

IDENTIFICAZIONE DI SISTEMI DINAMICI

(PARTE SECONDA)

GIORGIO PICCI

Dipartimento di Ingegneria dell'Informazione,
Università di Padova, Italy

Anno accademico 2012-2013

SEGNALI QUASI PERIODICI

Un processo QP è somma di v componenti armoniche elementari (in generale complesse),

$$\mathbf{z}(t) = \sum_{k=1}^v \mathbf{z}_k e^{j\omega_k t}, \quad t \in \mathbb{Z}$$

dove $\omega_k \in [-\pi, \pi]$ sono pulsazioni che si possono supporre diverse tra loro e le $\mathbf{z}_k, k = 1, \dots, v$ sono variabili aleatorie (complesse) a varianza finita.

Stazionarietà: Le correlazioni delle varie componenti armoniche

$$\mathbb{E} [\mathbf{z}_k \bar{\mathbf{z}}_h] e^{j\omega_k t - j\omega_h s} \quad k, h = 1, 2, \dots, v$$

debbono dipendere da $t - s$, il che può accadere solo per $k = h$ mentre per $k \neq h$ si deve necessariamente avere $\mathbb{E} [\mathbf{z}_k \bar{\mathbf{z}}_h] = 0 \Rightarrow$ **Le $\{\mathbf{z}_k\}$ debbono essere tra loro scorrelate.**

$$r(t, s) = \mathbb{E} \mathbf{z}(t) \bar{\mathbf{z}}(s) = \sum_{k=1}^v \sigma_k^2 e^{j\omega_k(t-s)}, \quad \sigma_k^2 = \mathbb{E} |\mathbf{z}_k|^2$$

da cui $r(t, s) = r(t - s)$ e \mathbf{z} è effettivamente un processo stazionario (in senso debole).

Per un processo *reale* le componenti armoniche sono a coppie complesse coniugate (frequenza zero è reale)

$$\mathbf{z}(t) = \frac{\mathbf{z}(t) + \bar{\mathbf{z}}(t)}{2} = \sum_{k=-v}^v \frac{1}{2} \mathbf{z}_k e^{j\omega_k t}, \quad \omega_{-k} = -\omega_k \quad \mathbf{z}_{-k} = \bar{\mathbf{z}}_k$$

dove le $\{\mathbf{z}_k\}$ sono variabili aleatorie tra loro scorrelate. Il termine corrispondente a $k = 0$ è un' eventuale componente continua a frequenza zero ($\omega_0 = 0$).

Scrivendo $\mathbf{z}_k = \mathbf{x}_k + iy_k$, per indice negativo ($-k$) si ha $\mathbf{z}_{-k} = \mathbf{x}_k - iy_k$; e quindi l'incorrelazione dei coefficienti a indice diverso implica che

$$\mathbb{E} \{ (\mathbf{x}_k + iy_k) \overline{(\mathbf{x}_k - iy_k)} \} = \mathbb{E} \{ (\mathbf{x}_k^2 - \mathbf{y}_k^2) + 2i\mathbf{x}_k\mathbf{y}_k \} = 0$$

da cui

$$\mathbb{E} \mathbf{x}_k^2 = \mathbb{E} \mathbf{y}_k^2, \quad \mathbb{E} \mathbf{x}_k \mathbf{y}_k = 0.$$

Ogni componente armonica elementare si può scrivere in forma reale come

$$\mathbf{z}_k(t) := \frac{1}{2} \{ \mathbf{z}_k e^{j\omega_k t} + \bar{\mathbf{z}}_k e^{-j\omega_k t} \} = \mathbf{x}_k \cos \omega_k t - \mathbf{y}_k \sin \omega_k t \quad k = 1, \dots, \nu.$$

ha potenza statistica $\sigma_k^2 := \mathbb{E} \mathbf{z}_k(t)^2 = \mathbb{E} \mathbf{z}_k(0)^2 = \mathbb{E} \mathbf{x}_k^2 = \mathbb{E} \mathbf{y}_k^2$ e ha una realizzazione di stato del tipo

$$\begin{aligned} \begin{bmatrix} \mathbf{x}_k(t+1) \\ \mathbf{y}_k(t+1) \end{bmatrix} &= \begin{bmatrix} \cos \omega_k & -\sin \omega_k \\ \sin \omega_k & \cos \omega_k \end{bmatrix} \begin{bmatrix} \mathbf{x}_k(t) \\ \mathbf{y}_k(t) \end{bmatrix} \\ \mathbf{z}_k(t) &= [1 \quad 0] \begin{bmatrix} \mathbf{x}_k(t) \\ \mathbf{y}_k(t) \end{bmatrix} \end{aligned}$$

con condizioni iniziali aleatorie scorrelate $\mathbf{x}_k(0) = \mathbf{x}_k$, $\mathbf{y}_k(0) = \mathbf{y}_k$ di uguale varianza σ_k^2 .

Il modello di stato complessivo per il segnale \mathbf{z} ha la forma

$$\mathbf{s}(t+1) = \mathbf{A} \mathbf{s}(t) \tag{1}$$

$$\mathbf{z}(t) = \mathbf{c}^\top \mathbf{s}(t) \tag{2}$$

$\mathbf{s}(t)$ è il vettore di stato di dimensione $2\nu + 1$ ottenuto incolonnando i vettori di stato elementari per $k = 0, 1, 2, \dots, \nu$.

La $\mathbf{z}_0(t) \equiv \mathbf{z}_0(0) \equiv \mathbf{z}_0$ è una componente continua costante. Togliendola, la matrice A ha una struttura diagonale a blocchi $A = \text{diag}\{A_1, \dots, A_v\}$ in cui tutti i blocchi A_k , di dimensione 2×2 , sono matrici ortogonali. La varianza di stato del modello è una matrice diagonale

$$P = \mathbb{E} \mathbf{s}(0) \mathbf{s}(0)^\top = \text{diag}\{\sigma_1^2 I_2, \dots, \sigma_v^2 I_2\}$$

e la funzione di covarianza è

$$\sigma(\tau) = c^\top A^\tau P c = \sum_{k=1}^v \sigma_k^2 \cos \omega_k \tau$$

da cui lo spettro di \mathbf{z} (che dev'essere una funzione pari) si può scrivere nella forma

$$\phi(\omega) = \sum_{k=-v}^v \frac{1}{2} \sigma_k^2 \delta(\omega - \omega_k), \quad \sigma_{-k}^2 = \sigma_k^2.$$

Osservazione 1 Ci si chiede se in un qualunque modello di stato di un processo scalare QP possano esserci autovalori multipli di A . La risposta, è no se si richiede l'osservabilità del modello. Questo fatto si può verificare usando il criterio di Hautus. Con lo stesso criterio si vede facilmente che le rappresentazioni di ordine 2 per le eventuali componenti a frequenze $\omega_k = \pm\pi$ sono ridondanti. Per queste componenti la rappresentazione minima ha ovviamente dimensione uno.

Al modello di stato corrisponde una descrizione “ingresso-uscita” del tipo

$$A(z^{-1})\mathbf{z}(t) = 0 \quad A(z^{-1}) = \prod_{k=1}^v (1 - 2\cos \omega_k z^{-1} + z^{-2}) \quad (3)$$

dove $A(z^{-1}) = z^{-n} \det(zI - A)$ è il polinomio caratteristico della matrice A . SE $\omega_k \neq \pm\pi$, questa descrizione ingresso-uscita è *minima* nel senso che non esiste polinomio di grado più basso di $A(z^{-1})$ che annulla il segnale $\mathbf{z}(t)$; i.e. l'equazione alle differenze ha l'ordine minimo possibile per descrivere $\mathbf{z}(t)$. Se ci sono zeri in $z = \pm 1$ bisogna moltiplicare $A(z^{-1})$ per i corrispondenti fattori $1 \pm z^{-1}$.

Notiamo che, introducendo le condizioni iniziali (aleatorie) la Z-trasformata della soluzione \mathbf{z} si può esprimere come una funzione razionale

$$\mathbf{z}(t) = \frac{N(z^{-1})}{A(z^{-1})}$$

dove $N(z^{-1})$ è un polinomio a coefficienti aleatori determinati dalle condizioni iniziali.

SEGNALI QUASI PERIODICI IN RUMORE BIANCO

Supponiamo che il segnale osservato $y(t)$ sia

$$\mathbf{y}(t) = \mathbf{z}(t) + \mathbf{e}(t)$$

dove $\mathbf{z}(t)$ è un segnale QP reale del tipo analizzato nella sezione precedente e $\mathbf{e}(t)$ è un segnale a spettro continuo. Spesso si assume che $\mathbf{e}(t)$ sia *rumore bianco* di varianza incognita σ^2 .

In questa ipotesi, possiamo descrivere il segnale mediante un modello ingresso-uscita del tipo

$$A(z^{-1})\mathbf{y}(t) = A(z^{-1})\mathbf{e}(t). \quad (4)$$

NB: la cancellazione del fattore comune $A(z^{-1})$ nei due termini non è lecita perchè il sistema parte da **condizioni iniziali non nulle** al tempo zero.

STIMA DI FREQUENZE COL METODO PEM

Descriveremo un metodo di stima di frequenze di un segnale QP immerso in rumore bianco basato su PEM. Il metodo è stato proposto da Nehorai nel 1985 e successivamente è stato riesaminato e affinato da vari autori.

Identificazione di modelli ARMA: si fissa, basandosi sull'informazione disponibile a priori, una classe parametrica di funzioni di trasferimento a fase minima $\{G_\theta(z); \theta \in \Theta\}$ e si calcola *l'errore di predizione del modello* G_θ che è definito come la differenza

$$\boldsymbol{\varepsilon}_\theta(t) := \mathbf{y}(t) - \hat{\mathbf{y}}_\theta(t | t-1) \quad \hat{\mathbf{y}}_\theta(t | t-1) = [G_\theta(z) - 1] G_\theta(z)^{-1} \mathbf{y}(t)$$

è il predittore “approssimato” costruito in base al modello G_θ .

L'errore di predizione associato al modello G_θ si ottiene semplicemente filtrando il processo con la funzione di trasferimento inversa G_θ^{-1}

$$\boldsymbol{\varepsilon}_\theta(t) = G_\theta^{-1}(z) \mathbf{y}(t)$$

Notiamo che questa operazione è possibile se G_θ non ha zeri sul cerchio unitario.

FILTRI NOTCH

Si vuole generalizzare la procedura PEM al caso di processi con una componente QP in rumore bianco, ai quali in senso stretto non sarebbe applicabile.

Il processo y è descritto dalla classe di modelli (4), a poli e zeri tutti sulla circonferenza unitaria. Il polinomio $A(z^{-1})$ è un polinomio **simmetrico**

$$\begin{aligned} A(z^{-1}) &= \prod_{k=1}^n (1 - 2 \cos \omega_k z^{-1} + z^{-2}) = \\ &= 1 + a_1 z^{-1} + a_2 z^{-2} + \dots + a_n z^{-n} + \dots + a_2 z^{-2n+2} + a_1 z^{-2n+1} + 1 z^{-2n} = \\ &= z^{-n} [1 z^n + a_1 z^{n-1} + a_2 z^{n-2} + \dots + a_n + a_n + \dots + a_2 z^{-n+2} + a_1 z^{-n+1} + 1 z^{-n}] \end{aligned}$$

che dipende solo da n parametri. Può essere parametrizzato mediante n coefficienti incogniti $\theta := [a_1, a_2, \dots, a_n]^\top$ che dipendono dai coseni delle frequenze $\omega_1, \omega_2, \dots, \omega_n$.

Per calcolare l'errore di predizione associato a questa classe di modelli gli zeri delle funzione di trasferimento debbono essere stabili; i.e. strettamente dentro il cerchio unit . Si approssima il polinomio a numeratore di $G_\theta(z)$ spostando gli zeri nei $2n$ punti $z_k = \rho e^{\pm j\omega_k}$ dove il parametro $0 < \rho < 1$   scelto prossimo a 1.

Se $z^{2n}A(z^{-1})$ ha uno zero in α allora $z^{2n}A(\rho z^{-1})$ ha uno zero in $\rho\alpha$.

L'errore di predizione si calcola col filtro approssimato:

$$\boldsymbol{\varepsilon}_\theta(t) \simeq \frac{A(z^{-1})}{A(\rho z^{-1})} \mathbf{y}(t)$$

che ha poli (di modulo ρ) interni al cerchio unit . L'equazione alle differenze corrispondente  

$$A_\theta(\rho z^{-1}) \boldsymbol{\varepsilon}_\theta(t) = A_\theta(z^{-1}) \mathbf{y}(t).$$

Questa equazione alle differenze può essere riscritta in modo da facilitare l'aggiornamento ricorsivo dei parametri, nella forma

$$\varepsilon_{\theta}(t) = \mathbf{y}(t) + \mathbf{y}(t - 2n) - \rho^{2n} \varepsilon_{\theta}(t - 2n) - \boldsymbol{\varphi}(t)^{\top} \boldsymbol{\theta} \quad (5)$$

dove $\boldsymbol{\varphi}_{\theta}(t) = [\varphi_1(t) \varphi_2(t) \dots \varphi_n(t)]^{\top}$ è definita dalle

$$\begin{aligned} \varphi_i(t) &= \begin{cases} -\mathbf{y}(t - i) - \mathbf{y}(t - 2n + i) + \rho^i \varepsilon_{\theta}(t - i) + \rho^{2n-i} \varepsilon_{\theta}(t - 2n + i) \\ \text{per } 1 \leq i \leq n - 1 \end{cases} \\ \varphi_i(t) &= -\mathbf{y}(t - n) + \rho^n \varepsilon_{\theta}(t - n) \quad \text{per } i = n. \end{aligned} \quad (6)$$

Per il calcolo del gradiente di $\varepsilon_{\theta}(t)$ cambiato di segno

$$\boldsymbol{\psi}_{\theta}(t) := -\frac{\partial \varepsilon_{\theta}(t)}{\partial \boldsymbol{\theta}}$$

si trova la relazione (vedere il lavoro di Nehorai (1985) per i dettagli)

$$\boldsymbol{\psi}_{\theta}(t) = \frac{1}{A_{\theta}(\rho z^{-1})} \boldsymbol{\varphi}_{\theta}(t) \quad (7)$$

che permette, una volta ottenuto il parametro corrente $\boldsymbol{\theta} = \boldsymbol{\theta}_k$, di calcolare il gradiente $\boldsymbol{\psi}_{\boldsymbol{\theta}_k}(t)$ allo stadio k -simo.

ALGORITMO DI OTTIMIZZAZIONE

Si usa l'algoritmo di Gauss-Newton nella forma standard dei metodi PEM per modelli ARMA

$$\theta_{k+1} = \theta_k + \left[\sum_{t=1}^N \psi_{\theta_k}(t) \psi_{\theta_k}(t)^\top \right]^{-1} \sum_{t=1}^N \psi_{\theta_k}(t) \varepsilon_{\theta_k}(t) \quad (8)$$

dove $\psi_{\theta_k}(t)$ è il gradiente.

Algoritmo 1 *Data la stringa dei dati di ingresso $\mathbf{y} = [\mathbf{y}(1) \dots \mathbf{y}(N)]^\top$, e la stima θ_k alla k -sima iterazione,*

1. *Si calcola la stringa degli errori di predizione $\boldsymbol{\varepsilon}_{\theta_k} = [\varepsilon_{\theta_k}(1) \dots \varepsilon_{\theta_k}(N)]^\top$ risolvendo l'equazione alle differenze (5). Lo schema esplicito richiede il calcolo dell'array $\boldsymbol{\varphi}_{\theta_k} = [\varphi_{\theta_k}(1) \dots \varphi_{\theta_k}(N)]$*
2. *Si calcola il gradiente $\boldsymbol{\Psi}_{\theta_k} = [\boldsymbol{\psi}_{\theta_k}(1) \dots \boldsymbol{\psi}_{\theta_k}(N)]$, usando l'equazione alle differenze (7) in cui $A_{\theta}(\rho z^{-1}) = A_{\theta_k}(\rho z^{-1})$.*
3. *si calcola la matrice pseudo-Hessiana*

$$H_{\theta_k} := \sum_{t=1}^N \boldsymbol{\psi}_{\theta_k}(t) \boldsymbol{\psi}_{\theta_k}(t)^\top = \boldsymbol{\Psi}_{\theta_k} \boldsymbol{\Psi}_{\theta_k}^\top$$

e la sua inversa $P_{\theta_k} := H_{\theta_k}^{-1}$. Questo calcolo si potrebbe anche organizzare in forma ricorsiva come visto nell'algoritmo generale PEM del capitolo precedente.

4. Si aggiorna θ_k usando la (8),

$$\theta_{k+1} = \theta_k + P_{\theta_k} \Psi_{\theta_k} \varepsilon_{\theta_k}$$

5. Si torna al passo 1) ponendo $\theta_k = \theta_{k+1}$.

Questo problema di stima è in genere mal condizionato. Nella funzione obiettivo si osserva un minimo molto pronunciato con una regione di attrazione che è tanto più piccola quanto più grande (prossimo a 1) si prende ρ . L'inizializzazione è quindi importante. Si può partire da delle stime iniziali ottenute da un periodogramma oppure per identificazione di un opportuno modello AR.

Fuori dalla regione di attrazione il gradiente è “piccolo” e la matrice Hessiana è mal condizionata per cui si possono avere esempi di convergenza estremamente lenta o di accumulo di errori di arrotondamento (con molte iterazioni). Per ovviare a questo problema si sceglie ρ variabile con il passo

di iterazione. Per k piccoli, quando le stime sono molto incerte, si prende ρ “piccolo” e poi lo si fa crescere con legge esponenziale, ad esempio

$$\rho(k+1) = \rho_0 \rho(k) + (1 - \rho_0) \rho(\infty) \quad (9)$$

dove $\rho(\infty)$ è il valore a regime desiderato (Nehorai suggerisce 0.995) e ρ_0 è la costante di tempo che determina il tasso di crescita di $\rho(k)$. Nehorai suggerisce di prendere $\rho_0 \simeq 0.99$.

Questa versione dell’algoritmo non è “ricorsiva” come quella dell’articolo ma usa tutti i dati disponibili in “batch”. Questo aggrava un pò i calcoli ma, in linea di principio, dovrebbe portare a prestazioni migliori. In ogni caso nel calcolo dell’errore di predizione e del gradiente i dati iniziali sono sempre male utilizzati e sarebbe opportuno usare un fattore d’oblio $\lambda(t)$ aggiornato con una relazione simile alla (9).

Considerazioni sulla stabilità dell'algoritmo

Come è facile intuire, con stime iniziali poco affidabili, il polinomio $A_{\theta_k}(z^{-1})$ potrebbe risultare instabile con conseguenze disastrose sul calcolo iterativo dell'errore di predizione e del gradiente. C'è però da notare che il polinomio che determina al passo k -simo la dinamica dell'errore di predizione e del gradiente (??), (7), non è $A_{\theta_k}(z^{-1})$ ma bensì il polinomio “scalato” $A_{\theta_k}(\rho_k z^{-1})$ in cui il fattore $\rho_k < 1$ ha un effetto stabilizzante e può riportare i poli a modulo leggermente maggiore di uno dentro il cerchio unitario. Questo è probabilmente il motivo per cui, a quanto afferma Nehorai, non si osserva praticamente mai il fenomeno dell'instabilità.

Considerazioni sul Bias

L'introduzione del fattore di scala ρ nel modello (4) porta ad una descrizione, a stretto rigore “non corretta” del segnale che si vuole identificare. Si può calcolare esplicitamente l'errore asintotico (*bias*) che si commette nella stima PEM del termine $\cos \omega_k$ utilizzando il modello “scalato”. Nel modello con una sola sinusoide gli autori di [?] trovano

$$\cos \hat{\omega} = \frac{(1 + \rho^2) \cos \omega}{2\rho}$$

che con $\rho = \rho(\infty) = .99$ diventa $1.00005 \cos \omega$. Come si vede si tratta di errori tollerabili.

In molte situazioni pratiche il modello “sinusoidi in rumore bianco” potrebbe essere poco realistico e il termine d'errore e sarebbe più accuratamente descrivibile come rumore colorato. Se il rumore additivo non è bianco la modellizzazione mediante un processo ARMA come in (4) non è più valida e i presupposti del metodo vengono a cadere. Inoltre il mal condizionamento rende l'algoritmo delicato e possono essere necessari aggiustamenti *ad hoc* (regolarizzazione etc.) per i casi problematici.

VERIFICA DI IPOTESI

Sia \mathbf{y} un vettore aleatorio di dimensione N e $\mathbf{y} \sim F_\theta$, con θ parametro incognito in $\Theta \subset \mathbb{R}^p$. Date $M + 1$ regioni disgiunte $\Theta_0, \Theta_1, \dots, \Theta_M$ di Θ . Chiameremo H_k la classe $\{F_\theta, \theta \in \Theta_k\}$, $k = 0, 1, \dots, M$.

Il problema della verifica delle $M + 1$ ipotesi

$$H_0, H_1, \dots, H_M$$

è quello di decidere in base ai dati osservati, y , se la distribuzione F_θ (incognita) di \mathbf{y} appartiene o no ad una delle classi $\{H_k\}$. Si tratta cioè di trovare una funzione di decisione

$$\phi : \mathbb{R}^N \rightarrow \{0, 1, \dots, M\}$$

che assegni ad ogni possibile risultato y dell'osservazione \mathbf{y} una ed una sola delle classi H_k . Si decide in sostanza che una delle H_k è “vera” in base all'osservazione $\mathbf{y} = y$.

Naturalmente questa decisione è incerta, ovvero l'assegnazione $y \mapsto H_k$ avviene sempre con una certa probabilità d'errore.

Anche per la verifica d'ipotesi si possono seguire due approcci, quello Fisheriano e quello Bayesiano. Seguiremo l'approccio Fisheriano in cui *nulla si sa a priori sull'appartenenza di F_θ alle classi $\{H_k\}$* .

Definizione 1 *Un'ipotesi H_k si dice semplice (o composta) se Θ_k contiene un solo punto di Θ (oppure se ne contiene più d'uno).*

Ad esempio, se y è il vettore dei risultati dell'osservazione successiva di N lanci indipendenti di una moneta ($y = 0$ se esce croce, $y = 1$ se esce testa), $\theta = p$, $p \in (0, 1)$; l'ipotesi H_0 corrispondente a $p = \frac{1}{2}$ è *semplice*, mentre l'ipotesi $H_1 : \{p; p > \frac{1}{2}\}$ è *composta*.

Supporremo, d'ora in avanti, di avere solo due ipotesi, H_0 e H_1 .

Definizione 2 *Si definisce test dell'ipotesi H_0 contro l'ipotesi alternativa H_1 una statistica a valori binari, $\phi : \mathbb{R}^N \rightarrow \{0, 1\}$. Diremo anche che la ϕ definisce il test.*

Tradizionalmente H_0 (ipotesi nulla) ha un ruolo **privilegiato**. Se $\phi(y) = 0$ si dice che si *accetta* H_0 , se $\phi(y) = 1$ si dice che si *rifiuta* H_0 . Naturalmente queste decisioni si prendono con una certa probabilità d'errore e la teoria si occupa di trovare statistiche per cui questa probabilità d'errore sia la più piccola possibile.

È tradizione definire ϕ dando la regione *di rifiuto di H_0* cioè

$$\mathcal{C} \subset \mathbb{R}^N : = \{y; \phi(y) = 1\}.$$

che si chiama **regione critica** del test ϕ . Un test è assegnato dando la sua regione critica \mathcal{C} .

In genere \mathcal{C} è un sottoinsieme dello spazio campionario definito da un sistema di disuguaglianze del tipo $\{y; \psi_k(y) \leq c_k, k = 1, \dots, m\}$ dove ψ è una statistica, in generale a valori vettoriali in \mathbb{R}^m .

Ogni statistica ψ la cui regione critica è la stessa di ϕ , *definisce lo stesso test*.

Nel decidere se accettare o no H_0 si possono ovviamente commettere degli errori. Il caso ideale sarebbe trovare una statistica ϕ che *discrimina esattamente* le due famiglie di probabilità H_0 e H_1 , cioè una ϕ per cui

$$\begin{cases} \phi(\mathbf{y}) = 0 & \text{se } \mathbf{y} \sim H_0 \\ \phi(\mathbf{y}) = 1 & \text{se } \mathbf{y} \sim H_1 \end{cases},$$

con ovvio significato delle notazioni. Questo accade solo in situazioni degenerate. Supponiamo per esempio che H_0 e H_1 siano entrambi *semplici* e che F_{θ_0} e F_{θ_1} ammettano densità $f_0(y)$ e $f_1(y)$.

Lemma 1 *Si ha discriminazione perfetta tra f_0 e f_1 se e solo se f_0 e f_1 sono ortogonali; i.e.*

$$\int_{\mathbb{R}^n} f_0(x) f_1(x) dx = 0.$$

In particolare f_1 è strettamente positiva negli insiemi in cui f_0 si annulla e, viceversa, dove f_1 si annulla f_0 è strettamente positiva.

Sia infatti \mathcal{C} la regione in cui $f_1(y) > 0$

$$\mathcal{C} = \{y; f_1(y) > 0\}.$$

Chiaramente, se $y \sim f_0$ allora y appartiene a \mathcal{C} con probabilità zero, cioè

$$\mathbb{P}_0(y \in \mathcal{C}) = 0$$

$$\mathbb{P}_1(y \in \mathcal{C}) = 1$$

Quindi basta prendere \mathcal{C} come regione critica e si ha un test perfetto.

In generale si ha invece la situazione descritta dalla seguente tabella

	Ipotesi vera	
	H_0	H_1
Decisione per H_0	O.K.	II
Decisione per H_1	I	O.K.

(10)

- Se è vera H_0 ma $y \in \mathcal{C}$ e quindi ($\phi(y) = 1$) *si decide di rifiutarla* si ha un cosiddetto errore di *prima specie*.

- Se H_0 è falsa (cioè $y \sim H_1$) ma accade che y appartenga al complementare $\bar{\mathcal{C}}$ di \mathcal{C} e quindi si decide che vale H_0 , si ha un errore *di seconda specie*.

Il calcolo delle probabilità d'errore è essenziale per valutare il comportamento di un test. Questo calcolo si presenta particolarmente facile se H_0 e H_1 sono ipotesi semplici.

