

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

Note Title

02/03/2011

Nicolodi + Morandin

↳ statistica 6 CFU

A Design of experiment
(Sleepers - Design for Six Sigma
Statistics - 59 tools for ...)

Capitolo 10

★ Middleton (Apogeo)

"Analisi statistiche con Excel"

★ Altre fonti:

- www.statsoft.com "statistics textbooks"

- dispense di Soliani ~3000 pagine

▣ Modalità di esame da definire

▣ Materiale → lea.unipr.it

scritti, lezioni (avi/pdf), forum, avvisi, ...

↳ anche vecchie, soprattutto lab

▣ Ricevimento: gio 13-14.30

▣ 10 laboratori di 2 ore

32 ore di lezione

★ Green belt, black belt, master black belt

Ross:

3-8 già fatti

7 intervalli di confidenza

8 test statistici★

9 regressione

10 analisi della varianza

11 i test del chi-quadro

★ Più curve OC e
simulazione MC

▣ REGRESSIONE

- R. lineare semplice

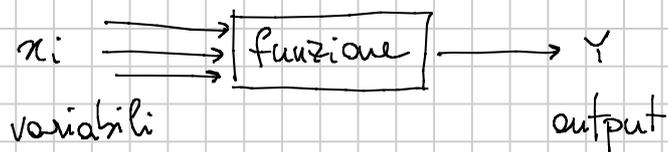
- R. lineare multiple

- R. nonlineare

- R. ... eteroschedastica

} omoschedastiche

+ DOF



dati ← esperimenti → dati

regressione → descrizione funzione

misura del rumore

ne tiene conto

REGRESSIONE LINEARE SEMPLICE

funzione (retta)

1 var di ingresso

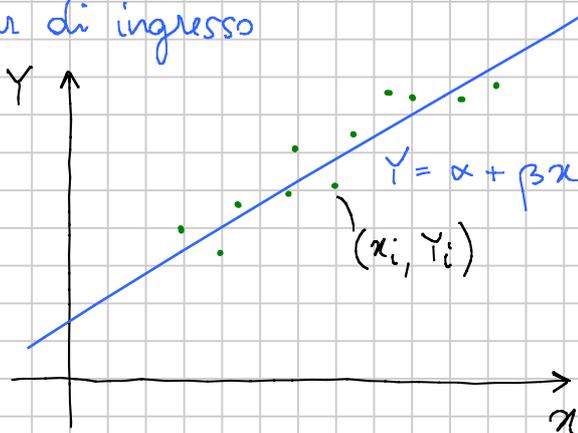
$(x_i, Y_i) \quad i=1, 2, \dots, n$

i dati su cui lavorare

per trovare la retta, ovvero α, β

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

errore



$$Y_1 = \alpha + \beta x_1 + e_1$$

$$Y_2 = \alpha + \beta x_2 + e_2$$

...

supponiamo che e_1, \dots, e_n siano

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

- indipendenti

σ non dipende da $i/x_i/Y_i$

OMOSCHEDASTICITÀ

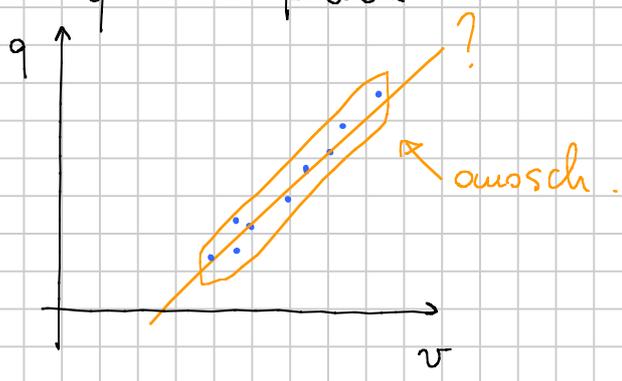
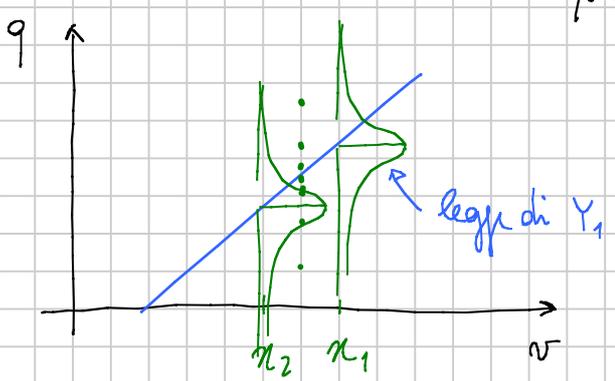
$$* Y_i = \alpha + \beta x_i + e_i \sim \mathcal{N}(\alpha + \beta x_i; \sigma^2)$$

Y_1, Y_2, \dots, Y_n sono indipendenti ma non con la stessa legge (tranne quando $\beta=0 \rightarrow$ si dice che non c'è regressione)

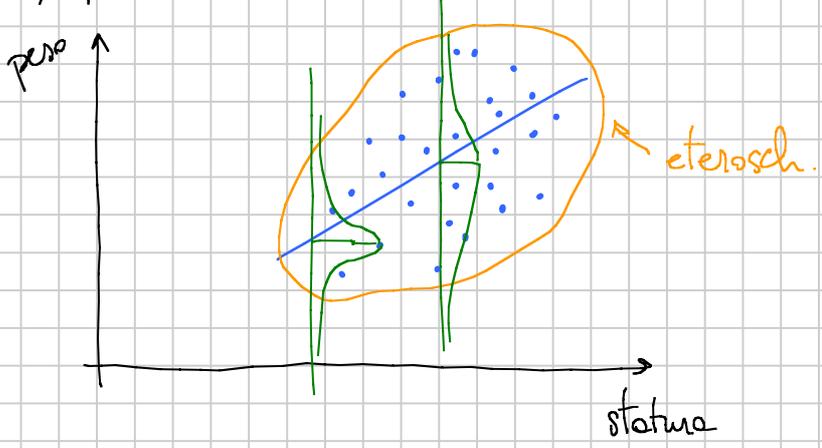
↳ quindi Y_i iid. e si usano tecniche elementari

Esempi

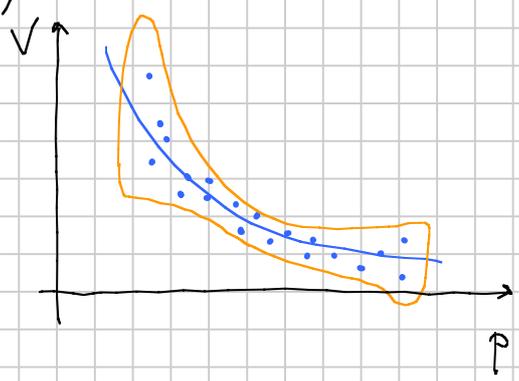
- 1) Ripartitore polveri - velocità v delle coclee
- quantità q della polvere



- 2) peso vs statura adulti italiani maschi



- 3) barometro elettronico - tensione V - pressione p



due variabili legate \rightarrow c'è regressione \rightarrow le var sono correlate

~~→~~ rapporto di causa / effetto

★ È possibile scambiare x e y $Y = \alpha + \beta x \Leftrightarrow x = -\frac{\alpha}{\beta} + \frac{1}{\beta} Y$

NB. Se si fanno le due regressioni si ottengono risultati (un po') diversi

★ Come scegliere?

- x causa, Y effetto
- x deterministica / precisa, Y casuale / rumorosa
- Y deterministica + rumore casuale

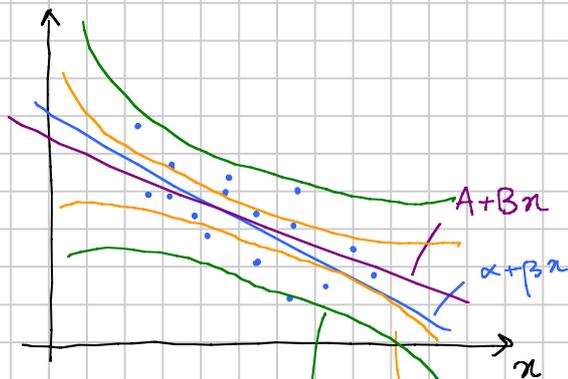
● Pratica:

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

α, β, σ incognite da stimare

$\downarrow \quad \downarrow \quad \downarrow$

$A \quad B \quad \text{se stimatori}$



$(x_i, Y_i) \quad i=1, \dots, n$

intervallo di confidenza per la risposta media

$$\alpha + \beta x \in A + Bx \pm \text{raggio}(x)$$

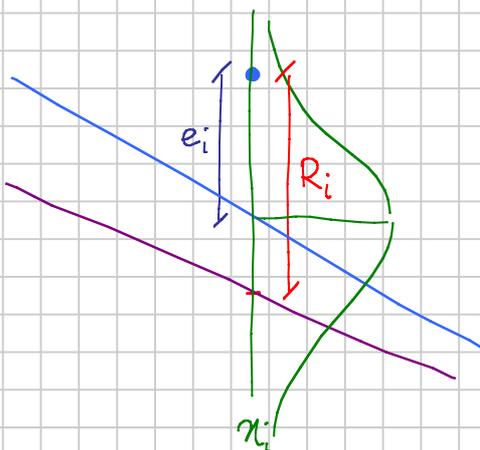
intervallo di predizione per una risposta futura

$$\alpha + \beta x + e \in A + Bx \pm \text{raggio_maggiore}(x)$$

● Come ritrovo la retta viola?

Si minimizza la somma dei quadrati dei residui

★ Residui o scarti sono l'errore misurato relativamente a $A+Bx$



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

può essere positivo o negativo

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

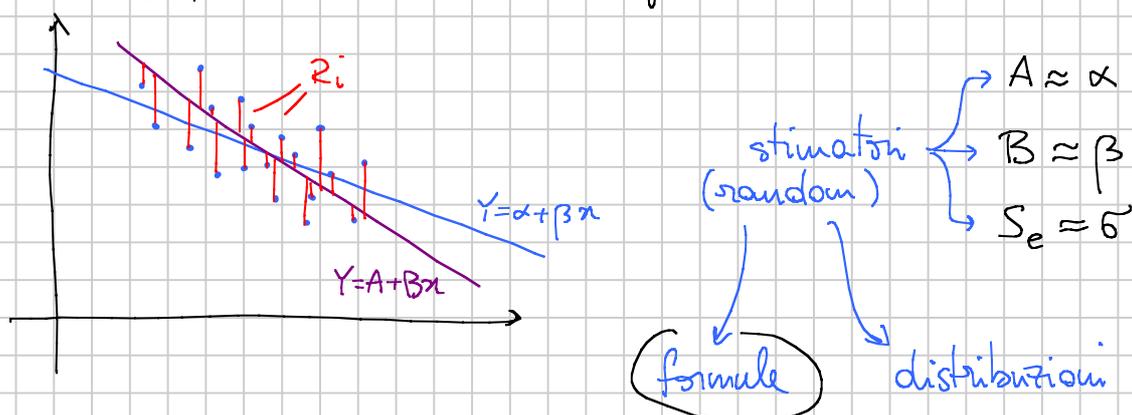
★ $SS_R = SS_R(A; B)$ funzione di A, B

HW: Trovare A, B che minimizzano SS_R (in funzione di x_i, Y_i)

REGR. LIN. SEMPLICE

$$Y_i = \alpha + \beta x_i + e_i \quad (x_i, Y_i), i=1, 2, \dots, n \quad e_i \sim \mathcal{N}(0, \sigma^2)$$

α, β, σ parametri incogniti



Troviamo A, B che rendono minima la somma dei quadrati

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

$$\begin{cases} 0 = \frac{\partial SS_R}{\partial A} = \sum_{i=1}^n 2(Y_i - A - Bx_i) \cdot (-1) \\ 0 = \frac{\partial SS_R}{\partial B} = \sum_{i=1}^n 2(Y_i - A - Bx_i) \cdot (-x_i) \end{cases} \begin{cases} \sum Y_i = \sum (A + Bx_i) \\ \sum x_i Y_i = \sum (A + Bx_i)x_i \end{cases}$$

$$\begin{cases} \sum Y_i = nA + B \sum x_i \\ \dots \end{cases} \begin{cases} \bar{Y} = A + B \bar{x} \\ n\bar{Y} = A n + B \sum x_i \end{cases} \text{ sistema delle equazioni normali}$$

$$\bar{xY} := \frac{1}{n} \sum_{i=1}^n x_i Y_i$$

$$\bar{x^2} := \frac{1}{n} \sum_{i=1}^n x_i^2$$

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

• Stimatore di σ

$$Y_i - \alpha - \beta x_i = e_i \sim \mathcal{N}(0, \sigma^2)$$

$$\frac{1}{n} \sum (Y_i - \alpha - \beta x_i)^2 \approx \sigma^2 \text{ ma non ho } \alpha, \beta$$

$$S_e^2 := \frac{1}{n-2} \sum_{i=1}^n (Y_i - A - Bx_i)^2 = \frac{SSR}{n-2}$$

errore standard

★ $\bar{Y} = A + B\bar{x}$: il punto (\bar{x}, \bar{Y}) sta sulla retta *viola*
 ↗ baricentro

★ $B = \frac{\overline{xY} - \bar{x}\bar{Y}}{\overline{x^2} - \bar{x}^2} = \frac{S_{xY}}{S_x^2}$ ← covarianza campionaria

$$S_x^2 = \frac{1}{n-1} \left[\sum x_i^2 - n\bar{x}^2 \right]$$

$$S_{xY} = \frac{1}{n-1} \left[\sum x_i Y_i - n\bar{x}\bar{Y} \right]$$

HW: cosa succede (e perché) se il den. è zero?

★ $Y_i - \alpha - \beta x_i \sim \mathcal{N}(0, \sigma^2)$ indipendenti

$R_i = Y_i - A - Bx_i$ NO (non hanno $\text{Var} = \sigma^2$ e non sono indep.)

$$S_R^2 = \frac{1}{n-1} \left[\sum R_i^2 - n\bar{R}^2 \right] = \frac{n-2}{n-1} S_e^2$$

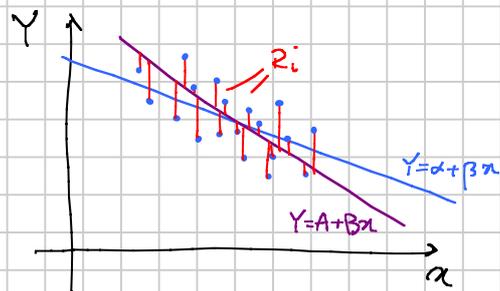
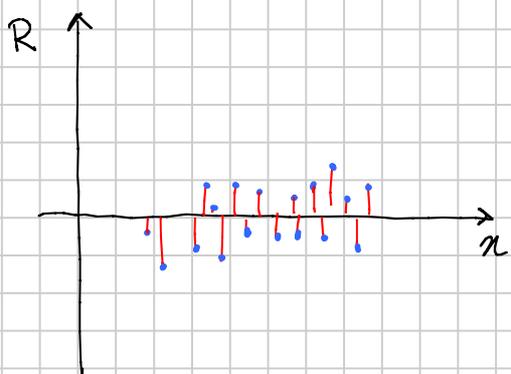
stimatore di σ^2
 che sottostima sistematicamente

↑ HW

★ Per come sono definiti i residui:

- $\bar{R} = 0$

- se faccio la regressione sui punti (x_i, R_i) trovo l'asse delle ascisse



DISTRIBUZIONE DI $A \in B$

$$B = \frac{\overline{\alpha Y} - \bar{\alpha} \bar{Y}}{\overline{\alpha^2} - \bar{\alpha}^2} = \frac{\frac{1}{n} \sum \alpha_i Y_i - \bar{\alpha} \frac{1}{n} \sum Y_i}{\overline{\alpha^2} - \bar{\alpha}^2} = \frac{1}{n(\overline{\alpha^2} - \bar{\alpha}^2)} \left[\sum \alpha_i Y_i - \bar{\alpha} \sum Y_i \right]$$

$$= \frac{1}{n(\overline{\alpha^2} - \bar{\alpha}^2)} \sum_{i=1}^n (\alpha_i - \bar{\alpha}) Y_i = \sum_{i=1}^n \underbrace{\frac{\alpha_i - \bar{\alpha}}{n(\overline{\alpha^2} - \bar{\alpha}^2)}}_{=: c_i} Y_i =: \sum_{i=1}^n c_i Y_i$$

i. B comb lineare di $Y_i \Rightarrow B \in \mathcal{N}$

ii. $E(B) = \beta$

(stimatore corretto)

iii. $\text{Var}(B) = \sigma^2 k_B$

(stimatore consistente)

$$E(X+Y) = E(X) + E(Y)$$

$$E(aX+bY) = aE(X) + bE(Y)$$

$$E(B) = E\left(\sum_{i=1}^n c_i Y_i\right) = \sum_{i=1}^n c_i E(Y_i)$$

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

$$= \sum c_i (\alpha + \beta x_i) = \frac{1}{n} \left[\sum (\alpha_i - \bar{\alpha}) \alpha + \sum (\alpha_i - \bar{\alpha}) \beta x_i \right]$$

$$\begin{aligned} \frac{1}{n} \sum \alpha_i &= \bar{\alpha} \\ \sum \alpha_i &= n\bar{\alpha} \\ \sum \bar{\alpha} &= n\bar{\alpha} \end{aligned}$$

$$= \beta \frac{\sum (\alpha_i^2 - \alpha_i \bar{\alpha})}{n(\overline{\alpha^2} - \bar{\alpha}^2)} = \beta \frac{\sum \alpha_i^2 - \bar{\alpha}^2 \cdot n}{n(\overline{\alpha^2} - \bar{\alpha}^2)} = \beta \frac{\overline{\alpha^2} - \bar{\alpha}^2}{\overline{\alpha^2} - \bar{\alpha}^2} = \beta$$

$$\begin{aligned} \text{Var}(X) &= \text{Cov}(X; X) \\ \text{Cov}\left(\sum a_i X_i; \sum b_j Y_j\right) &= \sum_i \sum_j a_i b_j \text{Cov}(X_i; Y_j) \end{aligned}$$

$$X, Y \text{ indipendenti} \Rightarrow \text{Cov}(X; Y) = 0$$

ora 4

$$\text{Var}(B) = \text{Cov}(B; B) = \text{Cov}\left(\sum_{i=1}^n c_i Y_i; \sum_{j=1}^n c_j Y_j\right) = \sum_i \sum_j c_i c_j \text{Cov}(Y_i; Y_j)$$

$$= \sum_{i=1}^n c_i^2 \text{Cov}(Y_i; Y_i) = \sum c_i^2 \text{Var}(Y_i) = \sigma^2 \sum_{i=1}^n c_i^2 =: \sigma^2 k_B$$

$$k_B := \sum_{i=1}^n c_i^2 = \sum_{i=1}^n \left[\frac{x_i - \bar{x}}{n(\bar{x}^2 - \bar{x}^2)} \right]^2 = \frac{1}{n^2(\bar{x}^2 - \bar{x}^2)^2} \sum (x_i - \bar{x})(x_i - \bar{x})$$

$$\frac{1}{n^2(\bar{x}^2 - \bar{x}^2)^2} \left[\underbrace{\sum (x_i - \bar{x}) x_i}_{= n(\bar{x}^2 - \bar{x}^2) \text{ come sopra}} - \underbrace{\sum (x_i - \bar{x}) \bar{x}}_{= 0 \text{ come sopra}} \right] = \frac{1}{n(\bar{x}^2 - \bar{x}^2)}$$

★ per $n \rightarrow \infty$ $k_B \rightarrow 0$

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

• Rifaccio per A

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

$$= \sum_{i=1}^n Y_i \left\{ \frac{1}{n} - \frac{x_i \bar{x} - \bar{x}^2}{n(\bar{x}^2 - \bar{x}^2)} \right\} = \sum_{i=1}^n Y_i \frac{\bar{x}^2 - \bar{x}^2 - x_i \bar{x} + \bar{x}^2}{n(\bar{x}^2 - \bar{x}^2)} = \sum_{i=1}^n d_i Y_i$$

$$d_i := \frac{\bar{x}^2 - x_i \bar{x}}{n(\bar{x}^2 - \bar{x}^2)}$$

HW: verificare che:

i. $A \sim \mathcal{N}$

ii. $E(A) = \alpha$

iii. $\text{Var}(A) = \sigma^2 k_A$

(stimatore corretto)

(stimatore consistente)

dove $k_A = \frac{\bar{x}^2}{n(\bar{x}^2 - \bar{x}^2)}$

★ $B \sim \mathcal{N}(\beta; \sigma^2 k_B)$

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

NB A, B non sono indipendenti

HW: $\text{Cov}(A; B) = -\sigma^2 \frac{\bar{x}}{n(\bar{x}^2 - \bar{x}^2)}$

$\text{Var}(A + B\bar{x}) = ?$

• DISTRIBUZIONE DI S_e

$$S_e^2 := \frac{SSR}{n-2} = \frac{1}{n-2} \sum_{i=1}^n (Y_i - A - Bx_i)^2$$

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

perché? Vedremo... per ora qualche analogia

a) $x_i \sim \mathcal{N}(\mu, \sigma^2)$ i.i.d.

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

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

$$\frac{n}{\sigma^2} \cdot \frac{1}{n} \sum (x_i - \mu)^2 = \frac{S_x^2}{\sigma^2} n \sim \chi^2(n)$$

$$\sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma} \right)^2 = \frac{n-1}{\sigma^2} \cdot \frac{1}{n-1} \sum (x_i - \bar{x})^2 = \frac{S_x^2}{\sigma^2} (n-1) \sim \chi^2(n-1)$$

b) $x_1 \dots x_n, y_1 \dots y_m$

$x_i \sim \mathcal{N}(\mu_x, \sigma^2)$ i.i.d.

$y_j \sim \mathcal{N}(\mu_y, \sigma^2)$ i.i.d.

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

$$S_y^2 = \frac{1}{m-1} \sum (y_i - \bar{y})^2$$

$$S_x^2 \approx \sigma^2$$

$$S_y^2 \approx \sigma^2$$

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

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

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

media pesata

$$= \frac{1}{n+m-2} \left[\sum_{i=1}^n (x_i - \bar{x})^2 + \sum_{j=1}^m (y_j - \bar{y})^2 \right]$$

$$\frac{S_p^2}{\sigma^2} (n+m-2) = \sum_i \left(\frac{x_i - \bar{x}}{\sigma} \right)^2 + \sum_j \left(\frac{y_j - \bar{y}}{\sigma} \right)^2 \sim \chi^2(n+m-2)$$

$$\sum_i \left(\frac{x_i - \mu_x}{\sigma} \right)^2 + \sum_j \left(\frac{y_j - \mu_y}{\sigma} \right)^2 \sim \chi^2(n+m)$$

(HW)

c) regressione

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

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

$$\frac{S_e^2}{\sigma^2} (n-2) = \frac{SSR}{\sigma^2} = \sum_{i=1}^n \left(\frac{Y_i - A - B x_i}{\sigma} \right)^2 \sim \chi^2(n-2)$$

★ A, Se indep

B, Se indep.

$$S_e := \sqrt{\frac{SSR}{n-2}}$$

$$S_e^2(n-2) = SSR$$

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

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

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

A, S_e indep. B, S_e indep

INFERENZA SUI PARAMETRI DI REGRESSIONE

funzione ancillare \rightarrow int. di conf.
 \rightarrow test statistici

Def. Si dice funz. ancillare di un parametro θ una v.a. che dipende da: θ , il campione, altre grandezze note di cui sia nota esattamente la distribuzione

\rightarrow Esempio: $X_i \sim \mathcal{N}(\mu, \sigma^2) \quad i=1, 2, \dots, n$ σ nota μ no

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

funz. anc. per μ ?

