

# STIMATORE BLUE (segue da pag. 209)

## Operazione

Se il parametro  $\beta$  è incognito si dimostra che

$$\hat{\beta} = \frac{E^T Y^{-1} E}{n - q}$$

è uno stimatore non polarizzato di  $\beta$ .

Ricordiamo che  $q$  = numero di parametri (di dimensione del vettore  $\beta$ ) e  $m$  = numero di dati e che valegono le ipotesi  $I_1$  e  $I'1$  ipotesi  $n > q$  (pag. 209).

## Nota

Questo affermato sullo stimatore BLUE vale anche con l'ipotesi di identificabilità, vale a dire le quali  $\text{rank}(E) = q$  (pag. 199)

L'osservazione fatta di questa che si può avere una conoscenza incompleta delle caratteristiche statistiche degli errori.

Si ricordi che  $\beta$  è la variante dell'errore di misura.

Nonostante ciò si assume di costruire il modello e con l'analisi dei residui si cercano di individuare le proprietà statistiche dell'errore di misura.

Se  $\Psi$  fosse la matrice identità lo stimatore si avrebbe la somma dei quadrati dei residui e si attenderebbe di dover dividere per  $n$  per avere l'ampio media del quadrato dei residui.

Al denominatore troviamo  $(n - q)$  data la polarizzazione, con come nella stessa ipotesi  $n - q$  compone  $(n - 1)$  (pag. 199).

Se si dividesse per  $n$  si avrebbe una stima polarizzata, ma si preferisce considerare uno stimatore non polarizzato.

## -) Operazione

213

Sempre nel caso di  $\theta^M$  incognita, essendo  $\Phi^T \Psi^{-1} \Phi$  si può scrivere anche la varianza del vettore dei parametri

$$\hat{\Sigma}_{\theta} = \hat{\sigma}^2 (\Phi^T \Psi^{-1} \Phi)^{-1}$$

La varianza var di  $\theta^M$  è:

$$\text{Var}[\theta^M] = \Sigma_{\theta} = \sigma^2 (\Phi^T \Psi^{-1} \Phi)^{-1}$$

Ovviamente, non conoscendo  $\sigma^2$  si usa la stima a disposizione.

Rimane aperto il problema degli vettori di confidenza per il vettore dei parametri.  $\theta^M$  è noto, si può risolvere, ma non si conosce la d.d.p. di  $\theta^M$ , cioè dello stimatore.

## -) Costa

$$\theta^M = cY = c(\Phi\theta^0 + V)$$

$$\text{costa } c = (\Phi^T \Psi^{-1} \Phi)^{-1} \Phi^T \Psi^{-1} \text{ (pag. 208)}$$

La matrice  $c$  è tutta deterministica, mentre  $Y$  è una v.c.

Quindi  $\theta^M$  risulta il prodotto di una matrice  $c$  deterministica, calcolabile, per un vettore di v.c.

Le proprietà statistiche di  $\theta^M$  dipendono dalle proprietà statistiche di  $Y$ .

A sua volta  $Y$  è la somma di un termine noto, essendo una deterministica  $(\Phi\theta^0)$ , ma caso è quello il vettore degli dei parametri, e di un termine casuale  $(V)$ .

Quindi, le proprietà statistiche, come la d.d.p. di  $\theta^M$  dipendono dalla d.d.p. di  $V$ , cioè dagli stessi di partenza.

Le proprietà degli errori di misura determinano le proprietà dello stimatore.

Per poter effettuare un'analisi per affidabilità, cioè ridurre un'ipotesi sulla d.d.p. di  $V$ , che indichiamo con  $I_V$ .

216 Ipotesi  $I_1$   

$$Y = \Phi \theta + V$$

con  $V$  distribuito gaussianamente con media nulla e matrice covarianza  $\Sigma_V$ .

$$V \sim N(0, \Sigma_V)$$

o.e.  

$$\Sigma_V = \sigma^2 \Psi$$

essendo  $\Psi$  nota e  $> 0$  e  $\sigma^2$  uno scalare eventuale da essere valutato.

Rispetto all'ipotesi  $I_2$  si sta considerando un set di misure di v.c. congiuntamente gaussiane, con mezzi e varianze distribuiti gaussianamente, anche se ciò non è sempre vero.

l'ipotesi  $I_1$  è identica alla  $I_2$  con l'aggiunta della considerazione della d.d.f.

Teorema sotto l'ipotesi  $I_1$  e rank  $(\Phi = q)$ ,  $\theta^M$  è distribuito gaussianamente:  

$$\theta^M \sim N(\theta^0, \sigma^2 (\Phi^T \Psi^{-1} \Phi)^{-1})$$

Dimostrazione

$$\theta^M = cY = c\Phi\theta^0 + cV$$

$\theta^M$  risulta una combinazione lineare di v.c. congiuntamente gaussiane e perciò è una v.c. gaussiane. (per  $\theta^0$ , essendo  $V$  un vettore di v.c. congiuntamente gaussiane e  $(c\Phi\theta^0)$  è c quantità (matrice) determinabile.

Intervallo di confidenza per  $\theta^0$

1) 1<sup>a</sup> nota

Definiamo la matrice covarianza di  $\theta^M$ .

$$\Sigma_{\theta^M} = \sigma^2 (\Phi^T \Psi^{-1} \Phi)^{-1}$$

e

$$\sigma_{\theta^M}^2 = [\Sigma_{\theta^M}]_{ii} \quad (\text{esimo elemento sulla diagonale principale di } \Sigma_{\theta^M})$$

Risultato

$I_{0,95}(\theta_i^0) = [\theta_i^M - 1,96 \hat{\sigma}_{\theta_i^M}, \theta_i^M + 1,96 \hat{\sigma}_{\theta_i^M}]$  215  
 essendo  $I_{0,95}(\theta_i^0)$  ~~essendo~~ l'intervallo di confiden-  
 za dell'elemento  $i$ -esimo di  $\theta_i^0$  con un valore di  $\alpha =$   
 fidabilità pari a 0,95.

Note: non confondere la  $t_{n-1}$  poiché non si sta costruendo  
 l'intervallo di confidenza della media campionaria.

negli intervalli di confidenza degli stimatori gaussiani  
 si compare sempre la deviazione standard dello stimatore.

2) 2<sup>a</sup> incognita

Definiamo

$$\hat{\Sigma}_{\theta^M} := \hat{\sigma}^2 (\Phi^T \Psi^{-1} \Phi) = \text{Var}[\theta^M]$$

$n^{\text{a}}$  "quadrato" varianza di  $\theta$  e una stima.

$$\hat{\sigma}_{\theta_i^M}^2 := [\hat{\Sigma}_{\theta^M}]_{ii}$$

Si dimostra che vale la seguente proprietà:

$$\left( \frac{\theta_i^M - \theta_i^0}{\hat{\sigma}_{\theta_i^M}} \right) \sim t_{n-1}$$

La quantità entro parentesi è distribuita come una  $t$  di Student con  $(n-1)$  gradi di libertà e la  $t$  presenta una pseudo-standardizzazione.

Da  $\theta_i^M$  si è tolto il valore vero  $\theta_i^0$  e si è diviso per la deviazione standard stimata, non conosciuta quella vera.