IPOTESI SEMPLICI

Siano $H_0 = \{F_0\}$ ed $H_1 = \{F_1\}$ due ipotesi *semplici*. Allora la probabilità, α , di commettere un errore di prima specie e quella, β , di commettere un errore di secondo specie, sono

$$\alpha = \int_{\mathcal{C}} dF_0(y) = \mathbb{P}_0(\mathcal{C}) \quad (11)$$

$$\beta = \int_{\overline{\mathcal{C}}} dF_1(y) = \mathbb{P}_1(\mathbb{R}^n \setminus \mathcal{C}) = 1 - \mathbb{P}_1(\mathcal{C}) \quad (12)$$

La probabilità

$$1 - \beta = \mathbb{P}_1(\mathcal{C}) = \mathbb{P}(y \in \mathcal{C}) \quad (13)$$

che, quando vale H_1 , i valori osservati cadano nella regione critica (cioè la probabilità di rifiutare H_0 quando H_0 è falsa) si chiama *potenza* (o potere discriminante) del test.

Notiamo subito che se H_1 non è semplice ovvero Θ_1 contiene più di un valore del parametro, *la potenza è una funzione di θ* .

La terminologia che si usa nelle comunicazioni elettriche è leggermente diversa. Si fa riferimento ad un problema di radar in cui H_0 e H_1 rappresentano rispettivamente *assenza* e *presenza di bersaglio*. Si definiscono

\mathbb{P}_F : probabilità di *falso allarme* (dire che c'è il bersaglio quando non c'è);

\mathbb{P}_M : probabilità di *perdere il bersaglio* (dire che non c'è, cioè accettare H_0 , quando invece il bersaglio c'è; M sta per *miss*, perdita);

\mathbb{P}_D : probabilità di discriminazione (dire che c'è il bersaglio quando in effetti è presente).

Ovviamente

$$\mathbb{P}_F = \alpha, \quad \mathbb{P}_M = \beta, \quad \mathbb{P}_D = 1 - \beta$$

cioè la probabilità di discriminazione coincide con la potenza.

Un test *ottimo* per discriminare tra H_0 e H_1 sarebbe evidentemente quello per cui α e β sono i più piccoli possibile. In realtà si tratta di due obiettivi contrastanti.

Siano ad esempio $H_0 \sim \mathcal{N}(\mu_0, \sigma^2)$ ed $H_1 \sim \mathcal{N}(\mu_1, \sigma^2)$ con $\mu_0 < \mu_1$. Supponiamo $N = 1$ e che la regione critica \mathcal{C} sia della forma

$$\mathcal{C} := \{y; y \geq c\}.$$

Dalla figura si vede che si può scegliere c in modo da avere α piccolo quanto si vuole,

$$\alpha = \frac{1}{\sqrt{2\pi}\sigma} \int_c^{+\infty} e^{-\frac{1}{2} \frac{(y-\mu_0)^2}{\sigma^2}} dy.$$

Più grande si prende c , però, più aumenta la probabilità β di classificare incorrettamente H_1 come H_0 . Infatti

$$\beta = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^c e^{-\frac{1}{2} \frac{(y-\mu_1)^2}{\sigma^2}} dy$$

cresce all'aumentare di c .

Il procedimento classico è allora quello *di fissare* α e di cercare la regione critica \mathcal{C} *che dà il* β *più piccolo possibile* (ovvero il massimo valore della potenza $1 - \beta$). Una siffatta regione critica si chiama la migliore regione critica di misura α (in inglese B.C.R. = *best critical region of size* α). Sussiste a questo proposito il fondamentale:

Lemma 2 (di Neyman – Pearson) *La miglior regione critica di misura α per verificare l'ipotesi $H_0 = \{f_0\}$ contro l'alternativa (semplice) $H_1 = \{f_1\}$ è l'insieme dei punti nello spazio campionario \mathbb{R}^n per cui*

$$\mathcal{C} = \{y; \Lambda(y) \geq k\} \quad (14)$$

dove

$$\Lambda(y) = \frac{f_1(y)}{f_0(y)} \quad (15)$$

e la costante k è scelta in modo tale da aversi

$$\int_{\mathcal{C}} f_0(y) dy = \alpha. \quad (16)$$

Ci sono varie versioni di questo Lemma che trattano del caso in cui H_0 e H_1 non si possono dare mediante densità e si preoccupano dell'eventualità in cui la frontiera di \mathcal{C} ha probabilità positiva. Rimandiamo ai testi di statistica per una trattazione completa.

Notiamo che per le (14) e (15), \mathcal{C} contiene i punti y in cui $f_0(y) = 0$, come intuitivamente ci si aspetta.

Prova: Sia $I_{\mathcal{C}}(y)$ la funzione indicatrice della regione critica \mathcal{C} .
 Dobbiamo massimizzare, rispetto a $\mathcal{C} \subseteq \mathbb{R}^n$

$$1 - \beta = \int_{\mathcal{C}} f_1(y) dy = \int_{\mathbb{R}^n} I_{\mathcal{C}}(y) \frac{f_1(y)}{f_0(y)} f_0(y) dy = \mathbb{E}_0[I_{\mathcal{C}}\Lambda] \quad (17)$$

con il vincolo

$$\alpha = \mathbb{E}_0[I_{\mathcal{C}}]. \quad (18)$$

Questo è un cosiddetto problema di *frontiera libera* del calcolo delle variazioni.

Introducendo il moltiplicatore di Lagrange λ , si tratta di massimizzare, rispetto a $I_{\mathcal{C}}$

$$J(\mathcal{C}) := \mathbb{E}_0[I_{\mathcal{C}}\Lambda] - \lambda \{ \mathbb{E}_0[I_{\mathcal{C}}] - \alpha \}.$$

Sia \mathcal{C}^* la miglior regione critica. Se perturbiamo \mathcal{C}^* di $\delta\mathcal{C}$ (infinitesima) dovremo avere $\delta J(\mathcal{C}^*) = 0$, ovvero

$$\mathbb{E}_0[I_{\delta\mathcal{C}}(\Lambda - \lambda)] = \int_{\delta\mathcal{C}} [\Lambda(y) - \lambda] f_0(y) dy = 0, \forall \delta\mathcal{C},$$

da cui si vede ($f_0(y) \geq 0$) che nei punti della frontiera di \mathcal{C}^* dev'essere

$$\Lambda(y) = \lambda \quad y \in \partial\mathcal{C}^*.$$

In particolare, il moltiplicatore λ dev'essere positivo visto che $\Lambda(y) \geq 0$.
Riscrivendo $J(\mathcal{C})$ come

$$J(\mathcal{C}) = \mathbb{E}_0[I_{\mathcal{C}}(\Lambda - \lambda)] + \lambda \alpha$$

si vede che \mathcal{C}^* dà un massimo di $J(\mathcal{C})$ solo se in \mathcal{C}^* si prendono i punti per cui $\Lambda(y) - \lambda \geq 0$. □

Esempio 1 Sia $H_0 = \{\mathcal{N}(\mu_0, \sigma^2)\}$, $H_1 = \{\mathcal{N}(\mu_1, \sigma^2)\}$ con la stessa varianza σ^2 e medie $\mu_0 \neq \mu_1$ assegnate. Sia y_1, \dots, y_N un campione casuale di numerosità N . Allora

$$\Lambda(y_1, \dots, y_N) = \exp -\frac{1}{2\sigma^2} \left\{ \sum_1^N (y_i - \mu_1)^2 - \sum_1^N (y_i - \mu_0)^2 \right\}.$$

Usando la classica decomposizione $\sum_1^N (y_t - \mu)^2 = \sum_1^N (y_t - \bar{y}_N)^2 + N(\bar{y}_N - \mu)^2$ si ottiene

$$\begin{aligned} \Lambda(y) &= \exp -\frac{1}{2\sigma^2} \{N(\bar{y}_N - \mu_1)^2 - N(\bar{y}_N - \mu_0)^2\} \\ &= \exp -\frac{N}{2\sigma^2} [(\bar{y}_N - \mu_0)^2 + 2(\bar{y}_N - \mu_0)(\mu_0 - \mu_1) + (\mu_0 - \mu_1)^2 - (\bar{y}_N - \mu_0)^2] \\ &= \exp -\frac{N}{2\sigma^2} [2(\bar{y}_N - \mu_0)(\mu_0 - \mu_1) + (\mu_0 - \mu_1)^2]. \end{aligned}$$

La disuguaglianza $\Lambda(y) \geq k$ che definisce la regione critica si può anche riscrivere $\log \Lambda(y) \geq \log k$ ovvero,

$$\frac{N}{2\sigma^2} [2\bar{y}_N(\mu_1 - \mu_0) + \mu_1^2 - \mu_0^2] \geq \log k$$

la quale, se $\mu_1 > \mu_0$ è equivalente alla

$$\bar{y}_N \geq \frac{1}{2}(\mu_1 + \mu_0) + \frac{\sigma^2}{N(\mu_1 - \mu_0)} \log k$$

Se viceversa fosse $\mu_0 > \mu_1$ la regione critica sarebbe definita dalla

$$\bar{y}_N \leq \frac{1}{2}(\mu_1 + \mu_0) - \frac{\sigma^2}{N(\mu_0 - \mu_1)} \log k.$$

Indichiamo il secondo membro di queste disuguaglianze con c_1 o c_2 ; la miglior regione critica è quindi definita dalle

$$\begin{aligned} \mathcal{C}_1 : \{y; \bar{y}_N \geq c_1\} & \quad \text{se} \quad \mu_0 < \mu_1 \\ \mathcal{C}_2 : \{y; \bar{y}_N \leq c_2\} & \quad \text{se} \quad \mu_1 < \mu_0. \end{aligned}$$

Regioni critiche per l'esempio 1

Notiamo che siamo riusciti ad esprimere la regione critica per mezzo della statistica \bar{y}_N . Quindi, assegnato α , e supponendo ad esempio $\mu_1 > \mu_0$ si

può ragionare come segue.

Sotto H_0 , $\bar{y}_N \sim \mathcal{N}(\mu_0, \frac{\sigma^2}{N})$. Normalizzando, $\frac{\sqrt{N}}{\sigma}(\bar{y}_N - \mu_0) \sim \mathcal{N}(0, 1)$, per cui esprimendo c_1 come

$$c_1 = a_1 \frac{\sigma}{\sqrt{N}} + \mu_0 \quad (19)$$

si ha

$$\mathcal{C}_1 = \{y \mid \frac{\sqrt{N}}{\sigma}(\bar{y}_N - \mu_0) \geq a_1\}$$

e, fissato α , si può andare sulle tavole della $\mathcal{N}(0, 1)$ e trovare a_1 in modo tale che

$$\mathbb{P}_0(\mathcal{C}_1) = \mathbb{P}_0\left(\frac{\sqrt{N}}{\sigma}(\bar{y}_N - \mu_0) \geq a_1\right) = \alpha.$$

Usando la (19) si ottiene c_1 e quindi \mathcal{C}_1 .

Questo esempio può essere usato per descrivere il ricevitore ottimo in un sistema di comunicazione digitale in cui la sorgente emette un segnale binario del tipo di figura,

Segnale binario campionato

Il periodo T del segnale trasmesso è esattamente n volte il periodo di campionamento T_c del ricevitore. Il ricevitore fornisce (dopo demodulazione) un segnale che è la somma di quello di figura più un rumore bianco Gaussiano $w(t)$ di media zero e varianza $T_c\sigma^2$ (nota). Questo segnale viene campionato con periodo di T_c secondi, in sincronismo con la sorgente, ottenendo così, ogni T secondi, n campioni descrivibili con lo schema,

$$\begin{cases} y_t = \mu_0 + w_t, & t = 1, \dots, n & \text{sotto } H_0 \\ y_t = \mu_1 + w_t, & t = 1, \dots, n & \text{sotto } H_1 \end{cases}$$

Alla fine di ogni periodo di durata T bisogna decidere in base alle n misure ricevute, y_1, \dots, y_n , se il segnale trasmesso era μ_0 o μ_1 . Lo schema del ricevitore che risulta dalla soluzione precedente è quello di figura

Ricevitore ottimo per la trasmissione numerica.

Questo ricevitore è progettato in modo tale che la probabilità di decidere erroneamente μ_1 quando il segnale effettivamente trasmesso è μ_0 è uguale al valore prefissato α .

È altrettanto importante però conoscere la potenza $1 - \beta = \mathbb{P}_D$. Questa si può in generale calcolare solo a posteriori. In questo caso si ha

$$1 - \beta = \mathbb{P}_1\{y; \bar{y} \geq c_1\}$$

dove c_1 ora è fissato. Basta allora notare che

$$\mathcal{L}_1 = \left\{y; \frac{\sqrt{n}}{\sigma}(\bar{y} - \mu_1) \geq \frac{\sqrt{n}}{\sigma}(c_1 - \mu_1)\right\}$$

dove $\frac{\sqrt{n}}{\sigma}(\mathbf{y} - \mu_1) \sim \mathcal{N}(0, 1)$, sotto H_1 .

In fase di progetto si hanno dei limiti inferiori su $1 - \beta$ in base ai quali si progetta il ricevitore. A questo scopo si può servirsi di grafici che danno $1 - \beta = \mathbb{P}_D$ in funzione di α e del rapporto

$$d = \frac{\sqrt{n} |\mu_1 - \mu_0|}{\sigma}.$$

Questi grafici vengono chiamati ROC (*Receiver Operating Characteristic*) e si possono trovare in letteratura, ad esempio nel libro di Van Trees [?] a pag. 38. Essi sono del tipo di figura

Caratteristiche del ricevitore (ROC).

Esempio 2 Siano $H_0 = \mathcal{N}(0, \sigma_0^2)$, $H_1 = \mathcal{N}(0, \sigma_1^2)$ due ipotesi da verificare in base all'osservazione di un campione casuale di numerosità N . Ovviamente

$$f_i(y_1 \dots y_N) = \frac{1}{(\sqrt{2\pi}\sigma_i)^N} \exp -\frac{1}{2} \frac{\sum_1^N y_t^2}{\sigma_i^2} \quad i = 0, 1$$

e

$$\Lambda(y) = \left(\frac{\sigma_0}{\sigma_1}\right)^N \exp \left\{ -\frac{1}{2} \left(\frac{1}{\sigma_1^2} - \frac{1}{\sigma_0^2} \right) \sum_1^N y_t^2 \right\}.$$

per cui la regione critica è definita dalla disuguaglianza

$$\frac{1}{2} \frac{\sigma_1^2 - \sigma_0^2}{\sigma_0^2 \sigma_1^2} \sum_1^N y_t^2 \geq N \log \frac{\sigma_1}{\sigma_0} + \log k.$$

Suponiamo che $\sigma_1^2 > \sigma_0^2$ allora la disuguaglianza $\Lambda(y) \geq k$ è equivalente alla

$$\sum_1^N y_t^2 \geq \frac{\sigma_0^1 \sigma_1^1}{\sigma_1^2 - \sigma_0^2} \left[N \log \frac{\sigma_1^2}{\sigma_0^2} + 2 \log k \right] := c^2$$

Regione critica per l'esempio 2

Per calcolare c^2 nella formula che dà la regione critica

$$\mathcal{C}^* = \left\{ y; \sum_1^N y_t^2 \geq c^2 \right\}$$

basta ricordare che, sotto H_0

$$\frac{\sum_1^N y_t^2}{\sigma_0^2} \sim \chi^2(N)$$

quindi, fissato α , si va sulle tabelle di $\chi^2(N)$ e si trova il valore a per cui $\mathbb{P}\left\{ \frac{\sum_1^N y_t^2}{\sigma_0^2} \geq a \right\} = \alpha$, e ovviamente si pone $c^2 = a\sigma_0^2$. Per calcolare la potenza,

ricordiamo che sotto H_1 si ha $\frac{\sum_1^N y_t^2}{\sigma_1^2} \sim \chi^2(N)$ e si va a vedere chi è

$$\mathbb{P}_1 \left[\frac{\sum_1^N y_t^2}{\sigma_1^2} \geq \frac{c^2}{\sigma_1^2} \right] = \mathbb{P} \left[\frac{\sum_1^N y_t^2}{\sigma_1^2} \geq \frac{a\sigma_0^2}{\sigma_1^2} \right] = 1 - \beta.$$

IL RICEVITORE A CORRELAZIONE

Nello studio dei sistemi di comunicazione si ha spesso a che fare con problemi di discriminazione di segnale che si formulano in modo naturale come problemi di verifica d'ipotesi. Supponiamo qui che t sia la variabile temporale continua. Tipicamente si deve decidere se vale una delle due ipotesi

$$\begin{cases} H_1 : \mathbf{y}(t) = s(t) + \mathbf{w}(t) & 0 \leq t \leq T \\ H_0 : \mathbf{y}(t) = \mathbf{w}(t) & 0 \leq t \leq T \end{cases} \quad (20)$$

dove $s(t)$ è il segnale utile, che può essere di forma completamente nota, nota a meno del valore di certi parametri, $s(t) = s(t, \theta)$, oppure ignota. Nell'ultimo caso si ha in genere una descrizione probabilistica di $s(t)$. Il rumore $w(t)$ è noto probabilisticamente (tipicamente è rumore bianco Gaussiano).

Supporremo $s(t)$ sia una *funzione nota* e $w(t)$ rumore *bianco Gaussiano* di media zero e varianza σ^2 nota.

Per risolvere il problema useremo un'idea di U. Grenander [?]. Prendiamo un sistema di funzioni ortonormali in $[0, T]$

$$\phi_1(t), \phi_2(t), \dots, \phi_n(t), \dots$$

con

$$\int_0^T \phi_i(t) \phi_j(t) dt = \delta_{ij}.$$

Come $\phi_1(t)$ possiamo sempre scegliere la funzione

$$\phi_1(t) = \frac{s(t)}{\|s(\cdot)\|} = \frac{s(t)}{\left[\int_0^T s^2(t) dt\right]^{\frac{1}{2}}} \quad (21)$$

dove il denominatore è la radice quadrata dell'energia, E , del segnale.

Calcolando la correlazione temporale di $\mathbf{y}(t)$ con $\phi_i(t)$ si trova

$$\mathbf{y}_i := \langle \mathbf{y}, \phi_i \rangle = \int_0^T \mathbf{y}(t) \phi_i(t) dt \quad i = 1, 2, \dots$$

per cui se vale H_1 si ha

$$\begin{aligned} \mathbf{y}_1 &= \frac{1}{\sqrt{E}} \int_0^T s(t) s(t) dt + \int_0^T \mathbf{w}(t) \frac{s(t)}{E} dt \\ &= \sqrt{E} + \mathbf{w}_1 \end{aligned}$$

e

$$\mathbf{y}_i = \int_0^T \mathbf{w}(t) \phi_i(t) dt \quad i = 2, 3, \dots$$

dato che $s(t)$ e $\phi_i(t)$ sono ortogonali per $i \geq 2$.

Se vale H_0 invece si ha

$$\mathbf{y}_1 = \mathbf{w}_1,$$

$$\mathbf{y}_i = \mathbf{w}_i \quad i = 1, 2, \dots$$

Notiamo che le variabili casuali $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_i, \dots\}$ sono indipendenti e Gaussiane.

Infatti

$$\begin{aligned}\mathbb{E}\mathbf{w}_i\mathbf{w}_j &= \mathbb{E} \int_0^T \mathbf{w}(t)\phi_i(t) dt \int_0^T \mathbf{w}(\tau)\phi_j(\tau) d\tau \\ &= \int_0^T \int_0^T \phi_i(t)\phi_j(\tau)\mathbb{E}[\mathbf{w}(t)\mathbf{w}(\tau)] dt d\tau \\ &= \int_0^T \int_0^T \phi_i(t)\phi_j(\tau)\sigma^2\delta(t-\tau) dt d\tau = \sigma^2 \int_0^T \phi_i(t)\phi_j(\tau) dt = \sigma^2\delta_{ij}.\end{aligned}$$

Ne segue che il problema può essere riformulato come segue.

Sotto H_1

$$\begin{cases} \mathbf{y}_1 = \sqrt{E} + \mathbf{w}_1 \\ \mathbf{y}_i = \mathbf{w}_i \end{cases} \quad i = 2, 3, \dots \quad (22)$$

Sotto H_0

$$y_i = \mathbf{w}_i, \quad i = 1, 2, 3, \dots \quad (23)$$

dove il *processo* discreto $\{\mathbf{w}_i\}$ è *bianco*, Gaussiano di media zero e varianza σ^2 . In altri termini le *osservazioni* $\{\mathbf{y}_i\}$ sono una famiglia di variabili Gaussiane indipendenti sotto entrambe le ipotesi, in particolare,

$$\begin{cases} \mathbf{y}_1 \sim \mathcal{N}(\sqrt{E}, \sigma^2) & \text{sotto } H_1 \\ \mathbf{y}_1 \sim \mathcal{N}(0, \sigma^2) & \text{sotto } H_0 \end{cases} \quad (24)$$

e, se $i \geq 2$,

$$\mathbf{y}_i \sim \mathcal{N}(0, \sigma^2) \quad (25)$$

sotto *entrambe* le ipotesi.

Calcolando $\Lambda(\mathbf{y}_1 \dots \mathbf{y}_n)$ per un qualunque n finito si vede subito che per l'indipendenza e per la (25), si ha

$$\Lambda(\mathbf{y}_1 \dots \mathbf{y}_n) = \Lambda(\mathbf{y}_1) \quad (26)$$

ovvero *la decisione ottima è funzione solo del valore assunto da \mathbf{y}_1* . Facendo

i conti,

$$\begin{aligned}\Lambda(\mathbf{y}_1) &= \exp -\frac{1}{2\sigma^2} [(\mathbf{y}_1 - \sqrt{E})^2 - \mathbf{y}_1^2] \\ &= \exp -\frac{1}{2\sigma^2} [-2\mathbf{y}_1\sqrt{E} + E] \\ &= \exp \frac{1}{\sigma^2} \left[\int_0^T \mathbf{y}(t)s(t) dt - \frac{1}{2} \int_0^T s^2(t) dt \right].\end{aligned}\quad (27)$$

Questa formula si ritrova in numerose varietà di problemi di discriminazione. É nota col nome di Likelihood Ratio formula.

La regione critica del test si ottiene imponendo che

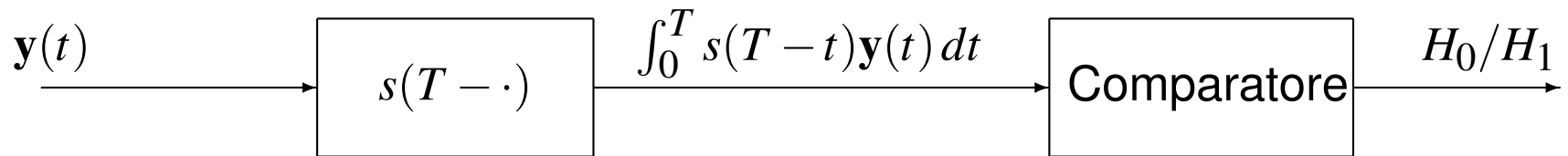
$$\log \Lambda(\mathbf{y}_1) \geq k$$

ovvero

$$\mathbf{y}_1 = \int_0^T \mathbf{y}(t)s(t) dt \geq \frac{1}{2} \int_0^T s^2(t) dt + \sigma^2 k = \frac{1}{2}E + \sigma^2 k = c.$$

Ricordiamo che $\mathbf{y}_1 \sim \mathcal{N}(\sqrt{E}, \sigma^2)$, sotto H_0 , per cui c si può calcolare dalla $c = \frac{a}{\sigma} - \sqrt{E}$, dove a è il punto della curva di $\mathcal{N}(0, 1)$ corrispondente a probabilità della coda superiore pari ad α .

Lo schema del ricevitore ottimo è riportato in figura



Struttura del ricevitore a correlazione

Questo schema viene anche chiamato “ricevitore a correlazione” o a “filtro adattato”. Il secondo nome deriva dal fatto che y_1 si può pensare come l’uscita all’istante T del filtro lineare di funzione di trasferimento

$$h(t) = s(T - t)$$

(che non è causale).

Sul problema della discriminazione di segnali sono stati scritti molti libri.

CRITICA ALL'APPROCCIO CLASSICO

Nella teoria statistica classica le ipotesi H_0 e H_1 giocano un ruolo non simmetrico; H_0 gioca un ruolo privilegiato, dato che fissando α ci si tutela solo sulla probabilità di commettere un errore di prima specie (rifiutando H_0 quando è vera). Questa libertà di scegliere α può condurre a paradossi. Accade normalmente che tutelandosi in modo molto conservativo per evitare un possibile rifiuto di H_0 , ovvero scegliendo α molto piccolo, si finisce coll'accettare H_0 quando invece la scelta di H_1 sarebbe più ragionevole.

Sebbene il Lemma di Neyman-Pearson garantisca che la probabilità dell'errore di seconda specie β venga minimizzata, in genere α non dice nulla sulla correttezza della scelta di H_0 quando in realtà vale H_1 . Occorrerebbe confrontare α con la probabilità dell'errore di seconda specie, β , di accettare H_0 quando invece vale H_1 . Nei problemi ingegneristici, spesso le due ipotesi giocano un ruolo completamente simmetrico e sarebbe forse più ragionevole considerare procedimenti che garantiscono uguali probabilità

d'errore α e β . Da questo punto di vista l'approccio Bayesiano, in cui si postula una distribuzione a priori $\{p_0, p_1\}$ che misura la verosimiglianza delle due ipotesi, è meno criticabile.

IPOTESI COMPOSTE

Molto spesso H_1 (e/o H_0) è un'ipotesi *composta*:

$$H_i = \{f(\cdot, \theta); \theta \in \Theta_i\} \quad i = 0, 1$$

dove Θ_i sono sottoinsiemi di Θ . Questo accade ad esempio nel problema di discriminazione di segnali se $s(t)$ è funzione di uno o più parametri incogniti. Nel rapporto

$$\Lambda(y, \theta) = \frac{f(y, \theta_1)}{f(y, \theta_0)} \quad \theta_1 \in \Theta_1, \theta_0 \in \Theta_0 \quad (*)$$

in cui $\theta = (\theta_1, \theta_0) \in \Theta_1 \times \Theta_0$, è naturale sostituire al parametro θ , uno *stimatore* $\hat{\theta}$, funzione di y , che tenga conto delle diverse regioni ammissibili del parametro sotto le due ipotesi.

Le cose funzionano bene, a patto di prendere $\hat{\theta}$ uguale allo **stimatore di massima verosimiglianza**, nel senso che se a θ_i nel rapporto (*) si sostituisce la statistica $\hat{\theta}_i, i = 0, 1$, che massimizza $f(y, \theta_i)$ nell'insieme $\Theta_i, i = 0, 1$, il test ha certe proprietà ottimali quali la consistenza, l'asintotica normalità, etc.