\bar{X} no perché la legge non è nota del tutto (μ)

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

Questa (volendo) è una funz. anc. per μ

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

Questa pure: è quella canonica

\rightarrow Esempio: stessa situazione, ma σ incognita

$\bar{X} - \mu$ no perché la legge non è nota del tutto (σ)

$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$ no, perché dipende anche da σ

$$\frac{\bar{X} - \mu}{S/\sqrt{n}}$$

campione noto

perché sia funz. anc. di μ occorre conoscere la distribuzione e che questa sia nota

Cos'è la t di Student?
 $Z \sim \mathcal{N}(0,1)$ $W \sim \chi^2(k)$ indip
 Allora $\frac{Z}{\sqrt{W/k}} \sim t(k)$
 (per definizione)

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

inoltre S e \bar{X} indipendenti

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

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

$$k := n-1$$

indipendenti:

$$Z \sqrt{\frac{k}{W}} = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \cdot \sqrt{\frac{n-1}{\frac{S^2}{\sigma^2}(n-1)}} = \frac{\bar{X} - \mu}{\cancel{\sigma}/\sqrt{n}} \cdot \frac{\cancel{\sigma}}{S}$$

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

• Torniamo alla regressione lineare semplice

- funz. ausiliarie per β e α :

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

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

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

$$k = n-2$$

indip.

HW: conto

$$\frac{B - \beta}{Se \sqrt{k_B}} \sim t(n-2)$$

funz. anc. per β

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

funz. anc. per α

HW: tutto

- funz. anc. per σ^2

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

- funz. ancillare per la retta azzerata: $\alpha + \beta x_0$
fissato un valore x_0 qualsiasi, $\alpha + \beta x_0$ è l'ordinata
della retta per quel valore di ascissa

i. Stimatore: $\alpha + \beta x_0 \approx A + B x_0$

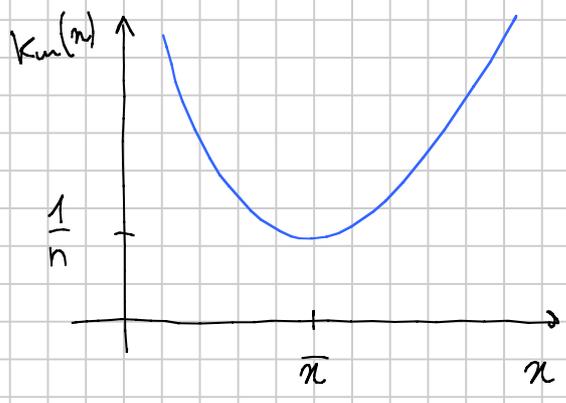
ii. Distribuzione $A + B x_0 \sim \mathcal{N}(\alpha + \beta x_0; \sigma^2 k_m(x_0))$

$$E(A + B x_0) = E(A) + x_0 E(B) = \alpha + \beta x_0$$

$$\begin{aligned} \text{Var}(A + B x_0) &= \text{Cov}(A + B x_0; A + B x_0) = \text{Cov}(A; A) + \\ &\quad \text{Cov}(A; B x_0) + \text{Cov}(B x_0; A) + \text{Cov}(B x_0; B x_0) \\ &= \text{Cov}(A; A) + x_0 \text{Cov}(A; B) + x_0 \text{Cov}(B; A) + x_0^2 \text{Cov}(B; B) \\ &= \text{Var}(A) + 2 x_0 \text{Cov}(A; B) + x_0^2 \text{Var}(B) \end{aligned}$$

$$\begin{aligned} &= \sigma^2 k_A - 2 x_0 \sigma^2 \frac{\bar{x}}{n(\bar{x}^2 - \bar{x}^2)} + x_0^2 \sigma^2 k_B = \sigma^2 \cdot \frac{\bar{x}^2 - 2 x_0 \bar{x} + x_0^2}{n(\bar{x}^2 - \bar{x}^2)} \\ &= \sigma^2 \frac{\bar{x}^2 - 2 x_0 \bar{x} + x_0^2 + \bar{x}^2 - \bar{x}^2}{n(\bar{x}^2 - \bar{x}^2)} = \sigma^2 \left(\frac{(x_0 - \bar{x})^2}{n(\bar{x}^2 - \bar{x}^2)} + \frac{1}{n} \right) =: \sigma^2 k_m(x_0) \end{aligned}$$

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



$$A + Bx_0 \sim \mathcal{N}(\alpha + \beta x_0; \sigma^2 k_m(x_0))$$

iii. funz. anc. $Z = \frac{A + Bx_0 - (\alpha + \beta x_0)}{\sigma \sqrt{k_m(x_0)}} \sim \mathcal{N}(0, 1)$

Se indep da A, da B e quindi da A + Bx_0 perciò

$$\frac{A + Bx_0 - (\alpha + \beta x_0)}{\sigma \sqrt{k_m(x_0)}} \sim t(n-2)$$

Funzione ancillare per la risposta media a livello di ingresso x_0

① INTERVALLI DI CONFIDENZA

- per la risposta media

i. fisso il livello di confidenza

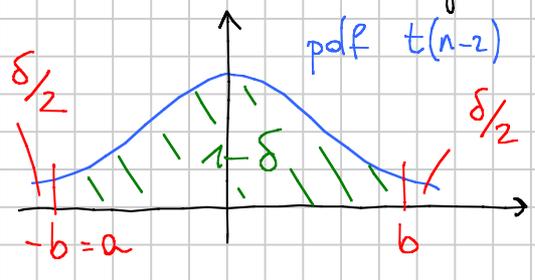
$$1 - \alpha$$

$1 - \delta$
0,95
0,90
0,99

ii. fisso la f. anc.

$$\frac{A + Bx_0 - (\alpha + \beta x_0)}{\sigma \sqrt{k_m(x_0)}} \sim t(n-2)$$

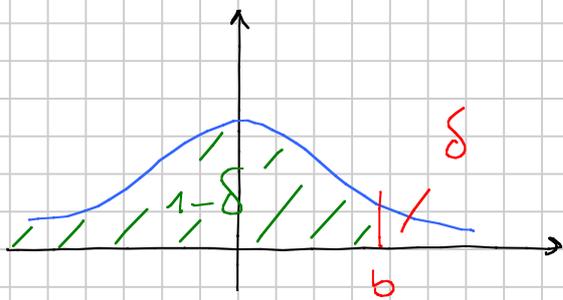
iii. trovo il/i quantili della legge relativa, per il livello di confidenza assegnato



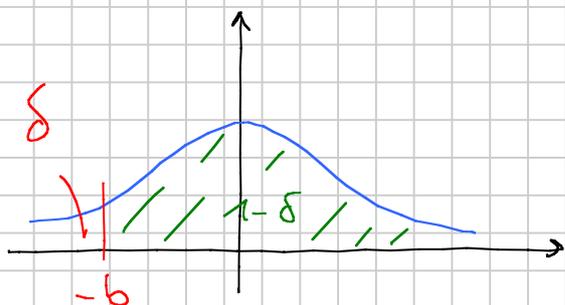
$$F_{t(n-2)}(b) := P(t(n-2) \leq b) = 1 - \frac{\delta}{2}$$

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

-b



$$b = F_{t(n-2)}^{-1}(1-\delta) = \text{INV.T}(2 \times \delta; n-2)$$



oppure $-b = -\text{INV.T}(2 \times \delta; n-2)$

iv. Scrivo $1-\delta = P(\dots)$ e ricavo il parametro $\alpha + \beta x_0$

$$1-\delta = P\left(-b \leq \frac{A + Bx_0 - (\alpha + \beta x_0)}{\underbrace{Se \sqrt{km(x_0)}}_{t(n-2)}} \leq b\right) = P\left(-b Se \sqrt{\quad} \leq A + Bx_0 - (\alpha + \beta x_0) \leq b Se \sqrt{\quad}\right)$$

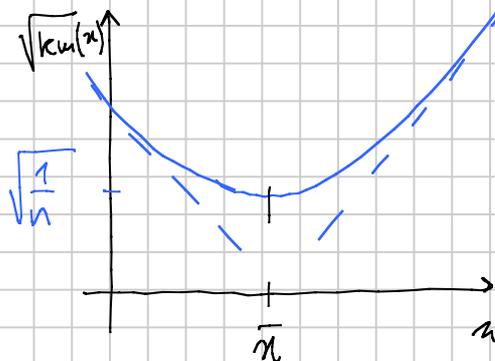
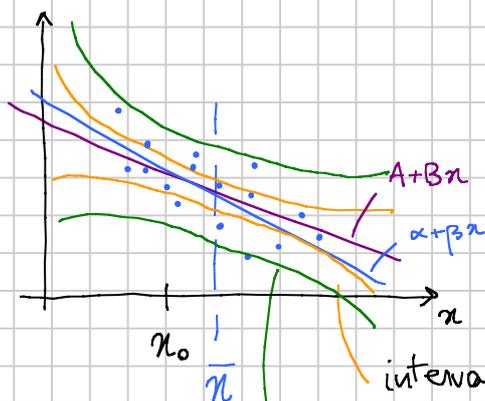
$$= P\left(- (A + Bx_0) - b Se \sqrt{\quad} \leq -(\alpha + \beta x_0) \leq - (A + Bx_0) + b Se \sqrt{\quad}\right)$$

$$= P\left(A + Bx_0 + b Se \sqrt{\quad} \geq \alpha + \beta x_0 \geq A + Bx_0 - b Se \sqrt{\quad}\right)$$

$$= P\left(A + Bx_0 - b Se \sqrt{\quad} \leq \alpha + \beta x_0 \leq A + Bx_0 + b Se \sqrt{\quad}\right)$$

$$\alpha + \beta x_0 \in \underbrace{A + Bx_0}_{\text{centro}} \pm \underbrace{b Se \sqrt{km(x_0)}}_{\text{raggio}} \quad \text{con lvl di conf } 1-\delta$$

v. Grafico



intervallo di confidenza per la risposta media
 $\alpha + \beta x \in A + Bx \pm \text{raggio}(x)$

intervallo di predizione per una risposta futura
 $\alpha + \beta x + e \in A + Bx \pm \text{raggio_maggiore}(x)$

INTERVALLI DI PREVISIONE

Simile all'intervallo di confidenza, ma invece che contenere un parametro incognito, contiene una v.a. futura

● Campione gaussiano classico

$$X_i \sim \mathcal{N}(\mu, \sigma^2) \quad \text{iid} \quad i = 1, 2, \dots, n \quad 1-\alpha \text{ liv di conf}$$

$$\rightarrow \text{int. conf medie} \quad \mu \in \bar{X} \pm q \frac{S}{\sqrt{n}} \quad q = F_{t(n-1)}^{-1} \left(1 - \frac{\alpha}{2} \right)$$

$$\rightarrow \text{int. di predizione per } X_{n+1} \quad X_{n+1} \in \bar{X} \pm q S \sqrt{\frac{1}{n} + 1}$$

$$\bar{X} - X_{n+1} \sim \mathcal{N}\left(0, \frac{\sigma^2}{n} + \sigma^2\right) \sim \mathcal{N}\left(0, \sigma^2 \left(\frac{1}{n} + 1\right)\right)$$

$\mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$ $\mathcal{N}(\mu, \sigma^2)$ HW: ripasso, approfondire
 indip perché X_{n+1} è nuovo

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

non è funz. ancillare ma si usa similmente

$$\frac{\bar{X} - X_{n+1}}{S \sqrt{\frac{1}{n} + 1}} \sim t(n-1) \quad \text{HW: verificare}$$

$$q \text{ quantile t. c.} \quad 1-\alpha = P(-q \leq t(n-1) \leq q)$$

$$1-\alpha = P\left(-q \leq \frac{\bar{X} - X_{n+1}}{S \sqrt{\frac{1}{n} + 1}} \leq q\right) \stackrel{\text{HW}}{=} P\left(\bar{X} - q S \sqrt{\frac{1}{n} + 1} \leq X_{n+1} \leq \bar{X} + q S \sqrt{\frac{1}{n} + 1}\right)$$

● Regressione

$$Y_i = \alpha + \beta x_i + e_i \quad e_i \sim \mathcal{N}(0, \sigma^2) \text{ iid } i = 1, 2, \dots, n$$

$A + Bx$ retta di regressione

Se $\approx \sigma$ errore standard

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

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

(x_{n+1}, Y_{n+1}) punto futuro

x_{n+1} noto
 ↑
 180 nell'esempio

$$Y_{n+1} = \alpha + \beta x_{n+1} + e_{n+1} \sim \mathcal{N}(\alpha + \beta x_{n+1}; \sigma^2)$$

$$A + Bx_{n+1} - Y_{n+1} \sim \mathcal{N}(0; \sigma^2 k_m(x_{n+1}) + \sigma^2) \sim \mathcal{N}(0; \sigma^2 (k_m(x_{n+1}) + 1))$$

$\underbrace{\mathcal{N}(\alpha + \beta x_{n+1}; \sigma^2 k_m(x_{n+1}))}_{\text{indip}}$
 $\underbrace{\mathcal{N}(\alpha + \beta x_{n+1}; \sigma^2)}_{\text{indip}}$

$$\frac{A + Bx_{n+1} - Y_{n+1}}{\sigma \sqrt{k_m(x_{n+1}) + 1}} \sim \mathcal{N}(0, 1)$$

1- δ livello di confidenza

$$\frac{A + Bx_{n+1} - Y_{n+1}}{S_e \sqrt{k_m(x_{n+1}) + 1}} \sim t(n-2)$$

HW: verificare

q quantile t.c. $1 - \delta = P(-q \leq t(n-2) \leq q)$

...

$$Y_{n+1} \in A + Bx_{n+1} \pm q S_e \sqrt{k_m(x_{n+1}) + 1}$$

● Approfondimento: int. di predizione quando i dati futuri sono tanti

→ domande di una bitte: prevedere tutto il mese prossimo non si possono fare intervalli singoli e "metterli assieme" in qualche modo.

$$\underbrace{X_1, X_2, \dots, X_n}_{\text{passato}} \quad \underbrace{X_{n+1}, \dots, X_{n+m}}_{\text{futuro}} \quad \text{iid } \mathcal{N}(\mu, \sigma^2)$$

$$T := \sum_{i=1}^m X_{n+i} \quad \text{somma}$$

$$U := \frac{1}{m} T \quad \text{media}$$

$$T \sim \mathcal{N}(m\mu; m\sigma^2)$$

$$U \sim \mathcal{N}\left(\mu; \frac{\sigma^2}{m}\right)$$

HW: check

HWC

$$T \in ? \quad m\mu \approx m\bar{X} \quad \bar{X} = \frac{1}{n} \sum_{i=1}^n x_i \quad m\bar{X} \sim \mathcal{N}\left(m\mu; \frac{m^2}{n}\sigma^2\right)$$

$$m\bar{X} - T \sim \mathcal{N}\left(0, \frac{m^2}{n}\sigma^2 + m\sigma^2\right) \sim \mathcal{N}\left(0, m^2\sigma^2\left(\frac{1}{n} + \frac{1}{m}\right)\right)$$

....

$$T \in m\bar{X} \pm q S m \sqrt{\frac{1}{n} + \frac{1}{m}}$$

HWC

quantile $t(n-1)$

$$\bar{X} - U \sim \mathcal{N}\left(0; \frac{\sigma^2}{n} + \frac{\sigma^2}{m}\right) \sim \mathcal{N}\left(0; \sigma^2\left(\frac{1}{n} + \frac{1}{m}\right)\right)$$

...

$$U \in \bar{X} \pm q S \sqrt{\frac{1}{n} + \frac{1}{m}}$$

HWC

ora 8

→ int. conf media

$$\mu \in \bar{X} \pm q \frac{S}{\sqrt{n}}$$

→ int. di predizione per X_{n+1}

$$X_{n+1} \in \bar{X} \pm q S \sqrt{\frac{1}{n} + 1}$$

L'int. di predizione per U generalizza entrambe ($m=1, \infty$)

★ L'int. di pred. per T è semplicemente quello di U "espanso" di un fattore m :

$$\rightarrow U \in [2,2; 2,5] \quad m=30 \text{ gg} \quad \rightarrow T \in [66; 75]$$

- Stessa cosa si può fare per la regressione (cfr i due anni precedenti per avere esempi di T e U con $x_{n+1} = x_{n+2} = \dots = x_{n+m}$ costante)

TEST STATISTICI SUI PARAMETRI DI REGRESSIONE

Esempio

Acceptable Quality Level

- β parametro incognito
- β_0 parametro di confronto (in altri casi target / AQL) di solito $\beta_0 = 0$ nella regressione
- ipotesi bilaterali o unilaterali (e nel I caso, H_0/H_1)

$$H_0: \beta = 0$$

"Y non dipende da x"
(non c'è regressione)

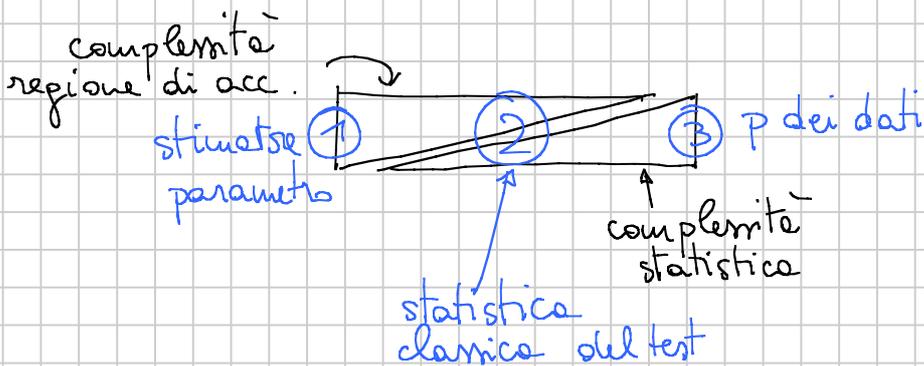
$$Y = \alpha + \beta x + e$$

$$H_1: \beta \neq 0$$

c'è regressione
conoscere x ci aiuta a prevedere Y

Nell'esempio, p dei dati 10^{-81} quindi nettamente H_1

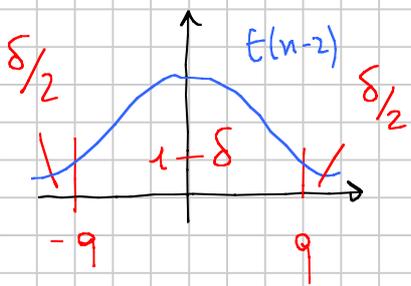
iv.



② stat. del test : funz. ancillare e cambio β con β_0

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

- Trovo i quantili bi- o uni-laterali per la distrib. $t(n-2)$ (per lvl di significatività δ assegnato)



$$\pm q$$

$$q = \text{INV.T}(\delta; n-2)$$

vi. Se H_0 è vera $T_B \sim t(n-2)$ quindi $T_B \in [-q; q]$ quasi sempre
 Allora, se $T_B \in [-q; q]$ dico H_0 , altrimenti dico H

$$RA_{T_B} := [-q; q] \quad \text{regione di accettazione}$$

① stimatore $T_B \rightarrow B$ statistica $RA_B = ?$

$$\text{vii. } T_B \in RA_{T_B} \Leftrightarrow \frac{B}{S_e \sqrt{k_B}} \in [-q; q] \Leftrightarrow B \in [-q S_e \sqrt{k_B}; q S_e \sqrt{k_B}]$$

$$RA_B := [-q S_e \sqrt{k_B}; q S_e \sqrt{k_B}]$$

③ p dei dati α^* (o δ^* ?)

$$T_B \in RA_{T_B} \Leftrightarrow T_B \in [-q; q] \Leftrightarrow |T_B| \leq q = F_{t(n-2)}^{-1}\left(1 - \frac{\delta}{2}\right)$$

$$\Leftrightarrow F_{t(n-2)}(|T_B|) \leq 1 - \frac{\delta}{2} \Leftrightarrow \frac{\delta}{2} \leq 1 - F_{t(n-2)}(|T_B|)$$

F monotona
crescente

$$\Leftrightarrow \delta \leq 2 - 2F_{t(n-2)}(|T_B|) =: \delta^*$$

p dei dati

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

$$RA_{\delta^*} := [\delta; 1]$$

★ Se non ho fissato il liv di significatività ha senso
 a maggior ragione calcolare i pdd: se è molto piccolo $\sim 10^{-4}$
 si sceglie H_1 ; se è molto grande $\sim 30\%$ si sceglie H_0 ;
 negli altri casi si discute

⊙ Altri test : $H_0: \alpha = 0$; b_1 /unilaterali in σ ; b_1 /unilaterali in $\alpha + \beta x_0$
 analoghi \rightarrow non approfondiamo

▣ VALUTAZIONE DEL MODELLO

- \rightarrow Ci sono errori nel modello o nelle ipotesi iniziali?
- \rightarrow E' un modello utile? Quanti € vale conoscere x ?