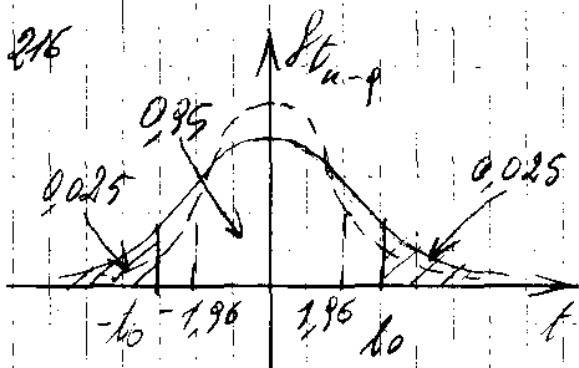
Risultato

$$I_{0,95}(\theta_i^0) = \left[ \theta_i^M - t_0 \hat{\sigma}_{\theta_i^M}, \theta_i^M + t_0 \hat{\sigma}_{\theta_i^M} \right]$$

con  $t_0$  di approssimazione

$$P(|t_{n-1}| \leq t_0) = 0,95$$

Se  $n \gg 1$ , la  $t_{n-1}$  diventa una gaussiana e  $t_0$  diventa 1,96.



dal grafico, la linea continua  
 rappresenta la  
 $t_{n-p}$ .  
 Le linee tratteggiate, indicano  
 la v.c. gaussiana standard.

### Applicazione: regressione lineare multiple

È un modello del tipo

$$f(t) = \theta_1 u_1(t) + \theta_2 u_2(t) + \dots + \theta_p u_p(t) + V(t)$$

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

Il vettore di

$$V \sim N(0, \sigma^2 I)$$

con

$$\Psi = I \quad (\text{per } \sigma^2 I)$$

e

$$V = [v(1) \ v(2) \ \dots \ v(n)]^T$$

Come esempio, possiamo immaginare di voler studiare  
 il valore del PIL in funzione di alcuni pa-  
 rametri quali investimenti in grandi opere, spesa per  
 la ricerca scientifica, ecc. Immagino peraltro il modello  
 applicato per un periodo storico.

Si può studiare come cambia il PIL al variare di  
 uno o più parametri.

Occorre scrivere le equazioni per tutti i casi

$$f(n) = \theta_1 u_1(n) + \theta_2 u_2(n) + \dots + \theta_p u_p(n) + V(n)$$

ovvero

In forma matriciale, si ha

$$Y = \begin{bmatrix} u_1(1) & u_2(1) & \dots & u_p(1) \\ u_1(2) & u_2(2) & \dots & u_p(2) \\ \vdots & \vdots & \ddots & \vdots \\ u_1(n) & u_2(n) & \dots & u_p(n) \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_p \end{bmatrix} + V$$

matrice  $\Phi$  ↗

Il modello è identificato attraverso i seguenti passi.

-) passo 1: calcolo di  $\hat{\theta}^M$  (vettore delle stime)

$$\hat{\theta}^M = (\Phi^T \Psi^{-1} \Phi)^{-1} \Phi^T \Psi^{-1} Y$$

-) passo 2: ricominciando che  $\hat{\theta}^2$  non è conosciuto, si calcola  $\hat{\sigma}^2$  (varianza dell'errore di misura)

$$\hat{\sigma}^2 = \frac{E^T \Psi^{-1} E}{n - q}$$

Se tutti gli errori hanno la stessa varianza e sono tra loro incorelati, allora

$$\Psi = I \quad (\text{diag. } 1 \text{ di } n)$$

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n E_i^2}{n - q}$$

$$\hat{\theta}^M = (\Phi^T \Phi)^{-1} \Phi^T Y = \hat{\theta}^LS \quad (\text{pag. 196})$$

-) passo 3: calcolo di  $\hat{\Sigma}_{\theta^M}$  (matrice covarianza dei parametri stimati)

$$\begin{aligned} \hat{\Sigma}_{\theta^M} &= \hat{\sigma}^2 (\Phi^T \Psi^{-1} \Phi)^{-1} = \\ &= \hat{\sigma}^2 (\Phi^T \Phi)^{-1} \quad \text{con } \Psi = I \end{aligned}$$

-) passo 4: intervalli di confidenza

Sono i risultati sono presentati nel paragrafo seguente:

$$f(t) = \hat{\theta}_1^M u_1(t) + \hat{\theta}_2^M u_2(t) + \dots + \hat{\theta}_q^M u_q(t) + V(t)$$

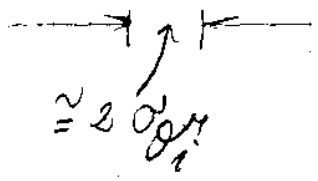
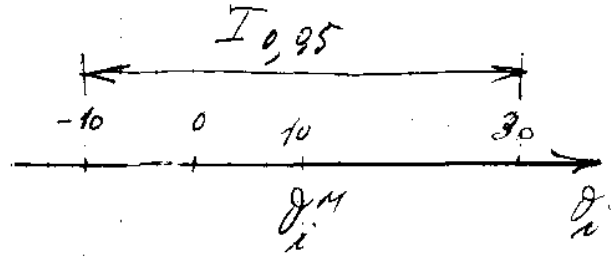
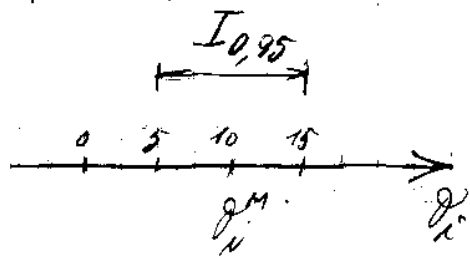
Al posto dei valori di  $\hat{\theta}$  si mettono i valori calcolati con la procedura di stima. Sotto ai  $\hat{\theta}_i^M$  e sotto al vettore  $V$  si collocano le deviazioni standard

$$f(t) = \hat{\theta}_1^M u_1(t) + \dots + V(t)$$

$$\begin{bmatrix} \hat{\sigma}_{\theta_1^M} \\ \hat{\sigma}_{\theta_2^M} \\ \vdots \end{bmatrix} \quad \begin{bmatrix} \hat{\sigma}_V \end{bmatrix}$$

218 Oltre al modello ed ai parametri stimati occorre fornire elementi che permettano di valutare il grado di affidabilità del risultato.

Conoscizione



In entrambi i grafici  $\theta_i^M$  è la stima del parametro  $\theta_i$ .

Supponiamo che  $\theta_i^M$  sia conosciuto oppure che il numero dei dati  $\theta_i$  sia molto grande.

Sotto questa ipotesi comprende a destra ed a sinistra di  $\theta_i^M$  due intervalli per  $\alpha \approx 2 \sigma_{\theta_i^M}$ .

Mentre l'intervallo di confidenza del caso a sinistra ~~non~~ va bene così, non accade per l'intervallo a destra, nel quale l'intervallo di confidenza comprende lo zero.

Di conseguenza non si sa come intervenire sul parametro quando il segno è positivo o negativo.

Regola pratica: l'intervallo di confidenza va bene se  $|\theta_i^M| > 2 \sigma_{\theta_i^M}$ .

Il valore di ogni parametro deve essere almeno il doppio della propria deviazione standard.