Definizione 3 Sia $H_0 = \{f(\cdot, \theta); \theta \in \Theta_0\}$ e $H_1 = \{f(\cdot, \theta); \theta \in \Theta_1\}$. Si chiama rapporto di massima verosimiglianza la quantità

$$L(y) := \frac{f(y, \hat{\theta}_1(y))}{f(y, \hat{\theta}_0(y))} \quad (28)$$

dove le statistiche $\hat{\theta}_i$, $i = 0, 1$, massimizzano nei rispettivi domini Θ_i la funzione di verosimiglianza $f(y, \cdot)$, ovvero,

$$\begin{aligned} \hat{\theta}_1(y) &:= \text{Arg} \max_{\{\theta \in \Theta_1\}} f(y, \theta) \\ \hat{\theta}_0(y) &:= \text{Arg} \max_{\{\theta \in \Theta_0\}} f(y, \theta). \end{aligned}$$

Si può allora pensare di definire la regione critica del test assumendo,

$$C := \{y; L(y) \geq k\}$$

Naturalmente per fissare la costante k bisogna usare una procedura un poco più complicata dato che ora $\alpha = \alpha(\theta)$ è funzione di $\theta \in \Theta_0$.

Nella grande maggioranza dei casi pratici, sotto H_0 , $L(\mathbf{y})$ è distribuita in modo indipendente da θ , cioè la d.d.p. $p_0(l)$ di $L(\mathbf{y})$, con $\mathbf{y} \sim \{f(y, \theta), \theta \in \Theta_0\}$, non dipende da θ .

In questo caso $p_0(\cdot)$ si può calcolare come se la d.d.p. di \mathbf{y} corrispondesse ad un *qualunque* valore $\bar{\theta}$ di Θ_0 , in particolare a quello, $\hat{\theta}_0$, che dà il massimo di $f(y, \theta)$ in Θ_0 . Ne segue che $\hat{p}_0 = p_0$ e quindi, fissato α , se si prende la regione critica $C = \{y; L(y) \geq k_\alpha\}$ dove k_α è determinata dalla

$$\int_{k_\alpha}^{\infty} p_0(l) dl = \alpha$$

si ha la probabilità α di commettere un errore di prima specie qualunque sia $\theta \in \Theta_0$.

In generale,

$$\alpha(\theta) = \int_C f(y, \theta) dy \quad \theta \in \Theta_0$$

e, se prendiamo il max rispetto a $\theta \in \Theta_0$ nei due membri e supponiamo che il max dell'integrale sia l'integrale del max rispetto a θ , si ha

$$\alpha_0 = \max_{\theta \in \Theta_0} \alpha(\theta) = \int_C f(y, \hat{\theta}_0(y)) dy.$$

Chiamiamo $\hat{f}_0(y)$ la densità $f(y, \hat{\theta}_0(y))$. Allora se si prende k in modo tale che

$$\int_C \hat{f}_0(y) dy = \int_{\{L(y) \geq k\}} \hat{f}_0(y) dy = \int_k^\infty \hat{p}_0(l) dl = \alpha_0,$$

dove $\hat{p}_0(l)$ è la distribuzione della v.a. $L = L(\mathbf{y})$ con $\mathbf{y} \sim \hat{f}_0(y)$, si commette un errore di prima specie $\leq \alpha_0$.

In generale *la distribuzione di $L(\mathbf{y})$ sotto H_1* dipende da $\theta \in \Theta_1$, per cui la *potenza del test*

$$[1 - \beta](\theta) = \int_{k_\alpha}^{\infty} p_1(l, \theta) dl$$

($p_1(\cdot, \theta)$ è la d.d.p. di $L(\mathbf{y})$ con $\mathbf{y} \sim \{f(\cdot, \theta); \theta \in \Theta_1\}$) è funzione di $\theta \in \Theta_1$.
Notiamo che, se si prende $\theta \in \Theta_0$ si ha

$$p_1(l, \theta) \equiv p_0(l)$$

e $[1 - \beta](\theta) = \alpha$, indipendentemente da θ .

Esempio 3 Sia $\mathbf{y} \sim \mathcal{N}(\mu, \sigma^2)$ in cui $\theta \equiv (\mu, \sigma^2)$ è incognito e

$$H_0 : \mu = \mu_0,$$

dove μ_0 è un valore fissato. Vogliamo verificare l'ipotesi H_0 , cioè che la media di \mathbf{y} sia proprio il valore assegnato, μ_0 , *contro tutte le possibili alternative*, sulla base di N osservazione indipendenti estratte da $N(\mu, \sigma^2)$. Ovviamente si ha

$$\begin{aligned}\Theta_0 &= \left\{ \theta; \mu = \mu_0, \sigma^2 > 0 \right\} \\ \Theta_1 &= \left\{ \theta; \mu \neq \mu_0, \sigma^2 > 0 \right\}\end{aligned}$$

e quindi Θ_1 corrisponde al semipiano aperto $\{\mu, \sigma^2 > 0\}$ privo della semiretta $\mu = \mu_0$. La funzione di verosimiglianza

$$f(y_1, \dots, y_N, \theta) = (2\pi\sigma^2)^{-\frac{N}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_1^N (y_t - \mu)^2 \right\} \quad (29)$$

va massimizzata separatamente su Θ_0 e su Θ_1 . Su Θ_0 , questo corrisponde a calcolare lo stimatore di M.V. per σ^2 *quando la media è nota e vale μ_0* .

Quindi

$$\hat{\theta}_0(y) = s_N^2(y) := \frac{1}{N} \sum_1^N (y_t - \mu_0)^2$$

per cui

$$f(y, \hat{\theta}_0(y)) = \left[2\pi s_N^2(y) \right]^{-\frac{N}{2}} \exp\left(-\frac{N}{2}\right).$$

Per massimizzare $f(y, \theta)$ su Θ_1 si può massimizzare su $\Theta = \{\mu, \sigma^2 > 0\}$ e controllare poi che i valori ottenuti di $\hat{\theta}_1$ non stanno sulla retta $\mu = \mu_0$. Si ha allora

$$\hat{\mu}(y) = \bar{y}_N = \frac{1}{N} \sum_1^N y_t$$

$$\hat{\sigma}_1^2(y) = \hat{\sigma}_N^2(y) = \frac{1}{N} \sum_1^N (y_t - \bar{y}_N)^2.$$

Ovviamente $\bar{y}_N = \mu_0$ con probabilità zero $\forall \theta \in \Theta$, per cui queste stime di

M.V. sono anche il massimo di $f(y, \theta)$ su Θ_1 . A conti fatti,

$$f(y, \hat{\theta}_1(y)) = \left[2\pi \hat{\sigma}_N^2(y) \right]^{-\frac{N}{2}} \exp\left(-\frac{N}{2}\right)$$

per cui, facendo il rapporto, si trova

$$\begin{aligned} L(y) &= \left[\frac{s_N^2(y)}{\hat{\sigma}_N^2(y)} \right]^{\frac{N}{2}} = \left[\frac{\hat{\sigma}_N^2(y) + (\bar{y}_N - \mu_0)^2}{\hat{\sigma}_N^2(y)} \right]^{\frac{N}{2}} \\ &= \left[1 + \frac{(\bar{y}_N - \mu_0)^2}{\hat{\sigma}_N^2(y)} \right]^{\frac{N}{2}}. \end{aligned}$$

Definiamo ora la variabile casuale

$$\mathbf{t} := \frac{\bar{y}_N - \mu_0}{\sqrt{\frac{\hat{\sigma}_N^2(\mathbf{y})}{N-1}}} \quad (30)$$

detta *t di Student*, che sotto l'ipotesi H_0 ha una distribuzione notevole che appare spesso in statistica.

LA DISTRIBUZIONE DI STUDENT

Siano $\mathbf{y} \sim \mathcal{N}(0, 1)$ e $\mathbf{x} \sim \chi^2(n)$, variabili indipendenti. Allora il rapporto

$$\mathbf{t} := \frac{\mathbf{y}}{\sqrt{\mathbf{x}/n}} \quad (31)$$

è distribuito secondo la densità di probabilità

$$p_n(t) = \frac{1}{\sqrt{n}B(1/2, n/2)} \left(1 + \frac{t^2}{n}\right)^{-\frac{n+1}{2}} \quad t \in \mathbb{R} \quad (32)$$

che si chiama *distribuzione di Student a n gradi di libertà* e si denota col simbolo $\mathcal{S}(n)$. Nella (32) B è la funzione Beta di Eulero, definita dalla

$$B(p, q) := \int_0^1 x^{p-1} (1-x)^{q-1} dx = \frac{\Gamma(p)\Gamma(q)}{\Gamma(p+q)}$$

dove la funzione Γ è la nota generalizzazione del fattoriale. Come è noto, per n intero maggiore di 1 la funzione Gamma vale $\Gamma(n) = (n-1)!$.

La distribuzione di Student è una delle distribuzioni notevoli della statistica classica. Essa è tabulata in varie forme in letteratura. Per $n = 1$ essa si

riduce alla distribuzione di Cauchy

$$\mathcal{S}(1) \equiv \frac{1}{\pi(1+t^2)}$$

e quindi si vede che alcuni momenti di $\mathcal{S}(n)$ possono non esistere. In effetti si dimostra che $\mathcal{S}(n)$ possiede momenti fino fino all' $n - 1$ -simo compreso e che questi valgono

$$\begin{aligned} \mu_r &= 0 \quad \text{se } r \text{ è dispari e } r < n \\ \mu_r &= \frac{\Gamma(\frac{1}{2}n - r)\Gamma(r + \frac{1}{2})}{\Gamma(\frac{1}{2}n)\Gamma(\frac{1}{2})} \quad \text{se } r \text{ è pari e } 2r < n. \end{aligned}$$

Si può poi mostrare direttamente che per $n \rightarrow \infty$ la $\mathcal{S}(n)$ converge ad una normale $\mathcal{N}(0, 1)$.

Da quanto esposto sopra si vede che la statistica \mathbf{t} in (30) è distribuita secondo la legge di Student, $\mathcal{S}(N - 1)$, ad $N - 1$ gradi di libertà. Dato che $L(\mathbf{y})$ si può scrivere come

$$L(\mathbf{y}) = \left[1 + \frac{1}{N-1} \mathbf{t}^2 \right]^{\frac{N}{2}} \quad (33)$$

si vede che $L(\mathbf{y})$ dipende da \mathbf{y} solo attraverso la statistica \mathbf{t} . Notiamo che la regione critica $C := \{\mathbf{y}; L(\mathbf{y}) \geq k\}$ può scriversi equivalentemente come

$$C := \left\{ \mathbf{y}; |\mathbf{t}(\mathbf{y})| \geq +\sqrt{(N-1)(k^{\frac{2}{N}} - 1)} \right\} \quad (34)$$

ovvero

$$C := \{\mathbf{y}; |\mathbf{t}(\mathbf{y})| \geq c\}, \quad c > 0.$$

Si vede che sotto H_0 , $L(\mathbf{y})$ ha distribuzione indipendente dal parametro libero, σ^2 e quindi la probabilità, α , di commettere un errore di prima specie, non dipende da σ^2 . Fissato α si può infatti trovare c_α tale che

$$\int_{c_\alpha}^{\infty} p_{N-1}(t) dt = \frac{\alpha}{2}$$

e questo valore c_α definisce una regione critica tale che la probabilità di rifiutare (quando è vera) l'ipotesi $\theta_1 = \mu$, è la stessa qualunque sia il valore della varianza incognita σ^2 .

Invece, *sotto* H_1 , la v.a. \mathbf{t} non è più distribuita come $\mathcal{S}(N-1)$ ed in generale la distribuzione di \mathbf{t} quando $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\theta}_1, \sigma^2)$ dipende da μ e da σ^2 attraverso il cosiddetto “parametro di non centralità”

$$\delta = \frac{\sqrt{N}}{\sigma}(\mu - \mu_0).$$

Il calcolo della potenza si presenta difficoltoso (si possono usare tavole della distribuzione \mathcal{S} “non centrale”). ◇

Osservazioni Vale la pena di estrapolare dall'esempio appena discusso alcune considerazioni di carattere generale. La prima è la seguente.

In molti casi l'ipotesi da verificare, H_0 , si presenta assegnando un valore fissato per alcune delle componenti $(\theta_1, \dots, \theta_p)$ del parametro p -dimensionale θ . Se supponiamo che queste componenti siano le prime $k(\geq 1)$ e scriviamo θ come

$$\theta = \begin{bmatrix} \beta \\ \eta \end{bmatrix} \quad \beta \in \mathbb{R}^k, \eta \in \mathbb{R}^{p-k} \quad (35)$$

allora possiamo supporre che H_0 sia della forma

$$H_0 := \{\theta; \beta = \beta_0\} \quad , \quad (36)$$

con β_0 un vettore *fissato* di \mathbb{R}^k . In questi casi H_1 è l'ipotesi alternativa

$$H_1 := \{\theta; \beta \neq \beta_0\} \quad . \quad (37)$$

Ne viene che Θ_0 , o si riduce ad un punto ($k = p$) oppure sta in un sottospazio di \mathbb{R}^p di dimensione inferiore a p per cui :

la massimizzazione di $f(y, \theta)$ su Θ_1 dà (con probabilità uno) lo stesso risultato della massimizzazione di $f(y, \theta)$ sull'intero spazio dei parametri Θ , questo a meno di situazioni patologiche (che hanno probabilità zero) in cui $f(y, \theta)$ è massimizzata da valori $\hat{\theta}_i$ dei parametri che *non dipendono dai dati osservati*. In generale, quindi

$$\max_{\theta \in \Theta_1} f(y, \theta) = \max_{\theta \in \Theta} f(y, \theta) \quad (38)$$

e cioè $\hat{\theta}_1(\mathbf{y})$ è *l'ordinario stimatore di M.V.*, $\hat{\theta}(\mathbf{y})$, di θ . Per questa ragione, quando H_0 e H_1 sono nella forma (36), (37) si può scrivere

$$L(\mathbf{y}) = \frac{f(\mathbf{y}, \hat{\theta}(\mathbf{y}))}{f(\mathbf{y}, \beta_0, \hat{\eta}(\mathbf{y}))} \quad (39)$$

dove $\hat{\eta}(\mathbf{y})$ è lo stimatore (“condizionato”) di M.V., ovvero, $\eta(\mathbf{y})$ massimizza $f(\mathbf{y}, \beta_0, \eta)$ rispetto ad η nella regione ammissibile Θ_0 , per l'ipotesi H_0 .

Notiamo che, siccome $\Theta_0 \subset \Theta$, si ha sempre

$$f(y, \hat{\theta}(y)) = \max_{\theta \in \Theta} f(y, \theta) \geq \max_{\theta \in \Theta_0} f(y, \theta) = f(y, \beta_0, \hat{\eta}(y))$$

e quindi $L(y) \geq 1 \forall y \in \mathbb{R}^N$ (notare che nella (34) bisogna supporre infatti $k \geq 1$).

Intuitivamente, tanto più grande è $L(y)$, tanto più “verosimile” è l’ipotesi H_1 , dato che $f(y, \hat{\theta}_1(y)) \gg f(y, \hat{\theta}_0(y))$, e quindi si *accetta* H_1 nella regione $\{L(y) \geq k\}$.

Nel seguito supporremo che H_0 e H_1 siano ipotesi del tipo (36), (37). La connessione con la stima di M.V., in particolare con il teorema di Wald permettono di dimostrare dei risultati molto precisi ed utili quando N è grande (per “grandi campioni”). Supporre che le osservazioni $(y_1 \dots y_N)$ siano un campione casuale.

Teorema 1 (Wald) *Si consideri la partizione (35) e sia H_0 definita come in (36), allora, sotto H_0 ,*

1. *la statistica*

$$l(\mathbf{y}) := 2 \log L(\mathbf{y}) \quad (40)$$

converge quando $N \rightarrow \infty$ con probabilità 1 verso

$$Q(\mathbf{y}) = \left[\hat{\beta}(\mathbf{y}) - \beta_0 \right]^\top I^{-1}(\beta_0) \left[\hat{\beta}(\mathbf{y}) - \beta_0 \right] \quad (41)$$

dove $\hat{\beta}(\mathbf{y})$ è il vettore delle prime k -componenti dello stimatore a M.V. di θ e $I(\beta)$ è la sottomatrice di informazione di Fisher, di dimensione $k \times k$, corrispondente al parametro β .

2. *Sotto H_0 , $\sqrt{N} \hat{\beta}_k(\mathbf{y}) \xrightarrow{L} \mathcal{N}(\beta_0, I^{-1}(\beta_0))$ e quindi, asintoticamente*

$$Q(\mathbf{y}) \sim \chi^2(k). \quad (42)$$

3. Sotto H_1 , ponendo $\theta_0 = (\beta_0, \eta)$ si ha

$$Q(\mathbf{y}) \sim \chi^2(k, \lambda)$$

(χ^2 non centrale) dove il parametro di non centralità λ è dato dalla formula

$$\lambda = [\boldsymbol{\theta} - \boldsymbol{\theta}_0]^\top I^{-1}(\boldsymbol{\theta}_0) [\boldsymbol{\theta} - \boldsymbol{\theta}_0]. \quad (43)$$

4. Inoltre il test di max verosimiglianza è consistente nel senso che se C_N è una regione critica di misura α fissata, basato su N campioni, allora,

$$\lim_{N \rightarrow \infty} \int_{C_N} f(y_1, \dots, y_N, \boldsymbol{\theta}) dy^N = 1 \quad \forall \boldsymbol{\theta} \in \Theta_1 \quad (44)$$

dove $C_N = \{y_1, \dots, y_N; L(y_1, \dots, y_N) \geq k_\alpha\}$.

IPOTESI SUL MODELLO LINEARE

Riscriviamo qui il modello lineare standardizzato N -dimensionale

$$\mathbf{y} = S\boldsymbol{\theta} + \sigma\mathbf{w}, \quad \mathbf{w} \sim \mathcal{N}(0, I_N)$$

ottenuto eventualmente col noto procedimento di normalizzazione ‘.

Si chiamano *ipotesi lineari* quelle esprimibili attraverso funzioni lineari di $\boldsymbol{\theta}$, nella fattispecie, H_0 è un’ ipotesi lineare, se per qualche $H \in \mathbb{R}^{k \times p}$, $k \leq p$ e β_0 un vettore assegnato in \mathbb{R}^k , si può scrivere

$$H_0 := \{\boldsymbol{\theta}; H\boldsymbol{\theta} = \beta_0\}. \quad (45)$$

È chiaro che vale la pena di considerare solo ipotesi espresse mediante matrici H di rango pieno, $k < p$.

Alla verifica di ipotesi lineari può essere ricondotta la diagnostica del modello lineare, che comprende l’insieme di verifiche a posteriori sulla significatività delle stime puntuali di $\boldsymbol{\theta}$, ottenute con i metodi della M.V. (o dei M.Q.) su un insieme assegnato di osservazioni. Tipici esempi sono:

1. Ipotesi di adeguatezza del modello lineare

$$H_0 := \{\theta = 0\} \quad (46)$$

(in questo caso $H = I$, $\beta_0 = 0$). Si tratta di verificare se esiste effettivamente un accoppiamento tra segnale e misure (almeno se ne esiste uno di tipo *lineare nei parametri*, come quello ipotizzato dal modello lineare). Se si accetta H_0 con livello di significatività α , si decide che con certezza statistica $1 - \alpha$, le misure y sono essenzialmente costituite da rumore.

2. Ipotesi sul numero di parametri significativi

$$H_0 := \{\theta_{k+1} = \theta_{k+2} = \dots = \theta_p = 0\} \quad . \quad (47)$$

Si tratta di verificare se il modello è *sovraparametrizzato*. In generale se si confrontano tra di loro modelli lineari con un numero diverso di parametri ci si può sempre ridurre a verificare un'ipotesi del tipo (47). Si voglia ad esempio decidere quale dei due modelli di regressione

$$y_t = \theta_0 + \theta_1 u_t + \varepsilon_t \quad (48)$$

$$y_t = \theta_0 + \theta_1 u_t + \theta_2 u_t^2 + \varepsilon_t \quad (49)$$

$t = 1, \dots, N$, si adatta meglio alle osservazioni (u_1, \dots, u_N) e (y_1, \dots, y_N) . Il modello (48) corrisponde evidentemente all'ipotesi $H_0 : \theta_2 = 0$.

3. Ipotesi sulla significatività delle stime puntuali.

Chiamiamo *regione di confidenza* (sotto l'ipotesi H_0) di misura $1 - \alpha$, la regione complementare a quella critica (di misura α). Se $\hat{\theta}(\bar{y})$ è la stima di θ corrispondente all'osservazione \bar{y} , si può pensare di validare la stima verificando l'ipotesi

$$H_0 : \{ \theta = \hat{\theta}(\bar{y}) \} \quad . \quad (50)$$

In questo caso cercare la regione critica del test significa sostanzialmente cercare una regione di confidenza per $\hat{\theta}(\bar{y})$ di coefficiente α uguale al livello di significatività del test.

Esaminiamo l'effetto del vincolo lineare $H\theta = \beta_0$ (che vale solo sotto l'ipotesi H_0) sulla stima di θ .

Ricordiamo che, data la normalità delle osservazioni, il metodo di M.V. si riduce a quello dei M.Q. (non pesati ovvero con $R^{-1} = I$, nel nostro caso).

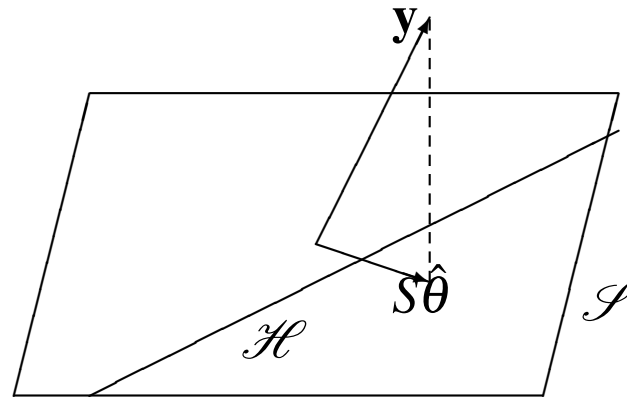
Sotto H_0 , lo stimatore $\hat{\theta}_0$ si trova minimizzando la distanza di y dal sottospazio colonne \mathcal{S} , di S , tenendo conto del vincolo (45); in altri termini, la combinazione lineare delle colonne di S che minimizza $\|y - S\theta\|^2$ non può più essere fatta con coefficienti $\theta \in \mathbb{R}^p$ *arbitrari*, ma deve invece essere costruita con coefficienti $(\theta_1, \dots, \theta_p)$ che soddisfano l'equazione $H\theta = \beta_0$.

$$\hat{\theta}_0(y) = \text{Arg} \min_{\theta \in \{\theta; H\theta = \beta_0\}} \|y - S\theta\|^2 \quad . \quad (51)$$

Se H ha rango k , la condizione $H\theta = \beta_0$ (k equazioni lineari in θ) fornisce solo $p - k$ “parametri liberi” tra le p componenti di θ e quindi solo $p - k$ combinazioni lineari indipendenti delle colonne di S . Questo significa che il minimo di $\|y - S\theta\|^2$ soggetto a $H\theta = \beta_0$ si trova proiettando y non più su tutto \mathcal{S} ma bensì su un opportuno sottospazio affine $p - k$ dimensionale, \mathcal{H} , di \mathcal{S} definito dalla

$$\mathcal{H} := \text{span} \{S\theta; H\theta = \beta_0\} \quad (52)$$

Tutto questo vale ovviamente nell'ipotesi che il vincolo (45) sussista effettivamente, cioè sotto H_0 . Se neghiamo H_0 , diciamo in sostanza che il vincolo non sussiste più e quindi (sotto H_1) la stima di M.V. $\hat{\theta}_1(y) = \hat{\theta}(y)$ si trova col metodo usuale, cioè proiettando y su \mathcal{S} . Si veda la figura



Proiezione ortogonale sul sottospazio \mathcal{H} .

Dato che $\mathcal{S} \supset \mathcal{H}$, la distanza del punto estremo del vettore y da \mathcal{S} dev'essere minore di quella da \mathcal{H} ; si vede quindi abbastanza chiaramente che la somma dei quadrati dei residui nelle due situazioni è *diversa*.

Se vale H_1 si ha un' "errore di approssimazione" di y mediante $S\hat{\theta}$ che è sempre *minore* del corrispondente errore $\|y - S\hat{\theta}_0(y)\|^2$ sotto H_0 . Definiamo allora la somma dei quadrati dei residui sotto le due ipotesi

$$\begin{aligned} H_0 : \quad R_0^2(y) &= \|y - S\hat{\theta}_0(y)\|^2 \\ H_1 : \quad R_1^2(y) &= \|y - S\hat{\theta}(y)\|^2 \end{aligned} \tag{53}$$

Tra R_0^2 ed R_1^2 sussiste una semplice relazione, ovvia se si guarda alla geometria del problema.

Lemma 3 $S\hat{\theta}_0(y)$ è la proiezione ortogonale di $S\hat{\theta}(y)$ su \mathcal{H} e quindi

$$R_0^2(y) = R_1^2(y) + \|S\hat{\theta}(y) - S\hat{\theta}_0(y)\|^2 \tag{54}$$

Prova: Introduciamo per comodità i simboli

$$\hat{\mu}_0 := S\hat{\theta}_0(y) \quad , \quad \hat{\mu}_1 := S\hat{\theta}(y).$$

Per provare la prima affermazione ricordiamo che $y - \hat{\mu}_1$ è ortogonale a \mathcal{S} e quindi in particolare a \mathcal{H} . Allora

$$R_0^2 = \|y - \hat{\mu}_1 + \hat{\mu}_1 - \hat{\mu}_0\|^2 = \|y - \hat{\mu}_1\|^2 + \|\hat{\mu}_1 - \hat{\mu}_0\|^2 + 2\langle y - \hat{\mu}_1, \hat{\mu}_1 - \hat{\mu}_0 \rangle \quad (55)$$

ma il prodotto scalare è nullo perchè $\hat{\mu}_0$ e $\hat{\mu}_1 \in \mathcal{S}$ (e quindi anche la loro differenza) e $y - \hat{\mu}_1$ è *ortogonale* a \mathcal{S} ($\hat{\mu}_1$ è la proiezione ortogonale su \mathcal{S} !). Quindi la (55) si riduce a

$$R_0^2 = R_1^2 + \|\hat{\mu}_1 - \hat{\mu}_0\|^2$$

(teorema di Pitagora) che è proprio la (54). □

Notiamo che se H_0 è vera gli scarti $\|\hat{\mu}_1 - \hat{\mu}_0\|^2$ (che dipendono dal campione osservato y) saranno in media *piccoli* dato che la stima, $S\hat{\theta}(y)$, anche se calcolata senza tener conto della (45) tende per n grande ad essere molto vicina a $S\theta_0$ (θ_0 è il valore vero) e $S\theta_0$ sta in \mathcal{H} per ipotesi.