▣ ERRORE STANDARD E COEFFICIENTE DI DETERMINAZIONE

★ $Se := \frac{\sqrt{SSR}}{\sqrt{n-2}} \approx \sigma$ margine di errore in Y nota x

\rightarrow ha le stesse unità di misura di Y
 9,3 kg nell'esempio è l'incertezza ineliminabile in Y conoscendo "solo x "

★ Coeff. di determinazione

$R^2 := 1 - \frac{SSR}{S_{YY}}$ colore percentuale di incertezza sulle Y_i che si ottiene "aggiungendo" le x_i

$SSR = \sum_i (Y_i - A - Bx_i)^2$ incertezza residua, note le x_i

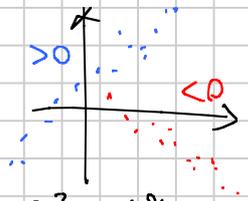
$S_{YY} = \sum_i (Y_i - \bar{Y})^2$ incertezza sulle Y , noto niente

$SSR \leq S_{YY}$

\rightarrow Più grande è migliore è la regressione, e più utili sono le x_i

\rightarrow Nel caso della regr. lin. semplice $R^2 = r^2$
 dove r è il coeff di correl. campionaria (cfr Cap 2)

$-1 \leq r \leq 1$ stimatore di $\rho(x; Y) = \frac{Cov(X; Y)}{\sqrt{Var(X) Var(Y)}}$



★ Se due var hanno $r = 0,4$ non è un gran che, perché $R^2 = 16\%$

ANALISI DEI RESIDUI

La regressione è un modello corretto se davvero:

$$Y_i = \alpha + \beta x_i + e_i$$

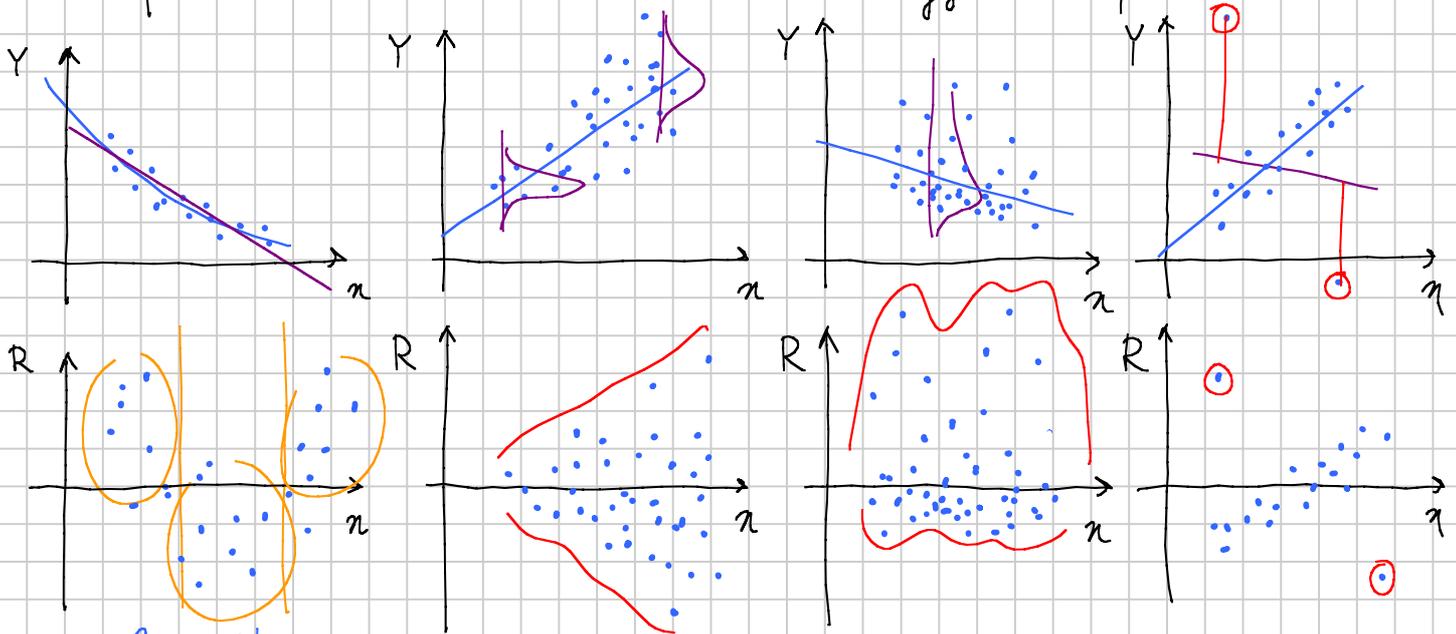
funzione lineare

$$e_i \sim \mathcal{N}(0, \sigma^2) \text{ indep.}$$

errori normali

omoschedasticità

★ Occorre diagnosticare se ci sono problemi rispetto a questi aspetti. Successivamente vedremo come correggerli in qualche caso



nonlinearità

$$Y = f(x) + e$$

f non lineare

eteroschedasticità

$$\sigma = \sigma(x)$$

$$\text{Var}(e_i) = \sigma^2(x_i)$$

errori non normali

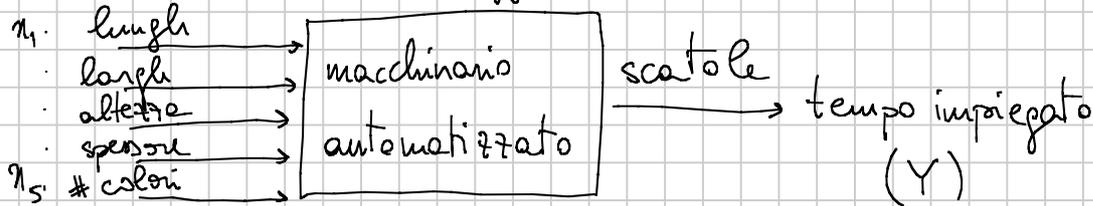
presenza di outliers

HW: Ripassare le matrici (prodotto riga / colonna ; inversa ; determinante ...)

REGRESSIONE LINEARE MULTIPLA

→ tecnica multivariata

→ esempio: imballaggi di cartone



$$Y = f(x_1, x_2, x_3, x_4, x_5) + e$$

* Non occorre a priori ragionare su quali variabili x_i potrebbero non incidere su Y : si fa a posteriori.

* Regressione: si fanno n esperimenti $i=1, 2, \dots, n$ con configurazioni variabili $x_j(i)$ $j=1, 2, \dots, p$ ($p=5$ nell'esempio)

$x_j(i) =: x_{ij}$: rettaggio della var. ingr. j nell'esperimento i

→ si ricavano n valori della var. di risposta

$Y(i) =: Y_i$: risposta nell'esperimento i

→ da tutti questi dati stimo il modello $Y = f(\vec{x}) + e$ da usare in futuro e per l'inferenza su qualunque parametro

* Ipotesi di Gauss:

i. modello lineare: $f(x_1, \dots, x_p) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$
 a volte suppongo di definire $x_0 \equiv 1$, con: $f(\vec{x}) = \sum_{j=0}^p \beta_j x_j$

ii. modello omoschedastico: $Y(i) = f(x_1(i), x_2(i), \dots, x_p(i)) + e(i)$
 $e(i) \sim \mathcal{N}(0, \sigma^2)$ σ^2 costante, non dipende da i

iii. $e(i)$ indipendenti

$$Y(i) = \sum_{j=0}^p \beta_j x_j(i) + e(i)$$

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

● Stimatori di β_j sono B_j $j=0, 1, \dots, p$

Residui : $R(i) = Y(i) - \hat{Y}(i) = Y(i) - \sum_{j=0}^p B_j x_j(i)$

Previsi : $\hat{Y}(i) = \sum_{j=0}^p B_j x_j(i) = B_0 + B_1 x_1(i) + B_2 x_2(i) + \dots$

$A + Bx(i)$ nella r. lin. semplice

★ Notazione matriciale

→ X matrice $n \times (p+1)$ $X_{ij} = x_{ij} = x_j(i)$
(occhio che $i=1, 2, \dots, n$ ma $j=0, 1, \dots, p$)

$$X = \begin{pmatrix} 1 & x_1(1) & x_2(1) & \dots & x_p(1) \\ 1 & x_1(2) & & & \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_1(n) & x_2(n) & & x_p(n) \end{pmatrix} \quad Y = \begin{pmatrix} Y(1) \\ \vdots \\ Y(n) \end{pmatrix} \quad e = \begin{pmatrix} e(1) \\ \vdots \\ e(n) \end{pmatrix}$$

→ Y, e vettori $n \times 1$

→ β e B vettori $(p+1) \times 1$

$$\beta = (\beta_0, \beta_1, \dots, \beta_p)^T$$

$$B = (B_0, B_1, \dots, B_p)^T$$

→ $Y = XB + e$

HWC

R vettore $n \times 1$ dei residui

$$R = Y - XB$$

→ $SS_R = \sum_{i=1}^n R_i^2 = \sum_{i=1}^n (Y(i) - XB(i))^2 = \|R\|^2$ scalare vettore

→ Equazioni normali (minimizzare SS_R in funzione di B)

$$\begin{aligned}
 [\nabla_B SS_R]_j &= \frac{\partial}{\partial B_j} \sum_{i=1}^n (Y_i - [XB]_i)^2 = 2 \sum_{i=1}^n (Y_i - [XB]_i) \frac{\partial}{\partial B_j} (Y_i - \sum_k X_{ik} B_k) \\
 &= -2 \sum_{i=1}^n (Y_i - [XB]_i) X_{ij} = -2 \left[\sum_i Y_i X_{ij} - \sum_i [XB]_i X_{ij} \right] \\
 &= -2 [Y^T X - [XB]^T X]_j = -2 [Y^T X - B^T X^T X]_j = 0
 \end{aligned}$$

$$B^T X^T X = Y^T X \quad (X^T X)^T B = X^T Y$$

$$\rightarrow (X^T X)^T = X^T X^{TT} = X^T X \quad \text{quindi } X^T X \text{ è simmetrica}$$

(p+1) × (p+1)

$$X^T X B = X^T Y \quad \text{equazioni normali}$$

$$B = (X^T X)^{-1} X^T Y \quad \text{se } X^T X \text{ è invertibile, posso ricavare } B$$

★ $X^T X$ è invertibile se X ha rango massimo (ovvero $p+1$)
 cosa che è quasi sempre vera, soprattutto se le x_{ij} sono misure del mondo reale

★ Se $\det(X^T X) = 0$ non si può procedere ma se $\det(X^T X)$ è anche solo piccolo la formula per B è malcondizionata, numericamente instabile, quindi errori approssimazioni e incertezze anche piccole su Y vengono amplificate su B

ora 11

● Funzioni da usare su Excel

- MATR.PRODOTTO (mat1; mat2)

- MATR.INVERSA (mat1)

- MATR.TRASPOSTA (mat1)

- MATR.SOMMAPRODOTTO (vet1; vet2)

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

(MMULT)

(MINVERSE)

(TRANSPOSE)

(SUMPRODUCT)

} Ctrl-shift-invio

prodotto scalare

STIMATORE PER σ^2

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

DISTRIBUZIONE STIMATORI

$$\begin{aligned} \bullet \quad \frac{S_e^2}{\sigma^2} (n-p-1) &= \frac{SSR}{\sigma^2} = \sum_{i=1}^n \left(\frac{Y_i - \sum_j B_j X_{ij}}{\sigma} \right)^2 \sim \chi^2(n-p-1) \\ &\quad \downarrow \mathcal{N}(0, \sigma^2) \text{ indep} \\ &= \sum_{i=1}^n \left(\frac{e(i)}{\sigma} \right)^2 \sim \chi^2(n) \end{aligned}$$

inoltre S_e indipendente da B
(tutto per la solita versione del thm di Cochran che vedremo più avanti)

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

S_e^2 stimatore corretto di σ^2
HWC

• Distribuzione di $B = (B_0, B_1, \dots, B_p)^T$

$$(X, Y) \quad f_{X,Y}(s, t)$$

$$B \quad f_B(s_0, s_1, \dots, s_p) = \left(2\pi \right)^{-\frac{p+1}{2}} \exp \left\{ \dots \right\}$$

Non occorre scrivere f_B . Ci basta osservare che
→ ciascuno dei B_j ha legge normale

$$B_j = \left[(X^T X)^{-1} X^T Y \right]_j = \sum_{k=1}^n \underbrace{\left[(X^T X)^{-1} X^T \right]_{jk}}_N Y_k = \sum_k N_{jk} Y_k$$

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

matrice $(p+1) \times n$

numeri

→ quindi B ha legge normale multivariata
per fissare la legge di B bisogna conoscere:

i. vettore delle medie

ii. matrice di covarianza

deterministico

$$\star E(B) = E((X^T X)^{-1} X^T Y) = E(NY) = NE(Y) = NE(X\beta + e)$$

$$= N(X\beta + 0) = NX\beta = (X^T X)^{-1} X^T X \beta = \beta$$

quindi B è uno stimatore corretto di β

\star Matrice di covarianza Σ matrice $(p+1) \times (p+1)$

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

i. Σ è simmetrica

$$\text{ii. } \Sigma_{ii} = \text{Cov}(B_i; B_i) = \text{Var}(B_i)$$

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

iv. Sia $V = \sum_j \alpha_j B_j = \alpha \cdot B$ una combinazione lineare dei B_j

$$(\alpha_0, \alpha_1, \dots, \alpha_p)^T = \alpha$$

Allora $\text{Var}(V) = \alpha^T \Sigma \alpha$, ovviamente $E(V) = \alpha \cdot E(B) = \alpha \cdot \beta$

$$(1) \times (p+1) \times (p+1) \times (p+1) \times (1) = (1, 1)$$

Per giustificare iii. e iv. diamo la proprietà più generale

X, Y vettori aleatori

$$X \in \mathbb{R}^m, Y \in \mathbb{R}^n$$

$$S = \text{Cov}(X; Y)$$

$$S_{ij} = \text{Cov}(X_i; Y_j) \quad S \text{ è } m \times n$$

M matrice $r \times m$

N matrice $s \times n$

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

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

$$= \sum_k \sum_l M_{ik} S_{kl} N_{lj}^T = [MSN^T]_{ij}$$

$$\text{Cov}(MX; NY) = MSN^T$$

$$\begin{aligned} \star \Sigma &= \text{Cov}(B; B) = \text{Cov}(NY; NY) \\ &= N \text{Cov}(Y; Y) N^T = N \sigma^2 I N^T \\ &= \sigma^2 N N^T = \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} \\ &= \sigma^2 (X^T X)^{-1} \end{aligned}$$

iii. verificato

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

Y_i indipendenti
 $Y_i \sim \mathcal{N}(X\beta; \sigma^2)$
 $\text{Cov}(Y; Y) = \sigma^2 I$

$$= \begin{pmatrix} \sigma^2 & & & \\ & \sigma^2 & & \\ & & \ddots & \\ 0 & & & \sigma^2 \end{pmatrix}$$

$$\star V = \alpha \cdot B = \alpha^T B \quad \alpha^T \text{ matrice } 1 \times (p+1)$$

$$\text{Var}(V) = \text{Cov}(V; V) = \text{Cov}(\alpha^T B; \alpha^T B) = \alpha^T \text{Cov}(B; B) \alpha = \alpha^T \Sigma \alpha$$

iv. verificato

ora 12

▣ TUTTE LE DISTRIBUZIONI DELLA REGRESSIONE LIN. MULTIPLA

$$Y_i \sim \mathcal{N}\left(\sum_j \beta_j x_{ij}, \sigma^2\right) \text{ indipendenti}$$

- Y vettore $\mathcal{N}(X\beta, \sigma^2 I)$

- B vettore $\mathcal{N}(\beta, \Sigma)$ stimatore β

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

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

- $\tilde{x} = (1, \tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_p)^T$ un punto fissato per le var di ingresso

$\tilde{Y} = B \cdot \tilde{x}$ valore previsto per la risposte Y a lvl di ingresso \tilde{x}

$\tilde{Y} \sim \mathcal{N}(\beta \cdot \tilde{x}; \tilde{x}^T \Sigma \tilde{x})$ si usa per stimare risposte medie e future
 è indipendente da Se

• Inferenza:

- su $\beta_j \approx B_j$ $\frac{B_j - \beta_j}{\sigma \sqrt{[(X^T X)^{-1}]_{jj}}} \sim \mathcal{N}(0,1)$

$$\frac{B_j - \beta_j}{S_e \sqrt{[(X^T X)^{-1}]_{jj}}} \sim t(n-p-1)$$

funz. ancillare per β_j

- su $\sigma \approx S_e$

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

f. anc. per σ

★ Permettono di eseguire tutti i test e gli intervalli di confidenza su β_j e σ

- Sulla risposta futura $\tilde{x} = (1, \tilde{x}_1, \dots, \tilde{x}_p)^T$

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

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

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

$$\frac{\tilde{Y} - \beta \cdot \tilde{x}}{\sigma \sqrt{1 + \tilde{x}^T (X^T X)^{-1} \tilde{x}}} \sim \mathcal{N}(0,1)$$

$$\frac{\tilde{Y} - \beta \cdot \tilde{x}}{S_e \sqrt{1 + \tilde{x}^T (X^T X)^{-1} \tilde{x}}} \sim t(n-p-1)$$

$$\rightarrow P(\tilde{Y} \in \beta \cdot \tilde{x} \pm q S_e \sqrt{1 + \tilde{x}^T (X^T X)^{-1} \tilde{x}}) = 1 - \delta \quad \text{dove } q = F_{t(n-p-1)}^{-1} \left(1 - \frac{\delta}{2} \right)$$

questo esempio è il caso bilaterale

- Sulla risposta media $\beta \cdot \tilde{x} \approx B \cdot \tilde{x}$

$$\frac{B \cdot \tilde{x} - \beta \cdot \tilde{x}}{S_e \sqrt{\tilde{x}^T (X^T X)^{-1} \tilde{x}}} \sim t(n-p-1)$$

→ Il p dei dati del test $\beta_j = 0$ dipende tantissimo da quali altre var sono nel modello

→ Questo accade in particolare quando le var in questione sono tra loro correlate.

★ Quando si usa Excel, la cosa più ragionevole è partire con tutte le variabili e andare per passi, togliendo di volta in volta una sola variabile tra quelle non significative. Di solito si sceglie quella con p dei dati maggiore, ma se c'è tempo non è male fare diversi tentativi

Questa tecnica si chiama stepwise backward e si arresta quando tutte le var residue risultano significative

Soglie ragionevoli per dire che α^* è significativo sono:

