{x}^T (X^T X)^{-1} \tilde{x})$$

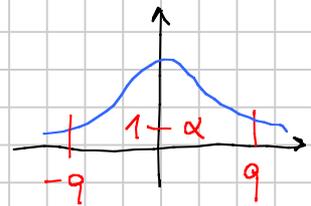
$$\frac{Y(n+1) - \tilde{x} \cdot B}{\sigma \sqrt{1 + \tilde{x}^T (X^T X)^{-1} \tilde{x}}} \sim \mathcal{N}(0, 1) \Rightarrow \frac{Y(n+1) - \tilde{x} \cdot B}{S_e \sqrt{1 + \tilde{x}^T (X^T X)^{-1} \tilde{x}}} \sim t(n-p-1)$$

(Non si tratta di una funz. ausiliaria)

$1-\alpha$ prob. fissata dell'intervallo di predizione

q quantile $t(n-p-1)$ bilaterale corrispondente a lvl conf $1-\alpha$

$$q = F_{t(n-p-1)}^{-1} \left(1 - \frac{\alpha}{2} \right) = \text{INV.T}(\alpha; n-p-1; 2)$$



$$1 - \alpha = P(-q \leq t(n-p-1) \leq q) =$$

$$= P\left(-q \leq \frac{Y(n+1) - \vec{x} \cdot \beta}{S_e \sqrt{1 + \vec{x}^T (X^T X)^{-1} \vec{x}}} \leq q\right) = \dots = P\left(Y(n+1) \in \vec{x} \cdot \beta \pm q S_e \sqrt{1 + \vec{x}^T (X^T X)^{-1} \vec{x}}\right)$$

$$Y(n+1) \in \vec{x} \cdot \beta \pm q S_e \sqrt{1 + \vec{x}^T (X^T X)^{-1} \vec{x}} \quad \text{intervallo di predizione}$$

↑
v.a. future

↑
+1

differenza rispetto int. di conf. $\vec{x} \cdot \beta$

HW: trovare intervallo di predizione per la somma e per la media delle prossime m risposte future, tutte con lvl di insp. \vec{x}

$$\begin{array}{l} \vec{x} \rightarrow Y(n+1) \\ \quad \searrow \\ \quad \quad \rightarrow Y(n+2) \\ \quad \quad \quad \vdots \\ \quad \quad \quad \rightarrow Y(n+m) \end{array} \left. \vphantom{\begin{array}{l} \vec{x} \\ \rightarrow \\ \rightarrow \\ \rightarrow \\ \rightarrow \end{array}} \right\} V := Y(n+1) + \dots + Y(n+m)$$

$$\left. \vphantom{\begin{array}{l} \vec{x} \\ \rightarrow \\ \rightarrow \\ \rightarrow \\ \rightarrow \end{array}} \right\} U := \frac{V}{m}$$

COEFFICIENTE DI DETERMINAZIONE

$$R^2 := 1 - \frac{SS}{S_{YY}} \quad \text{frazione di devianza spiegata}$$

$$\text{dove } S_{YY} = \sum_i Y_i^2 - n \bar{Y}^2 = (n-1) S_Y^2$$

$$S_Y^2 : \text{varianza campionaria} \quad S_Y^2 = \frac{1}{n-1} \left(\sum_i Y_i^2 - n \bar{Y}^2 \right) = \frac{1}{n-1} \sum_i (Y_i - \bar{Y})^2$$

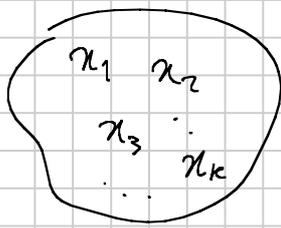
$$S_{YY} : \text{devianza campionaria} \quad S_{YY} = \sum_i (Y_i - \bar{Y})^2$$

$$SS : \text{devianza dei residui} \quad SS = \sum_i (Y_i - \hat{Y}_i)^2 \quad \hat{Y}_i = B \cdot x(i)$$

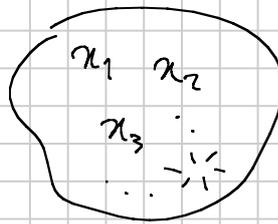
$$R^2 = 1 - \frac{SS}{S_{YY}} = 1 - \frac{(n-p-1) S_e^2}{(n-1) S_Y^2} = 1 - \frac{n-p-1}{n-1} \frac{S_e^2}{S_Y^2}$$

circa 1 se $p \ll n$ come auspicabile

- Se si effettuano due regressioni, con o senza una var di ingresso x_k il valore di R^2 è necessariamente migliore (più alto) in quella con la variabile in più



R^2 maggiore



R^2 minore

→ Si può definire un indice R_c^2 coeff di determinazione corretto che tiene conto anche di p :

$$R_c^2 = f(R^2, p) \quad R_c^2 \leq R^2$$

È fatto in modo che se si toglie una var x_k non signif. cattiva, R_c^2 non diminuisce, anzi, magari aumenta

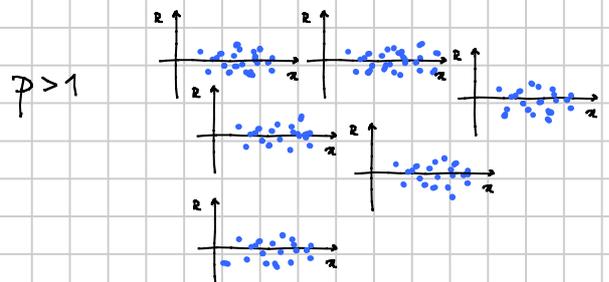
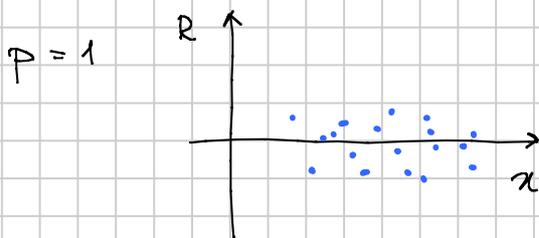
$$SS = \|Y\|^2 - B^T X^T X B = \|Y\|^2 - Y^T X B$$

ora 14

$$\|Y\|^2 = \sum_i Y(i)^2$$

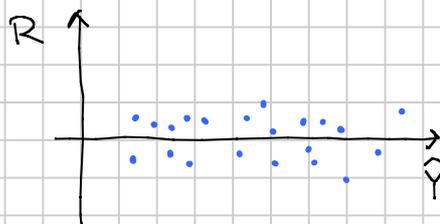
QUESTIONI PRATICHE SULLA REGRESSIONE

- Analisi dei residui



Possiamo anche farci dare il grafico previsti vs residui

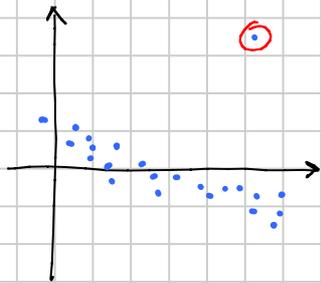
un grafico per ogni x_i
→ p grafici (tutti diversi)



$p=1$ è lo stesso grafico perché $\hat{Y} = \alpha + \beta x$ è una trasf. lineare di x | $p > 1$ non è nessuno di quelli precedenti, ma ne sintetizza alcuni aspetti

★ MINITAB ad esempio mostra di default questo grafico.

● Outliers : sono punti } → errati
 } → provenienti da un'altra distribuzione

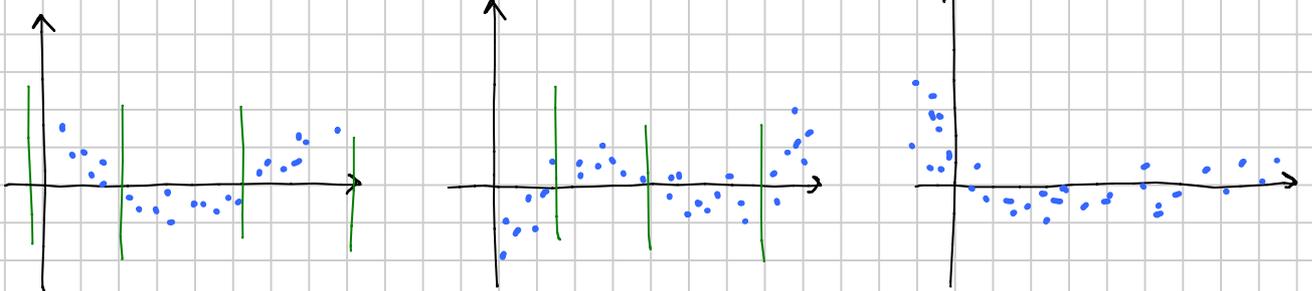


→ vanno capiti, studiati a parte e eliminati dal campione, poi si rifà

→ esempi : - stelle / seq principale
 - case di Philadelphia

● Nonlinearità : cerchiamo un andamento a parabola o simile

nei residui



→ se coinvolgono una variabile sola (o poche variabili) si corregge quella variabile, aggiungendo nel modello termini di grado più alto (vedi oltre)

→ in alternativa, soprattutto se le nonlinearità sono su tutte le var e si somigliano, si può tentare qualche trasformazione non lineare della var di risposta (Y)

■ Trasformazioni non lineari di Y

es: macchina che crea confezioni / scatole di cartone



$$V = \text{velocità} = f(x_1, x_2, \dots, x_p)$$

$$V \stackrel{?}{=} \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p \quad \text{lineare}$$

invece:

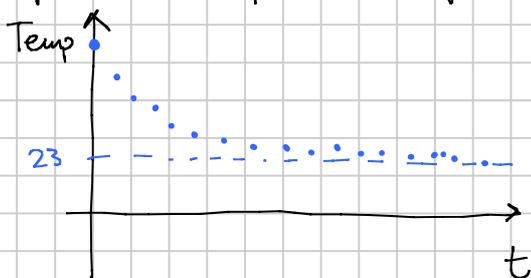
$$T = \text{tempo dedicato a ogni cartone} \stackrel{?}{=} \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

$$T = \beta_0 + \dots + \beta_p x_p + e$$

$$V = \frac{1}{T} = \frac{1}{\beta_0 + \dots + \beta_p x_p + e} \quad \text{non lineare}$$

In pratica si fa il reciproco di Y $Y \mapsto \frac{1}{Y} = T$
(una funzione non lineare)

→ esempio: tempo di raffreddamento di metallo



andamento esponenziale

$$\text{Temp} = 23 + e^{-\lambda t + \nu}$$

$$\log(\text{Temp} - 23) = \nu - \lambda t$$

$\log(\text{Temp} - 23)$ è una misura della Temp che potrebbe dipendere linearmente da t

$\text{Temp} \mapsto \log(\text{Temp} - 23)$ è una trasformazione non lineare

* Occhio all'effetto sui residui:

supponiamo che $\log(\text{Temp} - 23) = \nu - \lambda t$ sia giusto

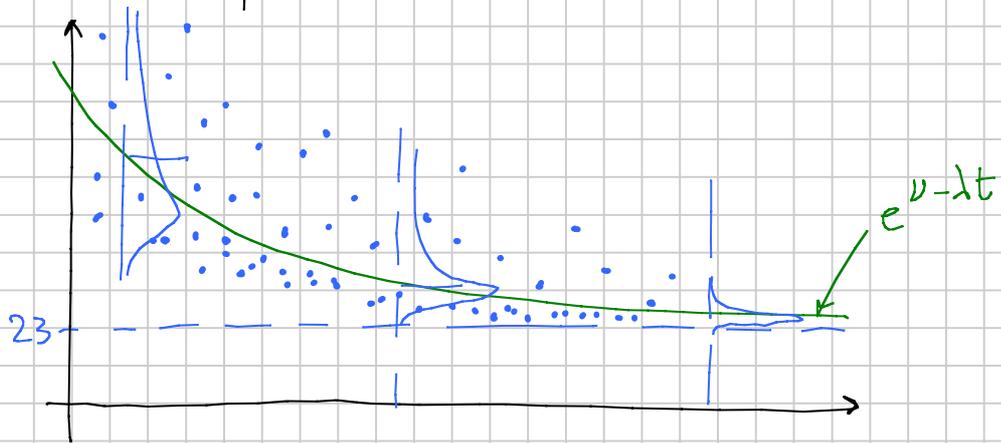
$$\log(\text{Temp} - 23) = \nu - \lambda t + e$$

Allora il modello originale è:

$$\text{Temp} = 23 + e^{\nu - \lambda t + e} \neq 23 + e^{\nu - \lambda t} + e$$

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

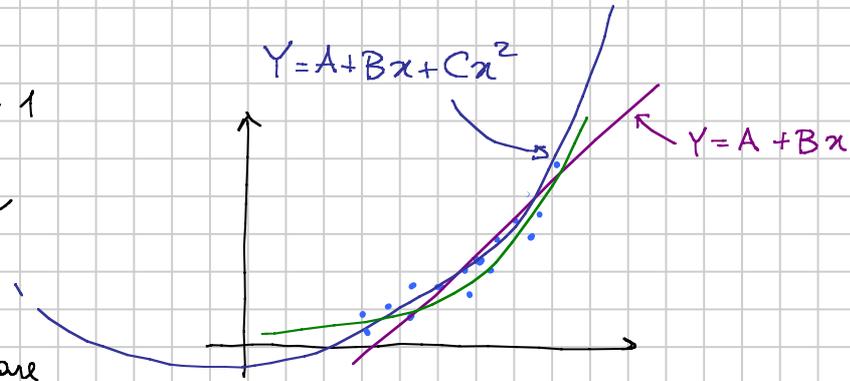
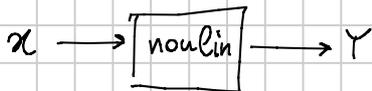
Nell'esempio



NONLINEARITÀ

- 1) Trasf. nonlineare di $Y \rightarrow$ può correggere
- 2) Aggiustare la nonlinearity di singole var di ingresso aggiungendo termini nonlineari

● Caso base : $p = 1$



Se il modello lineare

$Y = \alpha + \beta x$ non funziona, provo con polinomi di grado via via più alto, ogni volta controllando i residui

\rightarrow in pratica :

Y	x	x^2	x^3	...
$Y(1)$	$x(1)$	$x^2(1)$		
\vdots	\vdots	\vdots		
$Y(n)$	$x(n)$	$x^2(n)$		

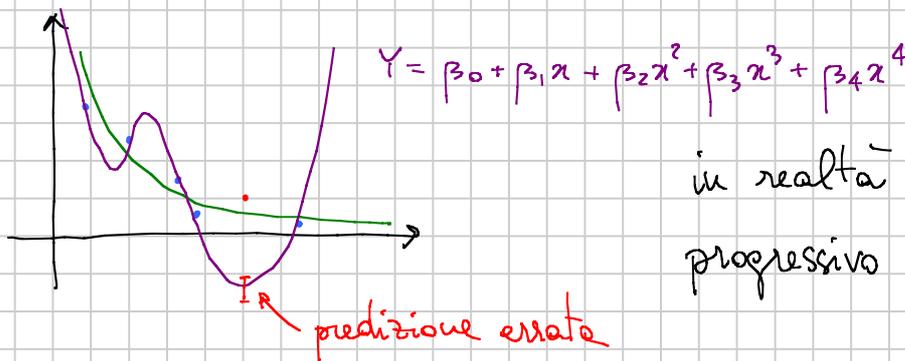
$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

cioè si aggiungono variabili fittizie (dummy)