Viceversa se H_0 è falsa $S\hat{\theta}(y)$ rimane fuori dal sottospazio \mathcal{H} anche se $N \rightarrow \infty$. Se ne ricava che il rapporto

$$\frac{\|\hat{\mu}_1 - \hat{\mu}_0\|^2}{R_1^2} = \frac{R_0^2 - R_1^2}{R_1^2} \quad (56)$$

è *piccolo* se vale H_0 e *grande* se vale H_1 .

Vedremo che questo intuitivo criterio di verifica di H_0 è in sostanza quello fornito dal test di massima verosimiglianza.

CALCOLO DEL RAPPORTO DI MAX VEROSIMIGLIANZA

Come abbiamo visto la densità $f(y, \theta, \sigma^2)$ del vettore aleatorio \mathbf{y} , si può scrivere come

$$f(y, \theta, \sigma^2) = (2\pi\sigma^2)^{-\frac{N}{2}} \exp -\frac{1}{2\sigma^2} \|\mathbf{y} - S\theta\|^2$$

e, dalla discussione precedente segue che *sotto* H_1 , lo stimatore $(\hat{\theta}_1, \hat{\sigma}_1^2)$ di (θ, σ^2) , può essere calcolato massimizzando su tutto lo spazio dei parametri. Quindi $\hat{\theta}_1$ e $\hat{\sigma}_1^2$ sono gli ordinari stimatori di M.V. di θ e σ^2 ,

$$\hat{\theta}_1 = \text{Arg min}_{\theta} \|\mathbf{y} - S\theta\|^2 \quad \hat{\sigma}_1^2 = \frac{1}{N} \|\mathbf{y} - S\hat{\theta}_1\|^2 = \frac{1}{N} R_1^2(y)$$

sostituendo si trova

$$f(y, \hat{\theta}_1(y), \hat{\sigma}_1^2(y)) = \left[2\pi \frac{R_1^2(y)}{N} \right]^{-\frac{N}{2}} \exp -\frac{N}{2}. \quad (57)$$

Sotto H_0 , lo stimatore $\hat{\theta}_0$ risolve il problema di minimo vincolato

$$\hat{\theta}_0(y) = \text{Arg} \min_{\theta \in \{\theta; H\theta = \beta_0\}} \|y - S\theta\|^2 \quad .$$

che per il momento non serve risolvere esplicitamente. Lo stimatore della varianza σ^2 , si trova eseguendo la massimizzazione della funzione di verosimiglianza rispetto a σ^2 , dopo aver sostituito a θ il suo massimo (vincolato). Si trova

$$\hat{\sigma}_0^2(y) = \frac{1}{N} \|y - S\hat{\theta}_0(y)\|^2 = \frac{1}{N} R_0^2(y)$$

e sostituendo nella funzione densità,

$$f(y, \hat{\theta}_0(y), \hat{\sigma}_0^2(y)) = \left[2\pi \frac{R_0^2(y)}{N} \right]^{-N/2} \exp(-N/2)$$

per cui,

$$L(y) = \left[\frac{R_0^2(y)}{R_1^2(y)} \right]^{N/2} = \left[\frac{R_0^2(y) - R_1^2(y)}{R_1^2(y)} + 1 \right]^{N/2} \quad (58)$$

e quindi $L(y)$ è una funzione biunivoca del rapporto (56).

Abbiamo così scoperto che il rapporto (56) è proprio la statistica prescritta dal test di massima verosimiglianza.

Dobbiamo ora studiare la distribuzione del rapporto $R_0^2(y) - R_1^2(y) / R_1^2(y)$. Finora abbiamo studiato il problema di minimo vincolato (51) solo da un punto di vista geometrico qualitativo senza però ottenere una soluzione esplicita. Dobbiamo calcolare la differenza $R_0^2(y) - R_1^2(y)$. Consideriamo la matrice H che definisce l'ipotesi H_0 e definiamo un parametro k -dimensionale β ponendo

$$\beta := H\theta.$$

Lo stimatore di M.V. di β è

$$\hat{\beta}(y) = H\hat{\theta}(y) = H \left[S^\top S \right]^{-1} S^\top y$$

e la sua varianza normalizzata è

$$\frac{1}{\sigma^2} \text{Var}(\hat{\beta}(y)) = H \left[S^\top S \right]^{-1} H^\top := D. \quad (59)$$

Notiamo che $H_0 \equiv \{\beta = \beta_0\}$. Se riusciamo a trasformare il problema (51) in uno in cui le prime variabili $(\theta_1 \dots \theta_k)$ sono sostituite da $(\beta_1 \dots \beta_k)$ il vincolo $H\theta = \beta_0$ si riduce ad aver imposti valori noti e prefissati $(\beta_{01}, \dots, \beta_{0k})$ alle prime k variabili.

Cerchiamo allora un cambiamento di base in \mathbb{R}^N nel quale le prime k equazioni del modello $\mathbf{y} = S\boldsymbol{\theta} + \sigma\mathbf{w}$ diventino del tipo $\mathbf{z} = H\boldsymbol{\theta} + \sigma\mathbf{e}$, con \mathbf{z} vettore k -dimensionale. Questo si traduce nel cercare una matrice Q tale per cui:

$$QS\boldsymbol{\theta} = H\boldsymbol{\theta}, \quad \forall \boldsymbol{\theta} \in \mathbb{R}^p$$

ovvero $QS = H$, dove ovviamente Q dovrà essere di dimensione $k \times N$. Il problema ha sicuramente soluzione, dato che le righe di H sono $k \leq p$ vettori di \mathbb{R}^p linearmente indipendenti che quindi stanno sempre dentro lo spazio righe di S , dato che quest'ultimo per ipotesi è \mathbb{R}^p . Proviamo a trovare una soluzione della forma $Q = CS^\top$, con $C \in \mathbb{R}^{k \times p}$. Dovrà essere $CS^\top S = H$ e quindi,

$$C = H(S^\top S)^{-1} \quad Q = H(S^\top S)^{-1}S^\top.$$

Nella nuova base allora,

$$H(S^\top S)^{-1}S^\top \mathbf{y} = H\boldsymbol{\theta} + \sigma H(S^\top S)^{-1}S^\top \mathbf{w} := \boldsymbol{\beta} + \sigma\mathbf{e}, \quad (60)$$

da cui si ricava il seguente utile risultato.

Lemma 4 Sotto H_0 lo stimatore di M.V. del parametro β è dato dalla formula

$$\hat{\beta}(\mathbf{y}) = H\hat{\theta}(\mathbf{y}) = \beta_0 + \sigma\mathbf{e}, \quad \mathbf{e} \sim \mathcal{N}(0, D) \quad (61)$$

dove $\mathbf{e} = H(S^\top S)^{-1}S^\top \mathbf{w}$, che ha matrice varianza D definita in (59).

Lemma 5 Si ha

$$\|S\hat{\theta}(\mathbf{y}) - S\hat{\theta}_0(\mathbf{y})\|^2 = \|\hat{\beta}(\mathbf{y}) - \beta_0\|_{D^{-1}}^2 \quad (62)$$

Il primo membro è la differenza $\|\hat{\mu}_1 - \hat{\mu}_0\|^2$ definita nella prova del lemma 3.

Prova: Dobbiamo finalmente risolvere il problema di minimo vincolato (51). Allo scopo, introduciamo il moltiplicatore di Lagrange $\lambda \in \mathbb{R}^k$ e consideriamo il problema

$$\min_{\theta} \{ \|y - S\theta\|^2 + \lambda^\top (H\theta - \beta_0) \}.$$

Calcolando il gradiente rispetto a θ della funzione Lagrangiana si trova la condizione

$$-2S^\top (y - S\theta) + H^\top \lambda = 0$$

da cui si ricava per l' estremo l'espressione

$$\hat{\theta}_0(y) = [S^\top S]^{-1} S^\top y - \frac{1}{2} [S^\top S]^{-1} H^\top \lambda, \quad (*)$$

che dipende da λ . Il moltiplicatore si ricava dalla condizione di vincolo che deve valere per $\hat{\theta}_0(y)$ e possiamo riscrivere come $H\hat{\theta}_0(y) = \beta_0$. Questa è equivalente alla

$$\frac{1}{2} H [S^\top S]^{-1} H^\top \lambda = H [S^\top S]^{-1} S^\top y - \beta_0$$

ovvero alla

$$\frac{1}{2} D \lambda = \hat{\beta}(y) - \beta_0.$$

Sostituendo in (*) si trova

$$\hat{\theta}_0(y) = \hat{\theta}(y) - [S^\top S]^{-1} H^\top D^{-1} [\hat{\beta}(y) - \beta_0]$$

ovvero

$$S [\hat{\theta}(y) - \hat{\theta}_0(y)] = S [S^\top S]^{-1} H^\top D^{-1} [\hat{\beta}(y) - \beta_0]$$

che conduce immediatamente a quanto si voleva provare. □

Teorema 2 *La decomposizione (54) si può riscrivere nella forma*

$$R_0^2(\mathbf{y}) = \|\hat{\beta}(\mathbf{y}) - \beta_0\|_{D^{-1}}^2 + R_1^2(\mathbf{y}) \quad (63)$$

e i due membri della somma a secondo membro $\|\hat{\beta}(\mathbf{y}) - \beta_0\|_{D^{-1}}^2$ e $R_1^2(\mathbf{y})$, sono variabili aleatorie indipendenti sotto entrambe le ipotesi.

Prova: Ricordiamo che $\hat{\beta}(\mathbf{y}) = H\hat{\theta}(\mathbf{y})$ e $R_1^2(\mathbf{y}) = \|\mathbf{y} - S\hat{\theta}(\mathbf{y})\|^2$; pertanto basta far vedere che $\hat{\theta}(\mathbf{y})$ e $\mathbf{y} - S\hat{\theta}(\mathbf{y})$ sono indipendenti. Usando il proiettore $P = S(S^\top S)^{-1}S^\top$, si trova

$$\begin{aligned} \text{Cov} [\hat{\theta}(\mathbf{y}), (\mathbf{y} - P\mathbf{y})] &= \mathbb{E} \left[\hat{\theta}(\mathbf{y}) (\mathbf{y} - S\hat{\theta}(\mathbf{y}))^\top \right] = \sigma(S^\top S)^{-1}S^\top \mathbb{E}(\mathbf{y}\mathbf{w}^\top)(I - P)^\top \\ &= \sigma^2(S^\top S)^{-1}S^\top(I - P) \\ &= \sigma^2 \left[(S^\top S)^{-1}S^\top - (S^\top S)^{-1}S^\top S(S^\top S)^{-1}S^\top \right] = 0. \end{aligned}$$

□

Notiamo adesso che sotto H_0 si ha $\hat{\beta}(\mathbf{y}) \sim \mathcal{N}(\beta_0, \sigma^2 D)$ (Lemma 4), e quindi

$$\frac{R_0^2(\mathbf{y}) - R_1^2(\mathbf{y})}{\sigma^2} = \frac{1}{\sigma^2} \|\hat{\beta}(\mathbf{y}) - \beta_0\|_{D^{-1}}^2 \sim \chi^2(k) \quad (64)$$

Si può dimostrare che sotto H_1 la (64) ha una distribuzione χ^2 non centrale, $\chi^2(k, \delta)$; δ essendo il parametro di non centralità *:

$$\delta = \frac{1}{\sigma^2} \|\hat{H}_0 \theta - \beta_0\|_{D^{-1}}^2.$$

Abbiamo poi ricordato piú volte che la somma dei residui "non vincolata" $R_1^2(\mathbf{y})$ è distribuita come $\chi^2(N - p)$:

$$\frac{R_1^2(\mathbf{y})}{\sigma^2} \approx \chi^2(N - p). \quad (65)$$

Possiamo adesso occuparci della distribuzione di probabilità del rapporto (58).

*Per la definizione e i dettagli il riferimento standard è il classico libro di Scheffè [?]

LA DISTRIBUZIONE F

Siano $\mathbf{x}_1 \sim \chi^2(n_1)$ e $\mathbf{x}_2 \sim \chi^2(n_2)$, variabili indipendenti. Allora il rapporto

$$\mathbf{z} := \frac{\mathbf{x}_1/n_1}{\mathbf{x}_2/n_2} \quad (66)$$

è distribuito secondo la densità di probabilità

$$p_{n_1, n_2}(z) = \left[\frac{\Gamma\left(\frac{n_1 + n_2}{2}\right)}{\Gamma\left(\frac{n_1}{2}\right) + \Gamma\left(\frac{n_2}{2}\right)} \right] \left(\frac{n_1}{n_2}\right)^{\frac{n_1}{2}} \frac{z^{\frac{n_1}{2} - 1}}{\left(1 + \frac{n_1}{n_2}z\right)^{\frac{n_1 + n_2}{2}}} \quad z \in \mathbb{R}_+ \quad (67)$$

che si chiama *distribuzione F di Snedecor* a n_1 e n_2 gradi di libertà e si denota col simbolo $\mathcal{F}(n_1, n_2)$.

Per la dimostrazione di questa formula si veda ad es. [Hogg and Craig]. Dopo la Gaussiana, la distribuzione F è forse una delle distribuzioni più importanti della statistica classica. Essa si trova tabulata in varie forme in

letteratura. Per $n_1 = 1$ essa è la distribuzione di probabilità di t^2 , la variabile di Student a n_2 gradi di libertà, elevata al quadrato.

La media μ_1 e la moda, m , di $\mathcal{F}(n_1, n_2)$ esistono se n_1 e n_2 sono strettamente maggiori di 1 e valgono:

$$\mu_1 = \frac{n_2}{n_2 - 2}, \quad m = \frac{n_2(n_1 - 2)}{n_1(n_2 + 2)}$$

Si dimostra che

$$L - \lim_{n_2 \rightarrow \infty} n_1 \mathbf{z} = \chi^2(n_1), \quad (68)$$

inoltre se $\mathbf{z} \sim \mathcal{F}(n_1, n_2)$ e $a := a(n_1, n_2)$ è definito dalla

$$\mathbb{P}(\mathbf{z} \geq a) = \alpha$$

il valore di b per cui

$$\mathbb{P}(\mathbf{z} \leq b) = \alpha$$

è uguale al reciproco di a calcolato in base alla distribuzione $\mathcal{F}(n_1, n_2)$ in cui i gradi di libertà sono scambiati; i.e.

$$b(n_1, n_2) = a(n_2, n_1).$$

Vedere ad esempio il sito web:

<http://econotools.com/jevons/java/Graphics2D/FDist.html>

Teorema 3 *Sotto H_0 , il rapporto*

$$\mathbf{z} := \frac{(N - p) \|\hat{\beta}(\mathbf{y}) - \beta_0\|_{D^{-1}}^2}{k R_1^2(\mathbf{y})} \quad (69)$$

è distribuito secondo la distribuzione F , $\mathcal{F}(k, N - p)$. La regione critica del test si può quindi scrivere

$$C := \{y; \mathbf{z} \geq k_\alpha\} \quad (70)$$

dove k_α corrisponde al valore assegnato di α secondo la distribuzione F .

Il test basato sul rapporto (69) si chiama **test F** ed è di impiego molto generale.

Qualche volta il calcolo di R_0^2 è semplice e conviene calcolare il numeratore di F senza passare attraverso l'espressione dello stimatore $\hat{\beta}$.

IPOSTESI LINEARI SU MODELLI DINAMICI

Proposizione 1 (Dalla teoria PEM) *Lo stimatore PEM $\hat{\theta}_N$, del parametro θ per un modello dinamico lineare che soddisfa le ipotesi di consistenza e normalità asintotica, e lo stimatore di massima verosimiglianza del parametro θ nel modello lineare statico Gaussiano $\mathbf{y} = S\theta + \sigma\mathbf{w}$, in cui*

$$S = \begin{bmatrix} \psi_{\theta_0}(1)^\top \\ \dots \\ \psi_{\theta_0}(N)^\top \end{bmatrix}, \quad \sigma\mathbf{w} = \begin{bmatrix} \mathbf{e}_0(1) \\ \vdots \\ \mathbf{e}_0(N) \end{bmatrix}, \quad t = 1, 2, \dots, N$$

*e i vettori $\{\psi_{\theta_0}(t); t = 1, 2, \dots, N\}$ sono quantità **deterministiche** (funzioni della traiettoria osservata dei processi di misura), **hanno lo stesso limite in legge.***

L'espressione limite per la varianza si ottiene sostituendo alla matrice $\frac{1}{N}S^\top S$ il suo limite per $N \rightarrow \infty$, uguale a $E_0\{\boldsymbol{\Psi}_{\theta_0}(t)\boldsymbol{\Psi}_{\theta_0}(t)^\top\}$.

TEORIA ASINTOTICA:

Sia θ una parametrizzazione identificabile di un modello lineare dinamico. Per $t \rightarrow \infty$ possiamo riferirci al modello lineare statico equivalente:

$$\mathbf{y}(t) = \boldsymbol{\psi}_{\theta_0}(t)^\top \boldsymbol{\theta} + \mathbf{e}_0(t) \quad (*)$$

nel senso che, sotto le ipotesi di normalità asintotica, lo stimatore statico di MV di θ converge in legge a quello PEM.

L'ipotesi lineare da testare *riguarda il modello vero*, i.e. la distribuzione di probabilità vera dei dati. Tipicamente, si prende per H una matrice che seleziona certe componenti del parametro θ ,

$$H\theta = \begin{bmatrix} \theta_{i_1} \\ \dots \\ \theta_{i_k} \end{bmatrix} \in \mathbb{R}^k$$

e si formula l'ipotesi $H\theta = 0$. Questa ipotesi dice che il modello vero ha k parametri nulli, quindi una struttura più semplice di quella ipotizzata in (*). In particolare si può formulare in questo modo un'ipotesi relativa all'ordine del modello.

Notiamo che in questo modo il modello vero appartiene sempre alla classe parametrica (*) ipotizzata all'inizio. Quindi valgono tutte le condizioni per la consistenza e normalità asintotica.

Esempio 4 *Vogliamo verificare l'ipotesi che nel modello AR*

$$\mathcal{M}_2 : \mathbf{y}(t) - a_1\mathbf{y}(t-1) - a_2\mathbf{y}(t-2) = \mathbf{e}(t), \quad \mathbf{e}(t) \sim \mathcal{N}(0, \lambda^2) \quad i.i.d.$$

sia $a_2 = 0$ e quindi i dati siano distribuiti secondo il modello

$$\mathcal{M}_1 : \mathbf{y}(t) - a_1\mathbf{y}(t-1) = \mathbf{e}(t).$$

Ovviamente in questo caso basta prendere $H = [0 \ 1]$.

IL TEST F SU MODELLI DINAMICI

Denotiamo con $\hat{\theta}_0(y)$ lo stimatore PEM sotto l'ipotesi H_0 definita dalla condizione

$$H\theta = 0$$

e con $\hat{\theta}(y)$ lo stimatore PEM sotto l'ipotesi H_1 che il modello sia nella classe generale (*).

Ricordiamo che (notazione: $\hat{\theta} \equiv \hat{\theta}(N)$ stimatore con N dati)

$$H_0 : R_0^2(y) = \|y - S\hat{\theta}_0(y)\|^2 = \sum_{t=1}^N \epsilon_{\hat{\theta}_0}^2 := NV_N^k(\hat{\theta}_0)$$
$$H_1 : R_1^2(y) = \|y - S\hat{\theta}(y)\|^2 = \sum_{t=1}^N \epsilon_{\hat{\theta}}^2 := NV_N^p(\hat{\theta})$$

Proposizione 2 *La statistica*

$$\hat{\mathbf{z}} := \frac{N-p}{k} \frac{V_N^k(\hat{\boldsymbol{\theta}}_0) - V_N^p(\hat{\boldsymbol{\theta}})}{V_N^p(\hat{\boldsymbol{\theta}})}$$

per $N \rightarrow \infty$ converge in legge alla distribuzione $\mathcal{F}(k, N-p)$.

Siccome se $\mathbf{z} := \frac{n_2}{n_1} \frac{\mathbf{x}_1}{\mathbf{x}_2}$ si ha

$$L - \lim_{n_2 \rightarrow \infty} n_1 \mathbf{z} = L - \lim_{n_2 \rightarrow \infty} n_2 \frac{\mathbf{x}_1}{\mathbf{x}_2} = \chi^2(n_1),$$

possiamo dire che asintoticamente,

$$\mathbf{z} := N-p \frac{V_N^k(\hat{\boldsymbol{\theta}}_0) - V_N^p(\hat{\boldsymbol{\theta}})}{V_N^p(\hat{\boldsymbol{\theta}})} \sim \chi^2(k)$$

quindi per N grande il test F si riduce ad un χ^2 .

ANALISI DELLA VARIANZA (ANOVA)

Supponiamo di avere p campioni di numerosità N_1, \dots, N_p estratti da p popolazioni $\mathcal{N}(\mu_i, \sigma^2)$ con $i = 1, \dots, p$, in cui le medie μ_1, \dots, μ_p e la varianza (che è la stessa) sono ignote. Ad esempio potrebbe trattarsi di N_1, \dots, N_p misure di una resistenza fatte su p ponti di Wheatstone fisicamente diversi, ma aventi le stesse caratteristiche di precisione. Indicando con \mathbf{y}_i l' i -esimo blocco di misure e con $\boldsymbol{\theta} = [\mu_1, \dots, \mu_p]^\top$ il vettore p dimensionale delle medie incognite, si può usare il modello lineare:

$$\mathbf{y} := \begin{bmatrix} \mathbf{y}_1 \\ \cdot \\ \cdot \\ \mathbf{y}_p \end{bmatrix} = \begin{bmatrix} e_{N_1} & 0 & 0 & 0 \\ 0 & \cdot & 0 & 0 \\ 0 & 0 & \cdot & 0 \\ 0 & 0 & 0 & e_{N_p} \end{bmatrix} \boldsymbol{\theta} + \boldsymbol{\sigma} \mathbf{w} \quad (71)$$

dove $e_{N_i} = [1 \dots 1]^\top \in \mathbb{R}^{N_i}$ e $\mathbf{w} \sim \mathcal{N}(0, I_N)$, è rumore bianco $N_1 + \dots + N_p = N$ dimensionale.

Si vuole verificare l'ipotesi:

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_p \quad (72)$$

contro l'alternativa che le medie siano diverse. Denotiamo con $\{y_{it}; i = 1, \dots, p; t = 1, \dots, N_i\}$ le p stringhe di osservazioni ottenute nei p procedimenti di misura e calcoliamo la somma dei residui sotto H_1 e sotto H_0 . Quando la (72) non vale (i.e. sotto H_1), si ha:

$$R_1^2(y) = \min_{\mu_1, \dots, \mu_p} \sum_{i=1}^p \sum_{t=1}^{N_i} (y_{it} - \mu_i)^2 = \sum_{i=1}^p \min_{\mu_i} \sum_{t=1}^{N_i} (y_{it} - \mu_i)^2 \quad (73)$$

e siccome il valore di μ_i che dà il minimo di ciascuna somma è :

$$\hat{\mu}_i = \bar{y}_{N_i} = \frac{1}{N_i} \sum_{t=1}^{N_i} y_{it},$$

si ottiene

$$R_1^2(y) = \sum_{i=1}^p \sum_{t=1}^{N_i} (y_{it} - \bar{y}_{N_i})^2 = \sum_{i=1}^p N_i \hat{\sigma}_{N_i}^2,$$

che è distribuita come $\sigma^2 \chi^2(N - p)$.

Sotto H_0 , la media è la stessa $\mu \equiv \mu_1 = \mu_2 = \dots = \mu_p$ e quindi

$$R_0^2(y) = \min_{\mu} \sum_{i=1}^p \sum_{t=1}^{N_i} (y_{it} - \mu)^2 = \sum_{i=1}^p \sum_{t=1}^{N_i} (y_{it} - \bar{y}_N)^2,$$

che sappiamo essere distribuita come $\sigma^2 \chi^2(N-1)$. La differenza $R_0^2 - R_1^2$ si può calcolare usando l'identità

$$\sum_{t=1}^{N_i} (y_{it} - \bar{y}_N)^2 = \sum_{t=1}^{N_i} [y_{it} - \bar{y}_{N_i} + (\bar{y}_{N_i} - \bar{y}_N)]^2 = \sum_{t=1}^{N_i} (y_{it} - \bar{y}_{N_i})^2 + N_i(\bar{y}_{N_i} - \bar{y}_N)^2$$

che permette di scrivere

$$R_0^2 - R_1^2 = \sum_1^p N_i(\bar{y}_{N_i} - \bar{y}_N)^2 \quad (74)$$

che è una somma pesata delle deviazioni delle medie per i singoli gruppi, dalla media totale \bar{y}_N . Notiamo che in questo problema $k = p - 1$ dato che la (72) si può riscrivere

$$\mu_1 - \mu_2 = 0; \dots; \mu_1 - \mu_p = 0,$$

e queste sono $p - 1$ equazioni indipendenti.

Esempio 5 Supponiamo di avere tre serie di misure indipendenti estratte da tre distribuzioni Gaussiane di ugual varianza con statistiche descritte nella tabella 5 seguente e di voler verificare l'ipotesi (H_0) che le medie delle tre distribuzioni siano uguali.

Serie	N_i	$\sum_t y_{it}$	\bar{y}_{N_i}
1	83	11.227	135.87
2	51	7.049	138.22
3	8	1.102	137.75

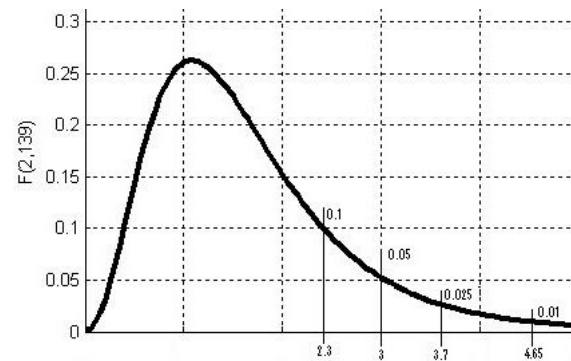
Con semplici calcoli si trova $R_0^2 = 4616.64$ e $R_0^2 - R_1^2 = 238.59$. Nel nostro caso $k = p - 1 = 2$ e $N - p = 142 - 3 = 139$ per cui

$$F = \frac{139}{2} \frac{238.59}{4616.64 + 238.59} = 3.79.$$

Andando nella tabella di $F(2, 139)$ con vari valori di α si trovano i valori critici:

α	0.10	0.05	0.025	0.01
k_α	2.30	3.00	3.70	4.65

per cui ci si trova nella regione critica a meno di non scegliere una probabilità d'errore di prima specie molto piccola.