~ 30%

~~~ 5%~~

$$\sim 1/100 = \frac{5\%}{p+1}$$

per modello predittivo (butta solo le var che non sembrano proprio dire niente)

per sapere quali var impattano di ricuso (tengo solo le var sicure)

## REGRESSIONE MULTIPLA : BONTÀ DEL MODELLO

| id | $x_0$ | $x_1$ | $x_2$ | ... | $x_p$ | $Y$ |
|----|-------|-------|-------|-----|-------|-----|
| 1  | 1     | ~     | ~     |     | ~     | ~   |
| 2  | 1     | ~     | ~     |     | ~     | ~   |
| ⋮  | ⋮     | ⋮     | ⋮     |     | ⋮     | ⋮   |
| ⋮  | ⋮     | ⋮     | ⋮     |     | ⋮     | ⋮   |
| ⋮  | ⋮     | ⋮     | ⋮     |     | ⋮     | ⋮   |
| n  | 1     | ~     | ~     |     | ~     | ~   |

$X$                        $Y$

$B$

|   |
|---|
| ~ |
| ~ |
| ⋮ |
| ⋮ |
| ~ |

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

$p+1$  coefficienti

$$(B_0, B_1, \dots, B_p)$$

### Selezione delle variabili

#### a) Stepwise backward :

- i. inizio con tutte,
- ii. faccio la regressione,
- iii. tolgo la peggiore variabile se è il caso
- iv. se l'ho tolta, torno a ii. altrimenti ho finito

ora 15

$$SSR = Y^T (Y - XB)$$

HWC

(c'è qualcosa di simile sul Ross)

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

va tolta se  $\beta_1 = 0$  o comunque molto piccolo

esiste un test  $H_0: \beta_1 = 0$      $H_1: \beta_1 \neq 0$

Excel dà automaticamente il  $p$  dei dati di ciascuna

$$\left( \begin{array}{l} \text{Oppure} \\ \text{stat} = \frac{B_1}{\text{Se} \sqrt{[X^T X^{-1}]_{1,1}}} \end{array} \right) \text{DISTRIB. T (ASS (stat); } n-p-1) \text{ HWC}$$

$\alpha_0^*, \alpha_1^*, \alpha_2^*, \dots, \alpha_p^*$  più è alto, più sono convinto di buttare una variabile

$\star$  Di solito si butta quella con  $\alpha^*$  maggiore, ma non è detto che sia ottimale come scelta

(Sotto  $H_0$   $\alpha^* \sim \text{unif}[0,1]$ )

\* Soglia per fermarsi : quando tutti gli  $\alpha_i^*$  residui sono inferiori a :

- se voglio un modello predittivo da usare su casi futuri

soglia intorno al 30%

- se voglio giudicare su quali variabili agire per ottenere effetti su  $Y$

soglia intorno a  $\frac{5\%}{p+1} \approx 1\% - 1\%$

### ◎ PARENTESI : LA CORREZIONE DI BONFERRONI

$$H_0: \beta_0 = 0 \quad H_1: \beta_0 \neq 0$$

$$H_0: \beta_1 = 0 \quad H_1: \beta_1 \neq 0$$

...

$$H_0: \beta_p = 0 \quad H_1: \beta_p \neq 0$$

$p+1$  test

che vanno fatti contemporaneamente

$\alpha_0^*$   
 $\alpha_1^*$   
 $\vdots$   
 $\alpha_p^*$

Voglio un lvl di sign del 5%

la prob di dire  $H_1$  quando non è vera  $\leq 5\%$

Se faccio i test al 5% ho sistematicamente dei falsi  $H_1$

Se abbasso la sign dei singoli test a  $\frac{5\%}{p+1}$ , quello complessivo risulta  $\leq 5\%$

b) Stepwise forward :

i. si fanno  $p$  regressioni lineari semplici  $Y = \beta_0 + \beta_i x_i + e$

ii. tengo la variabile risultata migliore ( $\alpha^*$  più piccolo) (sia  $k$ )

iii. si fanno  $p-1$  regr. con due variabili

$$Y = \beta_0 + \beta_k x_k + \beta_i x_i + e \quad i = 1, 2, \dots, p \quad \text{tranne } i = k$$

iv. tengo la var migliore come seconda dopo  $x_k$

v. con via fino a che  $\alpha^* < \text{soglia}$

★ Unrealistico farlo con Excel, tranne i panni i.e.ii. che si riescono ad automatizzare usando bene il comando REGR.LIN

c) Cosa succede a  $SS_R$ ,  $R^2$ ,  $Se$  quando cambio le variabili?

$$SS_R := \sum_{i=1}^n (Y_i - B_0 - B_1 x_{i1} - B_2 x_{i2} - B_3 x_{i3})^2$$

cerca  $(B_0, B_1, B_2, B_3)$  che minimizza  $SS_R$

test  $\beta_3 = 0$  dice  $H_0$   $\chi^2_3 = 87\%$ , toglgo  $x_3$

$$SS_R^* = \sum_{i=1}^n (Y_i - B_0 - B_1 x_{i1} - B_2 x_{i2})^2$$

cerca  $(B_0, B_1, B_2, 0)$  che minimizza  $SS_R$

cerca  $(B_0, B_1, B_2)$  che minimizza  $SS_R^*$

} stesse cose

★  $SS_R^* > SS_R$  sempre e comunque

★  $R^2 := 1 - \frac{SS_R}{S_{YY}}$  *quanto aumenta*  $S_{YY} = \sum_{i=1}^n (Y_i - \bar{Y})^2 = (n-1) S_Y^2$

$R_x^2 < R^2$  sempre e comunque

★  $Se^2 := \frac{SS_R}{n-p-1}$  *Se a volte aumenta, a volte diminuisce*

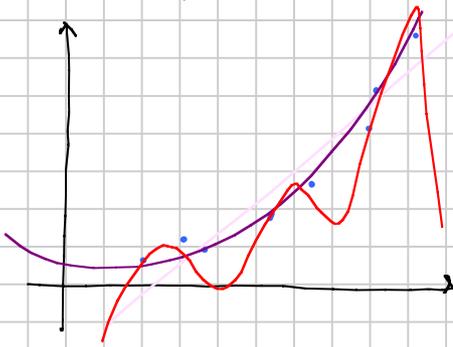
→  $Se$  è un indicatore ragionevole del fatto che la selezione delle var. abbia migliorato o no il modello

★  $\bar{R}^2 := 1 - \frac{Se^2}{S_Y^2}$  "  $R^2$  corretto " (adjusted)  
 si comporta come  $Se$  (al contrario)

● Quando si tolgono variabili, guardare se  $R^2$  aumenta (o se  $Se$  diminuisce) per capire se il modello è migliorato. Piccole variazioni non sono importanti. In ogni caso non è un metodo ottimale neanche questo.

## OVERFITTING

Esempio (rumoto)



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

$$Y = \alpha + \beta x + \gamma x^2$$

$$Y = \alpha + \beta x + \gamma x^2 + \delta x^3 + \dots + \zeta x^k$$

Nella regressione se i coefficienti stimati ( $p+1$ ) sono "troppi" rispetto alle osservazioni ( $n$ ) si ha lo stesso fenomeno: si trova un modello che si adatta fin troppo bene al campione in esame, ma che sarà pochissimo adattato alla popolazione, quindi alle osservazioni future

- $p+1 = n$  il modello ha residui nulli, passa esattamente per i punti  $R_i = 0$ ,  $SS_R = 0$ ,  $S_e = 0$ , DIV/0

- $p+1 > n$  errore

- $p+1$  poco minore di  $n$  il modello è sbagliatissimo

- $\frac{n}{p+1} \approx 5$  siamo al limite (grossi problemi, ma non da buttare)

★ in mezzo: buon senso

- $\frac{n}{p^2} \approx 2+$  va benissimo

★ Se occorre migliorare  $\frac{n}{p}$ , o si spende per aumentare  $n$ , oppure si usa la selezione delle var per diminuire  $p$

★ Grosso problema: un effetto naturale dell'overfitting è che conosciamo i coefficienti con precisione molto bassa

-  $\beta_i \in B_i \pm \text{roba}$   $\rightarrow$  è grande se c'è overfitting

-  $H_0: \beta_i = 0$   $H_1: \beta_i \neq 0$  ha potenza molto bassa

si dice quasi sempre  $H_0$  (a meno che  $|\beta_i|$  sia enorme)  
anche se è vera  $H_1$   
viene  $\alpha^*$  spesso grande

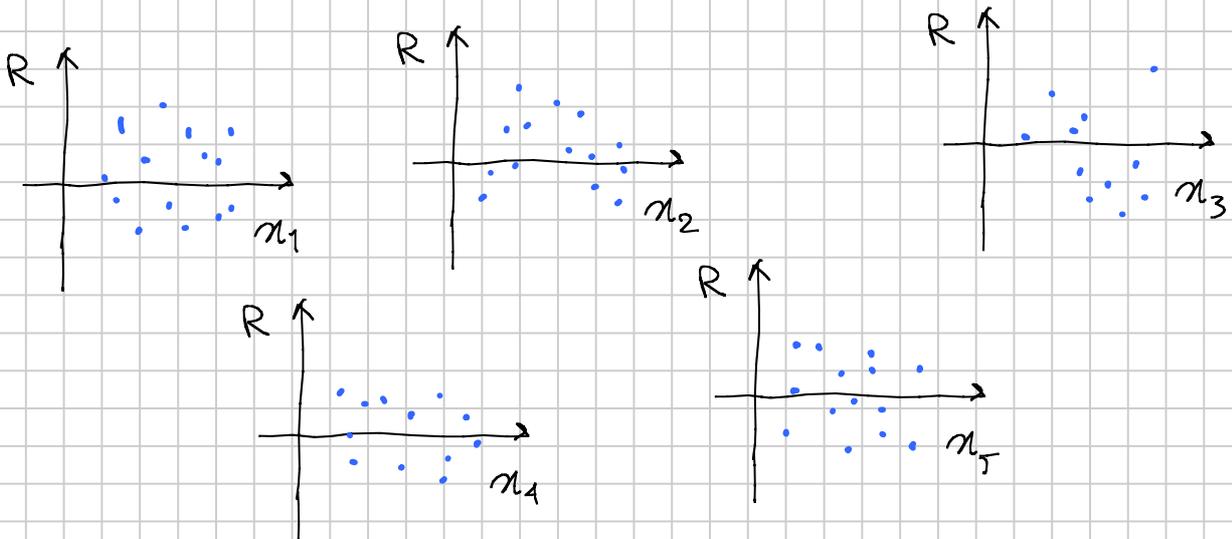
$\rightarrow$  Questo rende poco affidabile la stepwise backward  
(si buttano via all'inizio variabili forse buone)  
Usare in questi casi la forward, o almeno fare con la  
backward vari tentativi.

ora 15

## ANALISI DEI RESIDUI

$\rightarrow$  Come nella regr lin semplice, va fatta sempre  
Purtroppo non c'è un  $\alpha$  soltanto

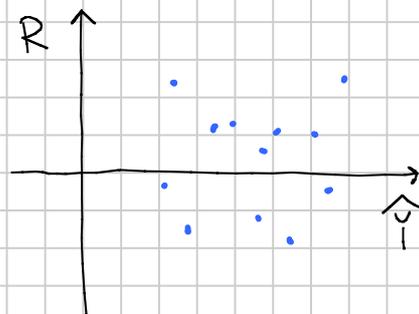
Si fanno tanti grafici: uno per ogni variabile di ingresso  
e si analizzano tutti



\* Spesso è utile aggiungerne uno:

quello dei residui vs  $\hat{y}$

$$\hat{y}_i = B_0 + B_1 x_{i1} + \dots + B_p x_{ip}$$



# ■ CURA DELLA REGRESSIONE (quando i residui mostrano stranezze)

■ OUTLIERS → correggibili  
 → non correggibili

|   |   |   |     |   |   |
|---|---|---|-----|---|---|
| 0 | 1 | 1 | 0   | 1 | 1 |
| 0 | 0 | 1 | 1   | 0 | 0 |
| 0 | 0 | 0 | 100 | 0 | 0 |
| 0 | 1 | 0 | 1   | 1 | 0 |

era un errore + era correggibile

Circoscrizioni di Philadelphia: le tre del centro andavano escluse perché troppo diverse

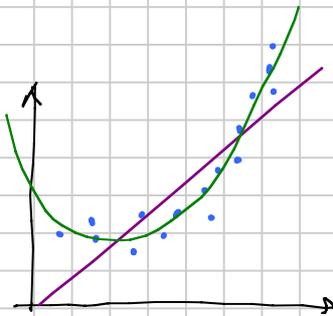
\* Vanno comunque eliminati o corretti

## ■ CORREGGERE LA NONLINEARITÀ

Metodo principale: aggiungere dummy variables

→ esempio 1

| id | x | Y |
|----|---|---|
| 1  | ~ | ~ |
| 2  | ~ | ~ |
| ⋮  | ⋮ | ⋮ |
| n  | ~ | ~ |



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



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

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

| id | x   | x <sup>2</sup>   | Y |
|----|-----|------------------|---|
| 1  | 2,1 | 2,1 <sup>2</sup> | ~ |
| 2  | 1,7 | 1,7 <sup>2</sup> | ~ |
| ⋮  | ⋮   | ⋮                | ⋮ |
| n  | ~   | ~ <sup>2</sup>   | ~ |

↑ dummy variable

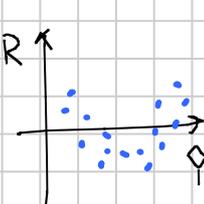
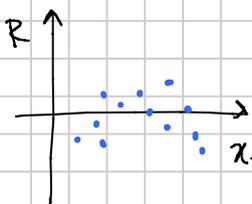
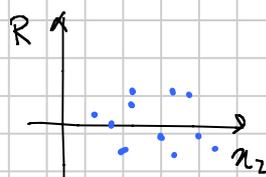
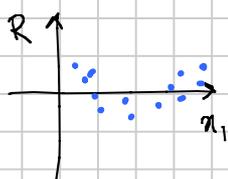
faccio la reg. lin. multiple sulle due variabili → i coeff che ottengo vanno interpretati come quelli del polinomio

\* Posso alzare il grado progressivamente aggiungendo dummy variables. (Occhio all'overfitting.)

→ Se si esagera con il grado del polinomio, tipicamente i p dei dati vengono tutti (quasi) uguali

→ esempio 2

| id | $x_1$ | $x_2$ | $x_3$ | $Y$ |
|----|-------|-------|-------|-----|
| 1  |       |       |       |     |
| 2  |       |       |       |     |
| ⋮  |       |       |       |     |
| n  |       |       |       |     |



| id | $x_1$ | $x_1^2$ | $x_2$ | $x_3$ | $x_3^2$ | $Y$ |
|----|-------|---------|-------|-------|---------|-----|
|    |       |         |       |       |         |     |
|    |       |         |       |       |         |     |
|    |       |         |       |       |         |     |
|    |       |         |       |       |         |     |
|    |       |         |       |       |         |     |

★ Nelle regre multivariate potrei dover aggiungere parecchie dummy variables. (Overfitting!)

$x_1 \quad x_2 \quad x_3 \quad x_1^2 \quad x_3^2 \quad x_1 x_3$

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_2 + \beta_4 x_3^2 + \beta_5 x_1 x_3$$

in un generico polinomio di II grado in più variabili ci possono essere questi prodotti che si chiamano termini di interazione.

★ Sono molto importanti. Approfondiremo sul DOE

p variabili  $x_1, x_2, \dots, x_p$

p quadrati  $x_1^2, x_2^2, \dots, x_p^2$

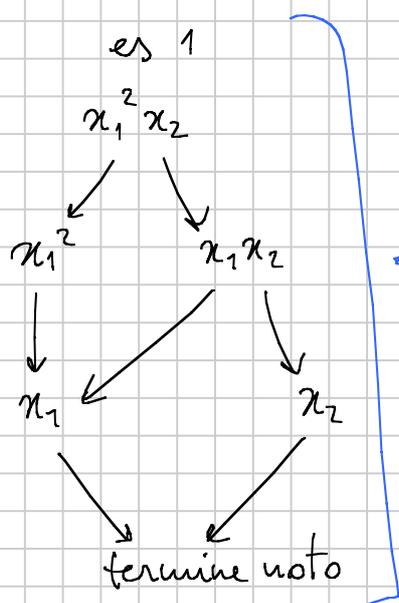
$\binom{p}{2}$  interazioni  $x_1 x_2, x_1 x_3, x_1 x_4, \dots, x_1 x_p, x_2 x_3, x_2 x_4, \dots, x_2 x_p, \dots, x_{p-1} x_p$

La regre. lin. multiple necessaria ha  $p' = p + p + \binom{p}{2} = \frac{p^2 + 3p}{2}$  variabili → overfitting sicuro

★ Vicererso, se sono dalla parte opposta rispetto all'overfitting (ad es  $n=10000$ ,  $p=3$ ) è legittimo e doveroso aggiungere un po' di termini (verranno significativi perché i test sono molto potenti)

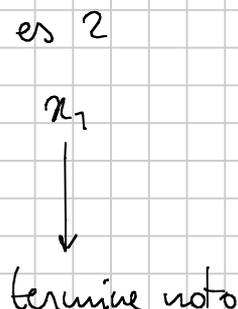
HW: capire!

● Regola gerarchica: se si inserisce un termine di grado  $\geq 1$ , vanno inseriti, anche se non significativi, tutti i termini "contenuti"



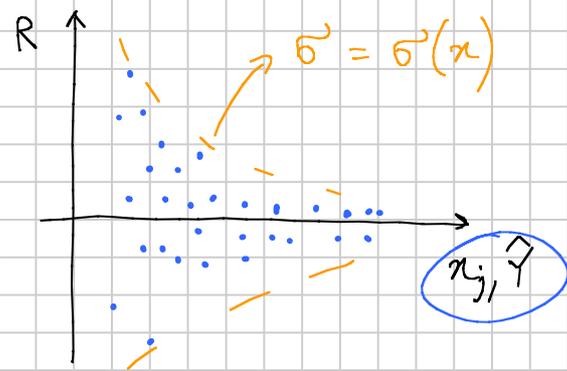
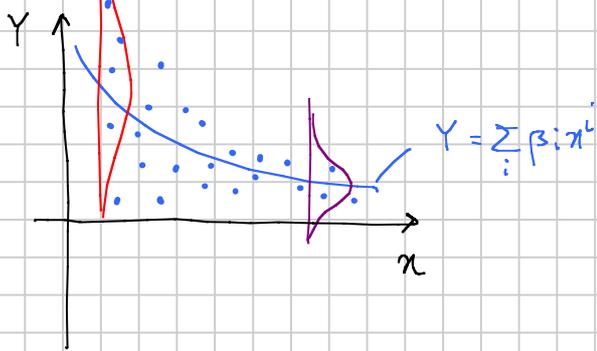
GENERALIZZAZIONE

già  
discusso  
in  $Y=A+Bx$



→ a meno che non ci siano ottimi motivi teorici a giustificare l'assente

## CURARE ETEROSCHEDASTICITÀ



$$Y(i) = \sum_j \beta_j x_j(i) + e_i$$

$$e_i \sim \mathcal{N}(0, \sigma_i^2)$$

$$\sigma_i^2 = \sigma^2$$

OMOSCHEDASTICITÀ

$$\sigma_i \neq \sigma_j \quad \sigma = \sigma(x)$$

ETEROSCHEDASTICITÀ

Per correggerla occorre conoscere come  $\sigma$  dipende da  $x$

★ Basta stabilire  $\sigma(x)$  a meno di una costante moltiplicativa

$$\sigma(x) = 7,2 \cdot x \quad \text{reale}$$

$$\sigma(x) \propto x \Rightarrow \sigma(x) = x$$

ipotizzate

funzione perfettamente  $\sigma(x) = x, \sigma(x) = 2x, \sigma(x) = 7,2x \dots$

★ Andamenti tipici :

a)  $\sigma \propto x_j$

b)  $\sigma^2 \propto x_j \Leftrightarrow \sigma \propto \sqrt{x}$

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

b)  $\sigma^2 \propto \hat{Y} \Leftrightarrow \sigma \propto \sqrt{\hat{Y}}$

valori di risposta previsti

~~$\hat{Y}$~~   $\rightarrow \sigma, \sigma^2 \propto E(Y) = \sum_j \beta_j x_j$

$$\hat{Y} = \sum_j \beta_j x_j$$

c) vedi oltre

\* Due situazioni tipiche

a) L'errore è "in percentuale" sul valore di  $Y$

$$\sigma \approx 25\% E(Y)$$

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

NB. anche  $\sigma \propto x$  è di questo tipo, infatti di solito accade quando  $Y \approx x$

$$Y = \alpha + \beta x \quad x \ll 1 \text{ o } x \approx 0$$

$$Y \approx \beta x \quad Y \approx x \quad \sigma \propto Y \Leftrightarrow \sigma \propto x$$

Rappresenta un sottocaso di a), ma se vi siamo conviene approfittarne:  $x$  è nota, " $Y$ " no (ovvero non occorre conoscere  $\beta$ )

b) L'errore è la somma di tanti contributi indipendenti, e quanti sono  $x \propto Y$  (o  $x \propto n$ )

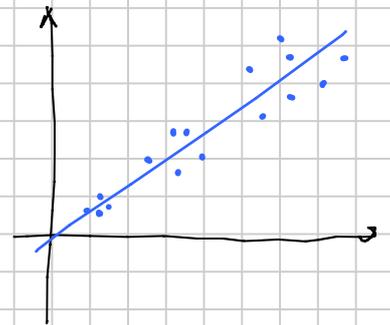
$x$ : km percorsi (tracciato cittadino)

$Y$ : tempo

$$T_1 + T_2 + \dots + T_n = Y$$

$\underbrace{\hspace{10em}}_{\Sigma \text{ iid}}$

$$\sigma^2 = \text{Var}(Y) = n \text{Var}(T_1) \propto n$$



NB. Come sopra

\* Difficile distinguere tra a), b) e tra NB si/no, anche guardando i residui. Comunque se non è chiaro, a parte i nomi i risultati saranno simili con qualunque modello

c) Situazione atipica: abbiamo stime di  $\sigma(x(i))$  per ogni  $x(i)$

- Esempio 1:  $x(i) \rightarrow \begin{cases} 5 \text{ risposte} \\ 4 \text{ risposte} \\ 3 \text{ risposte} \end{cases} \left. \vphantom{\begin{matrix} 5 \\ 4 \\ 3 \end{matrix}} \right\} Y(i) = \text{media}$

$$Y(i) = \frac{Y_{i,1} + Y_{i,2} + \dots + Y_{i,m}}{m(i)} \quad m=3,4,5 \text{ a seconda di } i$$

$$\text{Var}(\bar{X}) = \frac{\sigma^2}{n} = \frac{\text{Var}(X_1)}{n}$$

HWC

$$\sigma^2(Y(i)) = \text{Var}(Y(i)) = \frac{1}{m(i)} \text{Var}(Y_{i,1}) \propto \frac{1}{m(i)}$$

$\sigma^2 \propto \frac{1}{3}, \frac{1}{4}, \frac{1}{5}$  a seconda dell'esperimento

- Esempio 2:  $x_1, x_2, \dots, x_p \rightarrow$  probabilita'  $p = \sum_j \beta_j x_j$