## VALIDAZIONE E SCELTA DEL MODELLO

Contenuti:

- Validazione: Test  $\chi^2$
- Test  $F$
- Crossvalidazione
- Final Prediction Error (FPE)
- Akaike Information Criterion (AIC)
- Minimum Description Length (MDL)
- Un esempio semplice
- Conclusioni

**Motivazione:** Una volta nota la struttura del modello (cioè  $\Phi(U, \theta)$ ), è relativamente facile stimare  $\theta$ . Problemi aperti:

- come capire se un dato modello è "buono";
- come confrontare due o più modelli.

Statistica di base II

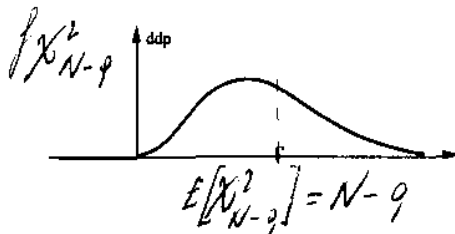
10

**Ipotesi I1:**  $Y = \Phi\theta + V$ ,  $V \sim N(0, \sigma^2 I)$

(Suppongo di aver azzeccato la struttura del "modello vero" che genera i dati e che gli errori siano gaussiani con varianza nota)

**Teorema:** Sotto l'ipotesi I1, dividendo per  $\sigma^2$  la somma dei quadrati dei residui (SSR: Sum of Squared Residuals) della stima LS si ottiene una V.C. di tipo  $\chi^2$  "con  $N-q$  gradi di libertà" ( $N = n^\circ$  di dati,  $q = n^\circ$  di parametri):

$$\frac{E^T E}{\sigma^2} \sim \chi^2(N-q)$$



**Nota:** Risulta  $E[\chi^2(N-q)] = N-q$ ,  $Var[\chi^2(N-q)] = 2(N-q)$ . Inoltre, per  $N-q$  "grande" la f.d.d. della V.C.  $\chi^2(N-q)$  è circa gaussiana.

**Idea:** Utilizzando la f.d.d. del  $\chi^2$  ho un criterio numerico per valutare quando la SSR è "anormale" ( $\Rightarrow$  modello sbagliato).

Statistica di base II

12

## VALIDAZIONE: TEST $\chi^2$

**Problema:** Avendo a disposizione  $N$  dati  $Y_1, \dots, Y_N$ , come faccio a sapere se il modello

$$Y = \Phi\theta + V, \quad Var[V] = \sigma^2 I$$

li descrive adeguatamente?

*Nota 1, pag. 218*

**Idea:**  $\theta LS \equiv \hat{\theta} \Rightarrow e = Y - \Phi\hat{\theta} LS \equiv V$  (il vettore dei residui è una "stima" del vettore degli errori di misura). Perciò mi aspetto che

$$\frac{1}{N} \sum_{i=1}^N \epsilon_i^2 \equiv \sigma^2$$

*2/218*

Se conosco  $\sigma^2$  è bene controllare che  $\sigma^2$  e la varianza campionaria dei residui abbiano lo stesso ordine di grandezza (idea ovvia: se ho errori di misura dell'ordine di  $10^{-3}$  e i residui sono dell'ordine di  $10^{-1}$ , il modello è sbagliato perché non riesce a spiegare i dati).

**Problema:** cosa significa "dello stesso ordine di grandezza"?

Statistica di base II

11

Tabella: Valori di  $\chi^2_{\alpha, n}$

n	$\alpha=0.975$	$\alpha=0.95$	$\alpha=0.05$	$\alpha=0.025$
1	0.00	0.00	5.02	5.02
2	0.05	0.10	4.99	7.38
3	0.22	0.35	7.81	9.35
4	0.48	0.71	9.49	11.14
5	0.83	1.15	11.07	12.38
6	1.24	1.64	12.59	14.45
7	1.69	2.17	14.07	16.01
8	2.18	2.73	15.51	17.53
9	2.70	3.33	16.92	19.02
10	3.25	3.94	18.31	20.48
11	3.82	4.57	19.68	21.92
12	4.40	5.23	21.03	23.34
13	5.01	5.89	22.36	24.74
14	5.63	6.57	23.68	26.12
15	6.27	7.26	25.00	27.49
16	6.91	7.96	26.30	28.85
17	7.56	8.67	27.59	30.19
18	8.23	9.39	28.87	31.53
19	8.91	10.12	30.14	32.85
20	9.59	10.85	31.41	34.17
25	13.12	14.61	37.65	40.65
30	16.79	18.49	43.77	46.98
40	24.43	26.51	55.76	59.34
50	32.36	34.76	67.50	71.42
60	40.48	43.19	79.08	83.30
70	48.76	51.74	90.53	95.02
80	57.15	60.39	101.88	106.63
90	65.65	69.13	113.14	118.14
100	74.22	77.93	124.34	129.56

*3/218*

Statistica di base II

13



Esempio: 10 dati, 3 parametri ( $N-q = 10-3 = 7$ ). Dalle tabelle della distribuzione  $\chi^2$  risulta  $P(\chi^2(7) < 14.07) = 0.95 \Rightarrow$  nel 95% dei casi  $\epsilon T \epsilon / \sigma^2 < 14.07$ .

Se accade che  $\epsilon T \epsilon / \sigma^2 > 14.07$ , sospetto che il modello non sia buono (sotto l'ipotesi che sia buono, accade raramente che la somma dei quadrati dei residui sia così grande).

Test  $\chi^2$ : Fissato il livello di significatività  $\alpha$  (tipicamente,  $\alpha = 0.05$ ) cerco sulle tabelle  $x_\alpha$  tale che  $P(\chi^2(N-q) < x_\alpha) = 0.95$ . Poi adotto la seguente regola:

- $\frac{\epsilon T \epsilon}{\sigma^2} < x_\alpha \Rightarrow$  accetto il modello
- $\frac{\epsilon T \epsilon}{\sigma^2} > x_\alpha \Rightarrow$  respingo il modello

4/219

Estensione:  $V \sim N(0, \Sigma_V)$ , dove  $\Sigma_V$  è una matrice nota. Basta usare  $\epsilon T \Sigma_V^{-1} \epsilon$  al posto di  $\epsilon T \epsilon / \sigma^2$ .

Punti deboli:

- Può essere difficile capire quali sono i motivi per cui viene scartato il modello. Almeno 4 possibilità:
  - a) la relazione  $Y = \Phi \theta^0 + V$  spiega male i dati;
  - b)  $V$  non è gaussiano;
  - c) il valore di  $\sigma^2$  è sbagliato per difetto;
  - d) gli errori di misura non hanno tutti la stessa varianza.
- Il test si basa sull'ipotesi che esista un "modello vero" di tipo lineare che genera i dati e che gli errori siano gaussiani (ipotesi molto semplificativa).

5/219

## TEST F

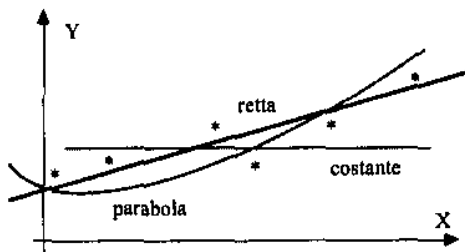
Confronto tra modelli: Spesso, non è chiaro quale sia il modello "giusto" e i considera un ventaglio di possibilità. Come si fa a scegliere il modello "ottimo"?

Modelli "matrioska": Una sequenza di classi di modelli in cui ciascuna classe comprende al suo interno (come casi particolari) le classi precedenti.