A questo punto ci si chiede se le medie μ_i delle tre popolazioni sono significativamente diverse oppure no. La risposta è che le medie si possono considerare diverse (si rifiuta H_0), se non si richiede a questa affermazione una certezza statistica troppo elevata. Diciamo che *con probabilità leggermente maggiore del 97,5 per cento le medie sono da considerarsi diverse*. Non c'è invece evidenza sufficiente per fare la stessa affermazione con

certezza statistica del 99 per cento. Se si volesse una certezza statistica del 99 per cento di non commettere errori rifiutando l'ipotesi, ci si troverebbe nella regione di accettazione e quindi si dovrebbe decidere di accettare H_0 . In questo caso però la probabilità α non direbbe nulla (o quasi) sulla correttezza della scelta quando vale H_1 .

Occorrerebbe a questo scopo calcolare la probabilità dell'errore di seconda specie, β , di accettare H_0 quando invece vale H_1 , il che è normalmente complicato perchè H_1 è *un'ipotesi composta*. Tener presente che di norma, β cresce al diminuire di α (vedere ad esempio la figure con le Gaussiane) per cui la decisione di **accettare H_0 con α molto piccoli si rivela in generale priva di senso**, dato che questo può comportare valori elevati di β , anche prossimi ad $1 - \alpha$. ◇

Sotto H_1 il rapporto F non è più distribuito come $F(k, N - p)$, ma bensì come una F non centrale dipendente dal cosiddetto parametro di non centralità λ , dato dalla:

$$\lambda^2 = \lambda^2(\theta, \sigma^2) = \frac{1}{\sigma^2} \|\beta - \beta_0\|_{D^{-1}}^2 = \frac{1}{\sigma^2} \|H\theta - \beta_0\|_{D^{-1}}^2.$$

La F non centrale può essere approssimata con una centrale. Una approssimazione sufficiente in molti casi si ottiene dalla relazione (che ovviamente va intesa tra variabili casuali)

$$F(n_1, n_2, \lambda) \cong \frac{n_1 + \lambda}{n_1} F(n_1^*, n_2)$$

con n_1^* dato da:

$$n_1^* = (n_1 + \lambda)^2 / (n_1 + 2\lambda).$$

In questo modo la potenza si può calcolare usando le tavole di $F(n_1^*, n_2)$. Chiaramente n_1^* non è più un intero e si usa l'intero più prossimo (ovvero si interpola). La potenza corrispondente a $\lambda = \lambda(\theta, \sigma^2)$ è allora:

$$1 - \beta(\theta, \sigma^2) = \int_{\frac{n_1}{n_1 + \lambda} a_\alpha}^{\infty} dF(n_1^*, n_2)$$

si vede che ponendo $\lambda = 0$ (ovvero $\beta = \beta_0$) in questa formula si ottiene α .
Bisogna però ricorrere a tavole complete della distribuzione F (in cui F è calcolata per valori grandi di α).

Esempio [Rao pag.227]

Si vuole stimare la capacità craniale C come funzione di 3 dimensioni lineari, L, B, H , mediante una formole del tipo

$$C \cong \alpha L^{\theta_1} B^{\theta_2} H^{\theta_3} \quad (75)$$

disponendo di una serie di 86 misure di C, L, B, H . Passando ai logaritmi e definendo:

$$y = \log C, \quad x_1 = \log L, \quad x_2 = \log B, \quad x_3 = \log H, \quad \theta_0 : \log \alpha,$$

la (75) si riscrive

$$y \simeq \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

ovvero

$$y_t = \theta_0 + \theta_1 x_{1t} + \theta_2 x_{2t} + \theta_3 x_{3t} + \varepsilon_t \quad (76)$$

$t = 1, \dots, 86$. Supporremo gli errori ε_t Gaussiani indipendenti a media nulla e di uguale varianza incognita σ^2 . Il problema è dunque ridotto ad un

problema di regressione lineare. Siano

$$\bar{x}_i = \frac{1}{86} \sum_1^{86} x_{it}, \quad i = 1, 2, 3$$

le medie campionarie delle variabili x_i , $i = 1, 2, 3$. Conviene sottrarre alla (76) l'equazione per le medie campionarie

$$\bar{y} = \theta_0 + \theta_1 \bar{x}_1 + \theta_2 \bar{x}_2 + \theta_3 \bar{x}_3 + \bar{\varepsilon}$$

ricavando

$$y_t - \bar{y} = \sum_1^3 \theta_i (x_{it} - \bar{x}_i) + (\varepsilon_t - \bar{\varepsilon}). \quad (77)$$

In questo modo si riducono i parametri a 3 e una volta ottenuta la stima $(\hat{\theta}_1, \hat{\theta}_2, \hat{\theta}_3)$; $\hat{\theta}_0$ si ricava dalla

$$\hat{\theta}_0 = \bar{y} - (\hat{\theta}_1 \bar{x}_1 + \hat{\theta}_2 \bar{x}_2 + \hat{\theta}_3 \bar{x}_3) \quad (78)$$

(Lo studente verifichi che se nel modello lineare $y = S\theta + \varepsilon$ la prima colonna è tutta di 1, lo stimatore della prima componente di θ ha proprio l'espressione

(79)). Riscriviamo la (77) in forma vettoriale

$$\Delta \mathbf{y} = S\boldsymbol{\theta} + \sigma \mathbf{w}, \quad \mathbf{w} \sim \mathcal{N}(\mathbf{0}, I).$$

Le medie campionarie tratte dalle 86 misure sono:

$$\bar{y} = 3.17, \quad \bar{x}_1 = 2.275, \quad \bar{x}_2 = 2.15, \quad \bar{x}_3 = 2.11$$

Inoltre si ha

$$S^T S = \begin{bmatrix} 0.0187 & 0.0085 & 0.0068 \\ 0.0085 & 0.029 & 0.0088 \\ 0.0068 & 0.0088 & 0.029 \end{bmatrix} \quad S^T \Delta \mathbf{y} = \begin{bmatrix} 0.030 \\ 0.044 \\ 0.036 \end{bmatrix}$$

Calcolando l'inversa, si trova

$$[S^T S]^{-1} = \begin{bmatrix} 64.21 & -15.57 & -10.49 \\ -15.57 & 41.71 & -9.00 \\ -10.49 & -9.00 & 39.88 \end{bmatrix}$$

da cui $\hat{\boldsymbol{\theta}} = [S^T S]^{-1} S^T \Delta \mathbf{y}$, ha i valori numerici

$$\hat{\theta}_1 = 0.88, \quad \hat{\theta}_2 = 1.04, \quad \hat{\theta}_3 = 0.73$$

e dalla (79)

$$\hat{\theta}_0 = -2.618$$

per cui la formula stimata è

$$C = 0.00241L^{0.88}B^{1.04}H^{0.73} .$$

La somma dei quadrati dei residui si calcola con la formula

$$R_1^2 = \|\Delta y\|^2 - \|S\hat{\theta}\|^2 = \|\Delta y\|^2 - \langle S\hat{\theta}, \Delta y \rangle = \|\Delta y\|^2 - \hat{\theta}^\top S^\top \Delta y \quad (79)$$

ovvero

$$R_1^2 = \sum_{t=1}^{86} (y_t - \bar{y})^2 - (\hat{\theta}_1 0.030 + \hat{\theta}_2 0.044 + \hat{\theta}_3 0.036) \quad (80)$$

$$= 0.127 - 0.099 = 0.028 \quad (81)$$

Per ottenere una stima corretta della varianza usiamo la formula:

$$\hat{\sigma}^2 = \frac{R_1^2}{N-4} = \frac{0.028}{82} = 0.00034$$

notiamo che il numero di parametri θ_i incogniti è sempre 4 anche se abbiamo usato il trucco di ridurre la dimensione del problema a 3. La matrice

di varianze e covarianze di $(\theta_1 \theta_2 \theta_3)$ si ottiene moltiplicando $[S^T S]^{-1}$ per σ^2 . Ad esempio, si trova

$$\text{var}\hat{\theta}_1 = 64.21 \times 3.4 \cdot 10^{-4} \cong 220 \cdot 10^{-4} = 0.022.$$

L'espressione per la varianza di $\hat{\theta}_0$ si può ricavare dalla (79). Lasciamo i facili calcoli al lettore.

Questi calcoli completano la fase di stima del modello. Occorre adesso passare alla fase di validazione della stima.

A: TEST DI ADEGUATEZZA DI UN MODELLO LINEARE

Verifichiamo l'ipotesi

$$H_0 : \theta_1 = \theta_2 = \theta_3 = 0$$

(Non ha molto senso verificare anche $\theta_0 = 0$, cioè $\alpha = 1$ nella (75), dato che questo corrisponderebbe a verificare $C = 1 + "$ rumore"').

Per verificare H_0 col test F basta procurarsi $R_0^2 - R_1^2 = \|\hat{\beta} - \beta_0\|_{D^{-1}}^2 = \|\hat{\beta}\|_{D^{-1}}^2$, dato che in questo caso $\beta = [\theta_1 \theta_2 \theta_3]^\top \equiv \theta$ e $\beta_0 = 0$. Allora

$$D = [S^\top S]^{-1}$$

e quindi

$$\|\hat{\beta}\|_{D^{-1}}^2 = \|\hat{\theta}\|_{D^{-1}}^2 = \hat{\theta}^\top S^\top S \hat{\theta} = \hat{\theta}^\top S^\top \Delta y$$

nell'ultimo passaggio si è sfruttata l'ortogonalità $S\hat{\theta} \perp \Delta y - S\theta$ per cui $\langle S\hat{\theta}, S\hat{\theta} \rangle = \langle S\hat{\theta}, \Delta y \rangle$. Ne segue che $R_0^2 - R_1^2$ è semplicemente l'ultimo addendo nella (79). A titolo di verifica notiamo che sotto H_0 ,

$$R_0^2 = \min_{\theta_0} \|\Delta y - S\theta\|^2 = \|\Delta y\|^2$$

che è proprio il primo addendo in (79).

Ne consegue che:

$$F_A = \frac{N - p}{k} \frac{R_0^2 - R_1^2}{R_1^2} = \frac{86 - 40.099}{3} \frac{0.099}{0.028} = 97.4$$

Sulla tabella di $F(3, 82)$ si trovano i seguenti valori critici k_α

α	0.10	0.05	0.025	0.01	0.005
k_α	2.15	2.70	3.90	4.00	4.60

il che porta sempre a rifiutare l'ipotesi (A), anche se si prende α estremamente piccolo.

B: UGUAGLIANZA DEI PARAMETRI

Verifichiamo l'ipotesi

$$H_0 = \theta_1 = \theta_2 = \theta_3$$

che in questo caso equivale a $H\theta = 0$ ovvero

$$\beta := \begin{bmatrix} 1 & -1 & 0 \\ 1 & 0 & -1 \end{bmatrix} \theta = 0$$

e si potrebbe usare la formula $\|\hat{\beta}\|_{D^{-1}}^2 = R_0^2 - R_1^2$, come fatto in precedenza. Però è più conveniente calcolare direttamente R_0^2 . Ponendo $\theta_1 = \theta_2 = \theta_3 = \eta$ si ha

$$\theta = \mathbf{1}\eta; \quad \mathbf{1} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

e, dato che $\Delta y = (S\mathbf{1})\eta + \sigma w$, le equazioni normali diventano

$$(\mathbf{1}^\top S^\top S \mathbf{1})\eta = \mathbf{1}^\top S^\top \Delta y.$$

Per calcolare $\hat{\eta}$ basta quindi sommare gli elementi di $S^\top S$ e di $S^\top \Delta y$, ottenendo

$$0.125\eta = 0.11 \Rightarrow \hat{\eta} = 0.887$$

per cui

$$R_0^2 = \|\Delta y - S\mathbf{1}\hat{\eta}\|^2 = \|\Delta y\|^2 - \hat{\eta}(\mathbf{1}^\top S^\top \Delta y)$$

e usando la (79),

$$R_0^2 - R_1^2 = 0.099 - \hat{\eta} \cdot 0.11 = 0.099 - 0.098 = 0.001 .$$

Si trova così

$$F_B = \frac{N - p}{k} \frac{0.001}{0.028} = \frac{84}{2} \frac{0.001}{0.028} = 1.4$$

Per $\alpha = 0.10$, il punto critico di $F(2, 82)$ è 2.77 per cui siamo nella regione di accettazione. Per $\alpha = 0.05$ siamo al limite della regione di accettazione. Caso dubbio.

C: ESPONENTI TUTTI UGUALI A UNO

Vogliamo ora considerare l'ipotesi

$$H_0 = \theta_1 = \theta_2 = \theta_3 = 1$$

equivalente a

$$I\theta = \mathbf{1}.$$

Si ha allora

$$R_0^2 - R_1^2 = \|\hat{\theta} - \mathbf{1}\|_{[S^\top S]}^2 = (\hat{\theta} - \mathbf{1})^\top (S^\top S) (\hat{\theta} - \mathbf{1}) = 0.0026$$

In questo caso $k = 3$ e $F_C = 2.5$. La tabella dei valori critici della distribuzione è

	0.10	0.05	0.025
	2.15	2.70	3.30

L'ipotesi si accetta con $\alpha = 0.05$ e si rifiuta con $\alpha = 0.10$. Il caso è un poco dubbio.

D: SOMMA DEGLI ESPONENTI UGUALE A 3

Vogliamo verificare se

$$\theta_1 + \theta_2 + \theta_3 = 3$$

Detto $\beta := \theta_1 + \theta_2 + \theta_3$, si ha $\beta_0 = 3$ e

$$\hat{\beta} = \hat{\theta}_1 + \hat{\theta}_2 + \hat{\theta}_3 = 2.65$$

mentre D è semplicemente la somma degli elementi di $[S^\top S]^{-1}$; i.e. $D = 75.7$, per cui

$$\|\hat{\beta} - 3\|^2/D = \frac{(2.65 - 3)^2}{75.7} = \frac{(0.35)^2}{75.7}$$

e si trova $F_D = 4.72$. I valori critici di $F(1, 82)$ sono

α	0.10	0.05	0.025	0.01
a_α	2.77	3.96	5.20	6.95

e con $\alpha = 0.05$ F_D è nella regione critica. È ragionevole rifiutare H_0 .

E: ESPONENTE DI H UGUALE A ZERO

Vediamo se la variabile H (quella che ha l'esponente più piccolo) è significativa. Saggiamo l'ipotesi:

$$\theta_3 = 0$$

In questo caso $\hat{\beta} = \hat{\theta}_3 = 0.73$, $D = 39.88$ e si calcola subito

$$\|\hat{\beta} - \beta_0\|_{D^{-1}}^2 = \frac{(0.73)^2}{39.9} = 0.0135$$

per cui

$$F_E = \frac{0.0135}{3.4 \cdot 10^{-4}} = 39.72$$

Dalla tabella dei valori critici di $F(1, 82)$ si vede che F_E cade nella regione critica anche per valori molto piccoli di α . Pertanto si rifiuta l'ipotesi $\theta_3 = 0$.

◇

STIMA DELLA COMPLESSITÀ DI UN MODELLO LINEARE

In molti casi il numero di parametri, p nel modello lineare $\mathbf{y} = S\theta + \sigma\mathbf{w}$ non è un dato del problema assegnato a priori, ma piuttosto un parametro che deve essere variato per confrontare l'adeguatezza di modelli più o meno complicati a descrivere i dati di misura.

Nella statistica classica Fisheriana, il problema della scelta ottima di p è un *problema di verifica d'ipotesi*: in base ai dati osservati decidere se il “modello vero” che li ha generati ha complessità p pari ad uno dei numeri naturali compresi in un certo intervallo $[p_{\min}, p_{\max}]$.

Però l'uso di test statistici richiede di fissare α in modo essenzialmente arbitrario, con le risultanti ambiguità che abbiamo visto..

Se la numerosità campionaria è fissa (cosa che da ora in avanti supporremo), è ovvio che all'aumentare di p si ottiene una descrizione sempre migliore dei dati, nel senso che l'errore quadratico medio

$$\hat{\sigma}^2(\mathbf{y}) = \frac{1}{N} \|y - S\hat{\theta}(\mathbf{y})\|_{\mathbb{R}^{-1}}^2$$

diminuisce all'aumentare di p **fino a diventare addirittura zero nel caso limite** $p = N$. Però, a parità di misure disponibili, la qualità delle stime ottenute, misurata ad esempio dalla varianza dei parametri stimati si deteriora all'aumentare di p . Al limite, per p molto grande, il “fit” perfetto ottenuto usando un elevatissimo numero di parametri è in pratica di nessuna utilità dato che la grande varianza delle stime rende inservibile il modello (il modello verrà usato poi per descrivere dati *diversi* da quelli usati in fase di stima).

Notiamo che il problema è adesso inquadrato in un'ottica diversa da quella Fisheriana, senza cioè assumere che esista necessariamente un modello vero di dimensione finita che ha generato i dati. In questo caso i modelli di diversa complessità sono da interpretare solo come **approssimazioni** usate per descrivere dei dati y che potrebbero anche non avere una descrizione “vera” del tipo ipotizzato ma richiedere anche un numero infinito di parametri. In questo caso (che è quello realistico) non esiste un modello stimato di dimensione finita che possa descrivere i dati in modo consistente. C'è quindi sempre un errore di modellizzazione (*bias*) che è diverso da zero anche con dati infiniti.

È quindi necessario procedere per tentativi successivi, aumentando p fino a che il compromesso raggiunto tra bontà del “fit” e varianza della stima sembra accettabile. Questo è il celebre

bias versus variance dilemma

CONFRONTO DI PARAMETRIZZAZIONI

Vogliamo confrontare le varianze degli stimatori di due diverse parametrizzazioni usate per descrivere gli stessi dati. Inizialmente formuleremo il problema in termini di confronto tra due sole alternative possibili. Consideriamo due modelli lineari statici in forma standard

$$\begin{aligned} M_1 : \quad \mathbf{y} &= S_1 \boldsymbol{\theta}_1 + \boldsymbol{\sigma} \mathbf{w} & \boldsymbol{\theta}_1 &\in \mathbb{R}^p \\ M_2 : \quad \mathbf{y} &= [S_1 S_2] \begin{bmatrix} \boldsymbol{\theta}_1 \\ \boldsymbol{\theta}_2 \end{bmatrix} + \boldsymbol{\sigma} \mathbf{w} & \boldsymbol{\theta}_2 &\in \mathbb{R}^k. \end{aligned} \quad (82)$$

In entrambi i casi \mathbf{w} è a media nulla e varianza I_N (matrice identità $N \times N$).

Nel modello più “semplice” M_1 , supporremo che $\text{rango } S_1 = p$.

Nel modello “complicato” M_2 , senza perdita di generalità assumiamo che

$$\text{rango } [S_1 S_2] = p + k \quad . \quad (83)$$

Se ciò non accade, il modello può facilmente essere riparametrizzato eliminando le colonne di S_2 che sono linearmente dipendenti e ridefinendo opportunamente $\boldsymbol{\theta}_2$.

REGRESSIONE LINEARE A STADI

Cercheremo di derivare delle formule per le stime dei parametri e per la varianza del modello M_2 che esprimano queste quantità come correzioni apportate alla stima e alla varianza del parametro θ_1 nel modello M_1 . Questo procedimento va sotto il nome di *regressione (ai M.Q.) a stadi*.

Indichiamo con \mathcal{S} lo spazio colonne della matrice $S := [S_1 S_2]$ e con θ il parametro $p+k$ dimensionale $[\theta_1^\top \theta_2^\top]^\top$ che compare nella (82). Naturalmente la stima ai M.Q. (di Markov) di θ è definita dalle solite formule,

$$\begin{aligned}\hat{\theta}(y) &= (S^\top S)^{-1} S^\top y \\ \text{Var } \hat{\theta} &= \sigma^2 (S^\top S)^{-1} \quad ,\end{aligned}$$

nelle quali però le matrici da invertire sono ora di dimensione $(p+k) \times (p+k)$. Vogliamo mettere in evidenza come si modifica la stima di θ_1 relativa al modello di ordine p per effetto dell'aggiunta dei k ulteriori parametri.

Per la (83) \mathcal{S} si può decomporre in somma diretta

$$\text{span}[S] = \text{span}[S_1 S_2] = \mathcal{S}_1 \oplus \mathcal{S}_2 = \text{span}[S_1] \oplus \text{span}[S_2] \quad (84)$$

e questa decomposizione può essere resa *ortogonale* se si introducono i due proiettori complementari

$$\begin{aligned} P_1 : \mathbb{R}^N &\rightarrow \mathcal{S}_1 & , & & P_1 &= S_1 (S_1^\top S_1)^{-1} S_1^\top & , \\ P_1^\perp : \mathbb{R}^N &\rightarrow \mathcal{S}_1^\perp & , & & P_1^\perp &= I - S_1 (S_1^\top S_1)^{-1} S_1^\top & . \end{aligned}$$

Per semplificare le notazioni in seguito denoteremo con Q_1 la matrice P_1^\perp . Dato che $P_1 + Q_1 = I$, si ha

$$S_2 = P_1 S_2 + Q_1 S_2$$

e siccome le colonne di $P_1 S_2$ stanno per definizione in \mathcal{S}_1 , l'ultimo addendo della (84) può venire sostituito da $\text{span}[Q_1 S_2]$. Quindi

$$\text{span}[S] = \text{span}[S_1] \overset{\perp}{\oplus} \text{span}[Q_1 S_2] \quad (85)$$

dove il simbolo $\overset{\perp}{\oplus}$ sta per somma diretta ortogonale.

Sia \hat{y} la proiezione ortogonale di y sullo spazio colonne, \mathcal{S} , della matrice S . Per l'indipendenza lineare delle colonne di S_1 e S_2 si può esprimere in modo unico \hat{y} nella forma

$$\hat{y} = S_1 \hat{\theta}_1 + S_2 \hat{\theta}_2 \quad , \quad (86)$$

dove $\hat{\theta}_1$ e $\hat{\theta}_2$ sono vettori che rappresentano i corrispondenti coefficienti nelle combinazioni lineari delle colonne di S_1 ed S_2 . Ovviamente $\hat{\theta}_1$ e $\hat{\theta}_2$ sono proprio le stime dei parametri θ_1 e θ_2 nel modello a $p + k$ parametri.

Per il principio di ortogonalità dovrà essere $y - \hat{y} \perp \mathcal{S}$ e quindi anche, separatamente,

$$y - \hat{y} \perp \mathcal{S}_1 \quad , \quad y - \hat{y} \perp Q_1 \mathcal{S}_2 \quad ,$$

che si riscrivono

$$\begin{aligned} S_1^\top (y - S_1 \hat{\theta}_1 - S_2 \hat{\theta}_2) &= 0 \quad , \\ S_2^\top Q_1 (y - S_1 \hat{\theta}_1 - S_2 \hat{\theta}_2) &= 0 \quad . \end{aligned}$$

Queste formule forniscono

$$\hat{\theta}_1 = (S_1^\top S_1)^{-1} S_1^\top [y - S_2 \hat{\theta}_2] \quad , \quad (87)$$

$$\hat{\theta}_2 = (S_2^\top Q_1 S_2)^{-1} S_2^\top Q_1 y. \quad (88)$$

Fatto : $S_2^\top Q_1 S_2$ è invertibile.

Questo si mostra facilmente ricordando che Q_1 è un proiettore. In effetti

$$a^\top S_2^\top Q_1 S_2 a = 0 \Rightarrow a^\top S_2^\top Q_1^\top Q_1 S_2 a = \|Q_1 S_2 a\|^2 = 0$$

e pertanto $S_2 a$ deve stare nello spazio nullo di $Q_1 = P_1^\perp$. Dato che $\text{Ker}(P_1^\perp) = \mathfrak{S}(P_1) = \mathcal{S}_1 = \text{span}[S_1]$, segue che $S_2 a \in \text{span}[S_1]$, ma questo può accadere solo se $a = 0$, dato che le colonne di S_1 ed S_2 sono indipendenti. \diamond

Se indichiamo con il simbolo $\bar{\theta}_1$ la stima di θ_1 ottenuta descrivendo i dati con un modello lineare a p parametri del tipo M_1 , la (87) può essere riscritta come

$$\hat{\theta}_1 = \bar{\theta}_1 - (S_1^\top S_1)^{-1} S_1^\top S_2 \hat{\theta}_2 \quad . \quad (89)$$

che esprime la stima di θ_1 ottenuta con il modello lineare a $p + k$ parametri, come la somma di $\bar{\theta}_1$ e di un termine di correzione dovuto all'introduzione del parametro ulteriore θ_2 .

INTERPRETAZIONE GEOMETRICA

Nella decomposizione (86) i due addendi $S_1 \hat{\theta}_1$ e $S_2 \hat{\theta}_2$ hanno il significato geometrico di *proiezioni oblique* rispettivamente di y su \mathcal{S}_1 lungo \mathcal{S}_2 e di y su S_2 lungo \mathcal{S}_1 .

Dalla formula (88) si vede in particolare che $\hat{\theta}_2$ si può ricavare dalla relazione di ortogonalità

$$Q_1 y - S_2 \theta_2 \perp Q_1 S_2$$

di modo che la proiezione obliqua di y su \mathcal{S}_2 lungo \mathcal{S}_1 , si può *calcolare* facendo prima la proiezione *ortogonale* di $Q_1 y = y - P_1 y$ sul sottospazio $(I - P_1)\mathcal{S}_2 = Q_1 \mathcal{S}_2$ (che è calcolabile risolvendo un problema di minimi quadrati ordinari) e poi moltiplicando per S_2 il parametro $\hat{\theta}_2$ trovato in questo modo*. La matrice di *proiezione obliqua su S_2 lungo \mathcal{S}_1* ha così la rappresentazione

$$P_{2\parallel 1} := S_2 (S_2^\top Q_1 S_2)^{-1} S_2^\top Q_1 \quad (90)$$

*Per evitare interpretazioni errate notiamo che $S_2 \hat{\theta}_2$ *non può essere la proiezione ortogonale di $Q_1 y = y - P_1 y$ sul sottospazio $Q_1 \mathcal{S}_2$* . In effetti quest'ultimo non è nemmeno un sottospazio di \mathcal{S}_2 .

usando la quale si controlla facilmente che in effetti $P_{2\parallel 1}^2 = P_{2\parallel 1}$, mentre

$$P_{2\parallel 1}^\top Q_1 = Q_1 P_{2\parallel 1}.$$

la quale, visto che Q_1 è un proiettore ortogonale e quindi $Q_1 = Q_1^\top$, si può riscrivere come $(Q_1 P_{2\parallel 1})^\top = P_{2\parallel 1}^\top Q_1^\top = Q_1 P_{2\parallel 1}$, i.e. $Q_1 P_{2\parallel 1}$ è simmetrica (e idempotente) e quindi è essa stessa un *proiettore ortogonale* che, per forza di cose, deve proiettare sul sottospazio $Q_1 \mathcal{S}_2$, che è il complemento ortogonale di \mathcal{S}_1 in \mathcal{S} . Infatti:

Proposizione 3 *Sia P la matrice proiezione ortogonale da \mathbb{R}^N sullo spazio \mathcal{S} e P_1 quella sul sottospazio $\mathcal{S}_1 \subset \mathcal{S}$. Allora $P - P_1$ è il proiettore ortogonale che proietta sul complemento ortogonale $\mathcal{S} \cap \mathcal{S}_1^\perp$ e che ha la rappresentazione*

$$P - P_1 = Q_1 P_{2\parallel 1} \tag{91}$$

dove $P_{2\parallel 1}$ è il proiettore obliquo definito in (90).