farò  $m(i)$  di esperimenti per ciascuna combinazione delle var di ingresso

$Y(i)$ : # di successi tra gli  $m(i)$  esperimenti di tipo  $i$

$$Y(i) \sim \text{bin}(m(i), p(i))$$

$$E(Y(i)) = m(i)p(i) = m(i) \sum_j \beta_j x_j(i) \quad \text{posso fare la regressione}$$

$$\hat{p}(i) = \frac{Y(i)}{m(i)} \quad E(\hat{p}(i)) = p(i) = \sum_j \beta_j x_j(i)$$

$$\text{Var}(\hat{p}(i)) \equiv \sigma$$

No

uso  $\hat{p}$  al posto di  $p$

$$\text{Var}(Y(i)) = m(i)p(i)(1-p(i)) \quad (\Rightarrow) \quad \text{Var}(\hat{p}(i)) = \frac{p(i)(1-p(i))}{m(i)}$$

$$\sigma^2(i) = \text{Var}(\hat{p}(i)) \propto \frac{1}{m(i)} \bar{p}(i)(1-\bar{p}(i))$$

dove  $\bar{p}(i) = \sum_j \beta_j x_j(i)$  e' il previsto

# ● REGRESSIONE ETEROSCHEDASTICA / PESATA (WEIGHTED)

$\sigma \propto$  qualcosa

$$SS_R = \sum_{i=1}^n R_i^2 = \sum_{i=1}^n \left( Y_i - \sum_j B_j x_{ij} \right)^2$$

questo non fa la cosa giusta se le varianze non sono uguali

nei

$$SS_R^W := \sum_{i=1}^n R_i^2 w_i = \sum_{i=1}^n \left( Y_i - \sum_j B_j x_{ij} \right)^2 w_i = \sum_{i=1}^n \left( \sqrt{w_i} Y_i - \sum_j B_j \sqrt{w_i} x_{ij} \right)^2$$

dove  $w_i$  è il peso dell'esperimento/dato  $i$   $w_i \propto \frac{1}{\sigma_i^2}$

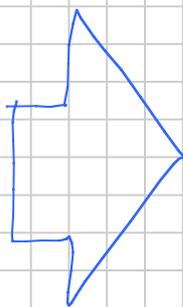
Se conosco  $\sigma \propto \dots$  posso determinare i pesi (a meno di una costante)

$$Y_i \sim \mathcal{N} \left( \sum_j \beta_j x_{ij}, \sigma_i^2 \right) \quad \sqrt{w_i} Y_i \sim \mathcal{N} \left( \sum_j \beta_j \sqrt{w_i} x_{ij}, \underline{w_i \sigma_i^2} \right)$$

↳ non dipende

★ Posso risolvere questa regressione, o (tecnica sconsigliata) <sup>da i</sup> minimizzando  $SS_R^W = SS_R^W(B_1, B_2, \dots, B_p)$ ; oppure trasformo tutti i dati moltiplicando per  $\sqrt{w_i}$

| Y     | $x_0$ | $x_1$      | $x_2$ | ... | $x_p$      |
|-------|-------|------------|-------|-----|------------|
| $Y_1$ | 1     | $x_{1(1)}$ | .     |     | $x_{p(1)}$ |
| $Y_2$ | 1     | $x_{1(2)}$ | .     |     | .          |
| ⋮     | ⋮     | ⋮          | ⋮     |     | ⋮          |
| $Y_n$ | 1     | $x_{1(n)}$ | .     |     | .          |



| Y'               | $x'_0$       | $x'_1$                | ... | $x'_p$                |
|------------------|--------------|-----------------------|-----|-----------------------|
| $\sqrt{w_1} Y_1$ | $\sqrt{w_1}$ | $\sqrt{w_1} x_{1(1)}$ |     | $\sqrt{w_1} x_{p(1)}$ |
| $\sqrt{w_2} Y_2$ | $\sqrt{w_2}$ | .                     |     | .                     |
| ⋮                | ⋮            | ⋮                     |     | ⋮                     |
| $\sqrt{w_n} Y_n$ | $\sqrt{w_n}$ | $\sqrt{w_n} x_{1(n)}$ |     | .                     |

$$Y_i' \sim \mathcal{N} \left( \sum_j \beta_j x_{ij}', \sigma_*^2 \right)$$

$$Y' = \sum \beta_j x_j' + e$$

è una regr. omoschedastica

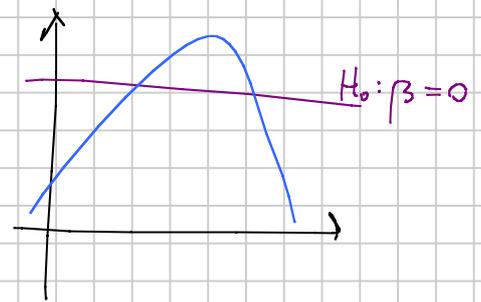
★ lancio la regr su  $Y', x'$  e trovo i coeff.  $B^W$  che minimizzano  $SS_R^W$

e quindi sono quelli giusti per la regr. eteroschedastica

- ★ La regr.  $Y', x'$  va fatta senza termine noto (spuntare "passe per l'origine")
- ★ Se servono intervalli di conf/predizione, si calcolano su  $Y', x'$  e si trasformano indietro, dividendo per  $\sqrt{w(\tilde{x})}$
- ★  $R^2$  non è confrontabile con quello della regressione precedente  
Se non ha senso
- ★ Per capire se il modello è migliorato si può solo osservare l'omoschedasticità dei nuovi residui.
- ★ In generale questo "trattamento" va fatto solo in casi eccezionali

## SCHEMA REGRESSIONE

- 1) Outliers
- 2) Selezione variabili
- 3) Nonlinearità
- 4) Eteroschedasticità



1, 3, 4) Trasformazioni nonlineari + 5) Residui aritmetici

## TRASFORMAZIONI NONLINEARI

Esempio: prodotto in bottiglia → lavaggio  
a che velocità faccio andare la macchina?

$$v = [\text{bottiglie/sec}]$$

$$v = v(x_1, x_2, \dots, x_p)$$

levello di pulizia richiesto (pointing to  $x_p$ )  
volume (pointing to  $x_1$ )  
profondità (pointing to  $x_2$ )  
tipo prodotto (pointing to  $x_p$ )



$$T: \text{tempo di lavaggio} \quad T = [\text{sec}]$$

$$T = T(x_1, x_2, \dots, x_p) = \sum_j \beta_j x_j + e$$

$$T = \sum_j \beta_j x_j + e \quad e \sim \mathcal{N}(0, \sigma^2)$$

$$v = \frac{3600}{T} = \frac{3600}{\sum_j \beta_j x_j + e}$$

→ relazione nonlineare  
 → residui non  $\mathcal{N}$  / simmetrici  
 → eteroschedasticità

$$\frac{a}{b + \mathcal{N}} \sim \frac{a}{\mathcal{N}'} \sim \frac{1}{\mathcal{N}''} \neq \mathcal{N} \rightarrow \text{finti outliers}$$

$$e \sim \pm 1$$

|                              |            |       |            |
|------------------------------|------------|-------|------------|
|                              |            |       | 360 ± 240  |
| $\sum \beta_j x_j \approx 4$ | den : 3-5  | $v$ : | 720 ~ 1200 |
| $\sum \beta_j x_j \approx 9$ | den : 8-10 | $v$ : | 360 ~ 450  |
|                              |            |       | 405 ± 45   |

• Se la variabile che dipende in maniera semplice dalle  $x_i$  non è  $Y$  ( $v$ ) ma una sua riscrittura o trasformazione allora posso migliorare tantissimo la regressione se al posto di  $Y$  lavoro su di essa ( $T$ )

★ Se faccio una trasformazione sbagliata posso anche peggiorare

★ Trasformazioni lineari non servono quasi a nulla

★ Anche trasformazioni nonlineari di alcune delle  $x_i$  possono in certi casi aiutare.

ora 18

## ■ TRASFORMAZIONI LINEARI NELLA REGRESSIONE

Non hanno quasi effetto!

$$Y' = \lambda + \mu Y, \quad x'_j = \lambda_j + \mu_j x_j$$

$$- \text{var } \in \rightarrow k \in \quad \lambda = 0 \quad \mu = \frac{1}{1000}$$

$$- \text{var } \% \rightarrow [0, 1] \quad \lambda = 0 \quad \mu = \frac{1}{100}$$

$$- T [\text{Fahr}] \rightarrow T [^{\circ}\text{C}] \quad \lambda, \mu$$

$$- \text{var } \{1, 2\} \rightarrow \{0, 1\} \quad \lambda = -1 \quad \mu = 1$$

$\alpha_j^*$   $j=1,2,\dots,P$  ,  $\mathbb{R}^2$  } non cambiano  
 grafici dei rendimenti

$S_e' = \mu S_e$   
 $B_j' = \frac{B_j}{\mu_j}$   $j=1,2,\dots,P$  } cambiano in modo prevedibile

★ Standardizzazione  $\mathcal{N}(\mu, \sigma^2) \rightarrow \mathcal{N}(0, 1)$   
 $X \rightarrow \frac{X - \mu}{\sigma}$

HWC

|                     |                                |
|---------------------|--------------------------------|
| $E(X) = \mu$        | $E(Y) = 0$                     |
| $Var(X) = \sigma^2$ | $Var(Y) = 1$                   |
| $X \rightarrow$     | $\frac{X - \mu}{\sigma} =: X'$ |

|                        |                                                                                                    |
|------------------------|----------------------------------------------------------------------------------------------------|
| $X_1, X_2, \dots, X_n$ | $X'_1, X'_2, \dots, X'_n$                                                                          |
| $\bar{X}, S_X$         | $\bar{X}' = 0, S_{X'} = 1$                                                                         |
| $X_i \rightarrow$      | $\frac{X_i - \bar{X}}{S_X} =: X'_i = -\frac{\bar{X}}{S_X} + \frac{1}{S_X} X_i = \lambda + \mu X_i$ |

operazione che si fa per colonne  
 → Per la regressione, questo significa  $B_0' = 0$   
 e inoltre i  $B_j'$  sono confrontabili

$B_3' = 2,4$      $B_2' = -0,7$   
 $\alpha_3$  o  $\alpha_3'$  ha un impatto su  $Y$  più forte di  $\alpha_2$  o  $\alpha_2'$   
 $\alpha_3^* \ll \alpha_2^*$

## VARIABILI CATEGORICHE

- la regressione per funzionare vuole tutte variabili numeriche "misurate"

- esempi di var non numeriche:

NOMINALI {  
 i. tipo di materiale {PET, PVC, VETRO}  
 ii. fornitore {6 fornitori}  
 iii. tipo di esca {4 categorie}

} non codificare

DICOTOMICA iv. sesso {m, f} → codificare come si vuole

ORDINALI {  
 v. titolo di studio {I, II media, maturità, laurea, dottorato}  
 vi. fascia di età {18-25, 25-35, ...}

} forse codificare

$x \xrightarrow{!} Y$

NON sono adatte alla regressione!

$T = T(x) + e$      $x$ : tipo di materiale {<sup>1</sup>PET, <sup>2</sup>PVC, <sup>3</sup>VETRO}

Esorse classico: codificare arbitrariamente le categorie con numeri

| $E(T)$ | $x$   |
|--------|-------|
| 9      | PET   |
| 8      | PVC   |
| 5      | VETRO |

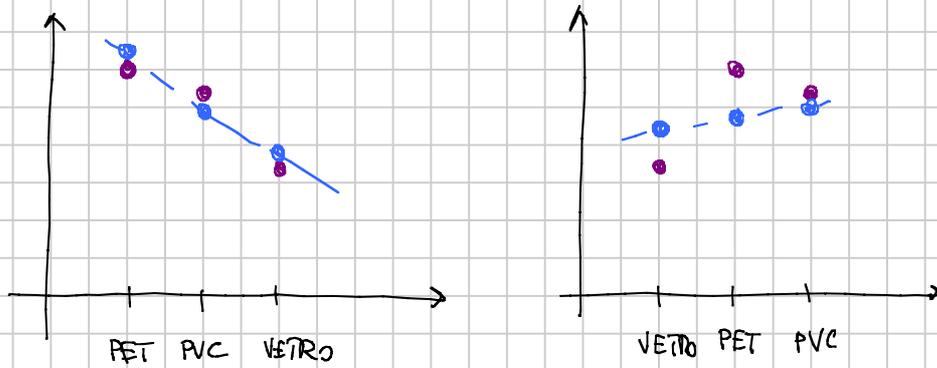
la dipendenza da  $x$   
 è una tabella, non una formula

$$T = \alpha + \beta x + e = \beta x + e \quad \text{dove} \begin{cases} \beta_{\text{PET}} = 9 \\ \beta_{\text{PVC}} = 8 \\ \beta_{\text{VETRO}} = 5 \end{cases}$$

c'è un coefficiente da stimare per ogni categoria

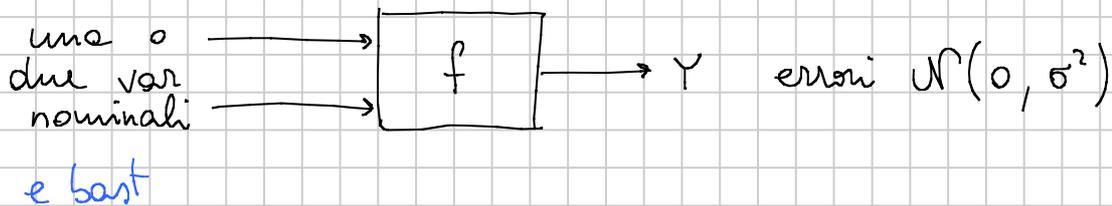
Esorse:  $T = \alpha + \beta x + e$   
 ↑  
 1, 2, 3

| $E(T)$            | $x$   |
|-------------------|-------|
| $\alpha + \beta$  | PET   |
| $\alpha + 2\beta$ | PVC   |
| $\alpha + 3\beta$ | VETRO |



• Variabili dicotomiche (e ordinali) si possono (forse) codificare quelle nominali NO

★ Il modo corretto per gestire:



è usare la analisi della varianza (Cap. 10)

• Negli altri casi c'è una tecnica un po' fragile che si può usare:

Le var nominali (e forse certe ordinali) si possono "esplodere" in tante dicotomiche quanto è il numero di categorie meno 1

Es: PVC, PET, VETRO → 2 dicotomiche

↑  
 sceglie la preferita/default/più comune/più simile alle altre

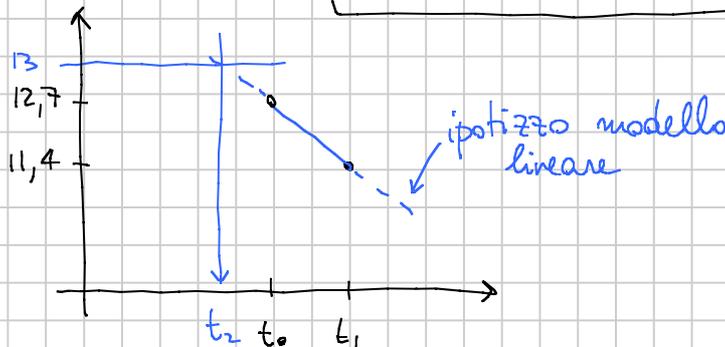
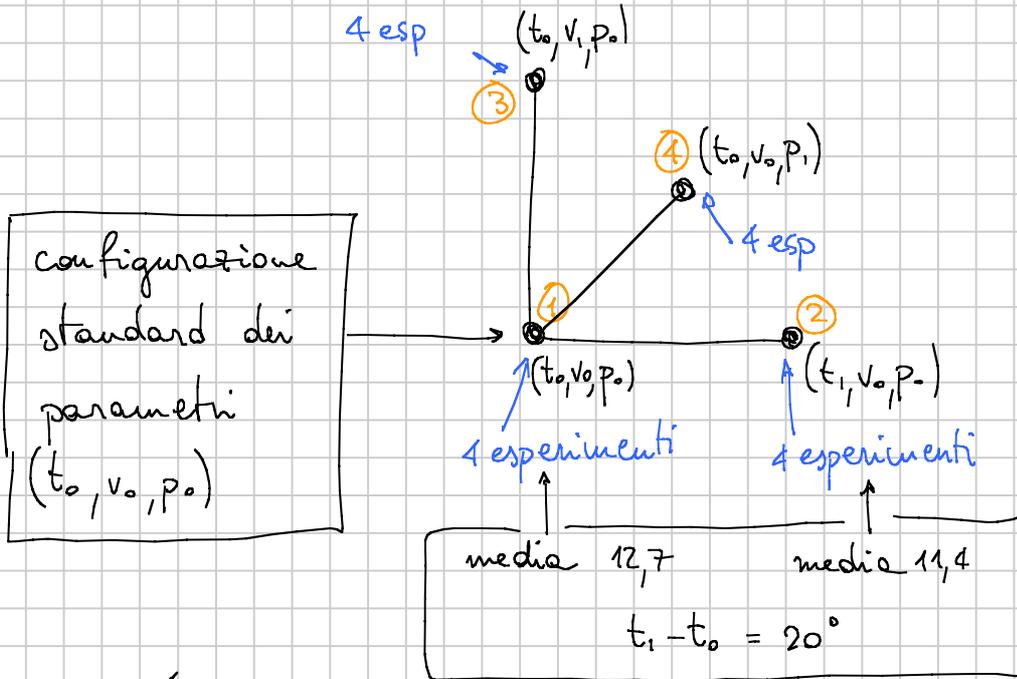
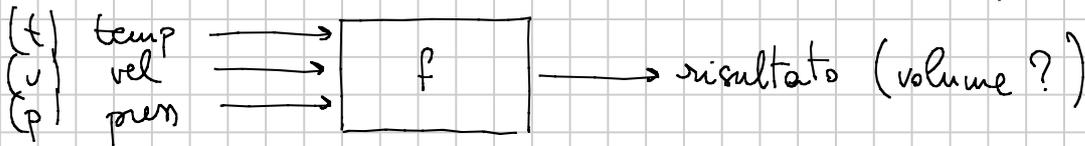
| id | $x$   | $x_{PVC}$ | $x_{VETRO}$ |       | $x_1$ | $x_2$ |
|----|-------|-----------|-------------|-------|-------|-------|
| 1  | PVC   | 1         | 0           | PET   | 0     | 0     |
| 2  | PVC   | 1         | 0           | PVC   | 1     | 0     |
| 3  | PET   | 0         | 0           | VETRO | 0     | 1     |
| ·  | PVC   | 1         | 0           |       |       |       |
| ·  | VETRO | 0         | 1           |       |       |       |
| ·  | PET   | 0         | 0           |       |       |       |
| ·  | VETRO | 0         | 1           |       |       |       |
| ·  | ·     |           |             |       |       |       |
| ·  | ·     |           |             |       |       |       |

★ Esplodere una dicotomica equivale a codificarla

## DESIGN OF EXPERIMENT (DoE) (Cap 10 Sleper)

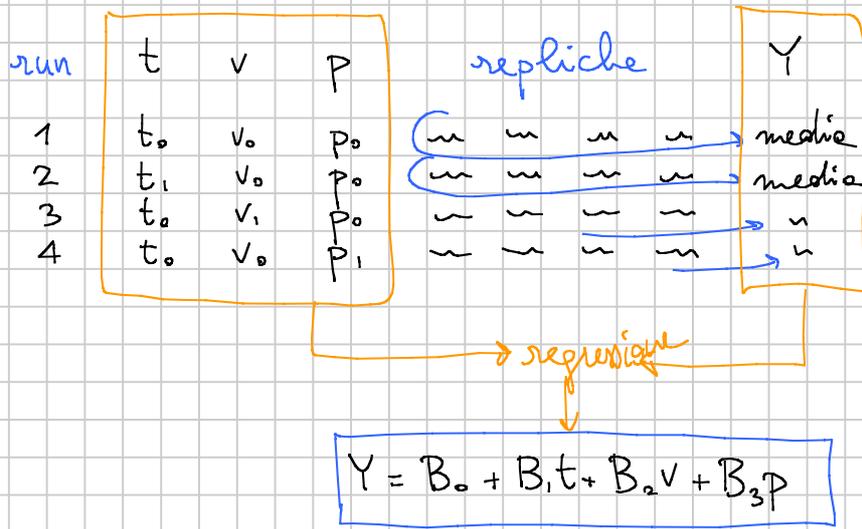
Se devo fare una regressione e devo ancora raccogliere i dati esiste un modo "migliore" per raccoglierti?

Esempio: stampo plastico: regolazioni  $\rightarrow$  temperature  
 $\rightarrow$  velocità  
 $\rightarrow$  pressione }  $\rightarrow$  risultato



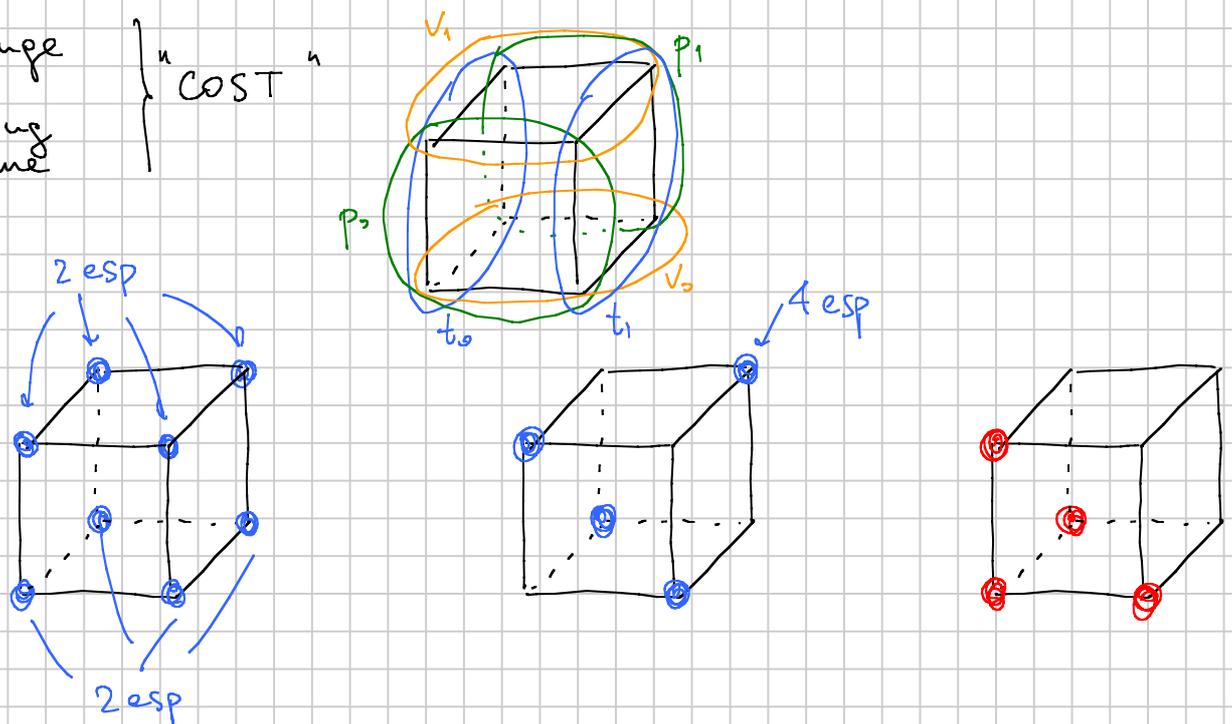
Allo fine riesco a ricavare

$$Y = B_0 + B_1 t + B_2 v + B_3 p$$



Change  
One  
Setting  
a Time

"COST"



design completo

design frazionario

### ☞ Come scegliere i valori di ingresso

- Due livelli soli per ogni variabile di ingresso vanno quasi sempre bene
- Si usano **3** livelli quando il modello deve essere fortemente nonlineare

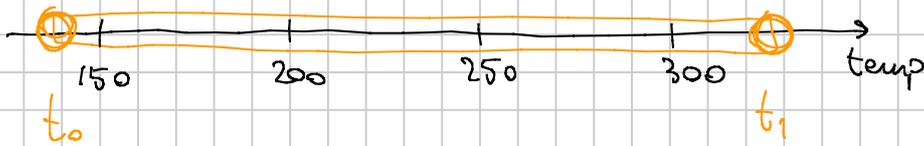
→ è possibile alla fine di una analisi su 2 livelli verificare con un ulteriore esperimento (center point) se sono sufficienti

→ con 2 livelli si prendono modelli lineari, anche con interazioni

$$Y = \beta_0 + \beta_1 t + \beta_2 v + \beta_3 p + \beta_4 tv + \beta_5 tp + \beta_6 vp + \beta_7 tvp$$

non si arriva ai "quadrati"

- La scelta dei due livelli è critica



occorre sempre decidere in brainstorming con ingegneri e tecnici tenendo conto degli aspetti industriali pratici

★ criteri statistici

- a) se l'intervallo è troppo largo il modello può essere non lineare

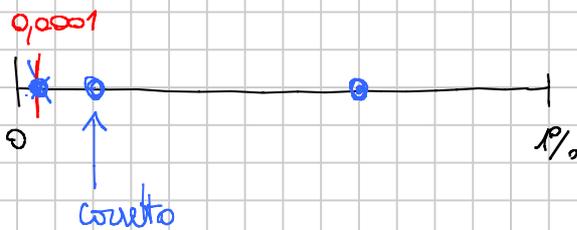


- b) se l'intervallo è troppo largo è plausibile che il "sistema" non funzioni come previsto

$t_0 = 200, t_1 = 250$  tutto ok

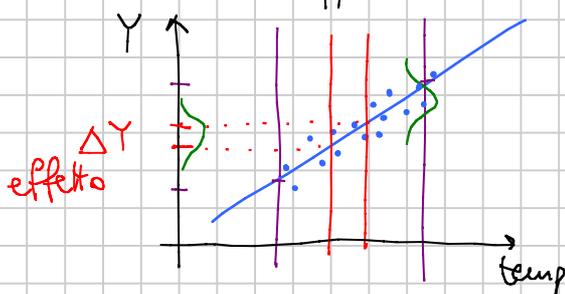
$t_0 = 150, t_1 = 300$

settaggi un po' al limite per la macchina qualcosa non funziona



(senza audio)

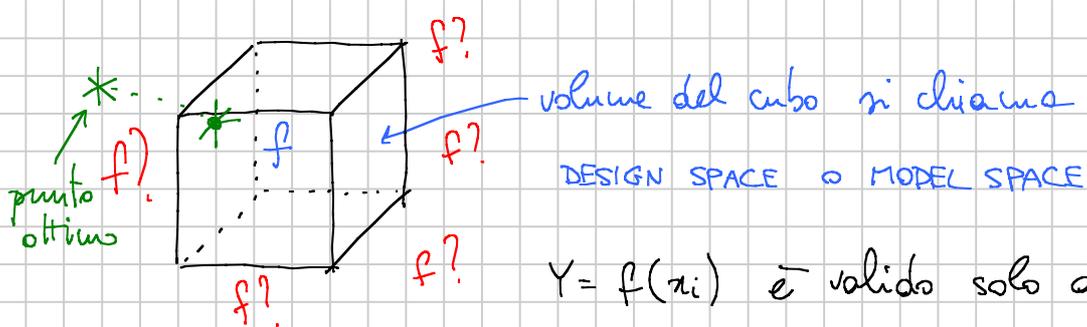
- c) se l'intervallo è troppo stretto può darsi che io non veda variazioni apprezzabili



rapporto segnale/rumore insufficiente

rapporto segnale/rumore adeguato

- d) se l'intervallo è troppo stretto il design space può essere inutile



$Y = f(x_i)$  è valido solo all'interno del D.S.

ORA 20

## ▣ Variabili di risposta

- Possono essere più di una (facendo varie regressioni)  
Questo può portare a problemi nella scelta delle var di ingresso e nei loro livelli
- Devono essere variabili numeriche misurabili con precisione su un range adeguato

esempio: riempimento buono/non buono è dicotomica

↳ devo trovare un modo di misurarla  
grammetre? scale di qualità?

→ Se  $Y$  dicotomica è l'unica soluzione allora saranno necessarie centinaia o migliaia di esperimenti

## ▣ Variabili categoriche

- tipo di plastica
- tipo di iniettore

★ solo come var di ingresso e solo con 2 categorie ciascuna (dicotomiche)

in ogni caso impediscono il center point

## ▣ Scelta delle var di ingresso e del design

esempio: 11 variabili: selezione delle var necessarie  $\frac{n}{p}$  piccolo

nel DOE siamo costretti comunque a  $n = n(p)$

→ fattoriale completo  $n = 2^p = 2^{11} = 2048$  run