Esempio:  $M_1 = \{\text{rette}\}$ ,  $M_2 = \{\text{parabole}\}$ ,  $M_3 = \{\text{cubiche}\}$ , ...

6/219

Problema tipico: Conosco  $N$  coppie  $(X_i, Y_i)$  e voglio identificare un modello  $Y = f(X)$ . Se mi limito ai polinomi di ordine  $k$ , cosa è meglio? Una retta, una parabola, una cubica, ...?



Idea (stupida): Considero i diversi modelli (retta, parabola, cubica, ...) e ne stimo i parametri mediante LS. Poi, tra i modelli stimati, scelgo quello che minimizza la SSR.

7/219

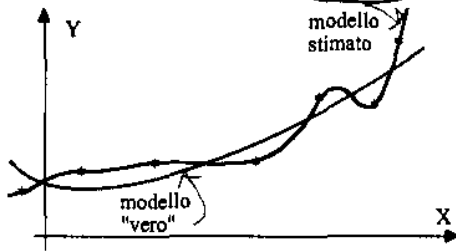
Fatto: Per i modelli matrioska, la SSR decresce sempre al crescere dell'ordine del modello (dato che la stima LS si basa sulla minimizzazione di SSR, non è possibile che la migliore parabola abbia una SSR maggiore di una retta)

8/219

Usare la minimizzazione di SSR per scegliere il modello ottimo conduce a scegliere sempre il modello più complesso (per assurdo,  $N = 100 \Rightarrow$  polinomio di ordine 99).

Che male c'è ad eccedere con il n° di parametri?

- Nessun problema se i dati fossero privi di rumore.  
Esempio: stimo con una parabola dei dati che stanno su una retta  $\Rightarrow$  il coefficiente del termine quadratico risulta = 0  $\Rightarrow$  non commetto errori.
- Se c'è rumore ed ho troppi parametri, il modello stimato tende ad essere influenzato dal rumore (riproduce oscillazioni che non hanno significato fisico ma che sono frutto degli errori di misura)



Principio di parsimonia: Non usare parametri addizionali per descrivere un fenomeno se essi non sono necessari.

Idea: Dati  $M_{k-1}$  e  $M_k$ , sappiamo che  $SSR_k < SSR_{k-1}$ . Sceglierò  $M_k$  solo se  $SSR_k$  è "molto più piccola" di  $SSR_{k-1}$ .

10/20

Problema: Cosa vuol dire "molto più piccola"?

Teorema: Sotto l'ipotesi II,

$$f = (N-k) \frac{SSR_{k-1} - SSR_k}{SSR_k}$$

11/20

è una V.C. distribuita come una F di Fisher con  $(1, N-k)$  gradi di libertà

Osservazioni:

- $f$  è un indice della riduzione % di  $SSR$  che ottengo passando dal modello  $M_{k-1}$  (meno complesso) al modello  $M_k$  (più complesso).
- Per  $N-k$  "grande", si ha  $F(1, N-k) = \chi^2(1)$

12/20

Test F: Fissato un livello di significatività  $\alpha$  (tipicamente,  $\alpha = 0.05$ ) cerco sulle tabelle  $f_\alpha$  tale che  $P(F(1, N-k) < f_\alpha) = 0.95$ . Poi adotto la seguente regola:

- $f < f_\alpha \Rightarrow$  scelgo il modello  $M_{k-1}$
- $f > f_\alpha \Rightarrow$  scelgo il modello  $M_k$

13/20

Osservazioni:

- Non è necessario conoscere  $\sigma^2$
- $V \sim N(0, \sigma^2 \Psi)$ , dove  $\Psi$  è una matrice nota  $\Rightarrow$  basta usare  $\epsilon^T \Psi^{-1} \epsilon$  al posto di  $SSR = \epsilon^T \epsilon$ .
- L'esito del test dipende dalla scelta del livello di significatività  $\alpha$  ( $\alpha$  "piccolo": aumenta la probabilità di sottostimare l'ordine;  $\alpha$  "grande": aumenta la probabilità di sovrastimare l'ordine)

14/20

Punti deboli:

- L'ipotesi II è restrittiva. Esiste un "modello vero"? Ammesso che esista, appartiene alla classe di modelli considerati?
- Il test si applica solo a classi di modelli "matrioska"

Approccio alternativo

15/20

Fatto: Supponiamo che valga II e che il modello vero è  $M_k$ . Allora  $\theta_j^0 = 0 \Rightarrow \theta_j^{LS} \sim G(0, \sigma^2) \Rightarrow \theta_j^{LS} / \sigma_0 \sim G(0, 1) \Rightarrow P(|\theta_j^{LS} / \sigma_0| \leq 1.96) = 0.95$ .



Se  $\theta_j^0 = 0$ , nel 95% dei casi risulta  $|\theta_j^{LS}| < 1.96 \sigma_0$

16/20

Morale: Se un parametro stimato è maggiore del doppio del valore della sua SD, è verosimile che il parametro sia  $\neq 0$ . In caso contrario, può essere conveniente azzerare quel parametro.

E' bene diffidare dei modelli in cui i parametri stimati non sono almeno 2-3 volte più grandi delle loro SD.

## CROSSVALIDAZIONE

*Idea:* Se ho abbastanza dati, li divido in due gruppi:

- 1) dati di identificazione  $Y$ ;
- 2) dati di validazione  $Y^v$ .

Uso il primo gruppo per identificare vari modelli (mediante LS, per es.). Poi uso il secondo gruppo per testare i modelli e capire quale è il migliore.

*Procedura:*

- Considero vari modelli (rette, parabole, cubiche, ...) e ne identifico i parametri mediante LS

$$\theta^{LS} = (\Phi^T \Phi)^{-1} \Phi^T Y$$

- Con i modelli identificati cerco di prevedere i dati di validazione e calcolo i relativi residui

$$SSR^v = \varepsilon^v T \varepsilon^v, \quad \varepsilon^v = Y^v - \hat{Y}, \quad \hat{Y} = \Phi^v \theta^{LS}$$

- Scelgo il modello che minimizza  $SSR^v$

*E se non ho abbastanza dati per formare due gruppi?*

*Ordinary Cross Validation (OCV):* Metto da parte il dato  $i$ -esimo e calcolo  $\theta^{LS(i)}$  usando tutti i dati rimanenti. Poi calcolo l'errore che commetto cercando di indovinare  $Y_i$  usando  $\theta^{LS(i)}$

$$\varepsilon^{(i)} = Y_i - \Phi^{(i)} \theta^{LS(i)} \quad (\Phi^{(i)}: i\text{-esima riga di } \Phi)$$

Ripeto la procedura per  $i = 1, \dots, N$  e uso come indice di bontà del modello:

$$OCV = \frac{1}{N} \sum_{i=1}^N \varepsilon^{(i)2}$$

*Problema:* Sembra necessario risolvere  $N$  problemi di stima LS (computazionalmente oneroso).