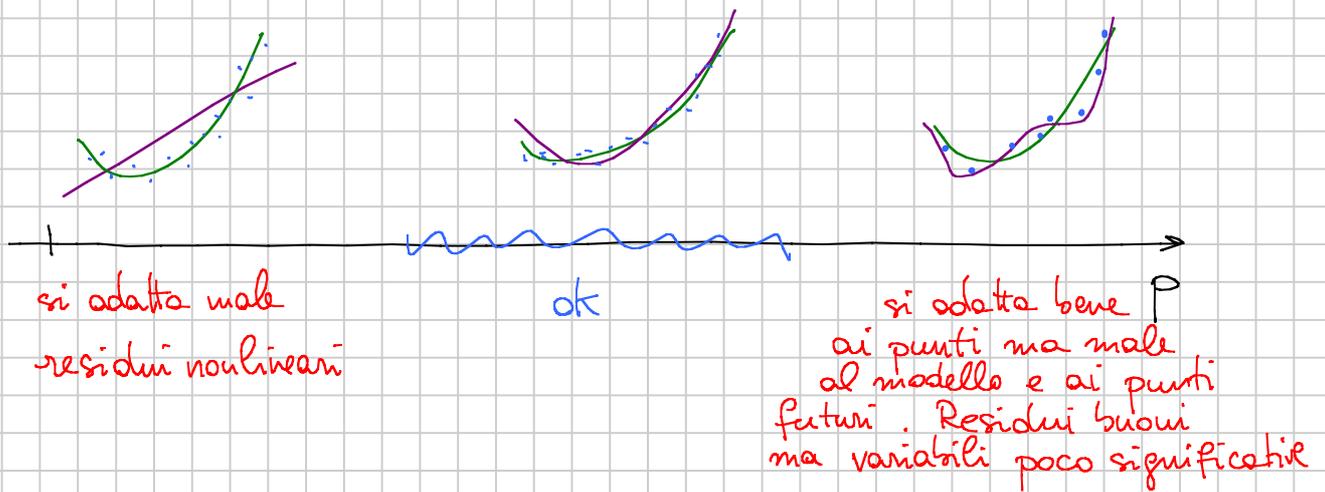
var ingresso di regr. lineare multiple

\rightarrow attenzione a non esagerare con il grado!

per n punti passa esattamente un solo polinomio di grado $n-1$



in realtà il fenomeno è progressivo : se il grado p



★ La situazione in cui p è troppo alto si chiama **overfitting**. In generale, per essere abbastanza sicuri che non si sia overfitting ci vuole $n \gg p$.

In letteratura si trovano diverse convenzioni, come $\frac{n}{p} \geq 5$, $\frac{n}{p^2} \geq 2$

★ Ci si può accorgere che p è eccessivo se le variabili risultano non significative al test $H_0: \beta_k = 0$

→ in particolare i p decisi vengono tutti alti e circa tutti uguali (spesso)

★ Inoltre maggiore è $\frac{n}{p}$, $\frac{n}{p^2}$, maggiore è la **potenza** dei test della regressione.

⊙ Nonlinearità nella regressione multiple $p > 1$
tutto simile, ma più laborioso:

→ il lavoro va fatto su ogni variabile di ingresso

$$\rightarrow Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 x_2^2 + \beta_5 x_1 x_2$$

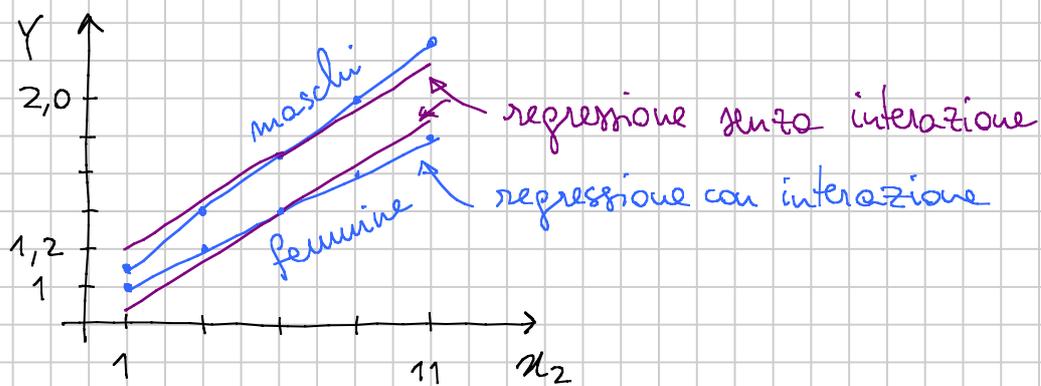
inoltre ci sono i termini di interazione!

	x_2	→ 1	3	5	7	11	— anni dalla laurea
maschi	0	1,1	1,4	1,7	2,0	2,3	
femmine	1	1	1,2	1,4	1,6	1,8	

↑
 x_1

interazione
↓

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



→ i termini di interazione sono necessari quando il modo in cui Y dipende da x_2 è legato a x_1

★

2	3	4	5	6	20
1	3	6	10	15		190

il numero di termini di interazione cresce rapidamente con p

$$\binom{p}{2} = \frac{p(p-1)}{2} \approx \frac{p^2}{2}$$

in questo modo, il "p effettivo", cioè il numero di coeff da stimare sul modello nonlineare sia troppo grosso (per n normali)

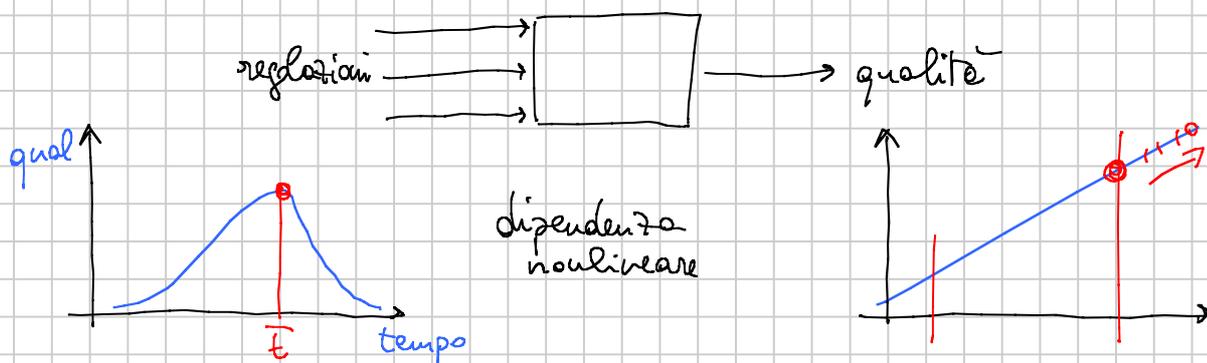
→ per questo motivo i termini di interazione da inserire nel modello vanno scelti con cura e parsimonia.

→ in ogni caso un termine inserito va testato ($H_0: \beta_k = 0$) per vedere se è significativo.

ora 16

● Regressione e ottimizzazione

Se lo scopo della regressione è massimizzare o minimizzare qualcosa...



... i modelli lineari si adattano male, perché trovano max e min sempre sul bordo.

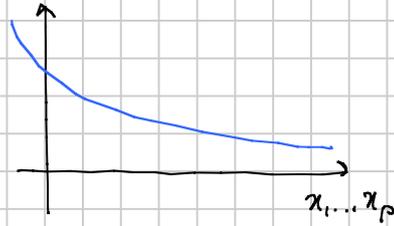
→ NB anche i termini di interazione mantengono questo problema

• Nel caso di metodo 1) : trasformazioni non lineari di Y
quali sono le trasformazioni più comuni?

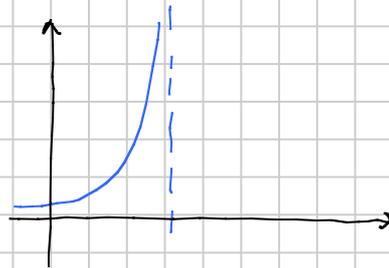
i. $\log(Y)$

$$Y = e^{\beta_0 + \beta_1 x_1 + \dots}$$

→ \sqrt{Y} , $\sqrt[3]{Y}$ sono simili

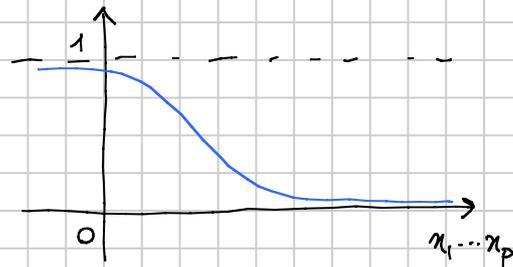


ii. Y^{-1} , Y^{-2} , ...



iii. $\log\left(\frac{Y}{1-Y}\right)$

"odds ratio"



forma sigmoide

"trasformazione logistica" = logit transform

$$Y = \frac{e^{\beta_0 + \beta_1 x_1 + \dots}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots}}$$

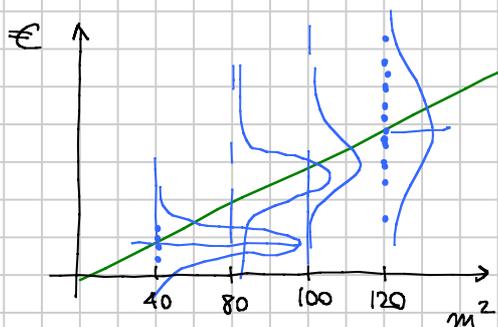
Se la sigmoide non va da 0 a 1, prima della trasformazione logistica, fare una trasformazione lineare in modo che Y vari circa tra 0 e 1.

MANCANZA DI OMOSCHEDASTICITÀ (problema 3)

Laborioso e spesso dagli effetti dubbi ...

= la varianza di $e \sim N(0, \sigma^2(x))$ dipende dalle var di ingresso (quindi anche le var dei residui)

Diagnosi



$\sigma^2 \propto x$ oppure $\sigma^2 \propto x^2$
 oppure $\sigma \propto Y$ oppure $\sigma^2 \propto Y$
 o altre dipendenze ancora ...

prezzi case

→ se la variabilità è $\pm 20\%$ del prezzo ($\Rightarrow \sigma \propto Y$)

$$Y = \beta x + e \quad \text{dev.st}(e) \approx 20\% \cdot \beta x$$

$$\Rightarrow \sigma \approx 0,2 \cdot \beta \cdot x \propto x, \text{ ma anche } \sigma \propto Y$$

(sia a x sia a Y perché $Y \propto x$)

$$Y = \alpha + \beta x + e \quad \text{dev.st}(e) \approx 20\% \cdot (\alpha + \beta x)$$

$$\Rightarrow \sigma \propto Y \text{ ma } \sigma \not\propto x$$

→ n : # incroci Y : tempo di percorrenza

$$Y = \beta x + e$$

$$Y = T_1 + T_2 + \dots + T_n \quad \text{modello mentale: il tempo complessivo}$$

è somma di tanti contributi indipendenti, uno per incrocio

$$\text{se } \text{Var}(T_i) = a^2 \text{ circa uguali } \Rightarrow \text{Var}(Y) = n a^2$$

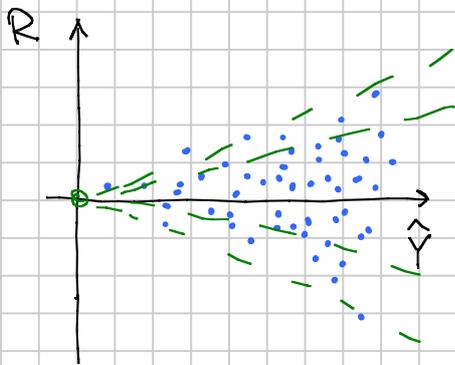
(si sommano per l'indipendenza)

$$\underbrace{T_1 + T_1 + T_1 + \dots + T_1}_n = n T_1 \Rightarrow \text{Var}(n T_1) = n^2 a^2$$

$$\sigma = \sqrt{\text{Var}(Y)} = \sqrt{n} a \propto \sqrt{n} \Rightarrow \sigma^2 \propto n \Rightarrow \sigma^2 \propto Y$$

$$Y = \alpha + \beta x + e \quad \Rightarrow \quad \boxed{\sigma^2 \propto x} \quad \boxed{\sigma^2 \propto Y}$$

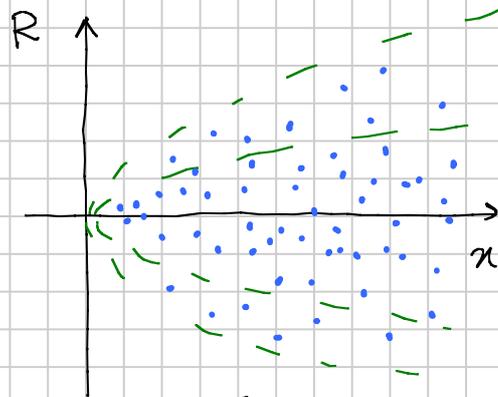
● Più frequentemente non abbiamo idea di motivi teorici per supporre una certa dipendenza ... comunque è bene guardare i residui



$$\sigma \propto Y$$

$$\sigma \propto x$$

stesso grafico x in ascissa

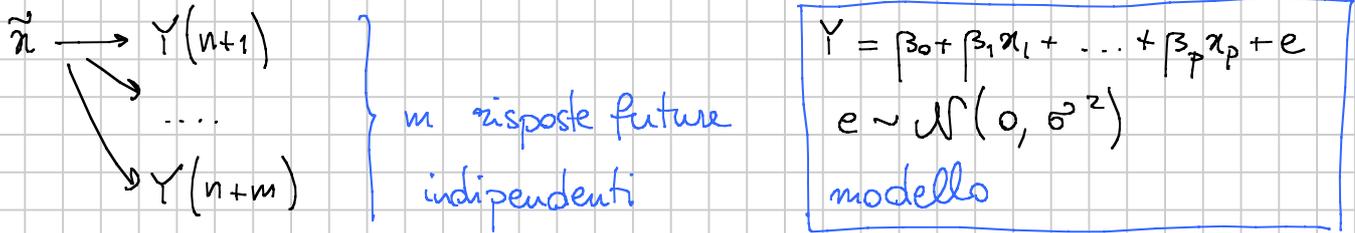


$$\sigma^2 \propto x \Rightarrow \sigma \propto \sqrt{x}$$

$$\sigma^2 \propto Y$$

stesso grafico Y in ascissa

laboratorio dalle sett prox spostato al mar 9.30 - 11.30



$$Y(n+i) \sim \mathcal{N}(\tilde{x} \cdot \beta; \sigma^2) \quad \forall i = 1, 2, \dots, m$$

$$V = \sum_{i=1}^m Y(n+i) \text{ la somma} \quad V \sim \mathcal{N}(m \tilde{x} \cdot \beta; m \sigma^2)$$

$$V - m \tilde{x} \cdot B \sim \mathcal{N}(0; m \sigma^2 + m^2 \underbrace{\sigma^2 \tilde{x}^T (X^T X)^{-1} \tilde{x}}_{\text{Var}(\tilde{x} \cdot B)})$$

$$\tilde{x} \cdot B \sim \mathcal{N}(\tilde{x} \cdot \beta; \underbrace{\tilde{x}^T \sigma^2 (X^T X)^{-1} \tilde{x}}_{\text{matrice di covarianze di B}})$$

$$\frac{V - m \tilde{x} \cdot B}{\sigma \sqrt{m + m^2 \tilde{x}^T (X^T X)^{-1} \tilde{x}}} \sim \mathcal{N}(0, 1) \quad \Rightarrow \quad \frac{V - m \tilde{x} \cdot B}{S_e \sqrt{m + m^2 \tilde{x}^T (X^T X)^{-1} \tilde{x}}} \sim t(n-p-1)$$