poi da moltiplicare per le repliche (almeno 2)

→ fattoriale frazionario  $2^3 \rightarrow 2^2 = 4$  run =  $2^{p-k}$   $k$  piccolo

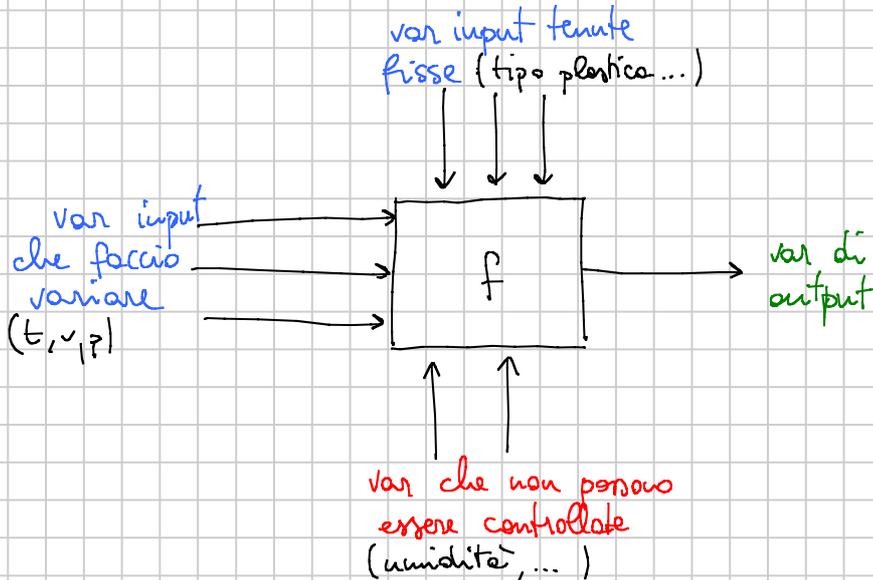
$k=1$  dimezzano le prove del completo

$k=2 \dots \frac{1}{4}$

$k=3 \dots \frac{1}{8}$

\* Le var di ingresso da far variare vanno scelte con cura, parlando con tutte le persone coinvolte

diagramma IPO



Le var da far variare vanno ridotte all'osso per tenere basso il numero degli esperimenti (vedi oltre esperimenti di screening/modeling)

### DESIGN FRAZIONARI & ALIASING

Esempio: esperimento con 3 variabili : completo L8 e fraz L4

| A  | B  | C  | <del>A<sup>2</sup></del> | AB | AC | BC | ABC |
|----|----|----|--------------------------|----|----|----|-----|
| -1 | -1 | -1 | +1                       | +1 | +1 | +1 | -1  |
| -1 | -1 | +1 | +1                       | +1 | -1 | -1 | +1  |
| -1 | +1 | -1 | +1                       | -1 | +1 | -1 | +1  |
| -1 | +1 | +1 | +1                       | -1 | -1 | +1 | -1  |
| +1 | -1 | -1 | +1                       | -1 | -1 | +1 | +1  |
| +1 | -1 | +1 | +1                       | -1 | +1 | -1 | -1  |
| +1 | +1 | -1 | +1                       | +1 | -1 | -1 | -1  |
| +1 | +1 | +1 | +1                       | +1 | +1 | +1 | +1  |

| I  | A  | B  | C  | AB | AC | BC | ABC |
|----|----|----|----|----|----|----|-----|
| +1 | -1 | -1 | -1 | +1 | +1 | +1 | -1  |
| +1 | -1 | -1 | +1 | +1 | -1 | -1 | +1  |
| +1 | -1 | +1 | -1 | -1 | +1 | -1 | +1  |
| +1 | -1 | +1 | +1 | -1 | -1 | +1 | -1  |
| +1 | +1 | -1 | -1 | -1 | -1 | +1 | +1  |
| +1 | +1 | -1 | +1 | -1 | +1 | -1 | -1  |
| +1 | +1 | +1 | -1 | +1 | -1 | -1 | -1  |
| +1 | +1 | +1 | +1 | +1 | +1 | +1 | +1  |

= X

design di Taguchi L8

↑ ↑ ↑  
 unità codificate:  $t_0 = -1$   $t_1 = +1$   $v_0 = -1$   $v_1 = +1$  ...

$X$  è ortogonale a meno di un fattore 8

$$X^T X = 8 I$$

$$(X^T X)^{-1} = \frac{1}{8} I$$

L4 a 3 fattori

| I  | A  | B  | C  | AB | AC | BC | ABC |
|----|----|----|----|----|----|----|-----|
| +1 | -1 | -1 | +1 | +1 | -1 | -1 | +1  |
| +1 | -1 | +1 | -1 | -1 | +1 | -1 | +1  |
| +1 | +1 | -1 | -1 | -1 | -1 | +1 | +1  |
| +1 | +1 | +1 | +1 | +1 | +1 | +1 | +1  |

- A è in alias con BC
- B è in alias con AC
- C è in alias con AB
- I è in alias con ABC

Alias structure:

- I + ABC
- A + BC
- B + AC
- C + AB

dell' L4 a 3 fattori

Se faccio la regressione con I, A, B, C e non gli altri ...

$$Y = \beta_0 + \beta_1 A + \beta_2 C + \beta_3 AC \rightarrow \beta_0 (1+ABC) + \beta_1 (A+BC) + \beta_2 (B+AC) + \beta_3 (C+AB)$$

$$Y = \beta_0 + \beta_1 A + \beta_2 B + \beta_3 C$$



trovo 4 coefficienti che sarei erroneamente portato ad interpretare come quelli di I, A, B, C mentre invece sono quelli di I, A, AC, C  
 $Y$  non dipende davvero da B, ma non ho modo di capirlo

• Nell'alias structure ho tante righe quante le run = coefficienti stimabili e tanti termini per riga quanto  $2^k$

"design di IV risoluzione"

Fractional Factorial Design

$2^6 = 64$  run il completo

Factors: 6 Base Design: 6, 16 Resolution: IV  
 Runs: 16 Replicates: 1 Fraction: 1/4  
 Blocks: 1 Center pts (total): 0

$k=2 \quad 16 = 2^{6-2} = 2^4$

$2^{4-k}$

Design Generators: E = ABC, F = BCD

Alias Structure

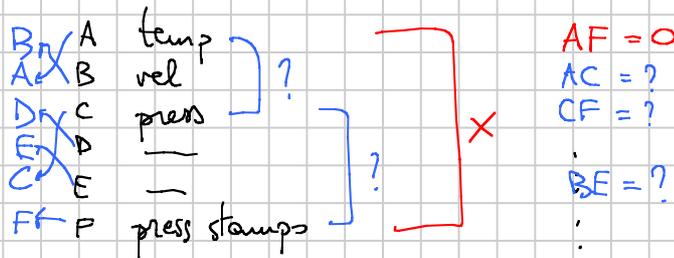
- 4 I + ABCE + ADEF + BCDF
- 4 A + BCE + DEF + ABCDF
- 4 B + ACE + CDF + ABDEF
- 4 C + ABE + BDF + ACDEF
- 4 D + AEF + BCF + ABCDE
- 4 E + ABC + ADF + BCDEF
- 4 F + ADE + BCD + ABCEF
- 4 AB + CE + ACDF + BDEF
- 4 AC + BE + ABDF + CDEF
- 4 AD + EF + ABCF + BCDE
- 4 AE + BC + DF + ABCDEF
- 4 ~~AF + DE~~ + ABCD + BCEF
- 4 BD + CF + ABEF + ACDE
- 4 BF + CD + ABDE + ACEF
- 6 ABD + ACF + BEF + CDE
- 6 ABF + ACD + BDE + CEF

L16  
 16 righe  
 16 run  
 16 coefficienti stimabili

2 coefficienti inutili

Adeguato (forse eccessivo) per determinare gli effetti singoli  
 solo in parte adeguato per determinare le interazioni

Figure 10-53 Alias Structure for an L16 Treatment Structure with Six Factors



Risoluzione: il minimo tra tutti i coeff della somma dei gradi dei due effetti di grado più basso in alias tra di loro

Available Factorial Designs (with Resolution)

| Runs | 2    | 3    | 4    | 5    | 6    | 7    | 8    | 9   | 10  | 11  | 12  | 13  | 14  | 15  |
|------|------|------|------|------|------|------|------|-----|-----|-----|-----|-----|-----|-----|
| 4    | Full | III  |      |      |      |      |      |     |     |     |     |     |     |     |
| 8    |      | Full | IV   | III  | III  | III  |      |     |     |     |     |     |     |     |
| 16   |      |      | Full | V    | IV   | IV   | IV   | III |
| 32   |      |      |      | Full | VI   | IV   | IV   | IV  | IV  | IV  | IV  | IV  | IV  | IV  |
| 64   |      |      |      |      | Full | VII  | V    | IV  |
| 128  |      |      |      |      |      | Full | VIII | VI  | V   | V   | IV  | IV  | IV  | IV  |

scuro = III ris  
 giallo = IV ris  
 verde  $\geq$  V ris



## ⊙ Esperimenti di screening

permettono di scremare un elevato numero di variabili con pochi esperimenti (un poco maggiore di  $p$ )

la III risoluzione è adeguata, gli alias tipo  $A+BC$  portano solo ad includere qualche var di troppo (falsi positivi)

## \* Esperimenti di modeling

cercano un modello preciso → servono (alcune) interazioni

↳ V risoluz o IV strutturata bene ... efficienza ...

sono fattibili se le variabili sono poche

\* Per esperimenti di screening, ancora meglio dei Taguchi di III ci sono i design di **PLACKETT-BURMAN** che non sono affetti da aliasing ma da **confounding**.  
(Meno falsi positivi.)

DOE → gestione variabili  
 → struttura design

## ● Repliche

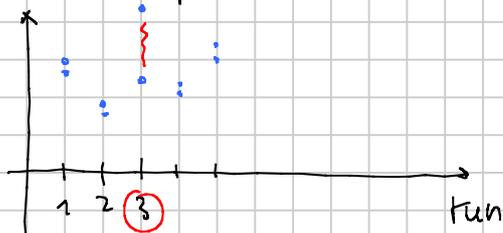
5 fattori → L16 : 16 run (m 32 combinazioni possibili)

16 run = 16, 32, 48, 64 ... esperimenti a seconda di quante repliche faccio per ogni run

| run | A  | B  | C  | D  | E  | rep1 | rep2 | rep3 | Y     |
|-----|----|----|----|----|----|------|------|------|-------|
| 1   | -1 | -1 | +1 | -1 | -1 | 11,4 | 11,7 | 10,8 | media |

→ Perché servono diverse repliche?

a) Per scoprire errori e outlier servono sempre almeno 2 repliche



qualcosa non va

varianza minore

b) Per alzare il rapporto segnale/rumore

$$\text{rep 1} \sim \mathcal{N}(\mu, \sigma^2) \quad \text{media (3 repliche)} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{3}\right)$$

Al diminuire del rumore diventano rilevabili via via effetti più piccoli

$$Y = 7,4 + 2,3A + e \quad e \sim \mathcal{N}(0, (3,1)^2)$$

→ Per determinare scientificamente il numero di repliche serve conoscere:

- $\sigma$  il rumore dell'esperimento
- ordine di grandezza degli effetti che voglio poter rilevare
- probabilità di rilevare un simile effetto (= potenza del test)

→ Siccome spesso non si conosce  $\sigma$ , esiste la regola empirica:

$$\left\lceil \frac{\text{run} + 32}{\text{run}} \right\rceil = \text{repliche per run}$$

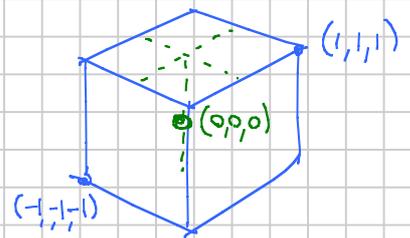
|     |   |          |
|-----|---|----------|
| L4  | 9 | repliche |
| L8  | 5 | "        |
| L16 | 3 | "        |
| L32 | 2 | "        |
| L64 | 2 | "        |
| ⋮   | ⋮ | "        |
| ⋮   | ⋮ | "        |

} per i Taguchi

### Center point

E' una run aggiuntiva

| run | A  | B | C | D | E | rep1 | rep2 | rep3 | media |
|-----|----|---|---|---|---|------|------|------|-------|
| ⋮   | ±1 | ⋮ | ⋮ | ⋮ | ⋮ | ~    | ~    | ~    | ~     |
| 17  | 0  | 0 | 0 | 0 | 0 | ~    | ~    | ~    | ~     |



→ stesso numero di repliche

→ si fa la regressione sulle altre run, lasciando fuori

→ si controlla se il valore osservato per il cp è compatibile con il modello di regressione trovato

$$\hat{Y} = B_0 + B_1 x_1 + \dots + B_p x_p \quad \text{valore previsto generico}$$

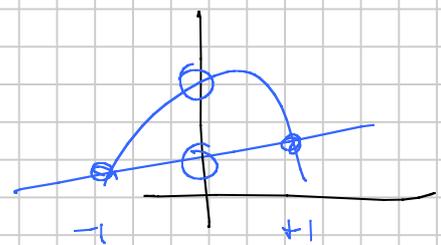
$$\hat{Y} = B_0 \quad \text{valore previsto per il c.p.}$$

confronto  $B_0 \pm Se$  o  $B_0 \pm 2Se$  con il valore misurato

→ se è compatibile, il modello va bene

se non è compatibile vuol dire che il fenomeno presenta forti non linearità

★ Va rifatto tutto il DOE passando a 3 livelli.



## ● Randomizzazione delle repliche

effetto temporale (tipo curve di apprendimento)

| run | A  | B  | C  | rep 1 | rep 2 |
|-----|----|----|----|-------|-------|
| 1   | -1 | -1 | +1 | 12    | 25    |
| 2   | -1 | +1 | -1 | 36    | 47    |
| 3   | +1 | -1 | -1 | 51    | 63    |
| 4   | +1 | +1 | +1 | 24    | 88    |

Lasso  
alto

Se faccio la regressione, A viene significativa, importante  
 Se faccio le repliche in ordine casuale, risolve il problema  
 → Occhio che la rand. costa, quindi si trovano resistenze

## ● Ortogonalità

Per evitare di perderla, la matrice delle risposte deve essere completa, senza buchi, anche dopo aver tolto gli outliers

- ↳ casomai si rifanno i singoli esperimenti problematici
- ↳ occhio ai difetti cronici del design

## ● Regressione vera e propria

↳ 1 punto per ogni run →  $\bar{Y}$  = media (repliche)

↳ oppure 1 punto per ogni replica →  $\bar{Y}$  = singole repliche

| run       | rep | A  | B  | C  | Y |
|-----------|-----|----|----|----|---|
| 1         | 1   | -1 | -1 | +1 | ~ |
| 1         | 2   | -1 | -1 | +1 | ~ |
| 2         | 1   | -1 | +1 | -1 | ~ |
| 2         | 2   | -1 | +1 | -1 | ~ |
| - - - - - |     |    |    |    |   |

\* Si trovano gli stessi identici valori per i coefficienti di regressione

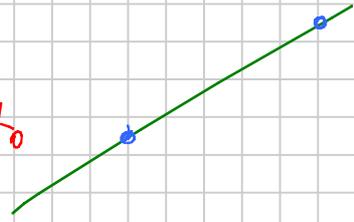
\* Cambia se e la potenza dei test

$\sim \frac{\sigma}{\sqrt{r}}$        $\sim \sigma$        $r$ : # di repliche

L16 5 fattori

- 1 I + ABCDE
- 2 A + BCDE
- 3 B + ACDE
- C + ABDE
- D + ...
- E + ...
- AB + CDE
- AC + BDE
- AD + ...
- AE + ...
- BC + ...
- BD + ...
- BE + ...
- CD + ...
- CE + ...
- 16 DE + ...

16 coefficienti con 16 punti  
residui nulli → se non definito



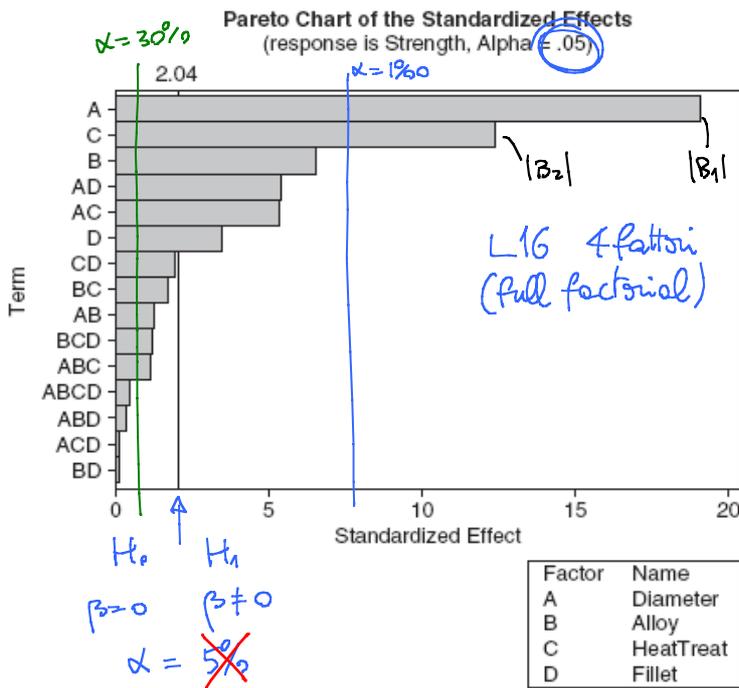
\* Basta togliere qualche variabile

ora 23

Selezione delle variabili

Si fa in modo completamente diverso dal solito.

Si fa il diagramma di Pareto dei valori assoluti dei coefficienti e si analizza qualitativamente



$$\frac{B_i}{\sqrt{\sum_{ii} E_{ii}}} = \text{stat} \quad H_0: \beta_i = 0$$

$|B_i|$

$$Y = B_0 + B_1 A + B_2 C + \cancel{B_3 B} + \dots$$

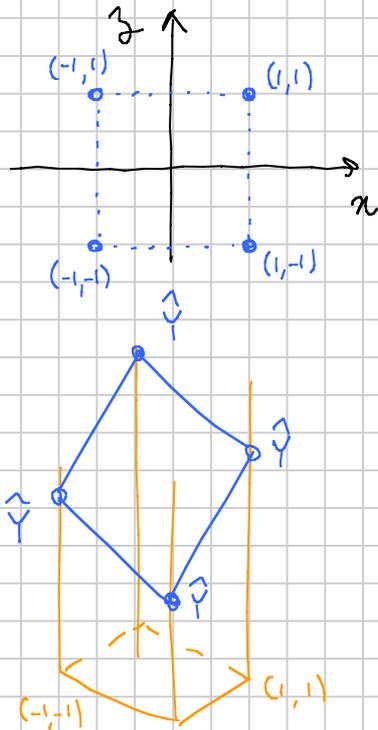
A, C    A, C, B    A, C, B, AD, D, AC

\* Posso togliere le variabili che non mi interessano *tutte insieme* perché il design è ortogonale e rilanciare la regressione. Ottengo solo una stima più giusta di  $\sigma$ , con un Se più preciso.

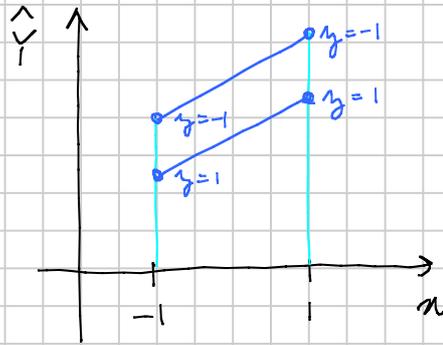
# Termini di interazione

$$\hat{Y} = B_0 + B_1 x + B_2 y + B_3 xy$$

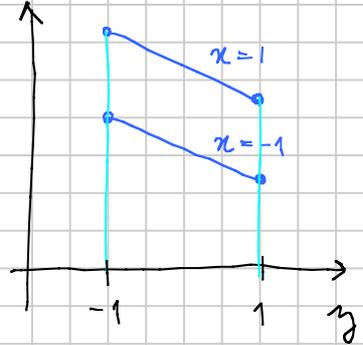
termini di interazione



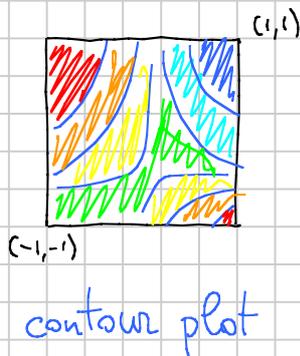
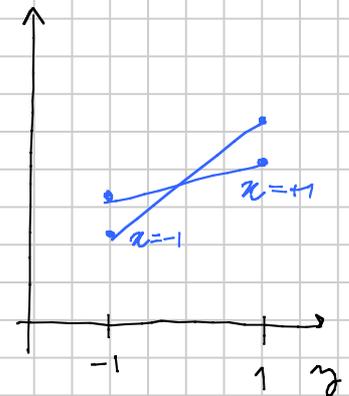
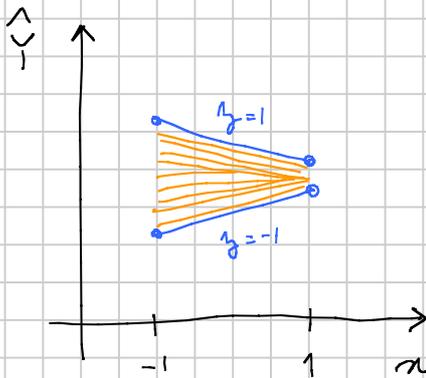
caso  $B_3 = 0$



$$\hat{Y} = B_0 + B_1 x + B_2 y$$

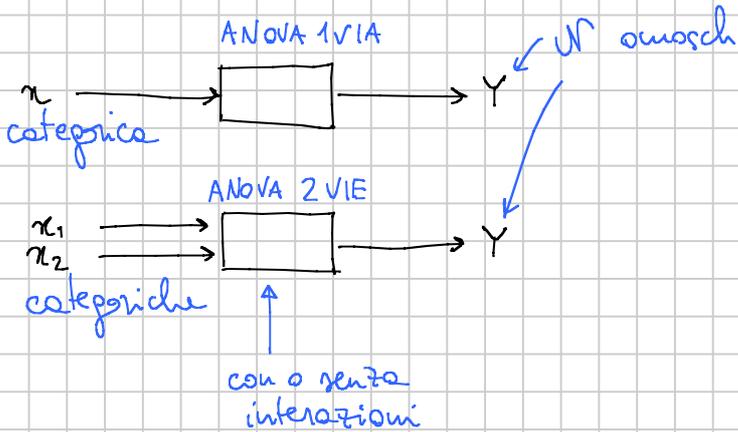


caso  $B_3 \neq 0$



## ANALISI DELLA VARIANZA

(Cap 10 Ross)

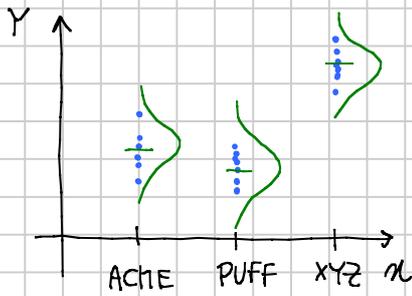


$Y$ : % di polvere  $< 5 \mu m$   
 $x$ : fornitore lattorio: 3 ditte

$Y$ : consumi di carburante  
 $x_1$ : diesel normale, 3 varianti sperimentali (4 carburanti)  
 $x_2$ : tipo di auto: utilitaria, berlina, mv

\* Si usa quando le categorie sono 2-5 (mai 10+)

# ANOVA A 1 VIA



$$Y = f(x) = a + bx + cx^2$$

|               |                    |                                               |
|---------------|--------------------|-----------------------------------------------|
| $\mu_{ACHIE}$ | $x = \text{ACHIE}$ | $+ e \sim \mathcal{N}(0, \sigma^2)$<br>omosch |
| $\mu_{PUFF}$  | $x = \text{PUFF}$  |                                               |
| $\mu_{XYZ}$   | $x = \text{XYZ}$   |                                               |

$$Y = \mu_x + e$$

• Struttura dei dati :  $m$  campioni di lunghezze variabili

$$\{\text{ACHIE, PUFF, XYZ}\} = \{1, 2, \dots, m\}$$

|                                      |                                         |                                                         |                                                                           |
|--------------------------------------|-----------------------------------------|---------------------------------------------------------|---------------------------------------------------------------------------|
| $Y_{1,1}, Y_{1,2}, \dots, Y_{1,n_1}$ | $\sim \mathcal{N}(\mu_1, \sigma^2)$ iid | $\mu_1 \approx Y_{1,*} := \frac{1}{n_1} \sum_i Y_{1,i}$ | $\sigma^2 \approx S_1^2 = \frac{1}{n_1 - 1} \sum_i (Y_{1,i} - Y_{1,*})^2$ |
| $Y_{2,1}, Y_{2,2}, \dots, Y_{2,n_2}$ | $\sim \mathcal{N}(\mu_2, \sigma^2)$ iid | $\mu_2 \approx Y_{2,*} := \dots$                        | $\sigma^2 \approx S_2^2 = \dots$                                          |
| .....                                |                                         |                                                         |                                                                           |