*Lemma del "lasciane-uno-fuori" ("leave-out-one" lemma):*

$$OCV = \frac{1}{N} \sum_{i=1}^N \frac{\varepsilon_i^2}{(1 - H_{ii})^2}$$

$$H = \Phi (\Phi^T \Phi)^{-1} \Phi^T$$

$$\varepsilon = Y - \Phi \theta^{LS}$$

*Osservazioni:*

- L'assegnamento di un'osservazione al set di identificazione o a quello di validazione deve essere casuale.
- Se uso modelli matriciali,  $SSR^v$  decresce al crescere dell'ordine. All'inizio anche  $SSR^v$  decresce. Ad un certo punto i parametri diventano troppi ed il modello identificato cerca di seguire troppo fedelmente i dati di identificazione  $\Rightarrow SSR^v$  comincia a salire.
- Talvolta la crossvalidazione suggerisce l'uso di modelli in cui alcuni parametri hanno  $SD$  elevata. Vale la pena di azzerare il valore di questi parametri e ricalcolare  $SSR^v$ . Se  $SSR^v$  aumenta di poco ( $1+2\%$ ) rispetto al minimo può essere conveniente adottare il modello semplificato.
- Non faccio nessuna ipotesi sul meccanismo "vero" di generazione dei dati. Non pretendo di trovare il modello "vero", ma solo il modello "migliore" (in termini di  $SSR^v$ ) entro una certa rosa di possibilità.
- Limitazione fondamentale: bisogna avere parecchi dati altrimenti sia l'identificazione che la validazione diventano poco affidabili.

*Osservazioni:*

- Grazie al "leave-out-one lemma" basta risolvere una sola stima LS e calcolare gli elementi sulla diag. principale di  $H$ .
- L'uso di  $OCV$  è in genere più oneroso che crossvalidare dividendo i dati in due gruppi (nel calcolo di  $\theta^{LS}$  in genere si evita di calcolare esplicitamente  $(\Phi^T \Phi)^{-1}$ , che è invece richiesta da  $OCV$ ).
- Punto debole: quando gli errori di misura non hanno tutti la stessa varianza.
- Per ridurre i calcoli, si può ricorrere ad una approssimazione di  $OCV$  (Generalized Cross Validation):

$$GCV = \frac{1}{N} \frac{\sum_{i=1}^N \varepsilon_i^2}{\left[ \frac{1}{N} \text{Tr}(I-H) \right]^2} = \frac{1}{N} \frac{\sum_{i=1}^N \varepsilon_i^2}{\left[ \frac{N-q}{N} \right]^2} = \frac{N}{(N-q)^2} SSR$$

## FINAL PREDICTION ERROR (FPE)

*Idea:* Se faccio delle ipotesi sul meccanismo di generazione dei dati posso cercare di minimizzare  $SSR^v$  senza doverla calcolare esplicitamente.

*Ipotesi I2:*  $Y = \Phi\theta^0 + V$ ,  $E[V] = 0$ ,  $Var[V] = \sigma^2 I$ .  
(rispetto a I1, non faccio ipotesi sulla gaussianità del rumore)

Supponiamo di considerare un vettore  $\theta$  e di sottoporlo a validazione. Ipotizzando  $\Phi^v = \Phi^t$ , si dimostra che, se estraggo a caso un campione  $Y^v$  di  $N$  dati di validazione (estrazione #1):

$$E[SSR^v] = N\sigma^2 + (\theta - \theta^0)^T \Phi^T \Phi (\theta - \theta^0)$$

*Osservazione (ovvia):*  $E[SSR^v]$  è minimizzata da  $\theta = \theta^0$ .

In molti casi  $\sigma^2$  non è nota. Uno stimatore non polarizzato di  $\sigma^2$  è fornito da:

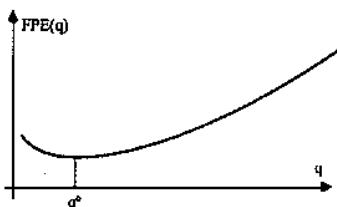
$$\hat{\sigma}^2 = \frac{SSR}{N-q}$$

*Criterio FPE:* Tra diversi modelli matriciali caratterizzati dal numero  $q$  di parametri, scelgo quello che minimizza la media (stimata) della  $SSR^v$ , ovvero il cosiddetto *Final Prediction Error*

$$FPE = \frac{N+q}{N-q} SSR$$

*Osservazioni:*

- Al crescere di  $q$  dapprima  $FPE$  diminuisce perché diminuisce  $SSR$ . Se però  $q$  cresce troppo,  $FPE$  cresce (c'è  $N-q$  al denominatore)  $\Rightarrow \exists$  un minimo di  $FPE$  al variare di  $q$ .
- Computazionalmente efficiente: per calcolare  $FPE$  basta solo conoscere  $\theta^{LS}$ .



Supponiamo ora che  $\theta = \theta^{LS}$  sia la stima ottenuta da un campione estratto a caso di  $N$  dati  $Y^t$  di identificazione (estrazione #2). Si dimostra che

$$E[E\{SSR^v\}] = \sigma^2(N+q), \quad q = \dim(\theta)$$

(ho due medie perché ho due estrazioni casuali)

*E' la formulazione matematica del principio di parsimonia: se il modello vero è una retta ( $q = 2$ ) e per identificare uso una parabola o una cubica ( $q = 3, q = 4$ ) peggioro inutilmente le prestazioni (medie) in validazione del mio modello.*

### Confronto criteri soggettivi/oggettivi

*Criteri soggettivi:* quelli basati sui test statistici (vedi test  $F$ ) che richiedono la scelta "soggettiva" del livello di significatività  $\alpha$ .

*Criteri oggettivi:* quelli basati sulla minimizzazione di una cifra di merito ( $OCV, GCV, FPE, AIC, MDL$ ). Non c'è da scegliere nessun livello di significatività.

*La differenza è meno profonda di quanto sembri.*

*Fatto (di facile ma noiosa dimostrazione):* Per  $N$  "grande", scegliere tra  $M_k$  e  $M_{k+1}$  basandosi su  $FPE$  equivale ad usare il Test  $F$  con  $\alpha = 0.157$

Per  $N \rightarrow \infty$ ,  $FPE$  ha il 15.7% di probabilità di scegliere (erroneamente) il modello più complicato anche se quello giusto è quello più semplice.

Per  $N \rightarrow \infty$ ,  $FPE$  sovrastima (in media) l'ordine del modello (non è uno stimatore consistente dell'ordine del modello)

## AKAIKE INFORMATION CRITERION (AIC)

Criterio ricavato in base alla "distanza" tra la d.d.p. vera dei dati e quella generata da un dato modello stimato (vale sotto l'ipotesi 11).

**Criterio AIC:** Tra diversi modelli matroska caratterizzati dal numero  $q$  di parametri, scelgo il modello che minimizza

$$AIC = \frac{2q}{N} + \ln(SSR)$$

**Osservazione:** Per  $N \rightarrow \infty$  si dimostra che  $FPE$  e  $AIC$  sono equivalenti (infatti, si vede che  $\lim_{N \rightarrow \infty} \ln FPE = AIC$ ).



Per  $N \rightarrow \infty$ , anche  $AIC$  sovrastima  
(in media) l'ordine del modello

## MINIMUM DESCRIPTION LENGTH (MDL)

**Idea:** Immagino di dover trasmettere i dati. Invece di trasmettere il vettore  $Y$ , posso trasmettere  $\theta^{LS}$  e il vettore dei residui  $e$ . Il ricevitore potrà ricostruire i dati calcolando  $Y = \Phi \theta^{LS} + e$ .

**Vantaggio:** se il modello è buono, l'errore di predizione è piccolo e, per una data precisione, bastano pochi bit per codificarlo.