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

Concludendo, se q è un quantile $t(n-p-1)$ con prob $1-\alpha$ (bilaterale):

$$V \in m \tilde{x} \cdot B \pm q S_e \sqrt{m + m^2 \tilde{x}^T (X^T X)^{-1} \tilde{x}} \quad \text{con prob } 1-\alpha$$

Analogamente per $U = \frac{V}{m}$ (media campionaria delle risposte future):

$$U \in \tilde{x} \cdot B \pm q S_e \sqrt{\frac{1}{m} + \tilde{x}^T (X^T X)^{-1} \tilde{x}} \quad \text{con prob } 1-\alpha$$

● Riguardo al lab 05

$$Y = f(\alpha_1, \alpha_2) + e$$

\uparrow \uparrow $c=1, 2$ canale
 $t = \text{temperatura}$

f nonlineare (di II grado)

Solo canale 1 : $Y = \beta_0 + \beta_1 t + \beta_2 t^2 + e$
 B_0, B_1, B_2 stimatori

Solo canale 2 : $\dots C_0, C_1, C_2$

Multiplo $Y = \alpha_0 + \alpha_1 t + \alpha_2 c + \alpha_3 t^2 + \alpha_4 ct + \alpha_5 ct^2$
 A_0, A_1, \dots, A_5 stimatori

Provo a considerare il modello multiplo per $c=1$

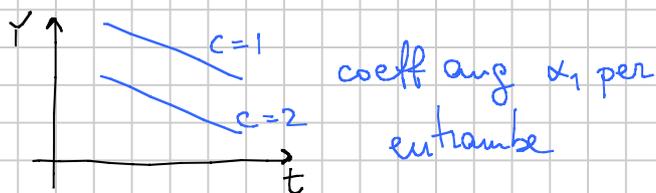
$$Y = \alpha_0 + \alpha_1 t + \alpha_2 + \alpha_3 t^2 + \alpha_4 t + \alpha_5 t^2$$

$$= (\alpha_0 + \alpha_2) + (\alpha_1 + \alpha_4)t + (\alpha_3 + \alpha_5)t^2$$

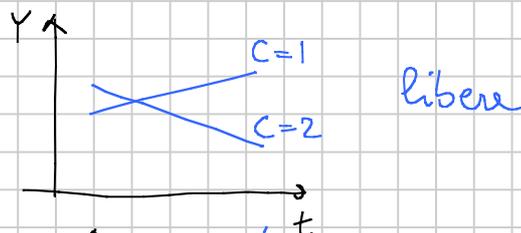
Si può verificare che $A_0 + A_2 = B_0$, $A_1 + A_4 = B_1$, $A_3 + A_5 = B_2$
 $A_0 + 2A_2 = C_0$, $A_1 + 2A_4 = C_1$, $A_3 + 2A_5 = C_2$



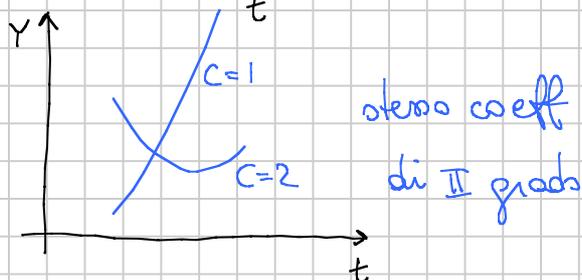
$$Y = \alpha_0 + \alpha_1 t + \alpha_2 c$$



$$Y = \alpha_0 + \alpha_1 t + \alpha_2 c + \alpha_3 ct$$



$$Y = \alpha_0 + \alpha_1 t + \alpha_2 c + \alpha_3 t^2 + \alpha_4 ct$$



Se si aggiunge $\alpha_5 ct^2$

si possono ottenere tutte le coppie di parabole possibili

★ L'unica differenza operativa tra il metodo in 2 casi e quello con regr. multiple è che nel primo stimo due distinti valori di σ , nel secondo uno solo: se i due gruppi (canali nell'esempio) mostrano varianze diverse, devo separare, se no posso tenere assieme.

(Tenerele assieme è leggermente più prudente.)

REGRESSIONE NON OMOSCHEDASTICA: CURA

$$Y = \beta \cdot x + e \quad e \sim \mathcal{N}(0, \sigma^2(x)) \quad \text{ad es } \sigma \propto x$$

$w(i)$ $i=1, 2, \dots, n$ dei pesi che vanno presi *inversamente proporzionali a σ^2*

$$w(i) \propto \frac{1}{\sigma^2(x(i))} \quad \text{nell'es } w(i) \propto \frac{1}{x(i)^2}$$

- Si rifà daccapo la regressione, minimizzando $SS^{(w)}$

$$SS^{(w)} := \sum_{i=1}^n (Y(i) - x(i) \cdot \beta)^2 \cdot w(i)$$

infatti, se prima minimizzavo $SS := \sum_{i=1}^n (Y(i) - x(i) \cdot \beta)^2$

$$e \quad Y(i) - x(i) \cdot \beta \sim \mathcal{N}(0, \sigma^2) \quad , \quad \text{ma}$$

$$Y(i) - x(i) \cdot \beta \sim \mathcal{N}(0, \sigma(i)^2) \quad \text{ma}$$

$$(Y(i) - x(i) \cdot \beta) \sqrt{w(i)} \sim \mathcal{N}(0, \underbrace{w(i) \sigma(i)^2}_{\text{costante}})$$

- In pratica, si "distribuisce" \sqrt{w} nell'equazione

$$SS^{(w)} := \sum_{i=1}^n (\sqrt{w(i)} Y(i) - \sqrt{w(i)} x(i) \cdot \beta)^2$$

ovvero si modificano X e Y moltiplicando ogni riga per il corrispondente $w(i)$ e si fa una normale regressione.

Y	x_0	x_1	\dots	x_p	Y'	x_0'	x_1'	\dots	x_p'
$Y(1)$	1	$x_{1(1)}$		$x_{p(1)}$	$Y(1)\sqrt{w(1)}$	$\sqrt{w(1)}$	$x_{1(1)}\sqrt{w(1)}$		$x_{p(1)}\sqrt{w(1)}$
$Y(2)$	1	$x_{1(2)}$	\dots	\vdots	$Y(2)\sqrt{w(2)}$	$\sqrt{w(2)}$	$x_{1(2)}\sqrt{w(2)}$	\dots	\vdots
\vdots	\vdots	\vdots		\vdots	\vdots	\vdots	\vdots		\vdots

↑ termine noto

var di ingresso

niente termine noto!

- ★ Lanciando la regressione su Y', X' , si trovano i coefficienti $B^{(w)}$ che minimizzano $SS^{(w)}$
- ★ Se servono interv di conf o di prediz, si fanno su Y', X' , poi si "torna" ai valori Y, X dividendo per $\sqrt{w(\bar{x})}$
- ★ Il coeff di determinazione R^2 del modello Y', X' non è confrontabile con quello di Y, X .
- ★ L'unica possibile verifica della bontà del procedimento è vedere se i residui di Y', X' ora sembrano omoschedastici.

FINE REGRESSIONE

ANALISI DELLA VARIANZA (ANOVA) (Cap 10)

● Un breve excursus su due test del capitolo 8

1] Test per il confronto delle medie di due popolazioni normali:

$$\begin{array}{l}
 X_1, X_2, \dots, X_m \sim \mathcal{N}(\mu_1, \sigma_1^2) \\
 Y_1, \dots, Y_n \sim \mathcal{N}(\mu_2, \sigma_2^2)
 \end{array}
 \quad
 \begin{array}{l}
 \bar{X} \approx \mu_1 \quad \bar{X} \sim \mathcal{N}(\mu_1; \frac{\sigma_1^2}{m}) \\
 \bar{Y} \approx \mu_2 \quad \bar{Y} \sim \mathcal{N}(\mu_2; \frac{\sigma_2^2}{n})
 \end{array}$$

$$H_0: \mu_1 = \mu_2 \quad H_1: \mu_1 \neq \mu_2 \quad \Rightarrow \quad H_0: \mu_1 - \mu_2 = 0 \quad H_1: \mu_1 - \mu_2 \neq 0$$

$$\bar{X} - \bar{Y} \approx \mu_1 - \mu_2 \quad \bar{X} - \bar{Y} \sim \mathcal{N}(\mu_1 - \mu_2; \frac{\sigma_1^2}{m} + \frac{\sigma_2^2}{n})$$

$$\frac{\bar{X} - \bar{Y} - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{m} + \frac{\sigma_2^2}{n}}} \sim \mathcal{N}(0,1)$$

se σ_1 e σ_2 note, questa è la funz. ausiliare

In alternativa, se $\sigma_1 = \sigma = \sigma_2$ sono uguali (ipotesi di omoschedasticità)

$$\frac{\bar{X} - \bar{Y} - (\mu_1 - \mu_2)}{\sigma \sqrt{\frac{1}{m} + \frac{1}{n}}} \sim \mathcal{N}(0,1) \quad \Rightarrow \quad \frac{\bar{X} - \bar{Y} - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{m} + \frac{1}{n}}} \sim t(m+n-2)$$

dove S_p^2 è lo stimatore pooled della varianza

$$S_p^2 = \frac{m-1}{m+n-2} S_x^2 + \frac{n-1}{m+n-2} S_y^2$$

media pesata di S_x^2 e S_y^2

coeff. positivi
a somma 1

$$\frac{S_p^2}{\sigma^2} (m+n-2) \sim \chi^2(m+n-2)$$

2] Test per il confronto delle varianze di due popolazioni normal:

$$X_1, X_2, \dots, X_m \sim \mathcal{N}(\mu_1, \sigma_1^2)$$

$$S_x^2 \approx \sigma_1^2$$

$$\frac{S_x^2}{\sigma_1^2} (m-1) \sim \chi^2(m-1)$$

$$Y_1, \dots, Y_n \sim \mathcal{N}(\mu_2, \sigma_2^2)$$

$$S_y^2 \approx \sigma_2^2$$

$$\frac{S_y^2}{\sigma_2^2} (n-1) \sim \chi^2(n-1)$$

$$H_0: \sigma_1^2 = \sigma_2^2$$

$$H_1: \sigma_1^2 \neq \sigma_2^2$$



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

$$H_1: \frac{\sigma_1^2}{\sigma_2^2} \neq 1$$

$$\frac{S_x^2}{S_y^2} \approx \frac{\sigma_1^2}{\sigma_2^2}$$

Def: Se W_1 e W_2 sono due v.a. chi-quadro, con a e b g.d.l. rispettivamente, e tra loro indipendenti, il loro rapporto ha (a meno di costanti) distribuzione **F di Fisher**:

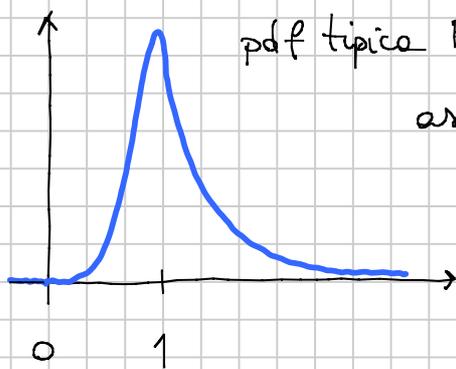
$$W_1 \sim \chi^2(a)$$

$$W_2 \sim \chi^2(b)$$

$$\frac{W_1/a}{W_2/b} \sim F(a; b)$$

"g.d.l. del denom."

"g.d.l. del numeratore"



pdf tipica F di Fisher
assume valori circa attorno a 1

$$\frac{S_x^2}{S_y^2} \approx \frac{\sigma_1^2}{\sigma_2^2} \quad W_1 = \frac{S_x^2}{\sigma_1^2} (m-1) \sim \chi^2(m-1) \quad W_2 = \frac{S_y^2}{\sigma_2^2} (n-1) \sim \chi^2(n-1)$$

$$\frac{W_1/a}{W_2/b} = \frac{W_1/(m-1)}{W_2/(n-1)} = \frac{S_x^2/\sigma_1^2}{S_y^2/\sigma_2^2} = \frac{S_x^2}{S_y^2} \cdot \frac{\sigma_2^2}{\sigma_1^2} \sim F(m-1; n-1)$$

$$\frac{\frac{S_x^2}{S_y^2}}{\frac{\sigma_1^2}{\sigma_2^2}} \sim F(m-1; n-1)$$

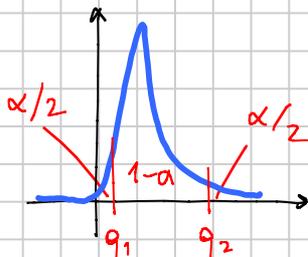
è una f. ausiliaria per $\frac{\sigma_1^2}{\sigma_2^2}$

→ La statistica del test si ottiene sostituendo 1 a $\frac{\sigma_1^2}{\sigma_2^2}$

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

$$V = \frac{S_x^2}{S_y^2} \stackrel{H_0}{\sim} F(m-1; n-1) \quad \text{statistica del test}$$

→ livello di sign α (test bilaterale)



$$q_1 = \text{INV.F}(1 - \alpha/2; m-1; n-1)$$

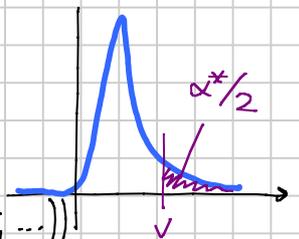
Hw: verificare i valori confrontando con le tabelle o gli esempi sul libro

$$q_2 = \text{INV.F}(\alpha/2; m-1; n-1)$$

... si guarda dove cade V

→ p dei dati

$$\alpha^* = 2 \min(\text{DISTRIB.F}(V; \dots); 1 - \text{DISTRIB.F}(V; \dots))$$



* Questo test è prezioso come preliminare a [1] per testare l'ipotesi di omoschedasticità

* Tutta l'ANOVA si basa su test simili a questo

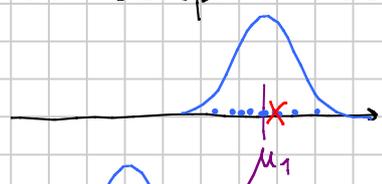
ANOVA A 1 VIA

Generalizzo il test [1] (nel caso omoschedastico) quando i campioni da confrontare sono un numero qualsiasi.