Prova Basta dimostrare la (91). Usando le formule (??) e (87) si ottiene

$$\begin{aligned}
 \hat{y} = Py &= S_1(S_1^\top S_1)^{-1} S_1^\top y - S_1(S_1^\top S_1)^{-1} S_1^\top S_2 \hat{\theta}_2(y) + S_2 \hat{\theta}_2(y) \\
 &= P_1 y + \left[I - S_1(S_1^\top S_1)^{-1} S_1^\top \right] S_2 \hat{\theta}_2(y) \\
 &= (P_1 + Q_1 P_{2\parallel 1}) y
 \end{aligned}$$

per cui effettivamente si ha $P - P_1 = Q_1 P_{2\parallel 1}$. La decomposizione $P = P_1 + Q_1 P_{2\parallel 1}$ è ovviamente ortogonale, stante che $P_1^\top (P - P_1) = P_1 Q_1 P_{2\parallel 1} = 0$. Notiamo che un'affermazione equivalente è la $\mathcal{S} = P_1 \mathcal{S} \oplus \mathcal{S} \cap \mathcal{S}_1^\perp$. \square

Problema 1 Verificare che $P_{2\parallel 1}$ è idempotente, il suo nucleo è S_1 e la sua immagine è lo spazio colonne di S_2 .

Si può dare una rappresentazione del tutto analoga della proiezione obliqua di y su \mathcal{S}_1 lungo \mathcal{S}_2 e arrivare ad una rappresentazione esplicita della decomposizione (86), del tipo

$$y = P_{1\parallel 2} y + P_{2\parallel 1} y = S_1(S_1^\top Q_2 S_1)^{-1} S_1^\top Q_2 y + S_2(S_2^\top Q_1 S_2)^{-1} S_2^\top Q_1 y \quad (92)$$

dove Q_2 ha un significato duale a Q_1 . Questa espressione è forse più semplice della decomposizione ortogonale che abbiamo illustrato sopra ma è meno comoda da usare perchè non è ortogonale.

Figura 5.3 (proiezione obliqua)

CONFRONTO DELLE VARIANZE

Concentriamoci ora sul calcolo delle varianze degli stimatori. Introduciamo allo scopo le seguenti notazioni:

$$\begin{aligned}\bar{\Sigma}_1 &:= [S_1^\top S_1]^{-1} \\ A_1 &:= [S_1^\top S_1]^{-1} S_1^\top \\ \Sigma_2 &:= [S_2^\top Q_1 S_2]^{-1} \quad ;\end{aligned}$$

ovviamente, $\bar{\theta}_1 = A_1 y$ e $\text{Var}_{\theta_1} \bar{\theta}_1 = \sigma^2 \bar{\Sigma}_1$. Nel seguito i pedici θ_1 e θ staranno ad indicare il modello rispetto a cui si calcola l'aspettazione (e quindi la varianza).

Proposizione 4 *Siano $\hat{\theta}_1(y)$ e $\hat{\theta}_2(y)$ gli stimatori di Markov definiti dalle formule (87) e (88). Si ha allora:*

$$\text{Var}_{\theta} \left\{ \begin{bmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \end{bmatrix} \right\} = \sigma^2 \begin{bmatrix} \bar{\Sigma}_1 + A_1 S_2 \Sigma_2 S_2^\top A_1^\top & -A_1 S_2 \Sigma_2 \\ -\Sigma_2 S_2^\top A_1^\top & \Sigma_2 \end{bmatrix} . \quad (93)$$

Prova: Incominciamo col dimostrare che $\text{Var}_\theta [\hat{\theta}_2] = \sigma^2 \Sigma_2$. Dalla (88) si ha

$$\text{Var}_\theta [\hat{\theta}_2] = \Sigma_2 S_2^\top Q_1 \text{Var}_\theta [\mathbf{y}] Q_1 S_2 \Sigma_2 = \sigma^2 \Sigma_2 S_2^\top Q_1 S_2 \Sigma_2 = \sigma^2 \Sigma_2 \quad ,$$

dato che $\text{Var}_\theta [\mathbf{y}] = \sigma^2 I$ ed Q_1 è idempotente.

Mostriamo ora che *i due stimatori $\bar{\theta}_1(\mathbf{y})$ e $\hat{\theta}_2(\mathbf{y})$ sono scorrelati*. Si ha infatti:

$$\text{Cov}_\theta [\bar{\theta}_1(\mathbf{y}), \hat{\theta}_2(\mathbf{y})] = \bar{\Sigma}_1 S_1^\top \text{Var}_\theta [\mathbf{y}] Q_1 S_2 \Sigma_2 = \sigma^2 \bar{\Sigma}_1 S_1^\top Q_1 S_2 \Sigma_2 = 0 \quad ,$$

perchè $S_1^\top Q_1 = Q_1 S_1 = 0$.

Usando ora la (89) si trova

$$\text{Cov}_\theta [\hat{\theta}_1(\mathbf{y}), \hat{\theta}_2(\mathbf{y})] = -A_1 S_2 \text{Var}_\theta [\hat{\theta}_2] = -\sigma^2 A_1 S_2 \Sigma_2 \quad .$$

Calcoliamo infine $\text{Var}_\theta [\hat{\theta}_1(\mathbf{y})]$. Dato che $\bar{\theta}_1(\mathbf{y})$ e $\hat{\theta}_2(\mathbf{y})$ sono scorrelati, si ha

$$\begin{aligned} \text{Var}_\theta [\bar{\theta}_1(\mathbf{y}) - A_1 S_2 \hat{\theta}_2(\mathbf{y})] &= \text{Var}_\theta [\bar{\theta}_1(\mathbf{y})] + A_1 S_2 \text{Var}_\theta [\hat{\theta}_2(\mathbf{y})] S_2^\top A_1^\top \\ &= \sigma^2 [\bar{\Sigma}_1 + A_1 S_2 \Sigma_2 S_2^\top A_1^\top] \quad . \end{aligned}$$

che conclude la dimostrazione della formula (93). □

La formula (93) descrive l'effetto dell'aumento del numero di parametri nel modello sulla varianza delle stime. In particolare mostra che la stima $\hat{\theta}_1$ di θ_1 nel modello maggiorato è generalmente “peggiore” della prima in termini di varianza. La varianza, Σ_1 , di $\hat{\theta}_1$ è in effetti *più grande* di quella di $\bar{\theta}_1$, essendo

$$\Sigma_1 = \bar{\Sigma}_1 + A_1 S_2 \Sigma_2 S_2^\top A_1^\top$$

e il termine che si somma a $\bar{\Sigma}_1$ è in generale non nullo.

Purtroppo però la varianza delle stime del parametro θ non è un criterio “oggettivo” per arrivare alla scelta dell'ordine del modello. Infatti, se le colonne di S_2 sono *ortogonali* a \mathcal{S}_1 , ovvero

$$S_1^\top S_2 = 0 \quad (S_2^\top S_1 = 0)$$

le formule si semplificano (dato che $Q_1 S_2 = S_2$) e i due stimatori $\hat{\theta}_1$ e $\hat{\theta}_2$ si possono calcolare indipendentemente l'uno dall'altro con le solite formule,

$$\hat{\theta}_i(\mathbf{y}) = (S_i^\top S_i)^{-1} S_i^\top \mathbf{y} \quad , \quad i = 1, 2 \quad .$$

In particolare si trova $\hat{\theta}_1 = \bar{\theta}_1$ e quindi anche $\Sigma_1 = \bar{\Sigma}_1$.

Per comprendere questo fenomeno (che a prima vista può sembrare sconcertante) basta pensare che ci sono molte parametrizzazioni del modello “ideale” $S\theta$ che sono assolutamente equivalenti agli effetti di descrivere i dati y . Per esempio, introducendo una fattorizzazione QR di S , vede facilmente che si può sempre fattorizzare S come prodotto di una matrice a colonne ortogonali (le prime $p+k$ colonne di Q) per una matrice quadrata $R \in \mathbb{R}^{p+k \times p+k}$ non singolare (a struttura triangolare inferiore). Definendo il nuovo parametro $\beta := R\theta$ si può riparametrizzare il modello in modo tale che le colonne di S siano ortogonali. In questo caso la varianza di $\hat{\beta}_1$ non aumenta aumentando la parametrizzazione del modello con k nuovi parametri.

La morale della storia è che la varianza delle stime dei parametri *dipende dal sistema di coordinate scelto per rappresentare il modello* (in breve, “dalla base”). I confronti dovrebbero essere quindi fatti solo tra quantità che sono *invarianti per cambio di base*. Quantità di questo genere sono ad esempio **gli errori residui di modellizzazione**. □

Dato che i modelli servono in ultima analisi a costruire predittori per dati “futuri” (non ancora osservati) si può allora porre un problema di scelta del modello che fornisce l'*approssimazione ottima dei dati* (non di un ipotetico modello vero). Si sceglierà così quel modello che dà *la migliore predizione dei dati futuri*. Beninteso l'errore di predizione dovrà qui tener conto anche dell' **incertezza introdotta nel modello usato per la predizione dal fatto che esso usa necessariamente un parametro stimato** che è, esso stesso, una variabile aleatoria.

Questa posizione del problema che verrà ripresa in modo più preciso più avanti, conduce alle soluzioni moderne del problema della stima dell'ordine.

IL CRITERIO FPE

Come abbiamo osservato **la bontà di un modello stimato non si può giudicare solo dall'accuratezza con cui esso esegue il *fit* dei dati usati per l'identificazione ma occorre in realtà valutare la bontà con cui il modello stimato riesce a descrivere dati *futuri*, non usati per l'identificazione del modello.**

Supponiamo allora di avere a disposizione due vettori di osservazioni $\mathbf{y} := [\mathbf{y}_1^\top \mathbf{y}_2^\top]^\top$ che per semplicità assumeremo di ugual dimensione N e di usare i primi N dati \mathbf{y}_1 per l'identificazione di un generico modello lineare standard di dimensione p . I dati \mathbf{y}_1 sono descritti dal modello lineare

$$\mathbf{y}_1 = S\boldsymbol{\theta} + \mathbf{w}_1, \quad \text{Var}[\mathbf{w}_1] = \sigma^2 I_N \quad (94)$$

ottenendo il classico stimatore $\hat{\boldsymbol{\theta}}(\mathbf{y}_1) = [S^\top S]^{-1} S^\top \mathbf{y}_1$.

Vogliamo ora valutare la “bontà statistica” del modello stimato, $S\hat{\boldsymbol{\theta}}(\mathbf{y}_1)$ per descrivere i dati \mathbf{y}_2 che abbiamo tenuto da parte.

Perchè questa operazione abbia senso dobbiamo supporre che i dati nei successivi N campioni siano stati *generati dallo stesso meccanismo che ha generato* \mathbf{y}_1 , il che si può esprimere dicendo che le d.d.p. (o almeno le statistiche del primo e secondo ordine) di \mathbf{y}_1 e \mathbf{y}_2 debbono essere le stesse.

Qui supporremo che le due componenti del vettore $[\mathbf{y}_1^\top \mathbf{y}_2^\top]^\top$ abbiano lo stesso vettore di media $\boldsymbol{\mu}$ (che potrebbe essere qualunque) e che la varianza complessiva di \mathbf{y} sia $\sigma^2 I_{2N}$. In questo modo \mathbf{y}_1 e \mathbf{y}_2 risultano scorrelati.

Consideriamo allora il vettore errore *finale* di predizione dei dati futuri

$$\boldsymbol{\varepsilon} := \mathbf{y}_2 - S\hat{\boldsymbol{\theta}}(\mathbf{y}_1) \quad (95)$$

che ha media $\boldsymbol{\mu} - S[S^\top S]^{-1}S^\top \boldsymbol{\mu}$ per cui sottraendo la media e calcolando la varianza di $\boldsymbol{\varepsilon}$ si trova

$$\text{Var}[\boldsymbol{\varepsilon}] = \sigma^2 I_N + S[S^\top S]^{-1}S^\top \sigma^2 I_N S[S^\top S]^{-1}S^\top = \sigma^2 \left[I_N + S[S^\top S]^{-1}S^\top \right].$$

Come misura dell'errore finale di predizione prendiamo la varianza scalare normalizzata che ha l'espressione

$$\begin{aligned} \frac{1}{N} \text{var}[\boldsymbol{\varepsilon}] &= \sigma^2 \frac{1}{N} \text{Tr} \{ I_N + S[S^\top S]^{-1} S^\top \} = \sigma^2 \left\{ 1 + \frac{1}{N} \text{Tr}([S^\top S]^{-1} S^\top S) \right\} \\ &= \sigma^2 \left(1 + \frac{p}{N} \right) \end{aligned} \quad (96)$$

dalla quale si vede che la varianza scalare dell'errore di predizione **cresce linearmente con p** . Per usare questo risultato per la stima della dimensione del modello, dobbiamo sostituire alla varianza σ^2 , che è un parametro incognito, una sua stima, naturalmente anch'essa basata su un modello a p parametri. Usando lo stimatore corretto della varianza

$$\frac{N}{N-p} \hat{\sigma}_p^2 = \frac{1}{N-p} \|\mathbf{y}_1 - S\hat{\boldsymbol{\theta}}(\mathbf{y}_1)\|^2 = \frac{1}{N-p} \|\hat{\boldsymbol{\varepsilon}}_p\|^2$$

dove $\hat{\boldsymbol{\varepsilon}}_p$ è il residuo di stima nel modello a p parametri, si arriva così a definire l'indice

$$FPE(p) := \frac{1}{N} \|\hat{\boldsymbol{\varepsilon}}_p\|^2 \frac{\left(1 + \frac{p}{N}\right)}{\left(1 - \frac{p}{N}\right)} := \hat{\sigma}_p^2 \frac{\left(1 + \frac{p}{N}\right)}{\left(1 - \frac{p}{N}\right)} \quad (97)$$

che viene anch'esso chiamato **errore finale di predizione** basato su un modello di dimensione p .

La stima dell'ordine del modello può essere basata sulla minimizzazione di questo indice. Naturalmente per effettuare la minimizzazione occorre preliminarmente identificare un certo numero di modelli di ordine crescente in un intervallo di valori plausibili di p e calcolare il relativo errore residuo quadratico medio. I calcoli si possono organizzare in modo efficiente usando algoritmi ricorsivi del tipo di quello illustrato nel seguente paragrafo.

IL CRITERIO FPE PER MODELLI DINAMICI

Supponiamo per il momento di avere una classe di modelli dinamici lineari a p parametri liberi e che lo stimatore PEM $\hat{\theta}_N$, del parametro θ soddisfi le ipotesi di consistenza e normalità asintotica. Abbiamo dati $\{\mathbf{y}^{N_1}, \mathbf{u}^{N_1}\}$ sull'intervallo $[1, N_1]$, con cui costruiamo lo stimatore PEM $\hat{\theta}_1$ e vogliamo descrivere col modello $M(\hat{\theta}_1)$ dei dati $\{\mathbf{y}^{N_2}, \mathbf{u}^{N_2}\}$ su un intervallo “lontano” $[t_0 + 1, t_0 + N_2]$.

Per l'ergodicità possiamo supporre che i due sets di dati siano approssimativamente **indipendenti**.

Calcoliamo la varianza dell'errore **finale** (asintotico) di predizione usando lo stimatore $\hat{\theta}_1$ supponendo $t_0 \rightarrow \infty$.

$$W_{N_1}(p) := \mathbb{E} \left[\boldsymbol{\varepsilon}^2(t, \hat{\theta}_1) \right], \quad t > t_0 \gg N_1$$

in cui $\boldsymbol{\varepsilon}^2(t, \hat{\theta}_1)$ dipende dai dati in $[1, N_1]$ attraverso lo stimatore $\hat{\theta}_1 = \hat{\theta}_1(\mathbf{y}^{N_1}, \mathbf{u}^{N_1})$ e dai dati $\{\mathbf{y}^{N_2}, \mathbf{u}^{N_2}\}$ usati per calcolare il predittore.

Per $N_1 \rightarrow \infty$ si ha $\hat{\boldsymbol{\theta}}_1 \rightarrow \boldsymbol{\theta}_0$ e si può approssimare

$$\begin{aligned} W_{N_1}(p) &\simeq \mathbb{E} \left[\boldsymbol{\varepsilon}(t, \boldsymbol{\theta}_0) + \frac{\partial \boldsymbol{\varepsilon}(t, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} (\hat{\boldsymbol{\theta}}_1 - \boldsymbol{\theta}_0) \right]^2 = \mathbb{E} \left[\mathbf{e}_0(t) + \boldsymbol{\psi}_{\boldsymbol{\theta}_0}^\top(t) (\hat{\boldsymbol{\theta}}_1 - \boldsymbol{\theta}_0) \right]^2 \\ &= \lambda_0^2 + \mathbb{E} \left\{ \text{Tr} \left[\boldsymbol{\psi}_{\boldsymbol{\theta}_0}(t) \boldsymbol{\psi}_{\boldsymbol{\theta}_0}(t)^\top (\hat{\boldsymbol{\theta}}_1 - \boldsymbol{\theta}_0) (\hat{\boldsymbol{\theta}}_1 - \boldsymbol{\theta}_0)^\top \right] \right\}. \end{aligned}$$

Scambiamo le operazioni di Traccia (che è lineare) e di aspettazione. Per l'indipendenza l'aspettazione è il prodotto di due aspettative. Quella rispetto ai dati in $[1, N_1]$, per $N_1 \rightarrow \infty$, è

$$\mathbb{E} (\hat{\boldsymbol{\theta}}_1 - \boldsymbol{\theta}_0) (\hat{\boldsymbol{\theta}}_1 - \boldsymbol{\theta}_0)^\top \simeq \frac{\lambda_0^2}{N_1} \left\{ \mathbb{E} \boldsymbol{\psi}_{\boldsymbol{\theta}_0}(t) \boldsymbol{\psi}_{\boldsymbol{\theta}_0}(t)^\top \right\}^{-1}$$

e quindi si trova

$$W_{N_1}(p) \simeq \lambda_0^2 \left(1 + \frac{p}{N_1} \right),$$

Prendendo come stimatore corretto di λ_0^2 , l'errore quadratico medio residuo

diviso per $N_1 - p$ si trova la formula asintotica (valida per $N_1 \rightarrow \infty$)

$$W_{N_1}(p) \simeq V_{N_1}(\hat{\boldsymbol{\theta}}_1) \frac{(1 + \frac{p}{N_1})}{(1 - \frac{p}{N_1})} \quad (FPE)$$

Si vede che al crescere di p l'errore quadratico medio di predizione $V_{N_1}(\hat{\boldsymbol{\theta}}_1)$ diminuisce ma il termine moltiplicativo aumenta (all'incirca linearmente) con p . Naturalmente per usare questo criterio di stima bisogna identificare una serie di modelli a diversa complessità p e andare poi a scegliere quello a minimo FPE.

IL CRITERIO AIC E L' MDL

Dato che normalmente $N \gg p$ si può approssimare

$$\frac{(1 + \frac{p}{N})}{(1 - \frac{p}{N})} \simeq (1 + \frac{2p}{N})$$

e prendendo i logaritmi, si trova il **critério AIC di Akaike**:

$$AIC := \log V_N(\hat{\boldsymbol{\theta}}_N) + \frac{2p}{N}.$$

Si dimostra che sia questo criterio che il FPE tendono a sovrastimare la dimensione p (non sono stimatori consistenti). Uno stimatore consistente di p è il valore di p che minimizza il **critério MDL (Minimum Description Length) di Rissanen**

$$MDL(p) := \log V_N(\hat{\boldsymbol{\theta}}_N) + \frac{\beta(p, N)}{N}$$

dove β è una funzione che asintoticamente cresce come $\log N$. Si prende normalmente $\beta(p, N) = p \log N$.

TESTS DI BIANCHEZZA DEI RESIDUI

La **struttura di un modello dinamico lineare** è una entità più articolata della dimensione p del parametro. Per un modello ARMAX a scatola nera ad es. la struttura è (n_a, n_b, n_c, n_k) , la terna di gradi dei polinomi A, B, C , più il ritardo con cui agisce l'ingresso.

Introducendo esplicitamente l'indice di struttura, v , un modello pilotato dal relativo errore di predizione, si può scrivere

$$\mathbf{y}(t) = F_{v,\theta}(z)\mathbf{u}(t) + G_{v,\theta}(z)\boldsymbol{\varepsilon}_{v,\theta}(t), \quad \theta \in \Theta_v. \quad (\dagger)$$

in questa rappresentazione $\boldsymbol{\varepsilon}_{v,\theta}(t)$ è il processo errore di predizione associato al predittore lineare a minima varianza costruito a partire dal modello.

Proposizione 5 *Se per qualche (v_0, θ_0) l'errore di predizione $\boldsymbol{\varepsilon}_{v_0, \theta_0}(t)$ è rumore bianco, allora il processo congiunto $\{\mathbf{y}, \mathbf{u}\}$ è descritto dal modello vero*

$$\mathbf{y}(t) = F_{v_0, \theta_0}(z)\mathbf{u}(t) + G_{v_0, \theta_0}(z)\mathbf{e}_0(t), \quad (98)$$

dove $\mathbf{e}_0(t) \equiv \boldsymbol{\varepsilon}_{v_0, \theta_0}(t)$. Viceversa, se il processo è descritto da un modello lineare (vero) di complessità finita come (\dagger) , allora l'errore di predizione del modello coincide con l'innovazione \mathbf{e}_0 .

Prova Dato che il modello d'innovazione, e in particolare, il processo d'innovazione della coppia $\{\mathbf{y}, \mathbf{u}\}$ sono unici, se $\boldsymbol{\varepsilon}_{v_0, \theta_0}(t)$ è rumore bianco, allora il modello (\dagger) corrispondente è il modello d'innovazione di $\{\mathbf{y}, \mathbf{u}\}$. Viceversa, se il processo è descritto dal modello vero (98), allora $\mathbf{e}_0(t)$ è anche l'errore di predizione corrispondente ai parametri v_0, θ_0 . \square

L' enunciato del teorema di consistenza della stima PEM si può interpretare nel presente contesto nel modo seguente.

Corollario 1 *Se il processo (vero) che genera i dati è ergodico del secondo ordine ed è descritto da un modello che appartiene alla classe parametrica \mathcal{M}_{v_0} e la classe parametrica dei modelli \mathcal{M}_{v_0} è identificabile localmente in $\theta = \theta_0$, allora lo stimatore PEM $\hat{\theta}_N$ è consistente e converge per $N \rightarrow \infty$, al parametro vero θ_0 con probabilità uno.*

Supponiamo ora che il modello $M_v(\theta)$ sia stato identificato con il metodo PEM sulla base di dati osservati di numerosità N e denotiamo con $\varepsilon_{v, \hat{\theta}_N}(t)$ l'errore residuo di predizione di un modello di struttura v basato su un campione di numerosità N .

Proposizione 6 Se $\hat{\theta}_N$ è uno stimatore consistente di θ per il modello (†) e se il limite, che esiste con probabilità uno,

$$\lim_{N \rightarrow \infty} \boldsymbol{\varepsilon}_{\mathbf{v}, \hat{\theta}_N}(t) := \boldsymbol{\varepsilon}_{\mathbf{v}, \theta_0}(t) \quad (99)$$

è rumore bianco, allora $\mathbf{v} \equiv \mathbf{v}_0$ è la struttura vera e $\boldsymbol{\varepsilon}_{\mathbf{v}, \theta_0}(t) = \mathbf{e}_0(t)$.

Prova Infatti, per la continuità dell'errore di predizione rispetto al parametro θ del modello, il processo limite $\boldsymbol{\varepsilon}_{\mathbf{v}, \theta_0}(t)$ è l'errore di predizione del modello $M_{\mathbf{v}}(\theta_0)$ e allora per concludere basta ricordare l'enunciato della proposizione 5. □

Quindi per verificare se il modello e la struttura identificati convergono al modello vero (e alla struttura vera) dobbiamo escogitare procedimenti per **verificare se il limite per $N \rightarrow \infty$ dell'errore residuo di predizione è un processo bianco.**

IL TEST DEL CORRELOGRAMMA

Questo è un test di bianchezza (asintotica) che usa la covarianza campionaria dell'errore residuo di predizione del modello stimato.

Denotando per semplicità $\boldsymbol{\varepsilon}_{\nu, \hat{\theta}_N}(t)$ come $\hat{\boldsymbol{\varepsilon}}(t)$, si costruisce la covarianza campionaria

$$\hat{\boldsymbol{\lambda}}(\tau) := \frac{1}{N} \sum_{t=\tau}^N \hat{\boldsymbol{\varepsilon}}(t) \hat{\boldsymbol{\varepsilon}}(t - \tau)$$

dove si prende $0 \leq \tau \leq \tau_{Max}$. Il valore di τ_{Max} è scelto opportunamente piccolo, tipicamente pari ad $1/20 \sim 1/50$ di N per evitare effetti ai bordi della varianza della stima, [?].