Se l'ordine  $q$  del modello aumenta i residui diventano più piccoli (occorrono meno bit per  $e$ ) ma aumenta il n° dei parametri da codificare (occorrono più bit per  $\theta^{LS}$ ).

**Criterio Minimum Description Length:** scelgo il modello che conduce alla codifica più compatta. Si dimostra che (sotto 11) ciò equivale a minimizzare la cifra di merito

$$MDL = \frac{\ln(N) q}{N} + \ln(SSR)$$

**Osservazione:** La penalità su  $q$  è più pesante che in  $AIC$ . Infatti  $MDL$  conduce a modelli più parsimoniosi. Anzi si dimostra che  $MDL$  è uno stimatore consistente dell'ordine del modello (per  $N \rightarrow \infty$  l'ordine indicato da  $MDL$  converge all'ordine vero).

## UN ESEMPIO SEMPLICE

Esempio tratto da (J.V. Beck e K.J. Arnold, "Parameter Estimation in Engineering and Science, Wiley 1977).

La conducibilità termica  $k$  di alcuni campioni di ferro è stata misurata a diverse temperature  $T$  (°F). I risultati sperimentali sono riportati nella seguente tabella.

$T$ (°F)	100	141	227	270	382	90	148	208	247	352
$k$	41.8	37.7878	38.4978	38.7884	63	42.346	39.6378	37.3628	66.38	24.818

Le prime cinque misure sono state prese in condizioni sperimentali diverse rispetto alle ultime cinque. In particolare, si sa che la varianza dell'errore di misura per i secondi cinque dati è quattro volte maggiore della varianza dell'errore per i primi cinque. Ci si pone l'obiettivo di identificare un modello che descriva la dipendenza di  $k$  nei confronti della temperatura. Si considerano i seguenti modelli:

1.  $k = \theta_1$
2.  $k = \theta_1 + \theta_2 T$
3.  $k = \theta_1 + \theta_2 T + \theta_3 T^2$
4.  $k = \theta_1 + \theta_2 T + \theta_3 T^2 + \theta_4 T^3$
5.  $k = \theta_1 + \theta_2 T + \theta_3 T^2 + \theta_4 T^3 + \theta_5 T^4$

**Problema 11:** Stima dei parametri

**Soluzione:** Minimi quadrati ponderati (WLS) (alcuni dati sono più affidabili di altri)

$$\theta = (\Phi^T Q \Phi)^{-1} \Phi^T Q Y$$

Vettore dei dati e delle variabili indipendenti:

$$Y = \begin{bmatrix} k(1) \\ k(2) \\ \dots \\ k(10) \end{bmatrix}, \quad U = \begin{bmatrix} T(1) \\ T(2) \\ \dots \\ T(10) \end{bmatrix}$$

Matrice  $\Phi(U)$  e vettore  $\theta$  nei 5 modelli:

$$1. \quad \Phi = \begin{bmatrix} 1 \\ 1 \\ \dots \\ 1 \end{bmatrix}, \quad \theta = \begin{bmatrix} \theta_1 \end{bmatrix}$$

$$2. \quad \Phi = \begin{bmatrix} 1 & T(1) \\ 1 & T(2) \\ \dots & \dots \\ 1 & T(10) \end{bmatrix}, \quad \theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

$$3. \quad \Phi = \begin{bmatrix} 1 & T(1) & T(1)^2 \\ 1 & T(2) & T(2)^2 \\ \dots & \dots & \dots \\ 1 & T(10) & T(10)^2 \end{bmatrix}, \quad \theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

$$4. \quad \Phi = \begin{bmatrix} 1 & T(1) & T(1)^2 & T(1)^3 \\ 1 & T(2) & T(2)^2 & T(2)^3 \\ \dots & \dots & \dots & \dots \\ 1 & T(10) & T(10)^2 & T(10)^3 \end{bmatrix}, \quad \theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \end{bmatrix}$$

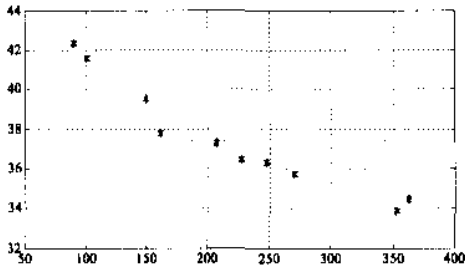
$$5. \quad \Phi = \begin{bmatrix} 1 & T(1) & T(1)^2 & T(1)^3 & T(1)^4 \\ 1 & T(2) & T(2)^2 & T(2)^3 & T(2)^4 \\ \dots & \dots & \dots & \dots & \dots \\ 1 & T(10) & T(10)^2 & T(10)^3 & T(10)^4 \end{bmatrix}, \quad \theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \\ \theta_5 \end{bmatrix}$$

Matrice  $Q$  di WLS

$$Q = \text{diag}(1 \ 1 \ 1 \ 1 \ 1 \ 1/4 \ 1/4 \ 1/4 \ 1/4 \ 1/4)$$

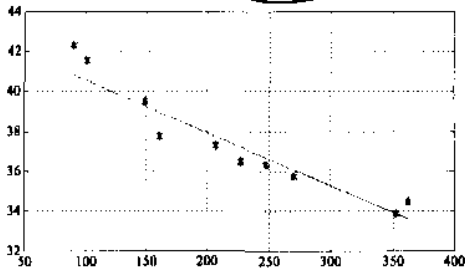
theta1 =  
3.7372e+01  
sigma\_theta1 =  
8.4316e-01

36/226



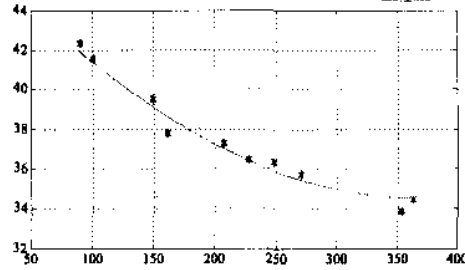
theta2 =  
4.3225e+01 -2.6486e-02  
sigma\_theta2 =  
7.9222e-01 3.3203e-03

37/226



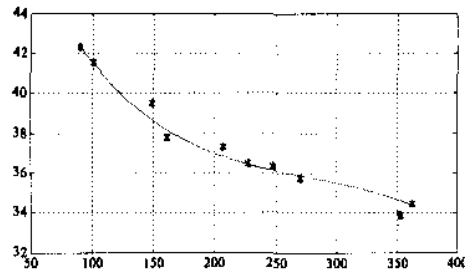
theta3 =  
4.7519e+01 -7.0826e-02 9.6665e-05  
sigma\_theta3 =  
1.0335e+00 9.8890e-03 2.1206e-05

38/226



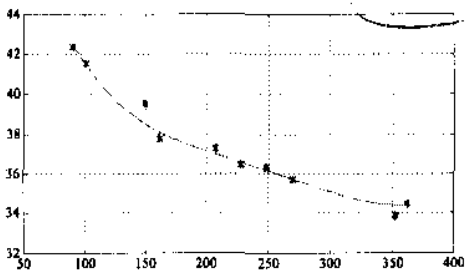
theta4 =  
5.3391e+01 -1.6786e-01 5.6508e-04 -6.8092e-07  
sigma\_theta4 =  
2.5573e+00 4.0883e-02 1.9459e-04 2.8187e-07

39/226



theta5 =  
6.2542e+01 -3.7495e-01 2.1788e-03 -5.8682e-06 5.8511e-09  
sigma\_theta5 =  
1.0209e+01 2.2728e-01 1.7526e-03 5.6051e-06 6.3143e-08

40/226



Problema #2: Scelta del modello ottimo

Osservazioni:

- Come previsto, all'aumentare dell'ordine del modello si ottiene una maggior aderenza ai dati sperimentali. Caso limite: con dieci parametri (polinomio di ordine 9) posso interpolare i dati.
- Nei modelli 1-3 lo deviazioni standard delle stime dei parametri sono accettabilmente inferiori ai valori dei parametri stimati.
- Nel modello 4 la SD di  $\theta_4^{ML}$  è di poco superiore al 40% del valore stimato del parametro. Nel modello 5 la SD di  $\theta_4^{ML}$  è circa uguale a  $\theta_4^{ML}$  o lo stesso accade per la SD di  $\theta_5^{ML}$ . Ciò potrebbe significare che il coefficiente del termine cubico (e a maggior ragione quello del termine di quarto grado) non è significativamente diverso da zero.