→ Nel caso i campioni siano 2, [1] e ANOVA danno sempre lo stesso risultato (α^* coincidente)

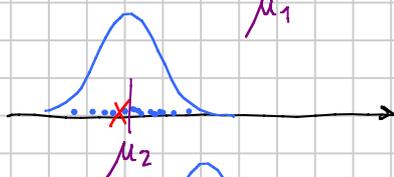
m campioni normali omoschedastici di medie $\mu_1, \mu_2, \dots, \mu_m$

ord 19



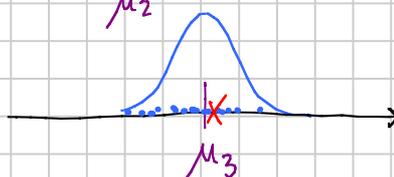
$X_{1,1}, X_{1,2}, \dots, X_{1,n_1}$

$\mathcal{N}(\mu_1, \sigma^2)$



$X_{2,1}, X_{2,2}, \dots, X_{2,n_2}$

$\mathcal{N}(\mu_2, \sigma^2)$



$X_{3,1}, X_{3,2}, \dots, X_{3,n_3}$

$\mathcal{N}(\mu_3, \sigma^2)$

.....

$X_{m,1}, \dots, X_{m,n_m}$

$\mathcal{N}(\mu_m, \sigma^2)$

* In questo contesto posso stimare σ^2 come facevo in [1]

$S_1^2, S_2^2, \dots, S_m^2$ varianze campionarie

Sono tutti stimatori indipendenti di σ^2 .

Per ciascuno di essi vale: $\frac{S_k^2}{\sigma^2} (n_k - 1) \sim \chi^2(n_k - 1)$

Def Si chiama devianza campionaria il prodotto di una qualunque varianza campionaria per i propri g.d.l. Di solito si denota con SS (sum of squares)

→ Ad es. nella regressione $\frac{SS}{n-2} = S_e^2 \Leftrightarrow SS = S_e^2(n-2)$

$$SS_1, SS_2, \dots, SS_m \quad SS_k := S_k^2(n_k - 1)$$

perciò $\frac{SS_k}{\sigma^2} \sim \chi^2(n_k - 1)$

Quindi: $\frac{SS_1 + SS_2 + \dots + SS_m}{\sigma^2} \sim \chi^2\left(\sum_{k=1}^m (n_k - 1)\right)$

Ad esempio nel caso $m=2$ $\frac{SS_1 + SS_2}{\sigma^2} \sim \chi^2(n_1 + n_2 - 2)$

$$SS_1 = (n_1 - 1) S_1^2 \quad SS_2 = (n_2 - 1) S_2^2$$

$$\frac{SS_1 + SS_2}{n_1 + n_2 - 2} = \frac{n_1 - 1}{n_1 + n_2 - 2} S_1^2 + \frac{n_2 - 1}{n_1 + n_2 - 2} S_2^2 = S_p^2$$

$$SS_p = SS_1 + SS_2 \quad \text{dove } SS_p = (n_1 + n_2 - 2) S_p^2$$

★ In generale per $m \geq 2$ il termine *pooled* viene sostituito con *within* (= entro i campioni ... poi ci sarà quella tra i campioni)

$$SS_w = SS_1 + SS_2 + \dots + SS_m \quad \text{devianza within}$$

$$\frac{SS_w}{\sigma^2} \sim \chi^2\left(\sum_k (n_k - 1)\right) \quad \text{perciò } \frac{SS_w}{\sum_k (n_k - 1)} \approx \sigma^2 \quad \text{stimatore corretto}$$

$$S_w^2 = \frac{SS_w}{\sum_k (n_k - 1)} \quad \text{varianza within} \quad S_w^2 \approx \sigma^2$$

$$S_w^2 = \sum_j \frac{SS_j}{\sum_k (n_k - 1)} = \sum_j \underbrace{\frac{n_j - 1}{\sum_k (n_k - 1)}}_{\text{coeff positivi a somma 1}} S_j^2 \quad \text{medie pesate delle } S_j^2$$

★ Per semplicità di notazione supponiamo che tutti i campioni abbiano la stessa numerosità n : $n_k = n \quad \forall k$

$$S_W^2 = \frac{SS_W}{mn - m}$$

* Sia $X_{k,*} = \bar{X}_k = \frac{1}{n} \sum_{i=1}^n X_{k,i}$ la media del k-esimo campione

Se i campioni hanno tutti la stessa media (vera) ...

$$\mu_1 = \mu_2 = \dots = \mu_m = \mu$$

$$X_{k,*} \sim \mathcal{N}\left(\mu_k, \frac{\sigma^2}{n}\right) \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

$X_{1,*}, X_{2,*}, \dots, X_{m,*}$ sono un campione iid di v.a. $\mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$

Sia S_M^2 la var campionaria di questo campione

$$S_M^2 = \frac{1}{m-1} \sum_{k=1}^m (X_{k,*} - X_{**})^2 \quad S_M^2 \approx \frac{\sigma^2}{n}$$

↑ media campionaria degli $X_{k,*}$

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

$$S_B^2 := n S_M^2 \quad \text{varianza between}$$

$$\frac{S_B^2}{\sigma^2} (m-1) \sim \chi^2(m-1) \quad S_B^2 \approx \sigma^2 \quad \text{stimatore corretto}$$

$$SS_B = (m-1) S_B^2 \quad \text{devianza between}$$

• Test dell'anova a 1 via (n_k tutti uguali)

$$H_0: \mu_1 = \mu_2 = \dots = \mu_m \quad H_1: \text{non tutte uguali}$$

→ Se è vera H_0 , S_W^2 e S_B^2 sono entrambi stimatori corretti di σ^2

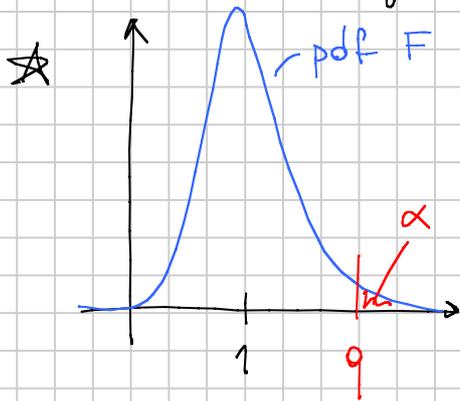
Quindi: $\frac{S_B^2}{S_W^2} \sim F(m-1; mn-m)$

infatti $\frac{S_B^2}{\sigma^2} (m-1) \sim \chi^2(m-1) \quad \frac{S_W^2}{\sigma^2} (mn-m) \sim \chi^2(mn-m)$

$$\frac{S_B^2 / \sigma^2}{S_W^2 / \sigma^2} = \frac{S_B^2}{S_W^2} \sim F(\dots)$$

→ Se è vero H_1 , S_w^2 è comunque uno stimatore corretto di σ^2 mentre S_B^2 è sempre grande (= sovrastima σ^2) perché le X_{k*} sono "sgreuate".

Perciò sotto H_1 la statistica $\frac{S_B^2}{S_w^2}$ sarà tipicamente (molto) più grande di 1



Si prende sempre un quantile unilaterale e sempre a destra

$$q = \text{INV.F}(\alpha; m-1; m \neq n - m)$$

DATE PROPOSTE PER GLI APPELLI

	Crocette & Excel	Scritto (Nicolodi)	Orale
1	19/01/2010	28/01/2010	08/02/2010
2	09/02/2010	17/02/2010	19/02/2010
3	08/06/2010	17/06/2010	28/06/2010
4	29/06/2010	08/07/2010	19/07/2010
5	01/09/2010	13/09/2010	16/09/2010
6	15/09/2010	27/09/2010	30/09/2010

ANOVA A 1 VIA

m campioni normali → varianze costante (omosched.)
 ↳ medie forse uguali o forse no (test)

$$X_{1,1}, X_{1,2}, \dots, X_{1,n_1} \sim \mathcal{N}(\mu_1, \sigma^2)$$

$$X_{i,1}, X_{i,2}, \dots, X_{i,n_i} \sim \mathcal{N}(\mu_i, \sigma^2) \quad i = 1, 2, \dots, m$$

$$n_1, n_2, \dots, n_m \text{ le numerosità} \quad N = n_1 + n_2 + \dots + n_m$$

Due stime della varianza

$$1) S_W^2 = \sum_{i=1}^m \frac{n_i - 1}{N - m} \cdot \frac{1}{n_i - 1} \sum_{j=1}^{n_i} (X_{i,j} - X_{i,*})^2$$

$$= \frac{1}{N - m} \sum_{i,j} (X_{i,j} - X_{i,*})^2$$

S_i^2 var. camp. del campione i

dove $X_{i,*} = \frac{1}{n_i} \sum_{j=1}^{n_i} X_{i,j}$

\bar{X}_i media camp. del camp. i

within
 (generalizzazione dello stimatore pooled)

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

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

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

media camp. totale

between

$$\neq \frac{1}{m} \sum_{i=1}^m X_{i,*}$$

$\frac{S_B^2}{\sigma^2} (m-1) \stackrel{H_0}{\sim} \chi^2(m-1)$ vera solo se è vera $H_0: \mu_1 = \mu_2 = \dots = \mu_m$
 se H_0 è falsa, S_B^2 in genere sovrastima σ^2

★ $x_{ij} \sim \mathcal{N}(\mu_i, \sigma^2)$ se H_0 è vera $x_{ij} \sim \mathcal{N}(\mu, \sigma^2)$

$$x_{i*} = \bar{x}_i \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n_i}\right)$$

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

$$\frac{1}{\sigma^2} \sum_{i=1}^m n_i (x_{i*} - \mu)^2 = \sum_{i=1}^m \left(\frac{x_{i*} - \mu}{\sigma/\sqrt{n_i}} \right)^2 \sim \chi^2(m) \quad \text{definizione di } \chi^2$$

Se al posto di μ metto $x_{**} \approx \mu$ per il solito t -test
 (verificare le ipotesi) $m \rightarrow m-1$

$$\frac{1}{\sigma^2} \sum_{i=1}^m n_i (x_{i*} - x_{**})^2 \sim \chi^2(m-1) \quad \text{CVD}$$

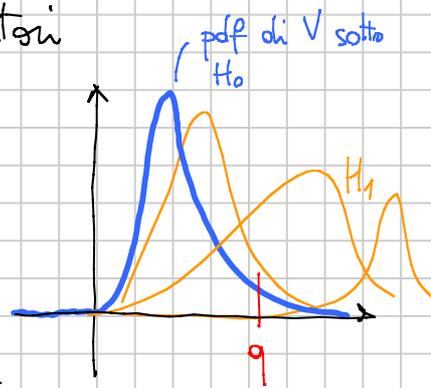
● L'ANOVA funziona confrontando i due stimatori

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

→ sotto H_0 , $V \sim F(m-1; N-m)$

→ sotto H_1 , V è in genere più grande

(quindi il test viene sempre fatto con quantili unilaterali)



● Identità algebraica delle devianze

$$S_W^2 (N-m) = SS_W$$

$$S_B^2 (m-1) = SS_B$$

$$S^2 (N-1) = SS_X$$

$$SS_X = SS_W + SS_B$$

$$S^2 = \frac{1}{N-1} \sum_{i,j} (x_{ij} - x_{**})^2$$

varianza campionaria globale

→ stima σ^2 sotto H_0 , sovrastima sotto H_1

→ Si usa spesso per ricavare $S_B^2 = \frac{1}{m-1} (SS_X - SS_W)$

$$\rightarrow \frac{SS_x}{\sigma^2} = \frac{SS_w}{\sigma^2} + \frac{SS_B}{\sigma^2} \quad (\text{dal punto di vista statistico})$$

$\xrightarrow{H_0} \chi^2(N-1)$ $\xrightarrow{H_0} \chi^2(m-1)$
 $\xrightarrow{H_0} \chi^2(N-m)$ $\xrightarrow{H_0} \chi^2(m-1)$
 sono indipendenti (H_0)

■ DUE COLLEGAMENTI CON LA REGRESSIONE

● Test globale di regressione

$$Y = \sum_{i=0}^p \beta_i x_i + e \quad x_0 = 1 \quad e \sim \mathcal{N}(0, \sigma^2)$$

C'è regressione (globalmente)?

$$H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0$$

H_1 : non tutti nulli

S_e^2 è l'analogo di S_w^2

$$SS_w = \sum_{ij} (x_{ij} - x_{ix})^2$$

$\xrightarrow{\mathcal{N}(x_i, \sigma^2)}$ $\approx x_i$

$$SS_e = SS = \sum_{i=1}^n \left[Y(i) - \underbrace{(B_0 + B_1 x_1(i) + B_2 x_2(i) + \dots + B_p x_p(i))}_{B \cdot x(i) \approx \beta \cdot x(i)} \right]^2$$

$\xrightarrow{\mathcal{N}(\beta \cdot x(i), \sigma^2)}$

S_Y^2 ha lo stesso ruolo di S^2

$$S_Y^2 = \frac{1}{n-1} \sum_{i=1}^n (Y(i) - \bar{Y})^2$$

$$SS_Y = S_Y^2 (n-1)$$

$$\frac{SS_Y}{\sigma^2} = \frac{SS_e}{\sigma^2} + ?$$

dev. totale dev. residui analogo di SS_B (dev. spiegata)

$$SS_D := SS_Y - SS_e$$

devianza spiegata

$$R^2 = 1 - \frac{SS_e}{SS_Y} = \frac{SS_Y - SS_e}{SS_Y}$$

dev. spiegata

$$\frac{SS_Y}{\sigma^2} = \frac{SS_e}{\sigma^2} + \frac{SS_D}{\sigma^2} \sim \chi^2(p)$$

$\xrightarrow{H_0} \chi^2(n-1)$ $\sim \chi^2(n-p-1)$ indipendenti (sotto H_0)

→ Posso effettuare il test per la regressione globale così:

$$V = \frac{S_D^2}{S_e^2} \text{ statistica :}$$

sotto H_0 ha legge $F(p; n-p-1)$

sotto H_1 in genere assume valori maggiori
(quindi per il test si usano quantili unilaterali)

ora 21

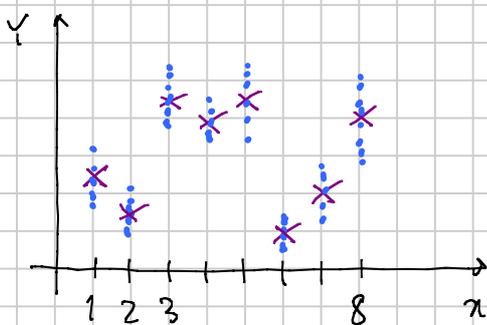
⊙ ANOVA come regressione con var di ingresso categorica.