| $Y_{m,1}, Y_{m,2}, \dots, Y_{m,n_m}$ | $\sim \mathcal{N}(\mu_m, \sigma^2)$ iid | $\mu_m \dots$                                           | $\sigma^2 \approx \dots$                                                  |

ora 24

$$Y_{i,*} \sim \mathcal{N}\left(\mu_i, \frac{\sigma^2}{n_i}\right) \quad \text{stimatore } \mu_i \text{ corretto e consistente}$$

↳ HWC

Ho poi  $m$  stimatori diversi e indipendenti per  $\sigma^2$

$$S_i^2 : \frac{S_i^2}{\sigma^2} (n_i - 1) \sim \chi^2(n_i - 1)$$

Voglio riarrangerli in uno stimatore solo per  $\sigma$ , migliore di tutti

$$S_M^2 = \theta_1 S_1^2 + \theta_2 S_2^2 + \dots + \theta_m S_m^2 \quad \theta_1 + \theta_2 + \dots + \theta_m = 1$$

media pesata

pesi

i.  $E(S_M^2) = \sigma^2$

ii.  $\text{Var}(S_M^2)$  dipende dai  $\theta_i \rightarrow$  il valore minimo si ha quando

$$\theta_i = \frac{gdi}{gdi_{\text{Tot}}} = \frac{n_i - 1}{\sum_j (n_j - 1)}$$

iii. Quando i  $\theta_i$  sono questi si usa il simbolo  $S_p^2$ , si

chiamo stimatore pooled della varianza e inoltre:

$$\frac{S_p^2}{\sigma^2} \cdot \text{gdl}_{\text{TOT}} \sim \chi^2(\text{gdl}_{\text{TOT}})$$

HWC

$$S_p^2 := \frac{n_1-1}{N-m} S_1^2 + \frac{n_2-1}{N-m} S_2^2 + \dots + \frac{n_m-1}{N-m} S_m^2$$

dove  $N = \sum_{i=1}^m n_i$  e quindi  $\text{gdl}_{\text{TOT}} = \sum_i (n_i - 1) = N - m$

★  $S_p^2 = S_w^2$  varianza "within" (=entro i campioni)

★  $S_B^2$  varianza "between" (=tra i campioni)

campione 1  $\rightarrow Y_{1,*}$   
 " 2  $\rightarrow Y_{2,*}$   
 .....  
 $\rightarrow Y_{m,*}$

} faccio la varianza campionaria

$$S_B^2 := \frac{1}{m-1} \sum_{j=1}^m n_j (Y_{j,*} - Y_{*,*})^2$$

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

media globale

★ Se  $Y$  non dipende da  $x$ , ovvero se  $\mu_1 = \mu_2 = \dots = \mu_m$ , allora si può verificare (HWC) che

$$\sigma^2 \approx S_B^2$$

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

Se  $Y$  dipende da  $x$ ,  $S_B^2$  è tipicamente più grande

$\rightarrow$  Inoltre  $S_B^2$  è sempre indipendente da  $S_w^2$

• L'anova a 1 via formalmente esegue il test:

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

( $Y$  non dipende da  $x$ )

$$H_1: \text{non tutte uguali}$$

( $Y$  dipende da  $x$ )

il test si basa su  $S_B^2$ : se è grande è un indizio per  $H_1$

in concreto calcola la statistica:

$$F := \frac{S_B^2}{S_W^2}$$

se è vera  $H_1$   $F$  tipicamente  $> 1$

se è vera  $H_0$   $F$  vicino a 1

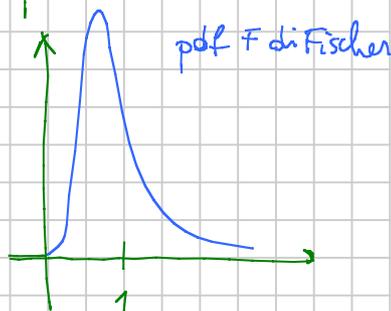
\* distribuzione sotto  $H_0$ : se  $S_B^2$  e  $S_W^2$  sono entrambi stimatori di  $\sigma^2$  basati su  $\chi^2$  con  $a$  e  $b$  gdl, allora

$$\frac{S_B^2}{S_W^2} \sim F(a; b)$$

$\uparrow$   $\uparrow$   $\uparrow$   
 F di Fischer    gdl denominatore    gdl numeratore

### F di Fischer

Classe di vaa continue, concentrate intorno a 1, positive, con 2 parametri

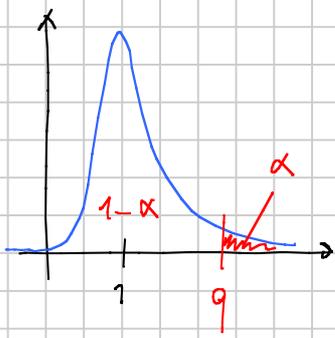


Se  $S_1^2, S_2^2$  sono indep con  $\frac{S_1^2}{\sigma^2} \sim \chi^2(a)$   $\frac{S_2^2}{\sigma^2} \sim \chi^2(b)$

Allora  $\frac{S_1^2}{S_2^2} \sim F(a; b)$

$$F := \frac{S_B^2}{S_W^2} \underset{H_0}{\sim} F(m-1; N-m)$$

$\uparrow$  sotto  $H_1$   $F$  è più grande



$F > q$  dico  $H_1$   
 $F \leq q$  dico  $H_0$

★ Altra inferenza:

$$\sigma^2 = S_w^2 \quad \frac{S_w^2}{\sigma^2} (N-m) \sim \chi^2(N-m)$$

$$\mu_i \approx Y_{i,x} \quad \frac{Y_{i,x} - \mu_i}{S_w / \sqrt{n_i}} \sim t(N-m)$$

⊙ Identità utili

$$S_w^2 = \frac{1}{N-m} \sum_{i=1}^m \sum_{j=1}^{n_i} (Y_{i,j} - Y_{i,x})^2 \quad (\text{HWC})$$

$$S_w^2, S_B^2, S_Y^2 = \frac{1}{N-1} \sum_{i,j} (Y_{i,j} - Y_{x,x})^2 \quad \text{varianze}$$

$SS_w, SS_B, SS_Y$  devianze

$$SS_w = S_w^2 (N-m) \quad SS_B = S_B^2 (m-1) \quad SS_Y = S_Y^2 (N-1)$$

$$SS_Y = SS_B + SS_w \quad \text{identità delle devianze:}$$

Info pratiche:

- oggi solo 2 ore
- domani scambio Nicolodi/Morandin  
(laboratorio 8:30-10:30, Nicolodi 14:30-16:30) ↗ aula boh!
- domani niente ricevimento (per cose brevi e urgenti, subito dopo lezione)
- mercoledì 11 maggio niente lezione

## ANOVA A 1 VIA

• Cfr audio

• Se  $m=2$   $Y_{1,1}, Y_{1,2}, \dots, Y_{1,n_1}; Y_{2,1}, Y_{2,2}, \dots, Y_{2,n_2}$   
 $\downarrow \qquad \qquad \qquad \downarrow$   
 $\mathcal{N}(\mu_1, \sigma^2) \qquad \qquad \mathcal{N}(\mu_2, \sigma^2)$

Anova a 1 via:  $H_0: \mu_1 = \mu_2$        $H_1: \mu_1 \neq \mu_2$

stesse ipotesi di lavoro del test per il confronto delle medie di due popolazioni normali

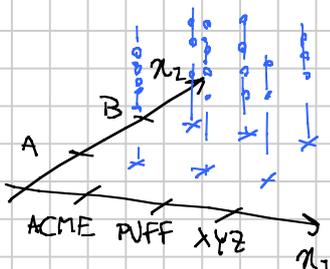
$$\frac{\bar{X} - \bar{Y}}{S_p \sqrt{\frac{1}{m} + \frac{1}{n}}}$$

$$\frac{Y_{1,x} - Y_{2,x}}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \stackrel{H_0}{\sim} t(n_1 + n_2 - 2)$$

★ Sono lo stesso test: il p-valore viene identico.

• Vedi ora 20 anno scorso: test globale di regressione

## ANOVA A DUE VIE



$x_2$ : tipo di mescolamento  
 $Y$ : % < 5µm dopo mescolato

| $\alpha_1$ | $\alpha_2$ | 1                         | 2                         |
|------------|------------|---------------------------|---------------------------|
| 1          |            | $X_{111}, X_{112}, \dots$ | $X_{121}, X_{122}, \dots$ |
| 2          |            | ...                       | ...                       |
| 3          |            |                           | $X_{321}, \dots$          |

$X_{ijk}$  dato  $k$ -esimo dell'esperimento fatto con  $\alpha_1 = i$  e  $\alpha_2 = j$

$X_{ij1}, X_{ij2}, \dots$  sono repliche

- In ogni casella deve lo stesso numero di repliche se no non funziona.
  - se è possibile ricondursi a numero di repliche uguali buttando via non troppi dati, è utile farlo (se no, regressione)
  - se una casella è vuota o ha comunque troppo pochi dati si può valutare se eliminare la riga o la colonna

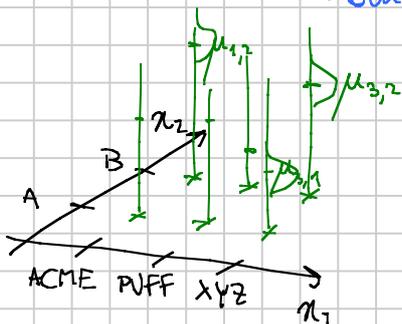
- $i = 1, 2, \dots, m$  valori di  $\alpha_1$
  - $j = 1, 2, \dots, n$  valori di  $\alpha_2$
  - $k = 1, 2, \dots, l$  repliche
- } var. categoriche
- $l = 1$  anova 2 vie classica
- $l > 1$  anova a 2 vie con repliche / con interazioni

→ NB Se  $l > 1$  posso ottenere info in più, ma comunque posso anche ricondursi a  $l = 1$  facendo la media su ogni casella.

### SENZA REPLICHE ( $l = 1$ )

$$X_{ij} \sim \mathcal{N}(\mu_{ij}, \sigma^2)$$

↑ omoschedastico



ipotesi fondamentale di linearità:

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

↑ media globale

↑ contributo dovuto ad  $\alpha_1$

↑ contributo dovuto ad  $\alpha_2$

|   |      |      |     |
|---|------|------|-----|
| B | 72   | 55   | 61  |
| A | 65   | 48   | 54  |
|   | ACME | PUFF | XYZ |

forzati dal modello lineare  
 se è troppo restrittivo, serve un  
 modello con interazioni :  $l > 1$

→ ipotesi tecnica (non restrittiva)

$$\sum_{i=1}^m \alpha_i = 0 \quad \sum_{j=1}^n \beta_j = 0$$

•  $X_{ij} \sim \mathcal{N}(\mu + \alpha_i + \beta_j; \sigma^2)$

si usa  
(HWC)

$$X_{i,*} = \frac{1}{n} \sum_{j=1}^n X_{ij} \sim \mathcal{N}(\mu + \alpha_i; \frac{\sigma^2}{n})$$

$$X_{*,j} = \frac{1}{m} \sum_{i=1}^m X_{ij} \sim \mathcal{N}(\mu + \beta_j; \frac{\sigma^2}{m})$$

$$X_{*,*} = \frac{1}{nm} \sum_{i,j} X_{ij} = \frac{1}{m} \sum_i X_{i,*} = \frac{1}{n} \sum_j X_{*,j} \sim \mathcal{N}(\mu, \frac{\sigma^2}{mn})$$

ora 26

Stimatori

$$\mu \approx X_{*,*} \sim \mathcal{N}(\mu, \frac{\sigma^2}{mn})$$

$$\alpha_i \approx X_{i,*} - X_{*,*} \sim \mathcal{N}(\alpha_i; \sigma^2 \frac{m-1}{mn})$$

HWC \*

effetto riga

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

$$\beta_j \approx X_{*,j} - X_{*,*} \sim \mathcal{N}(\beta_j; \sigma^2 \frac{n-1}{mn})$$

effetto colonna

$j = 1, 2, \dots, n$

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

risp. media riga

$i = 1, \dots, m$

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

risp. media colame

$j = 1, \dots, n$

$$\mu_{ij} = \mu + \alpha_i + \beta_j \approx X_{i,*} + X_{*,j} - X_{*,*} \sim \mathcal{N}(\mu_{ij}; \sigma^2 \frac{m+n-1}{mn})$$

$(i, j) \dots$

media che mi aspetto  $(i, j)$

previsto  $(i, j)$

HWC \*\*

$$R_{ij} := X_{ij} - (\mu_{ij}) \sim \mathcal{N}(0; \sigma^2 \frac{(m-1)(n-1)}{mn})$$

$(i, j) \dots$

residuo  $(i, j)$

ovvero osservato meno previsto

R crea confusione

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

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

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

$$\sum_{i,j} \frac{(X_{ij} - \mu_{ij})^2}{\sigma^2} \sim \chi^2(mn) \quad \text{per def di } \chi^2$$

$$\sum_{i,j} \left[ \frac{X_{ij} - (X_{i\cdot} + X_{\cdot j} - X_{\cdot\cdot})}{\sigma^2} \right]^2 = \sum_{i,j} \frac{R_{ij}^2}{\sigma^2} = \frac{SSe}{\sigma^2} \sim \chi^2(mn - (m+n-1))$$

gdl da togliere  
↓  
 $(m+n-1)$

$$\sim \chi^2((m-1)(n-1))$$

★  $\sigma^2 \approx Se^2$  stimatore che funziona sempre

● Test fondamentali dell'anova a 2 vie

→ Effetto riga:  $Y$  dipende da  $\alpha_1$ ?

→ Effetto colonna:  $Y$  dipende da  $\alpha_2$ ?

analogamente alla regressione (per la selezione delle variabili)

★ Eff riga:  $H_0: \alpha_1 = \alpha_2 = \dots = \alpha_m = 0$   $H_1$ : non tutti nulli  
( $Y$  non dipende da  $\alpha_1$ ) ( $Y$  dipende da  $\alpha_1$ )

$$F_R := \frac{S_R^2}{Se^2} \underset{H_0}{\sim} F(m-1; (m-1)(n-1))$$

→ piccola dico  $H_0$   
↳ grande dico  $H_1$

dove  $S_R^2 := \frac{SS_R}{m-1}$  varianza per righe

$$SS_R := \sum_{i=1}^m n (X_{i\cdot} - X_{\cdot\cdot})^2$$

devianza per righe

★ Eff colonna:  $H_0: \beta_1 = \dots = \beta_n = 0$   $H_1$ : non tutti nulli  
( $Y$  non dipende da  $\alpha_2$ ) ( $Y$  dipende da  $\alpha_2$ )

$$F_C := \frac{S_C^2}{Se^2} \underset{H_0}{\sim} F(n-1; (m-1)(n-1))$$

→ piccola dico  $H_0$   
↳ grande dico  $H_1$

dove  $S_C^2 := \frac{SS_C}{n-1}$  varianza per colonne

$$SS_C := \sum_{j=1}^n m (X_{\cdot j} - X_{\cdot\cdot})^2$$

devianza per colonne

- Se una delle due var risulta non significativa, può a volte essere rimosso (ovvero si rifà l'analisi come ANOVA a 1 via)



→ Mi interessa rifare il test perché spero di trovare un  $H_1$



→ Non mi interessa rifare il test (è già  $H_1$ ), ma posso trovare stime più precise dei parametri e dei previsti

guardare il test t-paired nel Cap 8 del Ross

\* Senti anche ora 27

★ Settimana prossima niente lezioni

Ultima volta teoria : 1 giugno

Ultima volta lab : da definire, ma quella settimana

● Legami ANOVA test classici

ANOVA 1 via con 2 campioni soli  $\Leftrightarrow$  test  $\frac{\bar{X} - \bar{Y}}{S_p \sqrt{\frac{1}{m} + \frac{1}{n}}}$

ANOVA 2 vie con  $m=2$ ,  $n$  qualsiasi  $\Leftrightarrow$  test t- Paired

$$\begin{array}{l} X_1, X_2, \dots, X_n \\ Y_1, Y_2, \dots, Y_n \end{array} \sim \mathcal{N}(\mu_x, \sigma^2)$$

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

$$X_i \sim \mathcal{N}(\mu_i, \sigma_x^2) \quad Y_i \sim \mathcal{N}(\mu_i + \delta, \sigma_y^2)$$

od esempio  $n$  persone pesate prima e dopo una dieta

$$Z_i := Y_i - X_i \sim \mathcal{N}(\delta, \sigma_z^2)$$

$$Z_1, Z_2, \dots, Z_n \sim \mathcal{N}(\delta, \sigma_z^2)$$

$$H_0: \delta = 0 \quad H_1: \delta \neq 0$$

$$\frac{\bar{Z} - \overset{0}{\mu_0}}{S_z / \sqrt{n}} \stackrel{H_0}{\sim} t(n-1)$$

★ In entrambi i casi i test sembrano differenti, ma si ottiene esattamente lo stesso  $p$  dei dati.

■ CON REPLICHE ( $l \geq 2$ )

$$X_{ij,k} \sim \mathcal{N}(\mu_{ij}, \sigma^2) \quad k=1, 2, \dots, l$$

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

$\mu_{ij}$  qualsiasi

NO REPLICHE

REPLICHE

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

scostamento dal modello lineare  
 = interazione (effetto simile ai termini  $\alpha_i \alpha_j$  nella regressione e nel DoE)

● Test che vengono effettuati:

a)  $H_0: \delta_{ij} \equiv 0 \quad \forall ij$      $H_1$ : non tutti nulli    test sulle interazioni

se dice  $H_0$  si può fare la media in ogni casellina e passare all'anova a 2 vie

$$X_{ij} := X_{ij*} = \frac{1}{e} \sum_{k=1}^e X_{ijik}$$

$$\sim \mathcal{N}\left(\mu + \alpha_i + \beta_j, \frac{\sigma^2}{e}\right)$$

→ ANOVA 2 VIE SENZA REPLICHE

attenzione che questo modello ha una varianza fittizia

→ Se dice  $H_1$ :

b)  $H_0: \alpha_i \equiv 0 \quad \forall i$      $H_1$ : non tutti nulli    test effetto riga

c)  $H_0: \beta_j \equiv 0 \quad \forall j$      $H_1$ : non tutti nulli    test effetto colonna

## ■ TEST DI ADATTAMENTO

→ Test del chi-quadro elementare

→ Test del chi-quadro con stima di parametri

} test di adattamento ad una distribuzione

→ Tabelle di contingenza → varie questioni

## ■ TEST $\chi^2$ ELEMENTARE

→ incluse le categoriche

esempio: va. discreta con una legge teorica nota

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

dati Istat →    4%    75%    7%    14%

Prendo un campione forse ha legge diversa: residenti a Treviso

$$n=300$$

$$6$$

$$183$$

$$60$$

$$51$$

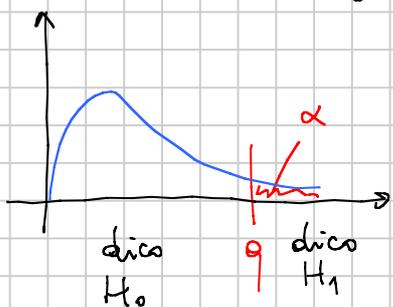
Q: questo campione può sensatamente venire da quella distribuzione ( $H_0$ ), oppure vi è evidenza che la distribuzione è diversa ( $H_1$ )?