Test F:

Per applicare il test F è utile costruire la seguente tabella, in cui  $f = (N-q)\Delta SSR/\Delta SSR_{mod}$  e  $F_{0.95}(1, N-q)$  indica il valore (ricavato dalla tabella della F di Fisher) in corrispondenza del quale la funzione di distribuzione della  $F \sim (1, N-q)$  gradi di libertà vale 0.95:

Mod.	q	N-q	SSR	$\Delta SSR$	f	$F_{0.95}(1, N-q)$
1	1	9	40.0078			
2	2	8	4.4680	35.5398	63.6341	5.32
3	3	7	1.1259	3.3421	20.7786	5.99
4	4	6	0.5708	0.5551	5.8356	5.99
5	5	5	0.4871	0.0837	0.8587	6.61

Confrontando la penultima e l'ultima colonna si vede che il modello 2 è significativamente migliore (in termini di riduzione dello scarto quadratico) del modello 1. Lo stesso vale per il modello 3 nei confronti del modello 2. Per il modello 4 si ha  $f < F_{0.95}(1, N-q)$  e pertanto la riduzione di scarto quadratico non è statisticamente significativa. Vista la vicinanza dei valori di  $f$  e  $F_{0.95}(1, N-q)$  è tuttavia opportuno non scartare a priori il modello 4. Infine, passando dal modello 4 al modello 5 la riduzione dello scarto quadratico è palesemente non significativa, cosicché il modello 5 è da scartare.

Criteri FPE, AIC, MDL:

Mod.	q	FPE	AIC	MDL
1	1	48.8984	3.8891	3.9193
2	2	6.7020	1.8969	1.9575
3	3	2.0910	0.7186	0.8094
4	4	1.3318	0.2392	0.3603
5	5	1.4613	0.2807	0.4320

È interessante notare che, in contrasto con il test F, tutti e tre i criteri "oggettivi" suggeriscono la scelta del modello 4. In conclusione, anche a causa dei pochi dati disponibili, è difficile scegliere con sicurezza tra il modello 3 e 4. Tuttavia, tenendo in considerazione le deviazioni standard dei parametri stimati, potrebbe essere più prudente orientarsi verso il modello 3.

## CONCLUSIONI

- L'unico metodo che non richiede ipotesi pesanti è la crossvalidazione.
- Nella pratica, i criteri *FPE*, *AIC*, *MDL* vengono usati senza preoccuparsi troppo delle ipotesi.
- *MDL* è meglio di *FPE* e *AIC*, ma solo asintoticamente.
- Suggerimento: usare più di un criterio. Anche l'esame delle *SD* dei parametri è importante e può aiutare a dirimere eventuali discordanze tra i criteri.
- Nella pratica non è essenziale trovare il miglior modello ma spesso basta un buon modello

# VALIDAZIONE E SCELTA DEL MODELLO

219

(note e commenti negli appunti disponibili in rete)

(1) Si fanno le seguenti ipotesi:

-  $I_1$  (pag. 214)

-  $Y = I \Rightarrow O^M = O^{15}$  (la stima di Cherkov coincide con la stima LS)

(2) Inizialmente si aveva stimato la stima ( $O^{10} \approx 80$ ).

In questa ipotesi il vettore dei residui ( $Y - \hat{Y}$ ) coincide praticamente con il vettore degli errori di misura.

Il vettore dei residui è considerato una stima del vettore degli errori di misura.

Conoscendo la varianza  $\sigma^2$  degli errori di misura, questa è utile alla ricerca dei residui, che si stima dividendo per  $n$  la somma dei quadrati.

(3) Le righe indicano i gradi di libertà.

Per  $p = 0,95$  si cerca nella tabella  $\alpha = 0,05$  ( $1 - p$ )

(4) Esiste la possibilità di respingere un modello lineare, e in questo caso, era solo il 5%.

Questo criterio di validazione molto semplice, richiede di conoscere  $\sigma^2$ , cioè le caratteristiche degli strumenti di misura.

(5)  $Y$  non è la stessa identica  $I$ .

(6) Dato un polinomio di grado 2, si rappresenta una parabola annullando il coefficiente di 2° grado la parabola degenera in una retta.

(7) Legando un modello la cui funzione pare vicino ai dati sperimentali, si cercano i valori dei parametri che minimizzano la somma dei quadrati dei residui.

Si può pensare di usare lo stesso criterio per valutare qual è il modello migliore, pensando che questo non quello che pare vicino ai dati sperimentali.

(8) Se si dispone avere come criterio quello dei minimi quadrati, si finirebbe sempre per scegliere il modello più complicato, si dovrebbe un modello che pare esattamente sopra i dati.



220 (9) si può dire che il modello ripropone il numero  
 (10)  $k$  e la complementa del modello.

Si può adottare il criterio di scegliere questo o quello  
 le ipotesi dei quadrati di residui e scegliere il  
 modello quando la differenza tra i due modelli è  
 la  $F$  di Fisher e scegliere il precedente  
 se la differenza è molto piccola, ma può essere dovuta  
 al numero di vari riprova.

(11) 
$$\frac{SSR_{k-1} - SSR_k}{SSR_k}$$
 rappresenta la differenza  
 relativa della somma dei  
 quadrati di residui.

La  $F$  di Fisher è una v.c. che ha due gradi di li-  
 bertà (parametri).

È una  $F$  indica che si sta confrontando un mo-  
 dello con quello immediatamente più complicato.

(12) Quando si hanno tutti i dati la  $F$  di Fisher tende  
 ad essere una  $\chi^2$  ad un grado di libertà.

(13) Dopo aver individuato sulla tabella il valore  $f_{\alpha}$  lo  
 si confronta con il valore di  $f$

$$f = (N-k) \frac{SSR_{k-1} - SSR_k}{SSR_k}$$

Se  $f < f_{\alpha}$  si sceglie il modello più semplice. Se  
 invece  $f > f_{\alpha}$  si sceglie il modello più complicato.  
 Il valore di  $f_{\alpha}$  si trova nella tabella  
 dei quadrati di residui. Si adotta il mo-  
 dello più complicato.

(14) È il caso in cui  $\Psi \neq I$

(15) È un approccio alternativo equivalente al Test  $F$   
 del caso  $\Psi = I$  sta decidendo tra una parte ed una  
 separata, anziché scegliere la  $F$  di Fisher si  
 considera il coefficiente del termine di secondo  
 grado.

del caso semplice la parte il coefficiente deve essere  
 nullo, altrimenti esso è diverso da zero.