Esempio → verbali di multe

id	data	persona	#di verbali
1	~	1	17
2	~	6	9
3	~	2	15
...	...	1	24
		8	32
		2	6
		⋮	⋮

Q: il # verbali (Y) dipende dalla persona (x)?

Si tratta di due var numeriche, quindi faccio la regressione... o no?



siccome 1, 2, ..., 8 (in questo ordine) è una codifica arbitraria la regressione non viene

★ la var "persona" infatti è una **variabile categorica**

1) Modo migliore: ANOVA a UNA VIA

si riorganizzano i dati:

1	2	3	4	5	6	7	8
17	15	-	-	-	9	14	32
24	6	-	-	-	~	-	~
~	~	-	-	-	-	-	-
~	~	-	-	-	-	-	-

ANOVA

H_1 : le medie non sono tutte uguali → il numero di verbali dipende dalla persona

H_0 : tutti gli operatori sono uguali.

Proseguendo l'analogia, il modello "di regressione" sarebbe

$$Y = \mu x + e$$

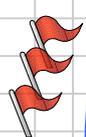
$$x = 1, 2, \dots, m$$

$$Y = \alpha + \beta x + e$$

$$Y = e^{\alpha + \beta x} + e$$

$$Y = \beta_0 + \beta_1 x + \beta_2 x^2 + e$$

- 2) C'è anche modo di recuperare la regressione
 E' meno potente e ha un sacco di difetti, però
 permette ad esempio di aggiungere altre variabili
 (in questo esempio la data, per vedere se c'è un trend)



Si "esplosa" la variabile categorica con k scelte in $k-1$ variabili dicotomiche

persona	d1	d2	d3	d4	d5	d6	d7
1	1	0	0	0	0	0	0
2	0	1	0	0	0	0	0
3	0	0	1	0	0	0	0
4	0	0	0	1	0	0	0
5	0	0	0	0	1	0	0
6	0	0	0	0	0	1	0
7	0	0	0	0	0	0	1
(8)	0	0	0	0	0	0	0

la scelgo come categoria "di default"

id	data	persona	# di verbali	d1	d2	d3	...	d6	...
1	~	1	17	1	0	0	...	0	0
2	~	6	9	0	0	0	...	1	0
...									

$$Y = \beta_0 + \beta_1 d_1 + \beta_2 d_2 + \dots + \beta_7 d_7 + e$$

$$\beta_0 = \mu_8 \quad \text{risposta della persona 8}$$

$$\beta_0 + \beta_3 = \mu_3 \quad \text{risposta della persona 3}$$

$$\beta_3 = \mu_3 - \mu_8 \quad \text{differenza tra la pers 3 e la 8}$$

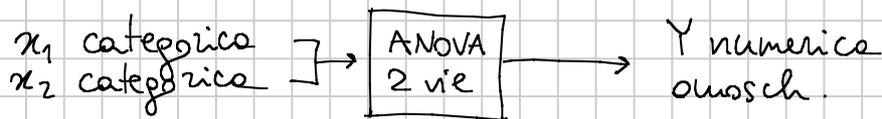
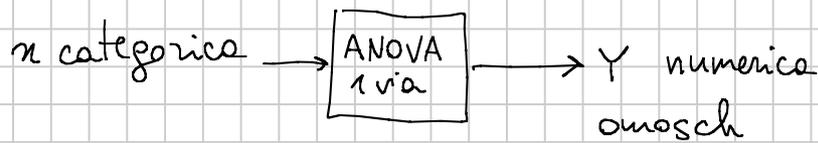
Ora posso inserire la var data (x)

$$Y = \beta_0 + \sum_1^7 \beta_i d_i + \beta_8 x + e$$

→ esempio di scelta del default : stato civile

	default	coniug. <input type="radio"/>	divor. <input type="radio"/>	vedovo <input type="radio"/>
- celibe/nubile				
- coniugato		1	0	0
- divorziato		0	1	0
- vedovo		0	0	1

ipotizzo dati sull'uso di internet



ANOVA a DUE VIE

esempio : Y : vendite di sigarette

x_1 : giorno della settimana (6 categ : lun ... sab)

x_2 : marca (4 categ : A, B, C, D)

id	x_1	x_2	Y
1	sa	B	12
2			
...			



m righe

	lu	ma	me	gi	ve	sa
A	<input type="checkbox"/>	x	x	x	x	x
B	x	x	x	x	x	12
C	x	<input type="checkbox"/>	<input type="checkbox"/>	x	x	x
D	x	x	x	<input type="checkbox"/>	x	x

n colonne

★ L'ANOVA a 2 VIE richiede esattamente un dato per ogni cella.

- i. se ci sono dei buchi al limite posso pensare di togliere tutta la riga o la colonna
- ii. se in ogni cella ci sono l dati, posso pensare di fare la media su ogni cella
- iii. se le numerosità cambiano, posso pensare di buttare via qualche dato (sempre scelto a caso)

● Ipotesi di lavoro :

- normalità + omosch. $X_{i,j} \sim \mathcal{N}(\mu_{i,j}, \sigma^2)$ dato cella i,j

- **effetto additivo** $\mu_{i,j} = \mu + \alpha_i + \beta_j$ ← effetto colonna
↑ "media globale" ↑ effetto riga

- $\alpha_1, \alpha_2, \dots, \alpha_m$ e $\beta_1, \beta_2, \dots, \beta_n$ sono un po' positivi e un po' negativi, e $\sum_{i=1}^m \alpha_i = 0 = \sum_{j=1}^n \beta_j$

★ Sull'effetto additivo

mol dire ad esempio che mediamente tutte le marche A, ..., D vendono 7 pacchetti in più il sabato, rispetto al mercoledì

	me	sa		me	sa
A	6	13	forse è più ragionevole :	6	9
B	4	11		4	6
C	1	8		1	1,5
D	2	9		2	3

+7
effetto additivo

$\mu_{i,j} = \mu + \alpha_i + \beta_j$
ok

+50%
x 1,5
effetto moltiplicativo

$\mu_{i,j} = \mu \times \alpha_i \times \beta_j$
KO

★ Nel caso si sospetti che gli effetti siano moltiplicativi si può pensare di fare il logaritmo di tutti i dati

$\log \mu_{i,j} = \log(\mu \cdot \alpha_i \cdot \beta_j) = \log \mu + \log \alpha_i + \log \beta_j$
 $\underbrace{\log \mu_{i,j}}_{\mu'_{i,j}} = \underbrace{\log \mu}_{\mu'} + \underbrace{\log \alpha_i}_{\alpha'_i} + \underbrace{\log \beta_j}_{\beta'_j}$

Funziona bene se i dati erano lognormali ...

● Che ipotesi verifica l'ANOVA a 2 vie ?

Vengono fatti due test distinti

a) Verifica se vi sia effetto riga

$$H_0: \alpha_1 = \alpha_2 = \dots = \alpha_m = 0 \quad H_1: \text{non tutti nulli}$$

b) Verifica se vi sia effetto colonna

$$H_0: \beta_1 = \beta_2 = \dots = \beta_n = 0 \quad H_1: \text{non tutti nulli}$$

● Come si fa l'ANOVA a 2 vie "a mano"?

→ Come si stimano α_i , β_j e μ ?

$$X_{i*} = \bar{X}_i = \frac{1}{n} \sum_{j=1}^n X_{ij} \quad \text{medie per riga} \quad i=1, 2, \dots, m$$

$$X_{*j} = \bar{X}^j = \frac{1}{m} \sum_{i=1}^m X_{ij} \quad \text{medie per colonna} \quad j=1, \dots, n$$

$$X_{**} = \bar{X} = \frac{1}{mn} \sum_{i,j} X_{ij} = \frac{1}{m} \sum_{i=1}^m X_{i*} = \frac{1}{n} \sum_{j=1}^n X_{*j} \quad \text{media globale}$$

i. X_{i*} , X_{*j} , X_{**} tutte normali

ii. Calcolo i valori attesi

$$E(X_{i*}) = \frac{1}{n} \sum_j E(X_{ij}) = \frac{1}{n} \sum_j (\mu + \alpha_i + \beta_j) = \mu + \alpha_i + 0 = \mu + \alpha_i$$

$$E(X_{*j}) = \mu + \beta_j \quad E(X_{**}) = \mu$$

iii. Varianze

$$\text{Var}(X_{i*}) = \frac{1}{n^2} \sum_{j=1}^n \text{Var}(X_{ij}) = \frac{\sigma^2}{n} \quad \text{Var}(X_{*j}) = \frac{\sigma^2}{m} \quad \text{Var}(X_{**}) = \frac{\sigma^2}{mn}$$

iv. $\alpha_i = (\mu + \alpha_i) - \mu \approx X_{i*} - X_{**}$

$$\left. \begin{aligned} \alpha_i &\approx X_{i*} - X_{**} \\ \beta_j &\approx X_{*j} - X_{**} \\ \mu &\approx X_{**} \end{aligned} \right\} \mu_{ij} = \mu + \alpha_i + \beta_j \approx X_{i*} + X_{*j} - X_{**}$$

★ Questi stimatori hanno legge normale con media ovvia (sono stimatori corretti) e varianze complicate (ad esempio X_{i*} e X_{**} non sono indipendenti)

→ Calcolo varianze e devianze campionarie delle medie per riga e per colonna

x_{1*} x_{2*} ... x_{m*} medie per riga x_{**} è la loro \bar{x}

$$S_{MR}^2 = \frac{1}{m-1} \sum_{i=1}^m (x_{i*} - x_{**})^2 \approx \frac{\sigma^2}{n} \quad m-1 \text{ g.d.l.}$$

↑
sotto H_0 a) ovvero se $\alpha_i \equiv 0$

ovvero se $x_{i*} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$

$$\sum_{i,j} \frac{(x_{ij} - \mu_{ij})^2}{\sigma^2} \sim \chi^2(mn) \quad \text{fatto ovvio}$$

$$\sum_{i,j} \frac{(x_{ij} - \underbrace{x_{i\cdot} - x_{\cdot j} + x_{\cdot\cdot}}_{\text{stimatore}})^2}{\sigma^2} \sim \chi^2(mn - \underbrace{m - n + 1}_{\text{gdl persi}}) \sim \chi^2((m-1)(n-1))$$

HW: Verificare ipotesi e applicazione del "teorema chi-quadrato"

$$SS_e = \sum_{i,j} (x_{ij} - x_{i\cdot} - x_{\cdot j} + x_{\cdot\cdot})^2 \quad \text{devianza errori}$$

$$S_e^2 = \frac{SS_e}{(m-1)(n-1)} \quad \text{varianza errori}$$

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

HW: Calcolare E e Var di S_e^2

$$\rightarrow \sigma^2 \approx S_e^2 \quad \text{stimatore corretto e consistente}$$

STIMATORI DI σ^2 BUONI SOTTO H_0

a) $H_0: \alpha_i \equiv 0 \quad \forall i$ $H_1: \text{non tutti nulli}$ test sull'effetto righe

b) $H_0: \beta_j \equiv 0 \quad \forall j$ $H_1: \text{non tutti nulli}$ test sull'effetto colonne

$$x_{i\cdot} \sim \mathcal{N}(\mu + \alpha_i, \frac{\sigma^2}{n})$$

sotto H_0 a)

$$x_{i\cdot} \sim \mathcal{N}(\mu, \frac{\sigma^2}{n}) \Rightarrow \sum_{i=1}^m n \frac{(x_{i\cdot} - \mu)^2}{\sigma^2} \sim \chi^2(m)$$

$x_{\cdot\cdot}$ al posto di μ e colonna di 1 i gdl

$$SS_R = \sum_{i=1}^m n (x_{i\cdot} - x_{\cdot\cdot})^2 \quad \text{devianza righe}$$

$$\frac{SS_R}{\sigma^2} \stackrel{H_0^a)}{\sim} \chi^2(m-1) \quad S_R^2 = \frac{SS_R}{m-1} \quad \text{varianza righe}$$

★ Sempre sotto H_0 a) S_e^2 e S_R^2 sono indipendenti, quindi:

$$V := \frac{S_R^2}{S_e^2} \stackrel{H_0^a)}{\sim} F(m-1; (m-1)(n-1)) \quad \text{statistica test a)}$$

★ Sotto H_1 a) V tende ad assumere valori più grandi

Analogamente ...

sotto $H_0^{(b)}$

$$SS_c = \sum_{j=1}^n m (x_{x,j} - x_{xx})^2 \quad \text{devianza colonne} \quad \frac{SS_c}{\sigma^2} \stackrel{H_0^{(b)}}{\sim} \chi^2(n-1)$$

$$S_c^2 = \frac{SS_c}{n-1} \quad \text{varianza colonne}$$

$$\star W := \frac{S_c^2}{S_e^2} \stackrel{H_0^{(b)}}{\sim} F(n-1; (m-1)(n-1)) \quad \text{statistica del test } b)$$

\star Sotto $H_1^{(b)}$ W tende ad assumere valori più grandi

ANALISI DEI RESIDUI

→ esempio: dieta delle banane - 10 persone per il test

pesi prima:

- 82 = $x_{1,1}$
- 97 = $x_{2,1}$
- 75
- 63
- 68
- 79
- ...
- 111 = $x_{10,1}$

pesi dopo:

- 85 = $x_{1,2}$
- 94 = $x_{2,2}$
- 74
- 68
- 67
- 80
- ...
- 109 = $x_{10,2}$

$$x_{1,1} \sim \mathcal{N}(\mu_{1,1}, \sigma^2)$$

$$x_{2,1} \sim \mathcal{N}(\mu_{2,1}, \sigma^2)$$

$$x_{1,2} \sim \mathcal{N}(\mu_{1,1} + \delta, \sigma^2)$$

$$x_{2,2} \sim \mathcal{N}(\mu_{2,1} + \delta, \sigma^2)$$

$$x_{i,2} \sim \mathcal{N}(\mu_{i,1} + \delta, \sigma^2)$$

effetto dieta

sto supponendo che l'effetto della dieta sia lo stesso per tutti

$$\mu := \frac{1}{10} \sum_{i=1}^{10} \mu_{i,1} + \frac{\delta}{2} \quad \alpha_i := \mu_{i,1} + \frac{\delta}{2} - \mu \quad \beta_1 := -\frac{\delta}{2} \quad \beta_2 := \frac{\delta}{2}$$