ora 28

$$\frac{(6 - \overset{4\% \cdot 300}{12})^2}{12} + \frac{(183 - 225)^2}{225} + \frac{(60 - 21)^2}{21} + \frac{(51 - 42)^2}{42} \rightarrow \text{statistica}$$

↳ se è grande dico  $H_1$ , se è piccola dico  $H_0$

↳ uso quantili  $\chi^2$  con 3 (categorie - 1) g.d.l.



$$q = \text{INV.CHI}(\alpha; 3)$$

• Il test del  $\chi^2$  elementare funziona su una v.a.  $X$  che assume solo i valori  $\{1, 2, \dots, k\}$  dove  $k$  è finito, piccolo e si suppone una certa legge  $\varphi_0$

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

$H_1$ : non tutte uguali

→ Si chiama test di adattamento a distribuzione discreta completamente specificata

→ Statistica del test

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

dove  $O_i$ : osservati → quanti del campione  $X_1, X_2, \dots, X_n$  sono  $i$

dove  $A_i$ : atenni  $\rightarrow A_i = n \cdot \varphi_0(i)$

$\hookrightarrow$  NB. Non sempre gli  $A_i$  sono interi, non occorre che lo siano

$\rightarrow W \sim \chi^2(k-1)$  sotto  $H_0$

$W$  assume valori maggiori sotto  $H_1$

★ L'approssimazione è buona quando gli atenni  $A_i$  sono tutti numeri grandi

$\hookrightarrow n$  grande

$\hookrightarrow k$  piccolo

$\hookrightarrow \varphi_0$  abbastanza uniforme

} tutti aspetti auspicabili

$\rightarrow$  "Rule of thumb": tutti gli  $A_i > 1$  e  
almeno  $i \frac{4}{5}$  degli  $A_i > 5$

★ Quando l'approssimazione non è buona, si dovrebbe usare la vera legge di  $W$  sotto  $H_0$ .  
È calcolabile con metodi sofisticati (distrib. multinomiali) oppure si può fare una simulazione Monte Carlo.

Ad esempio, genero  $N=1000$  volte le 300 variabili  $X_1, \dots, X_{300}$  con legge discreta  $\varphi_0 = (4\%, 75\%, 7\%, 14\%)$

Per 1000 volte calcolo  $W \rightarrow$  ottengo 1000 valori  $W_1, \dots, W_{1000}$  che hanno la distribuzione vera di  $W$  sotto  $H_0$

$\alpha = 5\%$   $q = \text{PERCENTILE}([W_1, W_2, \dots, W_{1000}]; 0,95)$

$95\%$   $\downarrow$   $5\%$   
 $950 W_i \leftarrow q \rightarrow 50 W_i$

★ NB  $(O_1, O_2, O_3, O_4)$  sotto  $H_0$  ha legge multinomiale di parametri  $n=300$ ,  $p_1=4\%$ ,  $p_2=75\%$ ,  $p_3=7\%$ ,  $p_4=14\%$

\* Ad esempio si può usare questo test per verificare se la distanza in mesi dal compleanno al decesso ha legge uniforme. Si scopre che (per i vip) non è così.

■ Generalizzazione a TANTI o INFINITI valori

→ Si raggruppa  $X$  in pochi valori e si fa il test elementare

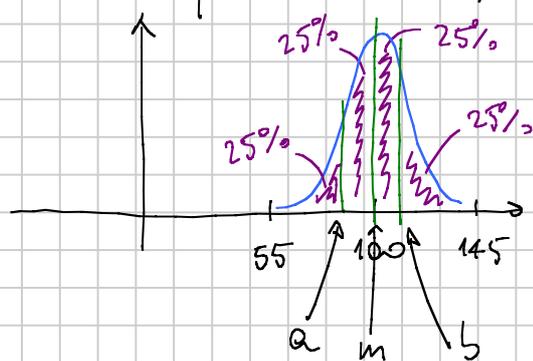
Ad esempio: vedi file Excel

ora 29

■ Generalizzazione a variabili CONTINUE

→ Si dividono i valori possibili in fasce (pochi e possibilmente equiprobabili) e si fa il test elementare

ad esempio:  $N(100, 15^2) \leftrightarrow$  quoziente di intelligenza.



$$a = \text{INV.NORM}(25\%; 100; 15)$$

$$b = \text{INV.NORM}(75\%; 100; 15)$$

$$m = 100 \quad (= \text{INV.NORM}(50\%; 100; 15))$$

\* Se i bin sono più o meno equiprobabili l'approssimazione è migliore e il test è più potente

↳ Se  $X$  è continua, conoscendo  $F^{-1}$  è possibile (e raccomandato) ottenere bin esattamente equiprobabili.

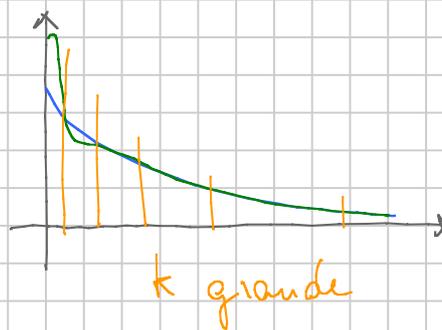
● Quanti bin bisogna fare :

i. maggiore  $n$ , maggiore  $k$

ii.  $k$  mai troppo grande (gli altri devono rispettare le "rule of th")

iii. più  $k$  è piccolo più il test è potente (differenze anche piccole, ma globali)

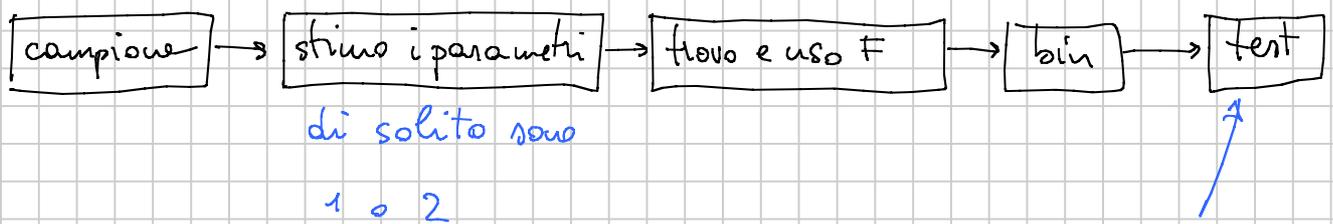
iv. più  $k$  è grande più il test è sensibile (differenze locali, non troppo piccole)



★ In mancanza di orientamenti precisi, consiglio  $k$  piccolo, dispari ( $k=5, 7, 9$ )

■ Generalizzazione a distribuzione specificate a meno di parametri

Q: Questo campione è normale?



$$gdl = k - 1 - \# \text{parametri stimati}$$

★ Scegliere  $k$  poco più grande per non trovarsi con troppo pochi gdl

# TABELLE DI CONTINGENZA

Esempio: parti cesarei / ospedale

|      | P.C. | P.N. |
|------|------|------|
| E.R. | 22   | 74   |
| H.   | 18   | 45   |
| G.R. | 33   | 71   |

- Q1: La propensione al P.C. è la stessa in tutti gli ospedali?
- Q2: La variabile "tipo di parto" è influenzata da quella "ospedale"?
- Q3: La variabile "ospedale" è influenzata da quella "tipo di parto"?
- Q4: Le variabili sono indipendenti?

→ C'è un solo test:

$H_1$ : No, sì, sì, No

$H_0$ : sì, No, No, sì

risponde a tutte le Q  
(perché in realtà sono una sola)

• Come si esegue il test:

|      | P.C. | P.N. |     |
|------|------|------|-----|
| E.R. | 22   | 74   | 96  |
| H.   | 18   | 45   | 63  |
| G.R. | 33   | 71   | 104 |
|      | 73   | 190  | 263 |

$27,8\% \cdot 96 = 26,6 = A_{ER,PC}$

otteni nell'ipotesi "di indipendenza"

$$A_{ij} = \frac{T_{i\cdot} \cdot T_{\cdot j}}{T_{xx}}$$

dove  $T_{i\cdot} = \sum_j O_{ij}$      $T_{\cdot j} = \sum_i O_{ij}$   
 $T_{xx} = \sum_i \sum_j O_{ij}$

$27,8\% = \frac{73}{263}$  stima dell'incidenza dei cesarei

|      | P.C. | P.N. |     |
|------|------|------|-----|
| E.R. | 22   | 74   | 96  |
| H.   | 18   | 45   | 63  |
| G.R. | 33   | 71   | 104 |
|      | 73   | 190  | 263 |

osservati

|      | P.C. | P.N. |     |
|------|------|------|-----|
| E.R. | 26,6 | 69,4 | 96  |
| H.   | 17,5 | 45,5 | 63  |
| G.R. | 28,9 | 75,1 | 104 |
|      | 73   | 190  | 263 |

otteni

$$W = \sum_{i,j} \frac{(O_{i,j} - A_{i,j})^2}{A_{i,j}}$$

↑  
somme su tutte le caselle

$H_0$

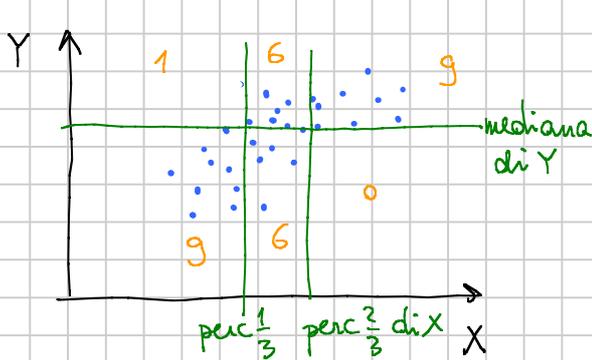
$$\sim \chi^2 [(m-1)(n-1)]$$

↑ # di righe  
↑ # di colonne

Esame : mar 7 18+ritardo round 0 (crocette)  
 mer 8 8.30 round 1 (Excel)  
 entro sab 4 8.30 inviare file desiderati al round 1  
 entro lun 6 li rimetto online corretti e ripuliti

## TABELLE DI CONTINGENZA

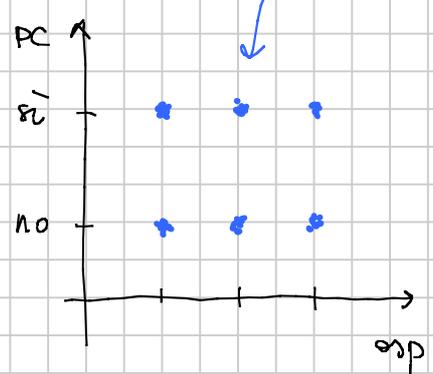
2 Variabili numeriche con tanti valori



|   | X |   |   |
|---|---|---|---|
| Y | 1 | 6 | 9 |
|   | 9 | 6 | 0 |

dira  $H_1$

var sono categoriche



\* Il test mi dice se X e Y sono indipendenti ( $H_0$ ) o correlate ( $H_1$ )

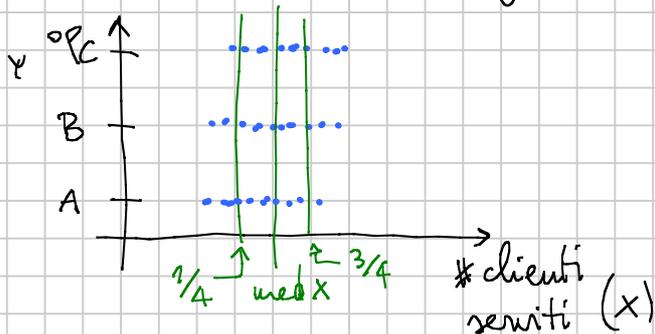
→ Test non parametrico analogo alla regr. semplice

non ha ipotesi (omosched, gauss, linearità, ...)

ma è meno potente e non fornisce un modello quantitativo

→ Conviene di solito conviene "affettare" ciascuna variabile in modo equiprobabile. Di solito  $2 \times 2$ ,  $2 \times 3$  o  $3 \times 3$  sono numeri adeguati

1 variabile categorica e una numerica



→ Test non parametrico analogo alla ANOVA a 1 via

★ Il test mi dice se la distrib di X è diversa per A, B, C ( $H_1$ )  
 (non solo la media ma "tutta" la distribuzione) o no ( $H_0$ )

↑  
 dipende da quante fette faccio

= ... se X dipende da Y ( $H_1$ ) o no ( $H_0$ )

N.B. Situazione interessante

|                 |                      | A    | B   | C   | D    |        |
|-----------------|----------------------|------|-----|-----|------|--------|
| X categorica    | : campione           | 17   | 6   | 2   | 12   | 37     |
| $p_0$ assegnate | $\chi^2$ sempl. 3gdl | 50%  | 10% | 10% | 30%  | 10%    |
|                 |                      | 18,5 | 3,7 | 3,7 | 11,1 | alteri |
|                 | campione 2           | 22   | 9   | 10  | 19   | 60     |

tab cont  
 $4 \times 2 \rightarrow 3 \text{gdl}$

■ Test esatto di Fisher o test di Fisher - Irwin

Permette di ottenere un p dei dati esatto nel caso di tabelle di contingenza  $2 \times 2$  (qualunque numerosità)

↳ Se le numerosità sono elevate non è pratico usarlo

|    |    |    |    |    |    |
|----|----|----|----|----|----|
| 7  | 2  | 9  | 4  | 5  | 9  |
| 4  | 12 | 16 | 7  | 9  | 16 |
| 11 | 14 | 25 | 11 | 14 | 25 |

alteri

legge  $V$

25 palline : 9 rosse, 16 blu

Ne pescò 11 a caso : quante rosse escano?

Questa legge è detta ipergeometrica, una legge discreta che ha 3 parametri interi

$$P(V=k) = \varphi_V(k) = \text{DISTRIB. IPERGEOM}(k; 11; 9; 25)$$

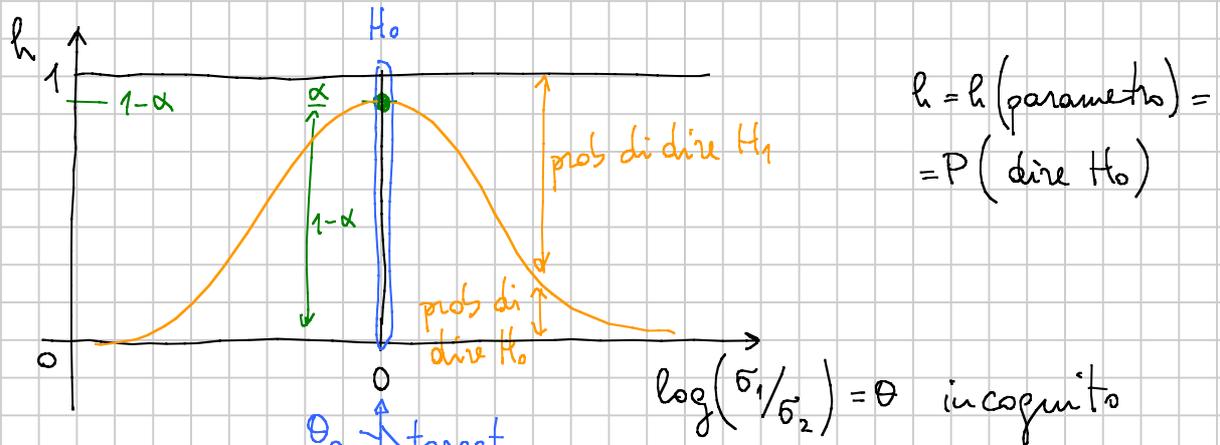
★ Sotto  $H_0$  (indipendenza delle 2 var)  $V$  ha legge ipergeom.  
 Se osservo valori di  $V$  troppo estremi ho il diritto di scegliere  $H_1$   
 → Si fa sempre con il  $p$  dei dati :

$\alpha^*$  : somma delle probabilità  $\varphi_V(k)$  per quei  $k$  che sono almeno altrettanto estremi/rari del valore effettivamente osservato

ora 31

TEST STATISTICI E CURVE OC

- |               |                               |                                                                                                                                                              |
|---------------|-------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------|
| ① stimatore   | } bilaterali<br>e unilaterali | } m parametri: $\mu, \sigma, p, \lambda, \mu_1, \mu_2$<br>$\beta, \alpha, \beta_i, \sigma, \frac{\sigma_1}{\sigma_2}$<br>(non m parametri: anova, $\chi^2$ ) |
| ② statistiche |                               |                                                                                                                                                              |
| ③ p dei dati  |                               |                                                                                                                                                              |



test bilaterale  $H_0: \sigma_1 = \sigma_2$   $H_1: \sigma_1 \neq \sigma_2$  ( $H_0: \frac{\sigma_1}{\sigma_2} = \eta$ )

②  $1-\alpha = P(\text{dire } H_0 \mid \text{vera } H_0) = \dots \Rightarrow RA$

$h(\theta) = P(\text{dire } H_0)$  :  $\beta$  prob en II specie (se  $\theta \notin H_0$ )  
 $1-h(\theta) = P(\text{dire } H_1)$  : potenza del test (se  $\theta \notin H_0$ )

★ Come si calcola  $h$ ? Prima devo costruire il test ② o ①

②  $\frac{S_1^2}{S_2^2} \stackrel{H_0}{\sim} F(m-1, n-1)$

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

$$h(\theta) = P(\text{dire } H_0; \theta) = P(a \leq \frac{S_1^2}{S_2^2} \leq b; \theta) = \begin{cases} \rightarrow 1-\alpha \text{ solo se } \theta = \theta_0 \Rightarrow \\ \rightarrow \text{e in generale!} \end{cases}$$

$$\frac{S_1^2}{\sigma_1^2} (m-1) \sim \chi^2(m-1) \quad \frac{S_2^2}{\sigma_2^2} (n-1) \sim \chi^2(n-1)$$

F di Fisher : rapporto tra due  $\chi^2$  indep, correte coi loro gdl

$$\frac{\frac{S_1^2}{\sigma_1^2} (m-1) / (m-1)}{\frac{S_2^2}{\sigma_2^2} (n-1) / (n-1)} = \frac{S_1^2 / \sigma_1^2}{S_2^2 / \sigma_2^2} = \boxed{\frac{S_1^2}{S_2^2} \cdot \frac{\sigma_2^2}{\sigma_1^2} \sim F(m-1; n-1)}$$

funz ancillare di  $\sigma_1/\sigma_2$

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

$$h(\theta) = P(a \leq \frac{S_1^2}{S_2^2} \leq b; \theta) = P\left(a \frac{\sigma_2^2}{\sigma_1^2} \leq \frac{S_1^2}{S_2^2} \frac{\sigma_2^2}{\sigma_1^2} \leq b \frac{\sigma_2^2}{\sigma_1^2}; \theta\right)$$

$$\theta = \log\left(\frac{\sigma_1}{\sigma_2}\right)$$

$$h(\theta) = \left( a e^{-2\theta} \leq F(m-1; n-1) \leq b e^{-2\theta} \right) = \text{DISTRIB. F}(b e^{-2\theta}; m-1; n-1) - \text{DISTRIB. F}(a e^{-2\theta}; m-1; n-1)$$

### • Caso unilaterale

campione di  $n$  oggetti,  $X$  il numero di difetti nel campione  
 $p$  la frequenza di difetti nella popolazione

N.B. di solito nel controllo  $Q$  popolazione = fornitura

$$p = \frac{\# \text{difetti fornitura}}{\# \text{fornitura}}$$

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

$$\text{test: } \begin{cases} X \leq a \\ X > a \end{cases}$$

accetto fornitura ( $H_0$ )

rifiuto fornitura ( $H_1$ )

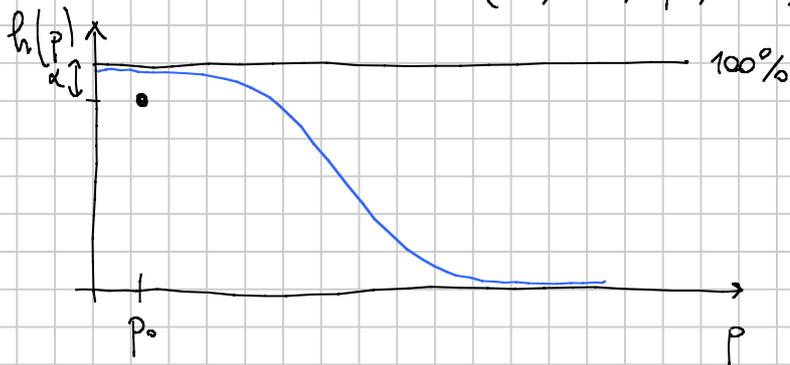
$$\rightarrow RA_X: \{0, 1, \dots, a\}$$

$$H_0: p \leq p_0$$

$$H_1: p > p_0$$

$$h(p) = P(\text{dire } H_0; p) = P(X \in RA_x; p) = P(X \leq a; p)$$

$$= \text{DISTRIB. BINOM} (a; n; p; 1)$$



of 32

Cos'è  $p_0$ ?

difetti nella  
fornitura

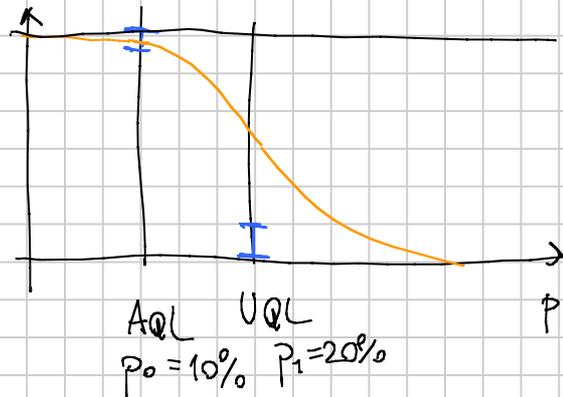
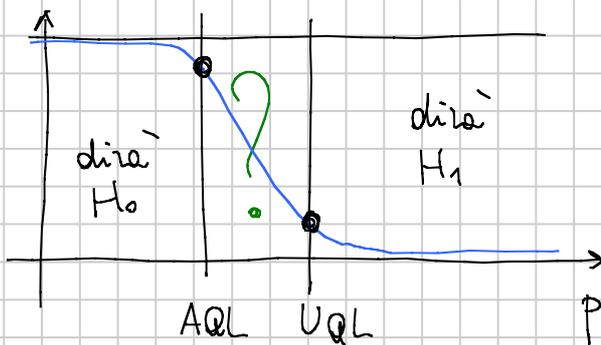
Voglio che se  $p \leq p_0$  allora accetterò la fornitura con prob elevata  $\geq 1 - \alpha$

$p_0$  : AQL *acceptable quality level*

11%  $\leq$  5%

$p_1$  : UQL *unacceptable quality level* (RQL)

30%  $\leq$  20%

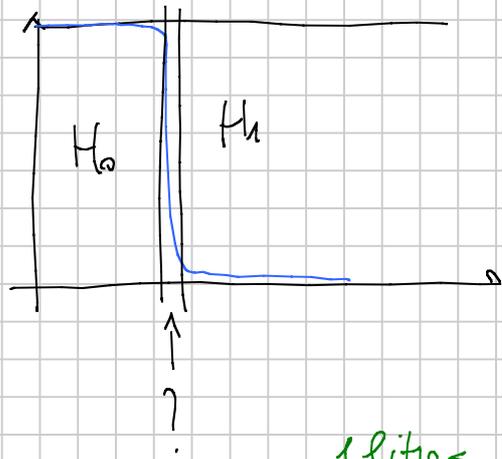
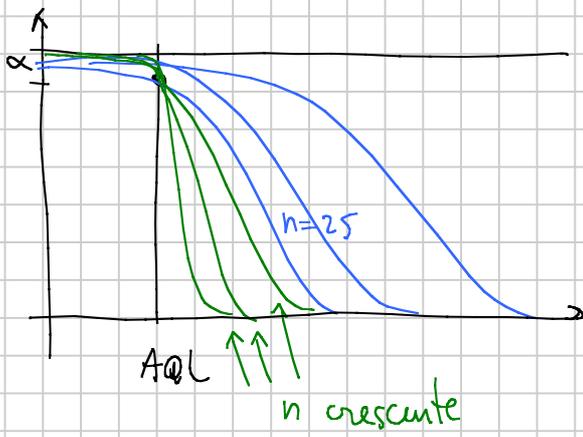


$$UQL = 20\%$$

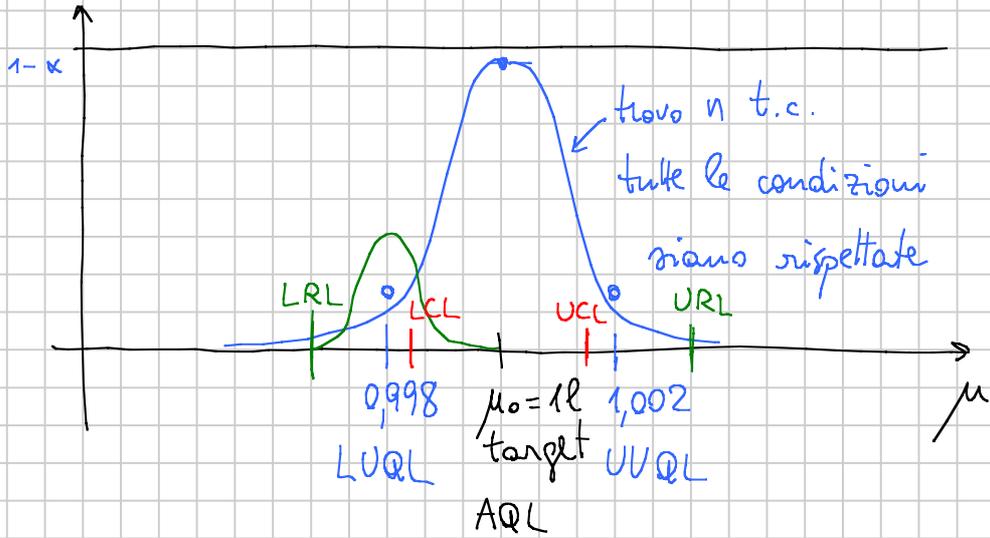
$$\beta = 10\%$$

$$AQL = 10\%$$

$$\alpha = 5\%$$



o UQL bilaterale



1 litro ← c'è  
 $X \sim \mathcal{N}(\mu, \sigma^2)$   
 riempimento cartoni di latte

5 note

$H_0: \mu = \mu_0$      $H_1: \mu \neq \mu_0$      $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$  iid  
 $\bar{X} \sim \mathcal{N}(\mu, \frac{\sigma^2}{n})$      $Z = \frac{\bar{X} - \mu_0}{\sigma / \sqrt{n}} \stackrel{H_0}{\sim} \mathcal{N}(0, 1)$      $\alpha = 0,0027$

②  $|Z| > q$  dico  $H_1$   
 $|Z| \leq q$  dico  $H_0$      $q = \text{INV.NORM.ST} \left( 1 - \frac{\alpha}{2} \right) = 3$

①  $\bar{X} \in \mu_0 \pm q \frac{\sigma}{\sqrt{n}}$     LCL, UCL

UQL determinati in funzione di quali livelli di riempimento sui singoli cartoni ritengo inaccettabili    LRL, URL

NB. Teorema di Cochran ora 6 anno scorso

FINE