L'ipotesi H_0 da testare è che $\nu = \nu_0$. Naturalmente manteniamo ferme le ipotesi che assicurano la consistenza dello stimatore PEM (Corollario 1). Sotto H_0 si ha allora, con probabilità uno:

$$\begin{aligned} \lim_{N \rightarrow \infty} \hat{\boldsymbol{\lambda}}(\tau) &= 0, & \text{for } \tau \neq 0 \\ \lim_{N \rightarrow \infty} \hat{\boldsymbol{\lambda}}(0) &= \boldsymbol{\lambda}_0^2, & \text{for } \tau = 0. \end{aligned}$$

Proposizione 7 *Si consideri la statistica a valori vettoriali*

$$\hat{\mathbf{I}} := \frac{1}{N} \sum_{t=m}^N \hat{\boldsymbol{\varepsilon}}(t) \begin{bmatrix} \hat{\boldsymbol{\varepsilon}}(t-1) \\ \vdots \\ \hat{\boldsymbol{\varepsilon}}(t-m) \end{bmatrix} = \begin{bmatrix} \hat{\boldsymbol{\lambda}}(1) \\ \vdots \\ \hat{\boldsymbol{\lambda}}(m) \end{bmatrix} \quad (100)$$

e la sua versione normalizzata $\hat{\mathbf{r}} := \frac{1}{\hat{\boldsymbol{\lambda}}(0)} \hat{\mathbf{I}}$. Si ha

$$N \hat{\mathbf{r}}^\top \hat{\mathbf{r}} = N \frac{1}{\hat{\boldsymbol{\lambda}}(0)^2} \hat{\mathbf{I}}^\top \hat{\mathbf{I}} \xrightarrow{L} \chi^2(m) \quad (101)$$

Prova: Notiamo che per il modello vero i prodotti $\mathbf{e}_0(t)\mathbf{e}_0(t-k)$; $k = 1, \dots, m$ sono d-martingale stazionarie per cui vale il teorema del limite centrale,

$$\sqrt{N} \frac{1}{N} \sum_{t=m}^N \mathbf{e}_0(t) \begin{bmatrix} \mathbf{e}_0(t-1) \\ \vdots \\ \mathbf{e}_0(t-m) \end{bmatrix} \xrightarrow{L} \mathcal{N}(0, P).$$

Dove

$$P = \mathbb{E} \left\{ \mathbf{e}_0(t) \begin{bmatrix} \mathbf{e}_0(t-1) \\ \vdots \\ \mathbf{e}_0(t-m) \end{bmatrix} \mathbf{e}_0(t) \begin{bmatrix} \mathbf{e}_0(t-1) \\ \vdots \\ \mathbf{e}_0(t-m) \end{bmatrix}^\top \right\} = \lambda_0^4 I_m,$$

dato che

$$\mathbb{E} (\mathbf{e}_0(t) \mathbf{e}_0(t-k))^2 = \lambda_0^4, \quad \mathbb{E} (\mathbf{e}_0(t) \mathbf{e}_0(t-k)) (\mathbf{e}_0(s) \mathbf{e}_0(s-k)) = 0 \text{ per } t \neq s.$$

Sotto l'ipotesi H_0 si ha per $N \rightarrow \infty$,

$$\hat{\boldsymbol{\varepsilon}}(t) \hat{\boldsymbol{\varepsilon}}(t-k) \rightarrow \mathbf{e}_0(t) \mathbf{e}_0(t-k)$$

con probabilità uno. Si ha così

$$\sqrt{N} \left\{ \frac{1}{N} \sum_{t=m}^N \hat{\boldsymbol{\varepsilon}}(t) \begin{bmatrix} \hat{\boldsymbol{\varepsilon}}(t-1) \\ \vdots \\ \hat{\boldsymbol{\varepsilon}}(t-m) \end{bmatrix} - \frac{1}{N} \sum_{t=m}^N \mathbf{e}_0(t) \begin{bmatrix} \mathbf{e}_0(t-1) \\ \vdots \\ \mathbf{e}_0(t-m) \end{bmatrix} \right\} \rightarrow 0$$

in legge, dato che la differenza dentro parentesi tende a zero con probabilità uno e il secondo termine moltiplicato per \sqrt{N} converge in legge a $\mathcal{N}(0, \lambda_0^4 I_m)$, per cui lo stesso deve accadere anche al primo termine.

Per il primo enunciato del teorema di Slutsky, la statistica (100) converge allora in legge alla distribuzione Gaussiana $\mathcal{N}(0, \lambda_0^4 I_m)$.

Dato che λ_0^2 è incognita conviene normalizzare la (100), introducendo il vettore dei coefficienti di correlazione campionari $\hat{\mathbf{r}}(k) = \hat{\boldsymbol{\lambda}}(k)/\hat{\boldsymbol{\lambda}}(0)$. Usando ancora il TLC e il teorema di Slutsky si vede che per $N \rightarrow \infty$,

$$\sqrt{N}\hat{\mathbf{r}}(k) \xrightarrow{L} \mathcal{N}(0, 1) \quad k \neq 0$$

Si ha quindi anche

$$\sqrt{N}\hat{\mathbf{r}} := \frac{1}{\hat{\boldsymbol{\lambda}}(0)}\hat{\mathbf{I}} \xrightarrow{L} \mathcal{N}(0, I_m). \quad (102)$$

Questo risultato implica che

$$N\hat{\mathbf{r}}^\top \hat{\mathbf{r}} = N \frac{1}{\hat{\boldsymbol{\lambda}}(0)^2} \hat{\mathbf{I}}^\top \hat{\mathbf{I}} \xrightarrow{L} \chi^2(m)$$

che è quanto si voleva dimostrare. □

La statistica $N\hat{\mathbf{r}}^T\hat{\mathbf{r}}$ può quindi essere usata come statistica di test. Si rifiuta l'ipotesi che $\mathbf{v} = \mathbf{v}_0$ con probabilità d'errore di prima specie α se i valori della statistica sono più grandi di k_α dove $P_{\chi^2(m)}\{x \geq k_\alpha\} = \alpha$. Il valore di m si prende normalmente pari a $m = 5 \sim 10$.

C'è una statistica simile che si usa per verificare *incorrelazione di $\boldsymbol{\varepsilon}(t)$ dagli ingressi passati*

IL PERIODOGRAMMA

Una cosa ovvia per verificare se un segnale è una traiettoria di rumore bianco stazionario è calcolarne lo spettro e vedere se questo ha ampiezza costante su tutto l'intervallo di frequenze $[-\pi, \pi]$. Naturalmente con un segnale di durata finita ci si aspetta che questo sia vero solo approssimativamente.

Supponiamo di prendere T campioni successivi dell'errore residuo di predizione, $\{\hat{\varepsilon}_N(t); t = 1, \dots, T\}$ basato su un campione di numerosità N . In quel che segue è opportuno distinguere l'ampiezza T della finestra temporale in cui consideriamo il processo e la numerosità campionaria N . Anche se per N finito potremmo prendere $T = N$, nell'analisi asintotica conviene far **tendere la numerosità campionaria N all'infinito tenendo T fisso**. Calcoliamo la stima dello spettro che si ottiene prendendo la norma quadrato della trasformata *Finita* di Fourier (FT) della sequenza (il pedice N è omissso per semplicità di notazioni):

$$\hat{\phi}_T(\omega) = \frac{1}{T} \left| \sum_{t=1}^T \hat{\varepsilon}(t) e^{-j\omega t} \right|^2 \quad (103)$$

Questa stima si chiama **Periodogramma** della sequenza $\{\hat{\varepsilon}(t); t = 1, 2, \dots, T\}$.

Il periodogramma (103) è una statistica che dipende anche dalla numerosità campionaria N del campione con cui è costrita $\hat{\varepsilon}$. Per chiarezza conviene per il momento studiare il periodogramma di un arbitrario processo stazionario \mathbf{y} avente funzione di correlazione $r(k)$.

Il motivo della normalizzazione (divisione per T) si può far risalire formalmente alla definizione della FT, ma è forse più intuitivo basarsi sulla classica relazione tra periodogramma e *correlazione campionaria*, che è definita da

$$\hat{r}_T(k) = \frac{1}{T} \sum_{t=1}^{T-k} \mathbf{y}(t+k)\mathbf{y}(t), \quad k = 0, 1, \dots, T-1$$

e da $\hat{r}_T(k) = \hat{r}_T(-k)$ per $k < 0$. Notare che $\hat{r}_T(k)$ è una sequenza finita definita solo per $|k| \leq T-1$.

A rigore questo stimatore di $r(k)$ non è corretto ma quello che si otterrebbe dividendo per $T-k$ anzichè T oltre ad avere varianza maggiore (specie per k vicino ad T) non gode della proprietà notevole descritta qui sotto.

Proposizione 8 *Si ha*

$$\hat{\phi}_T(\omega) = \sum_{k=-(T-1)}^{+(T-1)} \hat{r}_T(k) e^{-j\omega k}. \quad (104)$$

Prova: Facendo la sostituzione $t = s + k$ nella somma

$$\hat{\phi}_T(\omega) = \frac{1}{T} \sum_{t=1}^T \sum_{s=1}^T \mathbf{y}(t) e^{-j\omega t} \mathbf{y}(s) e^{j\omega s}$$

l'indice $k = t - s$ va da $-T + 1$ a $T - 1$ e, dato che si può porre $\mathbf{y}(s + k) = 0$ se $s > T - k$, si trova

$$\hat{\phi}_T(\omega) = \frac{1}{T} \sum_{k=-T+1}^{T-1} e^{-j\omega k} \sum_{s=1}^{T-k} \mathbf{y}(s+k) \mathbf{y}(s) = \sum_{k=-T+1}^{T-1} e^{-j\omega k} \hat{r}_T(k)$$

COMPORAMENTO ASINTOTICO DI $\hat{\phi}_T(\omega)$

Con qualche passaggio si ricava,

$$\mathbb{E} \hat{r}_T(k) = \left(1 - \frac{|k|}{T}\right) r(k), \quad -(T-1) \leq k \leq T-1$$

e si trova,

$$\begin{aligned} \mathbb{E} [\hat{\phi}_T(\omega)] &= \sum_{k=-(T-1)}^{+(T-1)} \mathbb{E} \hat{r}_T(k) e^{-j\omega k} = \sum_{k=-(T-1)}^{+(T-1)} \left(1 - \frac{|k|}{T}\right) r(k) e^{-j\omega k} \\ &:= \sum_{k=-\infty}^{+\infty} w_T(k) r(k) e^{-j\omega k} \end{aligned}$$

dove la funzione triangolare w_T (definita su tutto \mathbb{Z}) è chiamata *finestra di Bartlett*. Se assumiamo che $r(k)$ ammetta trasformata di Fourier, $\phi(\omega)$ (densità spettrale), nel dominio della frequenza questa relazione si scrive

$$\mathbb{E} [\hat{\phi}_T(\omega)] = \int_{-\pi}^{\pi} \phi(\lambda) W_T(\omega - \lambda) \frac{d\lambda}{2\pi} \quad (\dagger)$$

La trasformata di Fourier ordinaria, $W_T(\omega)$, di w_T è

$$W_T(\omega) = \frac{1}{T} \left[\frac{\sin(\frac{\omega T}{2})}{\sin(\frac{\omega}{2})} \right]^2.$$

che è una funzione positiva il cui integrale in frequenza vale $w_T(0) = 1$ qualunque sia T . Affinchè $\mathbb{E}[\hat{\phi}_T(\omega)]$ sia il più vicino possibile a $\phi(\omega)$, $W_T(\omega)$ dovrebbe essere una buona approssimazione dell'impulso di Dirac. Di fatto, dato che

$$W_T(\omega) = T \left[\frac{\sin(\frac{\omega T}{2})}{\frac{\omega T}{2}} \frac{\frac{\omega}{2}}{\sin(\frac{\omega}{2})} \right]^2$$

si ha $\lim_{T \rightarrow \infty} W_T(0) = \infty$ e quindi la successione $\{W_T\}$ approssima la delta di Dirac centrata in $\omega = 0$, per cui

$$\lim_{T \rightarrow +\infty} \mathbb{E}[\hat{\phi}_T(\omega)] = \phi(\omega)$$

quindi, se si prende $T = T(N)$ tendente all' ∞ con N , il periodogramma è uno *stimatore asintoticamente corretto* della densità spettrale di y .

Questo non accade per piccoli valori di T poichè la larghezza a 3 dB del lobo principale di $W_T(\omega)$ è $\frac{2\pi}{T}$ ($\frac{1}{T}$ in frequenza).

Questo causa il problema dello *smearing* per cui due picchi in $\phi(\omega)$ separati in frequenza per meno di $\frac{1}{T}$ non sono distinguibili; per questo $\frac{1}{T}$ è detto *limite di risoluzione spettrale del periodogramma*.

Questo fatto d'altro canto è ovvio se si interpreta la trasformata (104) come una trasformata *finita*, nel qual caso la frequenza ω della trasformata dev'essere interpretata come una variabile *discreta*, ad esempio $\omega \equiv \omega_k = \frac{2\pi}{T}k; k = 0, 1, 2, \dots, T - 1$.

LA TRASFORMATA FINITA DI UN PROCESSO STOCASTICO

La FT finita del processo $\mathbf{y} := \{\mathbf{y}(t) ; t = 1, 2, \dots, T\}$ è

$$\mathbf{Y}(\omega) = \sum_{t=1}^T \mathbf{y}(t) e^{-j\omega t} \equiv \mathfrak{F}(\mathbf{y}), \quad (FT)$$

Ovviamente $\mathbf{Y}(\omega) = \mathbf{Y}(\omega + 2\pi)$ e quindi $\mathbf{Y}(\omega)$ è una funzione periodica della variabile angolare (continua) ω . Il complesso-coniugato $\bar{\mathbf{Y}}(\omega)$ è $= \mathbf{Y}(-\omega)$. Dato che \mathfrak{F} è un operatore lineare su uno spazio di dimensione T , la funzione $\mathbf{Y}(\omega)$ è determinata univocamente da T valori discreti dell'angolo ω .

Si fissano i valori equispaziati $\omega \equiv \omega_k = \frac{2\pi}{T}k; k = 0, 1, 2, \dots, T - 1$.

Proposizione 9 *Il periodogramma è il modulo quadrato della FT finita del processo \mathbf{y} , diviso per T :*

$$\phi_T(\omega_k) = \frac{1}{T} |\mathbf{Y}(\omega_k)|^2 = \frac{1}{T} \mathbf{Y}(\omega_k) \bar{\mathbf{Y}}(\omega_k) \quad (105)$$

Il periodogramma “campionario” $\hat{\phi}_T(\omega_k)$ definito prima è semplicemente il valore campionario della variabile aleatoria $\phi_T(\omega_k)$.

LA TRASFORMATA FINITA SIMMETRIZZATA

Si può dare una formula “simmetrica” per la FT come la

$$\hat{\mathbf{Y}}(\omega) = \sum_{t=-M}^M \mathbf{y}(t)e^{-j\omega t}, \quad (FT')$$

che richiede che \mathbf{y} sia definito su un numero dispari di campioni, ad esempio supponendo T dispari, oppure T pari e facendo partire la somma (FT) da $t = 0$. Nel primo caso si può definire il punto medio M dell'intervallo $[1, T]$ come $M = (T - 1)/2 + 1$. Ponendo $t = \tau + M$ nella (FT) si trova

$$\mathbf{Y}(\omega) = e^{-j\omega M} \sum_{\tau=-M+1}^{M-1} \mathbf{y}(\tau + M)e^{-j\omega\tau} = e^{-j\omega M} \sum_{\tau=-(T-1)/2}^{(T-1)/2} \mathbf{y}(\tau + M)e^{-j\omega\tau}.$$

Se T è pari e $t_0 = 0$, con $M = T/2$ si ha una formula simmetrica. Adesso però $\omega \equiv \omega_h = \frac{2\pi}{T+1}h$; $h = 0, \pm 1, \pm 2, \dots, \pm T/2$.

Con la definizione simmetrica, $\hat{\mathbf{Y}}(\omega_{-h}) = \hat{\mathbf{Y}}(-\omega_h) = \hat{\mathbf{Y}}^*(\omega_h)$.

STAZIONARIETÀ E PERIODICITÀ

Il processo \mathbf{y} è **stazionario** se la matrice di correlazione $\mathbf{R}_T := \mathbb{E} \mathbf{y} \mathbf{y}^\top$ di dimensione $T \times T$ (che noi supponiamo definita positiva) ha struttura di Toeplitz, ovvero

$$\mathbf{R}_T = \begin{bmatrix} r(0) & r(1) & \dots & r(T-1) \\ r(1) & r(0) & r(1) & \dots \\ \dots & \dots & \dots & r(1) \\ r(T-1) & \dots & r(1) & r(0) \end{bmatrix}, \quad r(k) = \mathbb{E} \mathbf{y}(t+k) \mathbf{y}(t).$$

Definizione 4 *Un processo \mathbf{z} definito su tutto l'asse dei tempi \mathbb{Z} è **periodico di periodo T** se $\mathbf{z}(k+nT) := \mathbf{z}(k)$ (cp1) per $n \in \mathbb{Z}$.*

*Diciamo che il processo \mathbf{y} definito su $[1, T]$ è **periodico** se è la restrizione all'intervallo $[1, T]$ di un processo periodico su tutto \mathbb{Z} . Se \mathbf{y} è periodico può essere esteso fuori dell'intervallo di definizione, ponendo $\mathbf{y}(0) = \mathbf{y}(T), \mathbf{y}(-1) = \mathbf{y}(T-1), \dots$ ottenendo un processo periodico su \mathbb{Z} .*

È ovvio che la funzione di correlazione di un processo periodico è periodica di periodo T , ovvero $r(k+nT) = r(k)$.

MATRICI COVARIANZA CIRCOLANTI

Però la periodicità di $r(k)$ dev'essere compatibile con la parità $r(-k) = r(k)$.
Si deve avere

$$r(T - k) = r(-k) = r(k) \quad k = 0, 1, \dots, T - 1$$

si trova che r ha due rami speculari rispetto al punto medio dell'intervallo $[1, T]$:

$$r(T - 1) = r(1), r(T - 2) = r(2), \dots, r(T/2 + 1) = r(T/2 - 1), \quad (P)$$

Se T è dispari $r(T/2)$ non è definita ma con lo stesso ragionamento si trova

$$r\left(T - \frac{T - 1}{2}\right) = r\left(\frac{T + 1}{2}\right) = r\left(\frac{T - 1}{2}\right)$$

e vale la stessa relazione di simmetria.

Una matrice di Toeplitz che soddisfa alla relazione di simmetria (P) è una **Matrice Circolante**.

ESTENSIONE A UN PROCESSO PERIODICO

Proposizione 10 y è periodico se e solo se $r(0), r(1), \dots, r(T-1)$ sono campioni di una funzione di correlazione periodica di periodo T .

Un processo su $[1, T]$ è periodico se e solo se la sua matrice di correlazione \mathbf{R}_T è una matrice circolante simmetrica e (semi) definita positiva.

NB. In genere i campioni di correlazione $r(0), r(1), \dots, r(T-1)$ di un processo finito y non formano una sequenza periodica nel senso appena definito. La matrice (di Toeplitz) \mathbf{R}_T in generale non è circolante.

Però si dimostra che una arbitraria stringa di correlazioni $r(0), r(1), \dots, r(T-1)$ (con \mathbf{R}_T simmetrica e definita positiva) ha sempre un'estensione periodica ma il periodo dell'estensione è in generale più grande di T .

L'estensione periodica è definita come funzione della stringa originale e non richiede altri dati per la sua costruzione.

Quindi esiste sempre un'estensione periodica (finita) di y

Noi in seguito assumeremo che il processo (finito) y sia periodico.

Come abbiamo accennato, esiste un teorema di estensione di una arbitraria stringa stazionaria lunga T con una periodica su un intervallo di ampiezza $T' > T$. La correlazione dell'estensione è nota in base ai dati (noti) di correlazione iniziali e lo spettro relativo è una approssimazione dello spettro discreto originario con un errore che tende a zero per $T \rightarrow \infty$. È abbastanza intuitivo che per $T \rightarrow \infty$ un arbitrario processo stazionario possa essere approssimato (in senso opportuno) con un processo periodico di periodo T .

Una possibile dimostrazione si rifà alla teoria delle approssimazioni circolanti di matrici di Toeplitz. Vedere ad es il libro di R. GRAY: *Toeplitz and Circulant matrices* che si può scaricare dalla rete

<http://ee.stanford.edu/~gray/toeplitz.html>.

Teorema 4 Se y è periodico di periodo T , ovvero la stringa delle correlazioni

$$r^T := [r(0), r(1), \dots, r(T-1)] \quad (R)$$

del processo y è periodica di periodo T , le ampiezze della FT finita di y sono *scorrelate a frequenze diverse*; i.e.

$$\frac{1}{T} \mathbb{E} \mathbf{Y}(\omega_k) \bar{\mathbf{Y}}(\omega_h) = \phi(\omega_k) \delta_{k,h}$$

dove $\{\phi(\omega_k)\}$ è la FT finita della correlazione (R); i.e. lo spettro discreto del processo.

Prova: Scende da

$$\mathbb{E} \mathbf{Y}(\omega_k) \bar{\mathbf{Y}}(\omega_h) = \sum_{t=1}^T \sum_{s=1}^T r(t-s) e^{-j\omega_k t} e^{j\omega_h s} = \sum_{t=1}^T e^{-j(\omega_k - \omega_h)t} \sum_{\tau=t-T}^{t-1} r(\tau) e^{-j\omega_h \tau}$$

e dal fatto che l'ultima somma non dipende da t perchè sia $r(\tau)$ che $e^{-j\omega_h \tau}$ formano sequenze periodiche di periodo T .

Infatti con la sostituzione $u = \tau - (t - T)$ si trova

$$\sum_{\tau=t-T}^{t-1} r(\tau)e^{-j\omega_h\tau} = \sum_{u=0}^{T-1} r(u+t)e^{-j\omega_h(u+t)};$$

sviluppando l'ultima somma con un $t > 0$ arbitrario e sostituendo i valori di $r(u+t)e^{-j\omega_h(u+t)}$ esterni all'intervallo $[0, T-1]$ con i corrispondenti $\text{mod } T$, ad esempio $r(T-1+t)e^{-j\omega_h(T-1+t)}$ con $r(t-1)e^{-j\omega_h(t-1)}$ etc. si vede bene che l'ultima somma è proprio uguale a $\sum_{u=0}^{T-1} r(u)e^{-j\omega_h u} = \Phi(\omega_h)$.

Per concludere basta notare che le sequenze esponenziali $\{e^{-j\omega_k t}\}$ e $\{e^{-j\omega_h t}\}$, $t = 1, 2, \dots, T$, formano vettori ortogonali in \mathbb{C}^T e

$$\sum_{t=1}^T e^{-j\omega_k t} e^{j\omega_h t} = T \delta_{k,h}.$$

□

PROPRIETÀ STATISTICHE DI $\hat{\phi}_T(\omega)$

Dal teorema precedente il periodogramma $\phi_T(\omega_k) = \frac{1}{T} \mathbf{Y}(\omega_k) \mathbf{Y}^*(\omega_k)$, $k = 0, 1, \dots, T-1$ ha aspettazione $\phi(\omega_k)$ e quindi è uno stimatore corretto. Supponiamo che il processo \mathbf{y} sia Gaussiano a media zero. Allora, per il teorema 4, anche le variabili $1/\sqrt{T} \mathbf{Y}(\omega_k)$ sono Gaussianhe a media zero (tra loro indipendenti) e di varianza $\phi(\omega_k)$. Il loro modulo quadrato normalizzato ha quindi una distribuzione χ^2 ad un grado di libertà:

$$\left\{ \frac{\frac{1}{T} |\mathbf{Y}(\omega_k)|^2}{\phi(\omega_k)} \right\} \sim \chi^2(1)$$

Da cui scende la relazione

$$\boxed{\text{var}\{\hat{\phi}_T(\omega_k)\} = 2\phi(\omega_k)^2} \quad (106)$$

che vale qualunque siano T ed N . La relazione vale nel caso ideale di processi periodici ma per T sufficientemente grande vale per processi qualunque.

Teorema 5 Se per $N \rightarrow \infty$ y ha distribuzione asintotica Gaussiana e $T \rightarrow \infty$ con N , la varianza asintotica di $\hat{\phi}_T(\omega)$ è data dalla relazione seguente:

$$\lim_{N \rightarrow +\infty} E [\hat{\phi}_T(\omega) - \phi(\omega)]^2 = \begin{cases} 2\phi^2(\omega) & \omega = 0 \text{ oppure } \pm \pi \\ \phi^2(\omega) & 0 < \omega < \pi \end{cases} \quad (107)$$

mentre per $\omega_1 \neq \omega_2$ si ha

$$\lim_{N \rightarrow +\infty} E \left\{ [\hat{\phi}_T(\omega_1) - \phi(\omega_1)][\hat{\phi}_T(\omega_2) - \phi(\omega_2)] \right\} = 0 \quad (108)$$

Da questo risultato ricaviamo che i valori del periodogramma $\hat{\phi}_T(\omega)$ sono variabili aleatorie le cui deviazioni standard sono (anche per $N \rightarrow \infty$!) uguali ai corrispondenti valori dello spettro. Perciò il periodogramma è uno **stimatore spettrale inconsistente** (in media quadratica e in probabilità e anche con probabilità uno) che oscilla attorno al suo valore medio asintotico (il valore vero $\phi(\omega)$) con una varianza che si mantiene molto elevata anche per $N \rightarrow +\infty$. Inoltre le variabili aleatorie $\{\hat{\phi}_T(\omega); 0 \leq \omega \leq \pi\}$ sono asintoticamente scorrelate. In sostanza il periodogramma si comporta asintoticamente come un *rumore bianco* in frequenza.

Un fatto fondamentale dell'analisi statistica del periodogramma è che per $N \rightarrow \infty$ le deviazioni attorno alla media di $\hat{\phi}_T(\omega)$ a frequenze diverse diventano scorrelate. In sostanza, $\hat{\phi}_T(\omega) - \phi(\omega)$ diventa un **rumore bianco in frequenza**.

Notare che stiamo discutendo del comportamento asintotico per $N \rightarrow \infty$ ma con T fisso. Al tendere $T \rightarrow \infty$ le frequenze angolari ω_k si addensano sul cerchio unita e questo "rumore bianco limite" dovrebbe dipendere da una variabile continua (ω). Come è noto questo ipotetico oggetto a variabili continue scorrelate non può essere un processo stocastico nel senso standard del termine.

Di fatto, l'asintotica per $T \rightarrow \infty$ è molto sfuggente da maneggiare anche per altri motivi. Per esempio il limite della (FT') per $T \rightarrow \infty$:

$$\lim_{T \rightarrow \infty} \sum_{-M}^M \mathbf{y}(t) e^{-j\omega t}$$

dovrebbe dare la trasformata di Fourier ordinaria del processo \mathbf{y} . Sfortunatamente però questo limite non esiste nè in media quadratica nè con probabilità uno.