\star ho definito tutto in modo che:

$$\mu_{i,j} = \mu + \alpha_i + \beta_j \quad \forall i,j$$

$$\mu + \alpha_i + \beta_1 = \mu_{i,1} + \frac{\delta}{2} - \frac{\delta}{2} = \mu_{i,1}$$

$$\mu + \alpha_i + \beta_2 = \mu_{i,1} + \frac{\delta}{2} + \frac{\delta}{2} = \mu_{i,1} + \delta = \mu_{i,2}$$

$$\rightarrow \sum_{j=1}^2 \beta_j = 0$$

$$\rightarrow \sum_{i=1}^{10} \alpha_i = \sum_{i=1}^{10} \mu_{i,1} + 10 \frac{\delta}{2} - 10\mu = \sum_{i=1}^{10} \mu_{i,1} + 10 \frac{\delta}{2} - \sum_{i=1}^{10} \mu_{i,1} - 10 \frac{\delta}{2} = 0$$

ok

ora 24

82	85	83,5
97	94	95,5
75	74	74,5
...
82,2	83,4	82,8

$$\beta_1 \approx 82,2 - 82,8 = -0,6$$

$$\beta_2 \approx 83,4 - 82,8 = 0,6$$

$$\alpha_1 \approx 83,5 - 82,8 = 0,7$$

$$\alpha_2 \approx 12,7$$

$$\alpha_3 \approx -8,3$$

diff = effetto dieta + 1,2 $\approx \delta$

● Calcolo i previsti

$\mu_{1,1} \approx 82$ ma meglio $\mu_{1,1} \approx 82,8 + 0,7 + (-0,6) = 82,9$

previsti

è una stima che usa tutti i dati

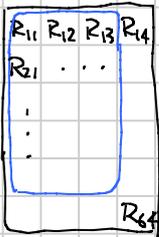
82,9	84,1	83,5
94,9	96,1	95,5
73,9	85,1	...
...
82,2	83,4	

$82,8 - 8,3 - 0,6 = 72,8 + 1,1 = 73,9$

● Calcolo i residui (osservato - previsto)

-0,9	0,9	○
2,1	-2,1	○
1,1	-1,1	○
...	...	○
○	○	○

i residui non possono per definizione avere effetti riga e colonna



i numeri liberi sono meno di $m \cdot n$

sono: $(m-1)(n-1)$

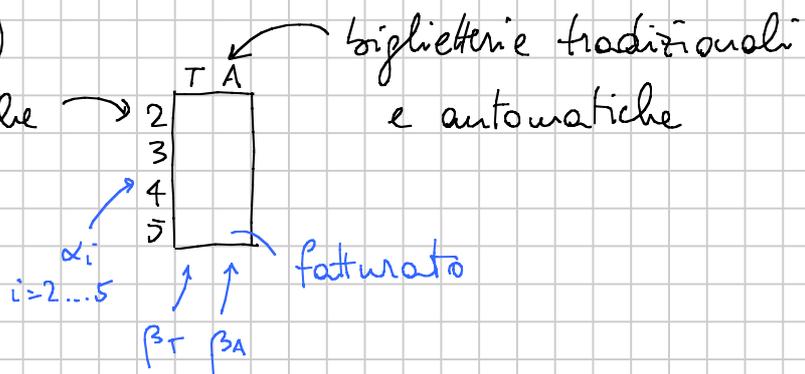
(ovvero i g d l dello stimatore S_e^2)

SS_e = somma quadrati tutti i residui

★ A parte il "bordo" i residui dovrebbero essere normali di media 0 e var circa S_e^2 , indipendenti e senza struttura

■ esempio 2: (sbagliato)

biglietti automatiche in funzione



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

$\beta_T = -\beta_A$ $\beta_A = \frac{\gamma}{2}$ dove γ è la differenza tra A e T

★ Non è naturale che γ sia lo stesso per tutte le righe

In realtà osservo:

65	12	38,5
64	14	39
62	17	39,5
60	20	40
62,75	15,75	78,5

[Vedi commento al file xlsx]

RELAZIONE TRA ANOVA A 1 VIA E A 2 VIE

● Analogie:

- ANOVA 1 via generalizza il test per il confronto delle medie di campioni normali omoschedastici a $n \geq 2$ campioni
- ANOVA 2 vie generalizza il test **t-paired** (vedi Cap 8 Ross) a $n \geq 2$ campioni; inoltre lo fa sia sulle righe, sia sulle colonne

● Ambiti di applicazione:

- Se i campioni hanno lunghezze diverse o se non c'è corrispondenza tra le righe di diverse colonne si può fare solo l'ANOVA a 1 via
- In caso contrario, capita spesso che uno dei due effetti sia evidente a priori che sarà **forte**: può sembrare inutile verificarlo esplicitamente e si ha la tentazione di fare ANOVA a 1 via sull'altro effetto.

In realtà questo è **sbagliatissimo**: proprio in questi casi la ANOVA a 2 VIE risulta **molto più potente** dell'ANOVA a 1 via nel cercare l'altro effetto (quello non ovvio), perché, tenendo conto che le righe (ad esempio) **sono diverse** riesce a schermare l'effetto e "vedere" un rumore minore (2 vs 14) sull'effetto delle colonne

- Se non è evidente se le due variabili hanno effetti significativi chiaramente uso l'ANOVA a 2 vie.

- Se entrambi i test dicono H_1 , fine
- Se un solo test dice H_1 , fine
- Se entrambi i test dicono H_0 : in questo caso può darsi (è frequente) che le due ANOVA a 1 via siano più potenti.

rifare con ANOVA 1 via

ANOVA A DUE VIE CON INTERAZIONI

Per ogni cella delle $m \times n$ combinazioni dei due fattori serve un campione di $l \geq 2$ dati

$$X_{i,j,k} \sim \mathcal{N}(\mu_{i,j}, \sigma^2) \quad \begin{matrix} \text{righe} \\ i=1, \dots, m \end{matrix} \quad \begin{matrix} \text{colonne} \\ j=1, \dots, n \end{matrix} \quad \begin{matrix} \text{campione} \\ k=1, \dots, l \end{matrix}$$

- ipotesi di omoschedasticità
- le medie non dipendono da k

$$\mu_{i,j} = \mu + \alpha_i + \beta_j \quad \text{ANOVA 2VIE} \quad \Leftrightarrow \quad \mu_{i,j} = \mu + \alpha_i + \beta_j + \delta_{i,j} \quad \text{ANOVA 2VIE CON INT.}$$

è sempre possibile scrivere $\mu_{i,j}$ come nell'espressione a destra e in modo tale che $\sum_i \alpha_i = 0 = \sum_j \beta_j = \sum_{i,j} \delta_{i,j}$

- $\delta_{i,j}$ si chiama termine di interazione

- Test che si fanno:

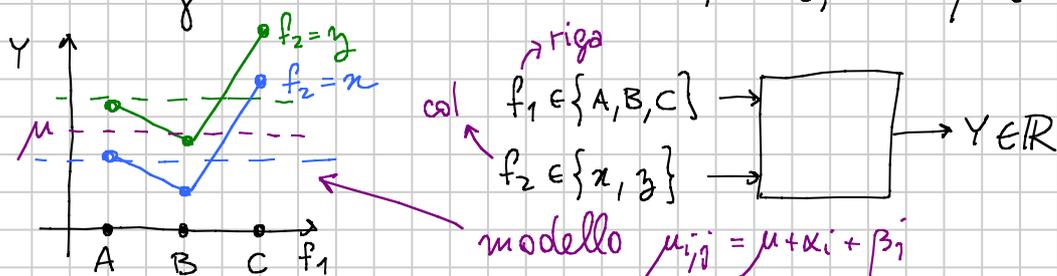
a) $H_0: \alpha_i \equiv 0 \quad \forall i$ $H_1: \text{non tutti nulli}$ test sull'effetto riga

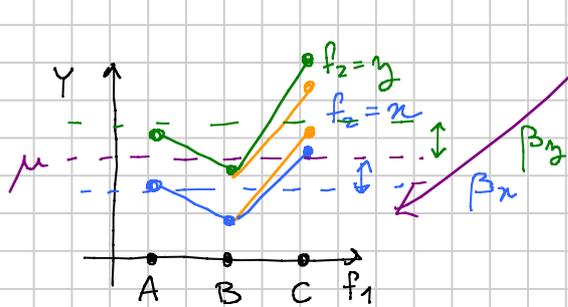
b) $H_0: \beta_j \equiv 0 \quad \forall j$ $H_1: \text{non tutti nulli}$ test sull'effetto colonna

c) $H_0: \delta_{i,j} \equiv 0 \quad \forall i,j$ $H_1: \text{non tutti nulli}$ test sulle interazioni

→ Se $H_0^{c)}$ allora il modello additivo è adeguato e si può usare la ANOVA a 2 VIE (+ potente, + robusto)

→ Se $H_1^{c)}$ allora vi sono interazioni non nulle; il modello additivo non è adeguato; la presenza/assenza di effetti riga e colonna globali vanno lette su a) e b) di questa ANOVA





modello con interazioni $\mu_{ij} = \mu + \alpha_i + \beta_j + \delta_{ij}$

i punti arancio sono i previsti di $\mu + \alpha_i + \beta_j$

i termini di interazione $\delta_{c,x}$ e $\delta_{c,y}$

dicono di quanto il modello osservato si discosta da quello semplificato additivo (arancione)

TEST DI ADATTAMENTO (Cap 11 Ross)

- 1) test del chi-quadro e derivati
- 2) test di Kolmogorov-Smirnov e parenti (leggere sul Ross)

TEST DEL CHI-QUADRO ELEMENTARE

esempio: v.a. discreta con una legge teorica nota

S: scolarizzazione

$S \in \{ \text{laurea}, \text{maturità}, \text{III media}, \text{meno} \}$

dati Istat
(distribuzione nota)

→ 4% 75% 7% 14%

Prendo un campione particolare: residenti a Treviso

$n = 100$ 2 61 20 17

Q: questo campione può sensatamente (H_0) venire da quella distribuzione o c'è evidenza (H_1) che la sua distrib sia diversa?

→ Il test del χ^2 elementare funziona su una v.a. X che assume solo i valori $\{1, 2, \dots, n\}$ dove n è finito, piccolo e si suppone una certa legge φ_0 .

$H_0: P(X=1) = \varphi_0(1), P(X=2) = \varphi_0(2), \dots, P(X=n) = \varphi_0(n)$

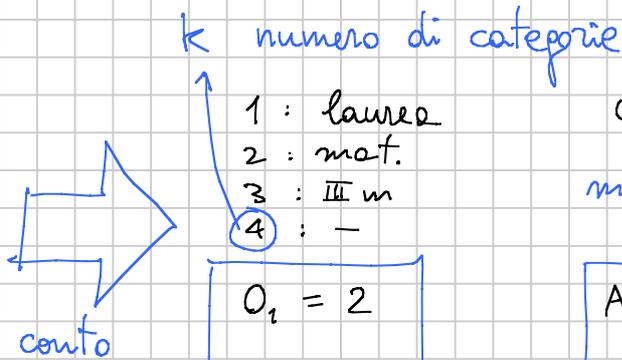
$H_1: \text{non tutte uguali}$

→ Si chiama test di adattamento a distribuzione discreta completamente specificata

• In pratica

id	sex	età	...	S
1				III m
2				la
⋮				III m
⋮				mat
⋮				⋮
100				⋮

n numero di dati



$O_1 = 2$
$O_2 = 61$
$O_3 = 20$
$O_4 = 17$

osservati

φ_0 nota
moltiplica per n

$A_1 = 4$
$A_2 = 75$
$A_3 = 7$
$A_4 = 14$

attesi
(normalmente non sono interi)

$$W := \sum_{i=1}^k \frac{(O_i - A_i)^2}{A_i}$$

statistica del test

- W è grande sotto H_1 e piccola sotto H_0
- Sotto H_0 , W ha una distribuzione calcolabile esattamente con metodi sofisticati (numerici o montecarlo) (in lab?)
- Fortunatamente, di solito, sotto H_0 : $W \approx \chi^2(k-1)$

ipotesi perché la approssimazione sia buona:

- Occorre che gli attesi siano grandi:

↳ "rule of thumb": tutti $A_i \geq 1$ e i $\frac{4}{5}$ degli $A_i \geq 5$

6 CFU → 8 h teoria
 ↳ 10 h pratica

1 CFU → 25 h

20 / 10 = 2 CFU }
 32 / 8 = 4 CFU } 6 CFU
 52 + 98 h di lavoro a casa

lun 21 dic 16.30-18.30
 Nicolodi + Morandini lab 2+3
 (da confermare)

TEST DEL CHI QUADRO

$A_i = n \cdot p_i$ affinché l'approssimazione sia buona
 bisogna che i p_i non siano piccoli

- in particolare se k è grande, questo non è possibile
- ↳ si può risolvere con metodi MC
- ↳ oppure raggruppando i valori :

Quando X ha tanti valori possibili o è continua

→ se X è discreta raggruppo valori simili in modo da ridurre il numero elevato di valori con prob piccole ad un numero basso di gruppi (bin) con prob non piccole

Vogliamo testare per X una distribuzione di Poiss

0	4,53999E-05	p0	
1	0,000453999	p1	
2	0,002269996	p2	
3	0,007566655		
4	0,018916637		50
5	0,037833275		
6	0,063055458	0,130141	6,507071
7	0,090079226	0,090079	4,503961
8	0,112599032	0,112599	5,629952
9	0,125110036	0,12511	6,255502
10	0,125110036	0,12511	6,255502
11	0,113736396	0,113736	5,68682
12	0,09478033	0,09478	4,739017
13	0,072907946	0,208444	10,42218
14	0,052077104	8 gruppi di valori poss	
15	0,03471807		
16	0,021698794		
17	0,012763996		

→ se X è continua, uguale legge esponenziale di media 10



4 bin di prob elevata

★ Principio generale: se i bin sono più o meno equiprobabili l'approx è migliore e il test è più potente

→ Se X è continua è possibile (e raccomandato) ottenere bin esattamente equiprobabili

→ voglio k bin equip. delle distrib. continua con pdf f e Cdf F

c_1, c_2, \dots, c_{k-1} punti di taglio sono tali che :

$$\frac{1}{k} = \int_{-\infty}^{c_1} f(t) dt = \int_{c_1}^{c_2} f(t) dt = \dots = \int_{c_{k-1}}^{\infty} f(t) dt$$

$$= F(c_1) = F(c_2) - F(c_1) \dots 1 - F(c_{k-1})$$

$$\Rightarrow F(c_1) = \frac{1}{k}, F(c_2) = \frac{2}{k}, \dots F(c_{k-1}) = \frac{k-1}{k}$$

$$\Rightarrow c_i = F^{-1}\left(\frac{i}{k}\right)$$

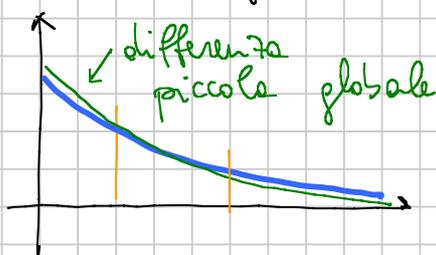
Vogliamo testare per X una distribuzione esponenziale di media 10

i	c_i	O_i / frequenza (array; ...)
1	1,541507	
2	3,364722	
3	5,596158	
4	8,472979	
5	12,52763	
6	19,4591	