Si utilizza l'intervallo di confidenza del coeff. qua-  
 drato e si guarda se è zero e compreso lo zero  
 (pag. 218). In questa ipotesi il modello più cor-  
 retto è quello con il coefficiente nullo. Se l'intervallo non contiene  
 lo zero, il coefficiente non è nullo e si sceglie il

modello più complicato.

221

Per questo, se il valore vero è un parametro e pari a zero, allora la stima è distribuita gaussianamente con una determinata varianza  $\sigma^2$ .

Inoltre il valore del parametro diviso per la sua deviazione standard  $\frac{\theta^{LS}}{\sigma}$

diventa una gaussiana standard.

$\frac{\theta^{LS}}{\sigma}$  rappresenta la standardizzazione made nel caso il valore medio è  $\theta^{LS}$ .

Avendo standardizzato abbiamo ottenuto una gaussiana standard che, nel 95% dei casi, ha

$$\left| \frac{\theta^{LS}}{\sigma} \right| \leq 1,96$$

(16) si conclude che, se un parametro vale zero, nel 95% dei casi la sua stima è minore di una volta la sua deviazione standard. Si veda anche quanto detto a pag. 218.

(17) si calcola  $\hat{Y} = \hat{\beta}^T \theta^{LS}$  che rappresenta la predizione dei dati di validazione.

Si calcolano, poi, i residui di validazione  $E^v = Y^v - \hat{Y}$

Successivamente si calcola la somma dei quadrati dei residui dei dati di validazione (SSR<sup>v</sup>), avendo di nuovo:

(18) occorre evitare di avere gruppi di dati relativi a caratteristiche diverse come nel caso di appartenere a due reticoli temporali ben distinti.

(19) se si cerca di prevedere i dati di validazione, si vuole il numero di dati tra i due gruppi di dati. Si può avere in validazione e una dipendenza di SSR, nel rivisto dati i modelli troppo "indimenticabili".

Quando il modello comincia ad imitare il suo nuovo SSR comincia a salire.

(20) non è stata riportata alcuna ipotesi statistica e si è scelto il modello che funziona meglio nel senso che è in grado di prevedere i dati.

(21) Mettendo a parte un dato si identifica il modello

222 e si cerca di prevedere il dato scontato, calcolando  
il residuo.

Si ripete l'operazione su tutti i dati.

(22) Il risultato che si otterrebbe con la procedura adica si  
può ottenere con un'unica procedura.

Si risolve il problema dell'identificazione con tutti i  
dati calcolando la somma dei quadrati dei residui.  
"operato" si ottiene lo stesso risultato che si otterrebbe  
con la procedura adica.

Si calcola matrice  $H$  e si prelevano gli elementi di  
matrice sulla diagonale principale.

(23) Nel calcolo di GCV si ipotizza che tutte le misure abbiano  
lo stesso varianza.

(24)

$$GCV = \frac{N}{(N-9)^2} SSR,$$

con

$$SSR = E^T E \quad (\text{somma dei quadrati dei residui})$$

Per ogni modello possibile candidato alla scelta si  
calcola il GCV e si sceglie il modello con il GCV  
più basso.

(25) Permette di effettuare la CROSSVALIDAZIONE, punto 23  
avere i dati

Con pochi dati, formulando alcune ipotesi statistiche  
si può prevedere il risultato di un'indagine estensiva  
con la non validazione.

Si minimizza la somma dei quadrati dei residui in  
valori di  $\beta$ . Punto calcolata

(26) L'ipotesi che la matrice  $F^V$  per i dati di validazione  
sia uguale alla matrice  $F^I$  (elaborati su dati di identificazione)  
dove

$$F^V = F^I$$

cio significa che i dati di validazione e quelli di  
identificazione sono sulle stesse  $x$ , essendo la matrice  
 $F$  basata sui valori delle  $x$  (pagg. 191 e 192).

Si sta ipotizzando che per ogni  $x$  si abbia un dato  
di validazione ed uno di identificazione.

Operando la validazione per il vettore  $D$ , il risultato  
è  $e$  - una v.c., avendo i dati di validazione un errore  
di misura.

(27) Il vettore  $D$  citato nella nota precedente è il risultato  
di una procedura di identificazione basata su dati  
 $y^I$  anche essi assenti poiché dipendenti dagli errori  
di misura.

Ci sono due diverse sorgenti annuali: una per i  
dati di validazione ed una per quelli di identificazione.

Prima si estraggono dai dati di identificazione (con  
mesi) la serie  $y^I$  calcolata  $D^I$  che viene poi applicata  
ai dati di validazione, pure annuali.

Si calcola il doppio errore atteso poiché ci sono due  
errori annuali.

Adoperando la somma dei quadrati dei residui in  
validazione vale

$$E[E[SSR^*]] = \sigma^2 (N + g)$$

con  $g = \dim(D^I)$

(28) Non conviene aumentare  $g$  poiché - occorrerebbe gli errori  
di validazione. Occorre prendere il minimo  
di  $g$  più piccolo.

Solo l'ipotesi  $I_2$ , perché la quale i dati sono

24  
 generati da un determinato  $\theta^0$  che ha un certo  
 valore di  $p$  (ad. esempio, un modello parabolico  
 indica che  $q = 3$ ).  
 Occorre assegnare il valore di  $q$  più piccolo che non  
 renda nullo il modello  $\theta^0$ .

- (29) Al posto di  $\theta^0$  si considera la sua stima.  
 (30) È il risultato della non validazione rispetto la sua ste-  
 tistica.  
 (31) Si rappresenta la probabilità di optare per il modello  
 più complicato, quando il modello più semplice è  
 quello corretto.  
 (32) Si basa sulla teoria della codifica automatica (complesso  
 scelto).

Se i dati non sono del tutto casuali ed hanno una  
 certa struttura, si può risparmiare nel numero di  
 codifica, prendendo in considerazione del modello di  
 i modelli, anche - considerabile tutti i possibili dati.

- (33) Rispetto al criterio AIC compare la  $(N)$  al posto della  
 costante  $2 \cdot \text{tercio}$ , verosimiglianza per il fatto =  
 numero di modelli più complicati.  
 (34) Lo stimatore considerato coincide con quello di Akaike  
 della formula presentata  $Q$  coincide con  $\Psi^{-1}$   
 $\Psi$  è il vettore delle conducibilità termiche.

U è il vettore delle temperature

- (35) È un esempio di modelli "matriciali", ogni volta  
 che si procede si aggiunge una colonna.  
 (36) Il modello "costante" è, palesemente cattivo, anche se  
 il valore di  $\theta$  (theta) è almeno, si differenzia del  
 la sua deviazione standard (negatività).  
 (37) Il modello "retto" appare graficamente decoroso e  
 la condizione di  $\theta$  (theta) è  $\theta = 0$  (vedi foglio).  
 (38) Il modello "parabola" appare più aderente ai dati e  
 la condizione è ancora soddisfatta.  
 (39) Anche il modello "cubico" è decoroso per la condizio-  
 ne relativa ai parametri, ma è rispettata sulla re-  
 a  $n$  sta avvicinando al limite.  
 (40) Sarebbe il modello migliore poiché pare quasi