★ Gli attesi sono sempre n per la prob di ciascun bin
 Gli osservati sono il numero di dati che cadono in ogni bin

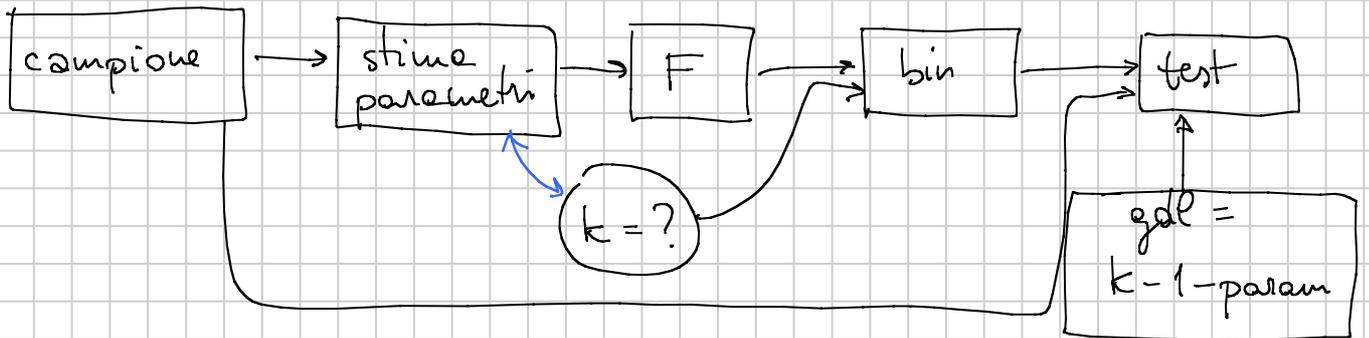
● Quanti bin bisogna fare :

- i. maggiore n , maggiore k
- ii. k mai troppo grande (gli attesi devono rispettare le "rule of thumb")
- iii. più k è piccolo più il test è potente (differenze piccole)
- iv. più k è grande più il test è sensibile (differenze locali)



Test di adattamento a distribuzione con parametri incogniti

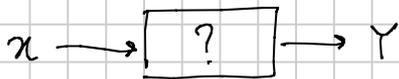
→ Si stimano puntualmente i parametri incogniti sui dati e poi si esegue il test usando quelli. In questo modo il campione si adatta però artificialmente bene alla distribuzione. Per correggere si devono calare i gdf della chi-quadro del numero di parametri stimati



ora 28

TABELLE DI CONTINGENZA

Primo scopo: vedere se due variabili sono indipendenti



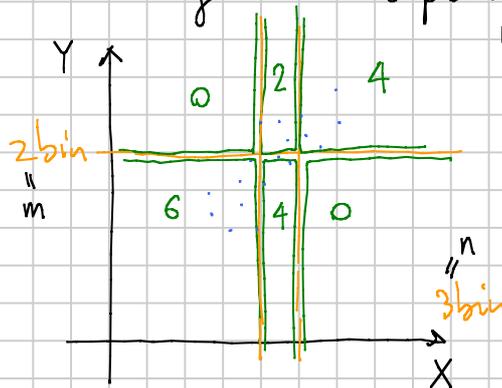
regressione: X, Y numeriche

anova: X categorica, Y numerica

t d cont: X, Y qualunque (ma è poco potente, quindi dice spesso H_0)

$X \leftrightarrow Y$ il legame di dipendenza cercato è simmetrico

id	X	Y
1	~	~
2	~	~
...
N	~	~



6 bin bidimensionali

	1	2	3	X
1	0	2	4	6
2	6	4	0	10
Y	6	6	4	16

caso entrambe numeriche

	1	2	3	X
1	0	2	4	6
2	6	4	0	10
Y	6	6	4	16

H_0 : le due var sono indipendenti

H_1 : no

sotto H_0 , ipotizzando che i totali verdi siano dati, che valori mi aspetterei

nelle 6 celle?

2,25	2,25	1,5	6
3,75	3,75	2,5	10
6	6	4	16

$$\frac{6}{16} \cdot \frac{6}{16} = \frac{36}{16^2} \text{ del totale } \frac{36}{16} = 2,25 \text{ attesi}$$

$$\frac{6}{16} \text{ della riga da } 10 = \frac{6}{16} \cdot 10 = \frac{15}{4} = 3,75$$

Questi sono gli attesi A_i che si avrebbero se ci fosse indipendenza

0	2	4
6	4	0

gli osservati

$$W = \sum_{i=1}^6 \frac{(A_i - O_i)^2}{A_i} = \frac{(0-2,25)^2}{2,25} + \frac{(2-2,25)^2}{2,25} + \frac{(4-1,5)^2}{1,5} + \frac{(6-3,75)^2}{3,75} + \frac{(4-3,75)^2}{3,75} + \frac{(0-2,5)^2}{2,5}$$

la statistica

sotto H_0 : $W \sim \chi^2(\text{gdl})$

richiede la solita "rule" sugli A_i , che qui è falsa

$\text{gdl} = (m-1)(n-1)$ dove m e n sono righe e colonne

★ Nelle scelte dei numeri di bin m, n occorre ricordare che i bin totali $m \cdot n$ devono essere adeguati alle numerosità N del campione

→ Se una variabile è **numerica**, converrà "tagliare" i bin ai percentili campionari

$$= \text{PERCENTILE}(\text{array}_x; 1/3)$$

$$= \text{PERCENTILE}(\text{array}_x; 2/3)$$

→ Se una var è **categorica**, converrà raggruppare le categorie fino ad ottenere il numero voluto di bin, che devono risultare \pm equifrequenti

stato famiglia	#
celibe	64
coniugato	112
convivente	45
divorziato	38
separato	35
vedovo	44
	<u>338</u>

3 bin ~ 113 l'uno
 cel+conv 109
 coniug 112
 altro 117

ora 29

TUTTI GLI SCOPI DELLE TAB. DI CONT.

- 1) Vedere se $H_0: X, Y$ sono indipendenti $H_1: \text{no}$
- 2) Vedere se $H_0: \text{no}$ $H_1: Y$ dipende da X
 (= X dipende da Y)

* alternativa **non-parametrica** a regressione semplice ($p=1$) e ANOVA a 1 via

↳ non ha ipotesi né di omoschedasticità, né di gaussianità, né di linearità.

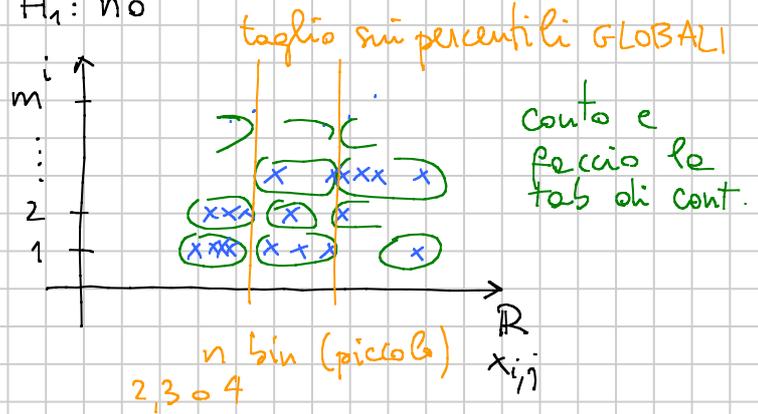
↳ è meno potente delle tecniche parametriche corrispondenti

3) Vedere se m campioni hanno la stessa distribuzione

$X_{1,1} X_{1,2} \dots X_{1,n_1} \sim D_1$
 $X_{2,1} \dots X_{2,n_2} \sim D_2$
 ...
 $X_{m,1} \dots X_{m,n_m} \sim D_m$

$H_0: D_1 = D_2 = \dots = D_m$
 $H_1: \text{no}$

X_{ij} dipende da i ? m bin
 ↑ numerica ma non necessariamente
 ↑ categoria



* Ricordiamo che il test del chi-quadrato standard permette simultaneamente di confrontare $D_1 \sim X_1 \dots X_n$ con una distribuzione assegnata

<table border="1" style="width: 100%; border-collapse: collapse;"> <tr> <td style="padding: 5px;">camp 1:</td> <td style="padding: 5px;">37</td> <td style="padding: 5px;">24</td> <td style="padding: 5px;">17</td> <td style="padding: 5px;">12</td> <td style="padding: 5px;">90</td> </tr> <tr> <td style="padding: 5px;">camp 2:</td> <td style="padding: 5px;">63</td> <td style="padding: 5px;">51</td> <td style="padding: 5px;">32</td> <td style="padding: 5px;">20</td> <td style="padding: 5px;">166</td> </tr> </table> <p style="text-align: center; margin-top: 10px;">test, 3 gdl hanno la stessa distrib?!</p>	camp 1:	37	24	17	12	90	camp 2:	63	51	32	20	166	<p style="margin-bottom: 5px;">distrib : 0,4 0,3 0,2 0,1</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <tr> <td style="padding: 5px;">camp 1:</td> <td style="padding: 5px;">37</td> <td style="padding: 5px;">24</td> <td style="padding: 5px;">17</td> <td style="padding: 5px;">12</td> <td style="padding: 5px;">90</td> </tr> <tr> <td style="padding: 5px;">alteri :</td> <td style="padding: 5px;">36</td> <td style="padding: 5px;">27</td> <td style="padding: 5px;">18</td> <td style="padding: 5px;">10</td> <td style="padding: 5px;">90</td> </tr> </table> <p style="text-align: center; margin-top: 10px;">test, 3 gdl hanno la stessa distrib?!</p>	camp 1:	37	24	17	12	90	alteri :	36	27	18	10	90
camp 1:	37	24	17	12	90																				
camp 2:	63	51	32	20	166																				
camp 1:	37	24	17	12	90																				
alteri :	36	27	18	10	90																				

4) Nel caso una delle variabili sia dicotomica, vedere se

H_0 : la prob di successo non dipende dall'altra variabile

H_1 : dipende

	ER	H	GR	...	
non ces.	74	45	~	~	osservati
cesarei	22	18	~	~	

	ok	ko	
farmaco	26	14	40
placebo	19	21	40
	45	35	80

▣ SUI TEST STATISTICI

(2008 da ore B.13)

→ parametro incognito all'interno di classe di distribuzioni note
(ad es μ per $N(\mu, \sigma^2)$)

→ due ipotesi contrapposte (ad es: $H_0: \mu = \mu_0$ $H_1: \mu \neq \mu_0$; $H_0: \mu \leq \mu_0$ $H_1: \mu > \mu_0$)
spesso con un valore target di confronto (ad es: μ_0) detto **AQL**
acceptable quality level

→ stimatore parametro ($\bar{X} \approx \mu$) a volte una statistica ($\frac{\bar{X} - \mu_0}{S/\sqrt{n}}$)

● ESEMPIO : misuro 50 pezzi $\bar{X} = 2,758$ il target doveva essere

$$\mu_0 = 2,735$$

la media con cui sto producendo questi pezzi (μ) è in target?



H_0 : in target $\mu = \mu_0$

H_1 : fuori target $\mu \neq \mu_0$

→ test dell'idiota :

se $\bar{X} = \text{target}$, ok se no regolo la macchina
dico sempre H_1 e mai H_0

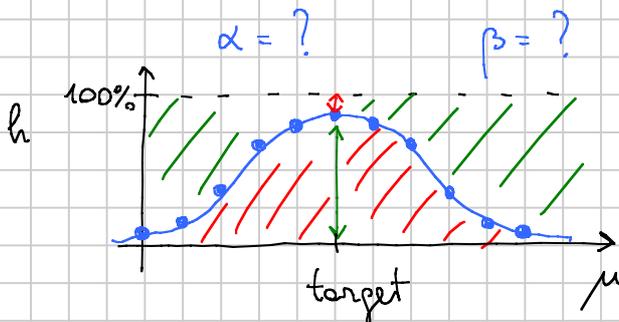
$$\alpha = 100\%$$

$$\beta = 0\%$$

→ test del dilettante : ci vuole un po' di tolleranza

se $\bar{X} \in \text{target} \pm 0,010$, ok se no regolo

↑ scelta a occhio



h misura la probs di dire H_0

$$h = h(\mu)$$

funzione o curva operativa caratteristica

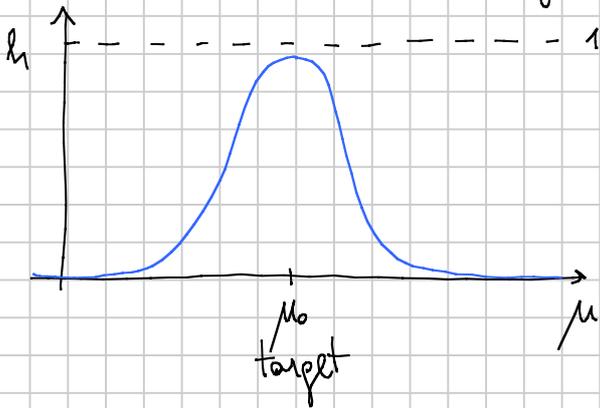
quando è vera H_0 , sono in target, dico H_1 con probs $1 - h(\mu_0)$

$$\alpha = 1 - h(\mu_0)$$

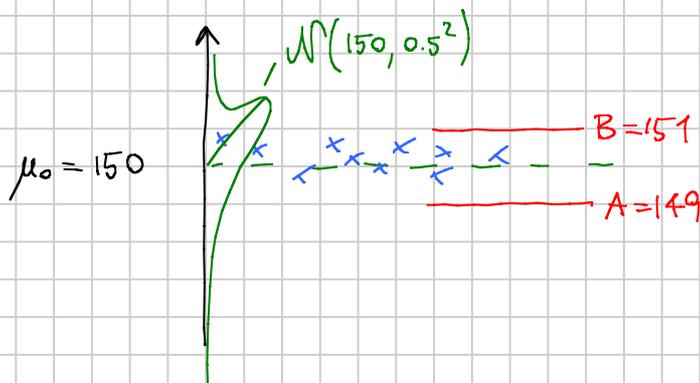
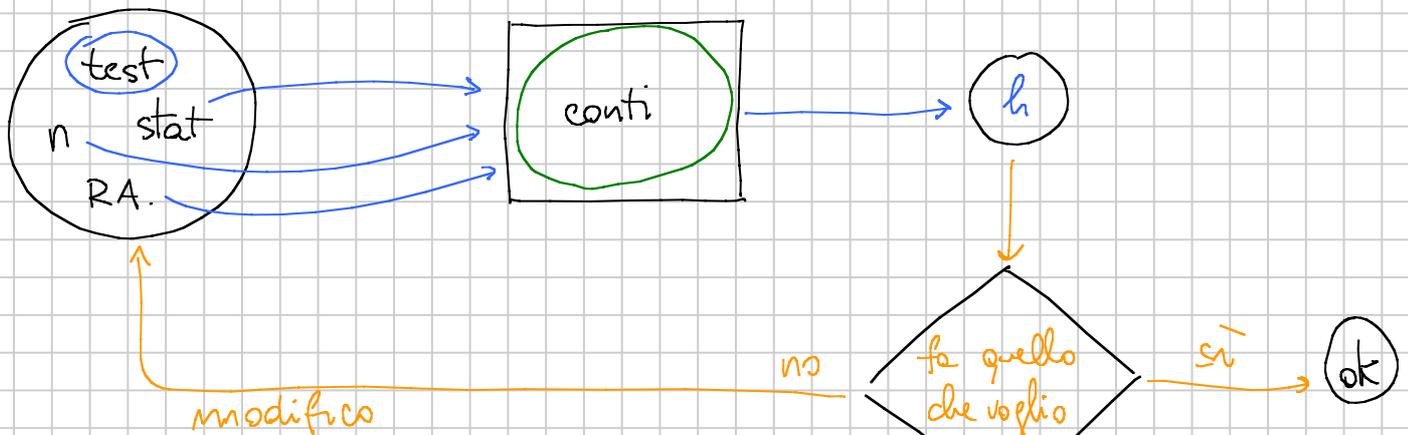
in tutti gli altri casi è vera H_1

$$\beta = h(\mu) \quad \text{cioè } \beta \text{ dipende da } \mu$$

Il "professionista" è in grado di calcolare la C.O.C.



$h(\mu)$: prob di dire H_0 quando il parametro vale μ



$n = 4$

$\sigma = 0,5 \text{ ml}$

$\bar{X} = \frac{1}{4} \sum_{i=1}^4 x_i$

$\mu_0 = 150 \text{ ml}$

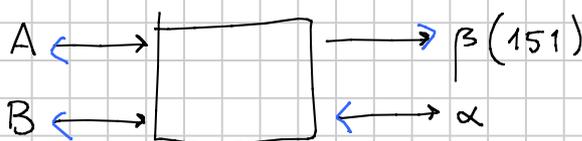
$RA_{\bar{X}} = [149; 151]$

$h = h(\mu) := \text{"prob di dire } H_0\text{"} = P(\bar{X} \in [A, B]) = P(A \leq \bar{X} \leq B)$

$\bar{X} \sim ? \quad F_{\bar{X}} = ?$

$\bar{X} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right) \quad F_{\bar{X}}(t) = \text{DISTRIB.NORM}\left(t; \mu; \sqrt{\frac{\sigma^2}{n}}\right)$

$h(\mu) = F_{\bar{X}}(B) - F_{\bar{X}}(A) = \text{DISTRIB.NORM}\left(B; \mu; \sqrt{\frac{\sigma^2}{n}}\right) - \text{DISTRIB.NORM}\left(A; \mu; \sqrt{\frac{\sigma^2}{n}}\right)$



$$\left. \begin{aligned} A &= \mu_0 - r \\ B &= \mu_0 + r \end{aligned} \right\} r \text{ diventa l'unica incognita}$$

$$r \longleftrightarrow \boxed{} \longleftrightarrow \alpha$$

$$\alpha = 1 - h(\mu_0) = 1 - (F_{\bar{X}}(\mu_0 + r) - F_{\bar{X}}(\mu_0 - r)) = 1 - \Phi\left(\frac{\mu_0 + r - \mu_0}{\sigma/\sqrt{n}}\right) + \Phi\left(\frac{\mu_0 - r - \mu_0}{\sigma/\sqrt{n}}\right)$$

$$\bar{X} \sim \mathcal{N}\left(\mu_0, \frac{\sigma^2}{n}\right) \quad F_{\bar{X}}(t) = \Phi\left(\frac{t - \mu_0}{\sigma/\sqrt{n}}\right)$$

$$= 1 - \Phi\left(\frac{r\sqrt{n}}{\sigma}\right) + \Phi\left(-\frac{r\sqrt{n}}{\sigma}\right) = 2\left(1 - \Phi\left(\frac{r\sqrt{n}}{\sigma}\right)\right)$$

$$\frac{\alpha}{2} = 1 - \Phi\left(\frac{r\sqrt{n}}{\sigma}\right)$$

$$1 - \frac{\alpha}{2} = \Phi\left(\frac{r\sqrt{n}}{\sigma}\right) \Leftrightarrow \frac{r\sqrt{n}}{\sigma} = \Phi^{-1}\left(1 - \frac{\alpha}{2}\right) \Leftrightarrow r = \frac{\sigma}{\sqrt{n}} \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)$$

ok 31

● Esempio 2 (unilaterale): mtst 0605.1

p : percentuale difettosi nelle forniture

↳ campione di n pezzi

X : numero di difettosi sul campione

$\hat{p} = \frac{X}{n} \approx p$ fraz difettosi sul campione

$p < 3\%$? $p > 3\%$?

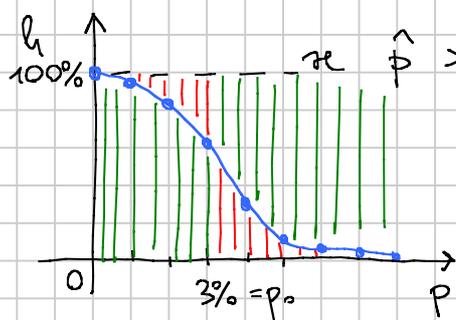
→ test dell'idiota

se $\hat{p} < 3\% \rightarrow p < 3\%$ tutto ok

se $\hat{p} > 3\% \rightarrow p > 3\%$ protesto col fornitore
(o rifiuto i pezzi)

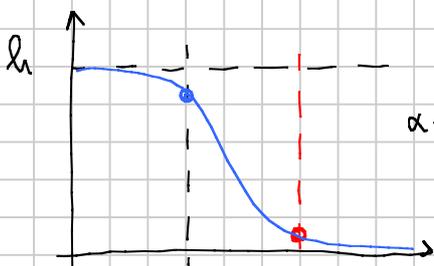
→ test del dilettante: un po' di tolleranza

se $\hat{p} \leq 3\% + 1,5\%$ → $p \leq 3\%$ tutto ok
↑
marginale arbitrario



se $\hat{p} > 3\% + 1,5\%$ → $p > 3\%$ protesto (H_1)

$h = h(p)$: "pds di dire H_0 "



α : "prob err di I specie" = $1-h(p)$ $p \leq 0,03$

$\alpha = \alpha(p) \leq \bar{\alpha} := 1-h(p_0) = 1-h(3\%)$

livello di significatività: max err I sp.

β : "prob err di II specie" = $h(p)$ $p > 0,03$

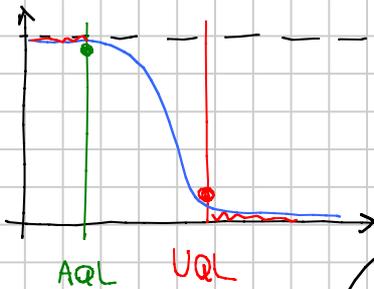
$\beta = \beta(p) \leq h(p_0) = 1-\bar{\alpha}$

AQL
acceptable quality level

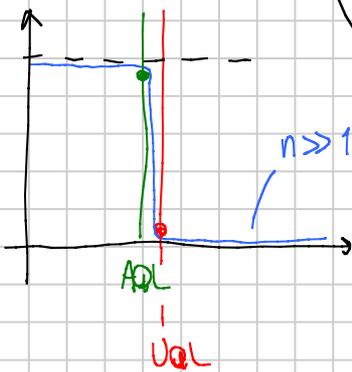
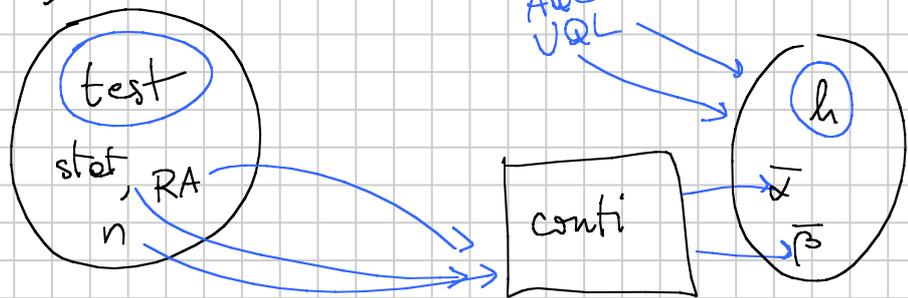
UQL
unacceptable quality level

se $p > p_1$ allora $\beta = \beta(p) \leq \bar{\beta} = h(p_1)$

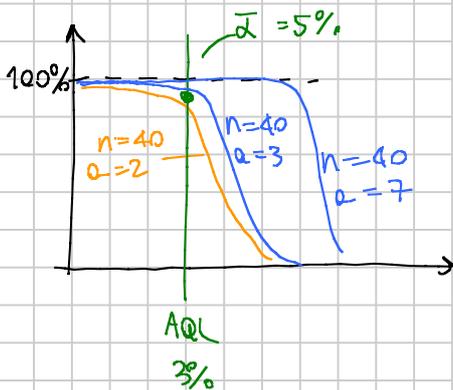
★ La definizione dell' UQL simmetrizza il test: scambiando le ipotesi e il ruolo di $\bar{\alpha}$ e $\bar{\beta}$ si ottiene il medesimo risultato



Si fissano UQL e AQL, $\bar{\alpha}$ e $\bar{\beta}$ e si ottiene un test con quelle caratteristiche

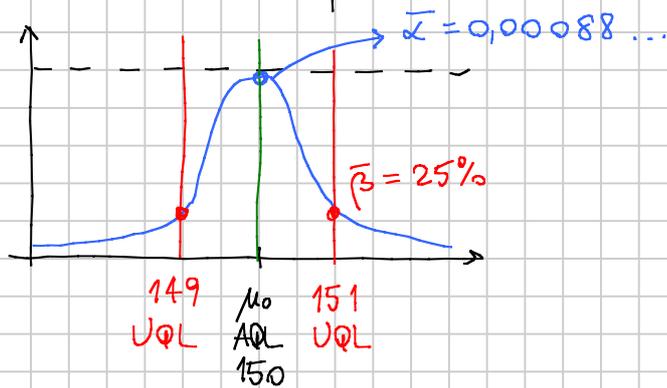


$h(p) = P(X \leq a) = F_X(a) = \text{DISTRIB. BINOM}(a; n; p; 1)$
bin(n, p)



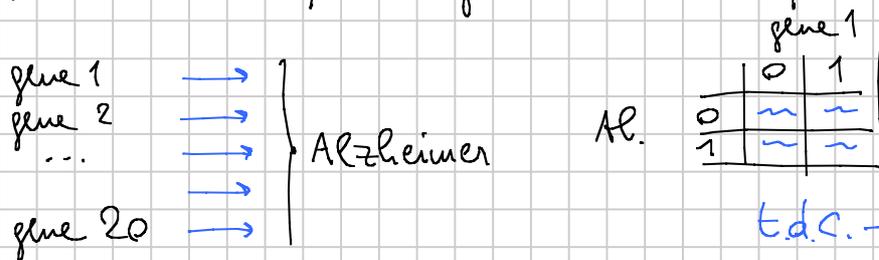
★ Può essere utile definire UQL anche nei casi bilaterali

Nel primo esempio:



CORREZIONE DI BONFERRONI

Se in un ambito singolo sono costretto a fare più test, le probabilità dei falsi positivi (α , per I specie) si sommano e quindi la frequenza globale di falsi positivi non è più piccola



t.d.c. → gene 1 è legato a Al.

$\alpha = 5\%$ per ciascun test ; ipotizziamo vere tutte H_0
ciascun test dice H_1 con prob 5% → la prob che dicano tutti H_0 è $0,95^{20} \approx 36\%$ → 64% di prob che almeno un gene risulti significativo *pur non essendolo*

★ Questo test : (H_0 : nessuno di questi geni è legato alla malattia vs H_1 : qualcuno sì) fatto così, ha livello di significatività del 64%. **INACCETTABILE**

★ Correzione di Bonferroni : se devo fare n test e voglio globalmente lvl di sign. α , ciascuno di questi test deve essere eseguito con lvl di sign. $\frac{\alpha}{n}$.

→ Ne escono test poco potenti

TEST DI FISHER-IRWIN (Cap 8)

tabelle di contingenza 2×2 quando i numeri sono piccoli per il χ^2

Al.

		gene 1		
		0	1	
0	6	7	13	
	1	14	15	
		7	21	28

3,25	9,75
3,75	11,25

gli altri hanno 2 numeri < 5
non uso il chi-quadrato

Mettiamoci in H_0 : le due var sono indipendenti

Y		13	
		15	
	7	21	28

è equivalente immaginare un'urna con 28 palline, di cui 15 nere (malattia) e 13 bianche che vengono pescate a caso da 28 persone: 7 con gene 0 e 21 con gene 1; alla fine l'associazione tra i due caratteri è indipendente.

Sia Y il numero di gene 0 - sani: che distribuzione ha Y ?

$Y \sim \text{hypergeom}(7; 13; 28)$ *ipergeometrica*

1	12	13
6	9	15
7	21	28

Sia y il valore osservato per Y

$$\alpha^* = 2 \cdot \min(F_Y(y); 1 - F_Y(y-1))$$

SULLA GENERAZIONE DI V.V.A.

= CASUALE() genera $\text{unif}(0,1)$

= $F^{-1}(\text{CASUALE}())$ genera una v.a con Cdf F

↳ = INV.NORM(CASUALE(); media; devst)

↳ = INV.GAMMA(CASUALE(); $\alpha \circ 1$; β)

...

↳ = CRIT.BINOM(CASUALE(); n ; p)

* Poisson: non c'è INV.POIS Pois(ν) ν la media

= CRIT.BINOM(CASUALE(); 1000; $\nu/1000$)

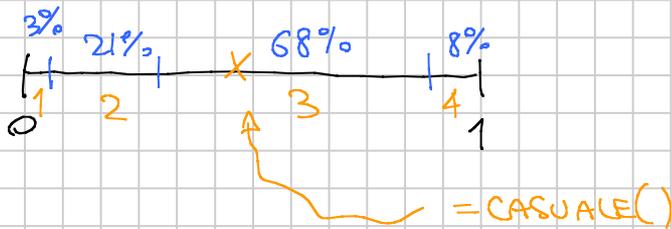
* Categoria generica

$n = 50$

1	no	3%
2	medie	21%
3	maturità	68%
4	lance	8%

2
11
30
7

$X_1, X_2, \dots, X_{50} \rightarrow \{1, 2, 3, 4\}$



= CASUALE() [A1]

= SE(A1 < 0,03; 1; SE(A1 < 0,03+0,21; 2; SE(A1 < 0,03+0,21+0,68; 3; 4)))

* Si ripete su 50 celle e si ottiene il campione X_i

* Si usa =CONTA.SE() per contare quante di ogni tipo nel campione

5/11