In realtà si può definire una rappresentazione di Fourier generalizzata di un processo stazionario su \mathbb{Z} (chiamata *rappresentazione spettrale*) che gode di una proprietà simile alla bianchezza dello spettro. L'interpretazione esatta di questa frase richiede però di entrare in dettagli matematici che non possiamo affrontare in questa sede. Come vedremo, si può invece mostrare con mezzi elementari che l'incorrelazione a diverse frequenze vale esattamente (e non solo asintoticamente) per il periodogramma costruito con i T campioni successivi di un processo (finito) periodico di periodo T .

Nella figura è riportato lo spettro vero (in nero) del processo descritto dal modello ARMA

$$\mathbf{y}(t) - 1.5\mathbf{y}(t - 1) + 0.7\mathbf{y}(t - 2) = \mathbf{e}(t) - \mathbf{e}(t - 1) + 0.2\mathbf{e}(t - 2) \quad (109)$$

mentre in blu è riportato il periodogramma con $T = 64$ e in verde quello con $T = 512$.

Spettro vero e periodogramma per l'esempio (109)

IL TEST DEL PERIODOGRAMMA CUMULATO

L'idea del test è di calcolare l'integrale discreto del periodogramma e confrontarlo con l'andamento ideale dell'integrale dello spettro di un rumore bianco. Come è noto l'operazione di integrazione riduce il rumore nei segnali, in particolare nella stima del periodogramma integrato.

Si assume che il numero di dati T sia un numero pari. Usando l'algoritmo FFT si calcola il periodogramma della sequenza $\{\hat{\varepsilon}(t); t = 1, \dots, T\}$ per $\omega_k = 2\pi \frac{k}{T} := 2\pi f_k, k = 1, 2, \dots, T$.

Il periodogramma cumulato è l'integrale indefinito discreto sull'intervallo $[0, \pi]$, con passo di discretizzazione $\Delta\omega = 2\pi/T$

$$\hat{I}(\omega_k) = 2 \sum_{i=0}^k \hat{\phi}_T(\omega_i) \frac{\Delta\omega}{2\pi}, \quad k = 1, 2, \dots, T - 1$$

che si può riscrivere in termini della frequenza normalizzata $f_k = \frac{k}{T}$

$$\hat{I}(f_k) = 2 \frac{1}{T} \sum_{i=0}^k \hat{\phi}_T(f_i) \quad k = 1, 2, \dots, T-1$$

Dalla correttezza asintotica di $\hat{\phi}_T$ segue che per $T \rightarrow \infty$

$$\mathbb{E} \hat{I}(\omega_k) \simeq 2 \frac{1}{T} \sum_{i=0}^k \phi(f_i), \quad k = 1, 2, \dots, T-1$$

per cui se $\hat{\epsilon}$ è rumore bianco (**ipotesi** H_0) si ha $\phi(f_i) = \lambda_0^2$ e

$$\frac{\mathbb{E} \hat{I}(f_k)}{\lambda_0^2} = 2f_k, \quad k = 1, 2, \dots, T-1$$

che è un segmento di retta per l'origine di pendenza 2.

Nella statistica del test si introduce la stima della varianza

$$\hat{\lambda}_T^2 = \sum_{i=0}^{T-1} \hat{\phi}_T(f_i)$$

e si tiene conto che la distribuzione asintotica della statistica è di tipo χ^2 .

STIMATORI PARAMETRICI DELLO SPETTRO

Gli stimatori parametrici di spettro sono **stimatori di densità spettrale** e assumono che il processo che genera i dati sia p.n.d.. Si riducono all'identificazione di modelli AR, e quindi sono semplici e facili da calcolare.

Supponiamo di avere stime di $n + 1$ campioni successivi della funzione di correlazione o di covarianza di un processo stazionario

$$\hat{r}(0), \hat{r}(1), \dots, \hat{r}(n); \quad \text{Toepl} \{ \hat{r}(0), \hat{r}(1), \dots, \hat{r}(n) \} > 0.$$

Vogliamo trovare una stima dello spettro del processo, che sia la più generale possibile, nel senso che non richieda altre informazioni sul processo oltre all'assegnazione dei suoi primi $n + 1$ momenti secondi.

Per risolvere questo problema ci si può rifare al **principio della massima entropia** (Burg). Si massimizza *l'entropia rate* di un processo stazionario di densità spettrale $\phi(\omega)$

$$H(\phi) = \frac{1}{2} \int_{-\pi}^{\pi} \log \det \phi(e^{j\omega}) \frac{d\omega}{2\pi} + \frac{m}{2} \log 2\pi e \quad ,$$

che ha i primi $n + 1$ campioni di covarianza assegnati.

Si dimostra che il processo stazionario con i momenti assegnati a massima entropia, **è un processo autoregressivo di ordine n** .

La densità spettrale (o lo stimatore spettrale) a massima entropia si ottiene costruendo il modello AR che è individuato dalle $n + 1$ stime di covarianza assegnate. Questo stimatore si può calcolare usando la matrice delle covarianze stimate (definita positiva), risolvendo un sistema di equazioni lineari, cosiddetto di *Yule-Walker*.

Se N è piccolo si ricorre alla cosiddetta tecnica di Burg, che allo stato è un procedimento semi-empirico che però dà in genere buoni risultati. Questo argomento è esposto nei libri ROBERTS MULLIS, SÖDERSTROM, etc.

STIMATORI DI FREQUENZE BASATI SULLA CORRELAZIONE CAMPIONARIA

Come abbiamo visto, l'ottimizzazione per la minimizzazione dell'errore quadratico medio di predizione deve essere effettuata con algoritmi iterativi di discesa che portano inevitabilmente a *minimi locali*. Per questo motivo sono stati proposti in letteratura dei metodi che non sono basati sulla minimizzazione di una cifra di merito ma sono classificabili come varianti dei *metodi dei momenti* e sono sostanzialmente basati su operazioni di fattorizzazione (eventualmente di tipo SVD) della matrice correlazione (o covarianza) campionaria del segnale di uscita. Questi metodi hanno il grosso vantaggio di essere robusti e numericamente stabili e di fornire il risultato in “un colpo”, senza bisogno di iterazioni.

Espressione matriciale della correlazione Vogliamo dare una formula semplice e compatta per la funzione di correlazione di un processo p.d. descritto da un modello di stato del tipo (1). Allo scopo possiamo anche pensare che $\mathbf{z}(t)$ sia a valori vettoriali (ad esempio in \mathbb{R}^m) e si ottenga mediante

una matrice stato-uscita generale $C \in \mathbb{R}^{m \times n}$; ($n = 2\nu$) invece di una matrice riga come c^\top in (1). Dalla rappresentazione di stato si ottiene facilmente la seguente espressione,

$$R(\tau) := \mathbb{E} \mathbf{z}(t + \tau) \mathbf{z}(t)^\top = \mathbb{E} \mathbf{z}(\tau) \mathbf{z}(0)^\top = CA^\tau PC^\top \quad \tau \geq 0 \quad (110)$$

dove P è la matrice varianza dello stato iniziale: $P := \mathbb{E} \mathbf{s}(0) \mathbf{s}(0)^\top$. Nella particolare base in cui è rappresentato il modello (1), P è diagonale a blocchi con blocchi diagonali 2×2 del tipo

$$P_k = \begin{bmatrix} \sigma_k^2 & 0 \\ 0 & \sigma_k^2 \end{bmatrix} \quad k = 1, 2, \dots, n.$$

Espressioni esplicite per la matrice di covarianza di vettori costruiti con traslazioni temporali di $\mathbf{y}(t)$. Da (1) si ottiene

$$\begin{aligned} \mathbf{y}^m &:= \begin{bmatrix} \mathbf{y}(t) \\ \mathbf{y}(t+1) \\ \vdots \\ \mathbf{y}(t+m-1) \end{bmatrix} = \begin{bmatrix} c^\top \\ c^\top A \\ \vdots \\ c^\top A^{m-1} \end{bmatrix} \mathbf{s}(t) + \begin{bmatrix} \mathbf{e}(t) \\ \mathbf{e}(t+1) \\ \vdots \\ \mathbf{e}(t+m-1) \end{bmatrix} \\ &:= \Omega_m \mathbf{s}(t) + \mathbf{e}^m \end{aligned} \quad (111)$$

dove Ω_m (Ω quando non c'è pericolo di confusione) è la *matrice di osservabilità* del sistema (1) che ha la seguente struttura,

$$\Omega_m = \begin{bmatrix} 1 & 0 & \dots & 1 & 0 \\ \cos \omega_1 & -\sin \omega_1 & \dots & \cos \omega_v & -\sin \omega_v \\ \dots & & & & \dots \\ \cos(m-1)\omega_1 & -\sin(m-1)\omega_1 & \dots & \cos(m-1)\omega_v & -\sin(m-1)\omega_v \end{bmatrix} \quad (112)$$

Nell'ipotesi che le frequenze ω_k siano tutte diverse tra loro, se $m \geq 2v$ la matrice ha rango $n = 2v$ e il sistema è osservabile. Calcolando la covarianza del vettore \mathbf{y}^m , nell'ipotesi che \mathbf{e} sia rumore bianco, si trova

$$R := E\{\mathbf{y}^m (\mathbf{y}^m)^\top\} = \Omega P \Omega^\top + \sigma^2 I_{2v} \quad (113)$$

dove

$$P = E\{\mathbf{s}(t)\mathbf{s}(t)^\top\} = \text{diag}\{P_1, \dots, P_V\}$$

in cui P_k è la varianza di stato del modello elementare di indice k , che ha la forma

$$P_k = \begin{bmatrix} \sigma_k^2 & 0 \\ 0 & \sigma_k^2 \end{bmatrix}$$

Si può verificare che questa matrice risolve l'equazione di Lyapunov $P_k = A_k P_k A_k^\top$.

Nel seguito useremo anche una formula per la covarianza di due vettori, costruiti analogamente a \mathbf{y}^m in (111) ma relativi a due intervalli disgiunti.

Tenedo conto del fatto che $\mathbf{s}(t-k) = (A^\top)^k \mathbf{s}(t)$, si trova

$$\Gamma := E \left\{ \begin{bmatrix} \mathbf{y}(t+1) \\ \mathbf{y}(t+2) \\ \vdots \\ \mathbf{y}(t+m) \end{bmatrix} [\mathbf{y}(t) \quad \mathbf{y}(t-1) \quad \dots \quad \mathbf{y}(t-l)] \right\} \quad (114)$$

$$= \begin{bmatrix} r(1) & r(2) & \dots & r(l+1) \\ r(2) & r(3) & \dots & r(l+2) \\ \dots & \dots & \dots & \dots \\ r(m) & r(m+1) & \dots & r(l+m) \end{bmatrix} \quad (115)$$

$$= E \{ \Omega_m \mathbf{s}(t) \mathbf{s}(t)^\top \Omega_l^\top \} = \Omega_m P \Omega_l^\top \quad (116)$$

dove

$$\Omega_m := \begin{bmatrix} c^\top A \\ \vdots \\ c^\top A^m \end{bmatrix}, \quad \Omega_l^\top := [c \quad Ac \quad \dots \quad A^l c].$$

Se m e l sono maggiori di $n = 2v$ entrambe queste matrici hanno rango n e quindi anche Γ ha rango n .

Γ è una matrice di Hankel, similmente a quanto accade per le matrici di

covarianza. costruite con le correlazioni dell'uscita di un sistema stocastico lineare che rappresenta un processo p.n.d. Di fatto questa struttura sta alla base dei metodi cosiddetti "a sottospazi" che verranno illustrati in seguito.

RASSEGNA DI ALCUNI METODI STANDARD

Passeremo in rassegna alcuni di questi metodi che sono sempre basati sull'ipotesi di osservare la somma di un segnale quasi periodico con sovrapposto del rumore bianco. Se il rumore additivo non è bianco si incorre in difficoltà e sembra che il metodo d'elezione da usare in queste circostanze sia l' *identificazione a sottospazi* che studieremo nel capitolo ???. I metodi di stima che passeremo in rassegna più sotto si possono vedere come rudimentali precursori dell'identificazione a sottospazi.

YULE WALKER STESO

Il metodo di Yule-Walker esteso, per il quale si usa l'acronimo inglese HOYW (High-Order Yule-Walker), si basa sul seguente ragionamento. Dato che in genere non si conosce il numero di componenti armoniche del segnale, si può sempre supporre che esso sia descritto da un modello ARMA del tipo (4), di ordine l sufficientemente elevato, superiore all'ordine $n = 2v$ del modello minimo "vero":

$$\mathbf{y}(t) + b_1\mathbf{y}(t-1) + \dots + b_l\mathbf{y}(t-l) = \mathbf{e}(t) + b_1\mathbf{e}(t-1) + \dots + b_l\mathbf{e}(t-l)$$

che si può scrivere anche come

$$B(z^{-1})\mathbf{y}(t) = B(z^{-1})\mathbf{e}(t) \quad (117)$$

dove il polinomio $B(z^{-1})$ si può pensare ottenuto moltiplicando $A(z^{-1})$ in (4) per un certo altro polinomio $\bar{A}(z^{-1})$ di grado pari a $(l-n)$.

Riscriviamo l'equazione (??) nella seguente forma piú concisa

$$\begin{bmatrix} \mathbf{y}(t) & \mathbf{y}(t-1) & \dots & \mathbf{y}(t-l) \end{bmatrix} \begin{bmatrix} 1 \\ b \\ \vdots \\ b_l \end{bmatrix} = \mathbf{e}(t) + \dots + b_l\mathbf{e}(t-l)$$

premultiplichiamo per $[\mathbf{y}(t-l-1) \dots \mathbf{y}(t-l-m)]^T$, dove m é un intero positivo che sar  specificato in seguito e calcoliamo le matrici covarianza a primo e secondo membro. Tenedo conto del fatto che per $k > 0$ si ha $E[\mathbf{y}(t-k)\mathbf{e}(t)] = 0$ e sostituendo $E[\mathbf{y}(t-k)\mathbf{y}(t-j)] = r(j-k) = r(k-j)$, si ottiene la:

$$\begin{bmatrix} r(l+1) & r(l) & \dots & r(1) \\ r(l+2) & r(l+1) & \dots & r(2) \\ \dots & & & \dots \\ r(l+m) & r(l+m-1) & \dots & r(m) \end{bmatrix} \begin{bmatrix} 1 \\ b \end{bmatrix} := \Gamma^c \begin{bmatrix} 1 \\ b \end{bmatrix} = 0. \quad (118)$$

dove la matrice Γ^c é la matrice Γ definita in (115) con le colonne in ordine opposto. La (118) scritta in forma scalare, é simile al classico sistema di equazioni di Yule-Walker [?, pp. 288]:

$$r(k) + \sum_{i=1}^l b_i r(k-i) = 0, \quad k = l+1, \dots, l+m.$$

Notiamo che scambiando l'ordine dei coefficienti, i.e. ponendo $\bar{b}_k := b_{l-k}$, il sistema di equazioni lineari (118) si potrebbe scrivere nella forma equiv-

alente

$$\Gamma \begin{bmatrix} \bar{b} \\ 1 \end{bmatrix} = 0.$$

In ogni caso il vettore dei parametri soddisfa il sistema di equazioni:

$$\begin{bmatrix} r(l) & \dots & r(1) \\ \vdots & \dots & \vdots \\ r(l+m-1) & \dots & r(m) \end{bmatrix} b = - \begin{bmatrix} r(l+1) \\ \vdots \\ r(l+m) \end{bmatrix} \quad (119)$$

Il metodo di HOYW per stimare le frequenze del segnale in esame utilizza il risultato precedente sostituendo però le covarianze teoriche $\{r(k)\}$ con le covarianze campionarie $\{\hat{r}(k)\}_{k=1}^{l+m}$ ricavate a partire dai campioni del segnale a disposizione.

Ovviamente a causa degli errori di stima in $\{\hat{r}(k)\}$ invece di una uguaglianza si ha in realtà solo una relazione approssimata del tipo:

$$\begin{bmatrix} \hat{r}(l) & \dots & \hat{r}(1) \\ \vdots & \dots & \vdots \\ \hat{r}(l+m-1) & \dots & \hat{r}(m) \end{bmatrix} \hat{b} \approx - \begin{bmatrix} \hat{r}(l+1) \\ \vdots \\ \hat{r}(l+m) \end{bmatrix} \quad (120)$$

Il passo fondamentale del metodo consiste nel risolvere (120) in \hat{b} con il metodo ai minimi quadrati. A questo proposito bisogna però fare alcune precisazioni.

Denotiamo con $\hat{\Gamma}$ la matrice $m \times l$ delle covarianze campionarie in (120). Anche se per $m, l \geq n$, il rango di Γ è in teoria uguale ad n , quello di $\hat{\Gamma}$ sarà sempre pieno (uguale al minimo tra m ed l).

Dato che $\text{rango} \hat{\Gamma} \simeq n$ il sistema (120) è mal condizionato dal punto di vista numerico. Infatti si può dimostrare che ogni metodo ai minimi quadrati che stima \hat{b} direttamente da (120) ha una scarsa accuratezza. Per far fronte a queste difficoltà si usa la *Decomposizione ai Valori Singolari (SVD)* della matrice $\hat{\Gamma}$.

Sia

$$\hat{\Gamma} = U \Sigma V^{\top} = [U_1 \quad U_2] \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{bmatrix} \begin{bmatrix} V_1^{\top} \\ V_2^{\top} \end{bmatrix} \quad (121)$$

la SVD di $\hat{\Gamma}$, dove U, V sono matrici ortogonali di dimensioni rispettivamente $m \times m$ e $l \times l$ e Σ é una matrice diagonale $m \times l$ a elementi positivi (i *valori singolari*) ordinati in senso decrescente.

Si dimostra che la matrice $\hat{\Gamma}_n$, ottenuta scartando in (121) la sottomatrice Σ_2 che contiene i valori singolari a indice maggiore di n ,

$$\hat{\Gamma}_n := U_1 \Sigma_1 V_1^\top$$

é la migliore approssimazione di rango n di $\hat{\Gamma}$ in una varietà di possibili norme. Usando $\hat{\Gamma}_n$ in (120) al posto di $\hat{\Gamma}$ si ottiene il sistema di equazioni di HOYW di rango troncato:

$$\hat{\Gamma}_n \hat{b} \approx - \begin{bmatrix} \hat{r}(l+1) \\ \vdots \\ \hat{r}(l+m) \end{bmatrix} \quad (122)$$

che può essere risolta con un metodo ai minimi quadrati ottenendo:

$$\hat{b} = -V_1 \Sigma_1^{-1} U_1^\top \begin{bmatrix} \hat{r}(l+1) \\ \vdots \\ \hat{r}(l+m) \end{bmatrix} \quad (123)$$

dove $V_1 \Sigma_1^{-1} U_1^\top$ é la pseudoinversa di $\hat{\Gamma}_n$. Una volta ottenuta la stima di \hat{b} , si considera il polinomio

$$1 + \sum_{k=1}^l \hat{b}_k z^{-k}$$

e le stime delle frequenze $\{\hat{\omega}_k\}$ del segnale si fanno coincidere con le posizioni angolari delle n radici del polinomio che si trovano piú vicine al cerchio di raggio unitario. Si assume cosí che le "radici del segnale", ovvero le radici di $A(z)$, siano sempre piú vicine al cerchio di raggio unitario delle "radici del rumore" o di $\bar{A}(z)$.

Per la stazionarietà, quando $N \rightarrow \infty$ tutte le radici di $B(z)$ debbono trovarsi all'interno del cerchio unitario (chiuso) ma quando si ha a disposizione un numero di campioni finito questa propriet  non pu  sempre essere garantita e di conseguenza il metodo HOYW produce delle stime di frequenza distorte. Questo é un problema comune a tutti i metodi che stimano le frequenze dalle radici di un polinomio di grado superiore a n , come vedremo in seguito.

Per concludere discutiamo brevemente la scelta dei parametri m e l . In pratica é consigliabile usare $l \approx m$ e scegliere i valori di questi interi in modo in modo che $l + m$ sia molto minore del numero di campioni osservati.

ANALISI ASINTOTICA

Supponiamo che il processo $\mathbf{y}(t)$ si possa esprimere come l'uscita di un modello di stato:

$$\mathbf{x}(t+1) = A\mathbf{x}(t) \quad (124)$$

$$\mathbf{y}(t) = c\mathbf{x}(t) + \mathbf{e}(t) \quad (125)$$

in cui possiamo anche supporre che \mathbf{e} sia un processo p.n.d. generale, scorrelato da \mathbf{x} . Ovviamente $c\mathbf{x}$ è la componente p.d. di \mathbf{y} . Come già visto nei paragrafi precedenti, si può supporre che la matrice A abbia una struttura diagonale a blocchi $A = \text{diag}\{A_1, \dots, A_v\}$ in cui i blocchi A_k sono matrici oscillatorie di dimensione 2×2 , mentre

$$c = [c_1^\top \quad c_2^\top \quad \dots \quad c_v^\top]$$

dove ciascun blocco riga c_k ha la forma $c_k^\top = [1 \quad 0]$. Si può quindi esprimere $\mathbf{y}(t)$ nel seguente modo:

$$\mathbf{y}(t) = \sum_{k=1}^v c_k^\top \mathbf{x}_k(t) + \mathbf{e}(t) = \sum_{k=1}^v c_k^\top A_k^t \mathbf{x}_k(0) + \mathbf{e}(t)$$

con condizioni iniziali aleatorie $\mathbf{x}_k(0)$, $k = 1, 2, \dots, v$ tra loro scorrelate. Consideriamo ora per semplicità di notazione solo la componente k -sima della parte p.d.:

$$\mathbf{y}_k(t) := \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} \cos \omega_k t & -\sin \omega_k t \\ \sin \omega_k t & \cos \omega_k t \end{bmatrix} \begin{bmatrix} \mathbf{x}_{k1}(0) \\ \mathbf{x}_{k2}(0) \end{bmatrix} \quad (126)$$

$$= \mathbf{x}_{k1}(0) \cos \omega_k t - \mathbf{x}_{k2}(0) \sin \omega_k t \quad (127)$$

$$:= \mathbf{a}_k \sin(\omega_k t + \varphi_k) \quad (128)$$

la componente QP del segnale si può esprimere come somma di v componenti sinusoidali ciascuna delle quali con ampiezza \mathbf{a}_k , frequenza ω_k e fase φ_k . L'ampiezza aleatoria \mathbf{a}_k è legata al vettore delle condizioni iniziali $\mathbf{x}_k(0)$ dalla relazione:

$$\mathbf{a}_k^2 = \mathbf{x}_{k1}^2(0) + \mathbf{x}_{k2}^2(0).$$

Supponiamo di osservare una singola traiettoria del processo \mathbf{y} e studiamo il limite della covarianza campionaria quando la numerosità campionaria T tende all'infinito. Come si è già visto ([?, p. 108]) il limite per $T \rightarrow \infty$ della

correlazione campionaria di una somma di segnali sinusoidali (deterministici) del tipo $y_k(t) = a_k \sin(\omega_k t + \varphi_k)$, $k = 1, 2, \dots, v$ con $\omega_k \neq \omega_j$, é :

$$\hat{r}(\tau) = \sum_{k=1}^v \frac{a_k^2}{2} \cos \omega_k \tau \quad (129)$$

mentre la correlazione “vera” è

$$r(\tau) = \sum_{k=1}^v \sigma_k^2 \cos \omega_k \tau \quad \sigma_k^2 = 1/2 E a_k^2 = E x_{k1}^2(0) = E x_{k2}^2(0). \quad (130)$$

Inoltre, si può dimostrare che la correlazione campionaria tra un segnale sinusoidale e una traiettoria $\{e(t)\}$ di un processo p.n.d. ergodico a media nulla è sempre uguale a zero

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{t=1}^N y_k(t + \tau) e(t) = 0$$

(si veda [?]) per cui il limite della covarianza campionaria del segnale complessivo y ha la forma

$$\hat{r}(\tau) = \sum_{k=1}^v \frac{a_k^2}{2} \cos \omega_k \tau + r_e(\tau) \quad (131)$$

dove, per l'ergodicità di e , $r_e(\tau) = E e(t + \tau) e(t)^\top$ è la covarianza “vera” di e .

Notiamo adesso che considerando ancora una volta solamente il blocco k -simo si può rappresentare (il limite del-) la correlazione campionaria come

$$\hat{r}_k(\tau) = \frac{a_k^2}{2} \cos \omega_k \tau = \quad (132)$$

$$= \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} \cos \omega_k \tau & \sin \omega_k \tau \\ -\sin \omega_k \tau & \cos \omega_k \tau \end{bmatrix} \begin{bmatrix} \frac{x_{k1}^2 + x_{k2}^2}{2} & 0 \\ 0 & \frac{x_{k1}^2 + x_{k2}^2}{2} \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad (133)$$

$$= c_k A_k^\tau \hat{P}_k c_k^\top \quad (134)$$

Nel caso generale di segnali che sono somma di v componenti sinusoidali e di una componente p.n.d. la correlazione campionaria ha pertanto la

seguinte espressione

$$\hat{r}(\tau) = \sum_{k=1}^n \frac{a_k^2}{2} \cos \omega_k \tau + r_e(\tau) = \quad (135)$$

$$= cA^\tau \hat{P} c^\top + r_e(\tau) \quad (136)$$

in cui le matrici A e P hanno strutture diagonali a blocchi. Si ottiene così un'espressione formalmente analoga alla (129) in cui però la matrice P è sostituita da una matrice diagonale a blocchi \hat{P} che dipende dalla particolare traiettoria del segnale. Notiamo che \hat{P} è simmetrica e non singolare con probabilità uno.

PISARENKO E MUSIC

Il metodo MUSIC (*MUltiple Signal Classification*) e quello di Pisarenko, che é un caso speciale del primo, come sar a spiegato in seguito, si basano sul modello di covarianza introdotto in (113) con $m > 2\nu$ e in particolare sulla matrice di covarianza R che per comodit a riscriviamo qui sotto:

$$R = \Omega P \Omega^\top + \sigma^2 I \quad (137)$$

Poich e la matrice $\Omega P \Omega^\top$ ha rango 2ν , essa possiede 2ν autovalori, che denoteremo $\{\tilde{\lambda}_k, k = 1, 2, \dots, 2\nu\}$, strettamente positivi e i rimanenti $(m - 2\nu)$ tutti uguali a zero. Denotiamo con $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$ gli autovalori di R e siano inoltre $\{s_1, \dots, s_{2\nu}\}$ gli *autovettori ortonormali* associati a $\{\lambda_1, \dots, \lambda_{2\nu}\}$ e $\{g_1, \dots, g_{m-2\nu}\}$ quelli corrispondenti a $\{\lambda_{2\nu+1}, \dots, \lambda_m\}$.

Vale il seguente utile risultato.

Lemma 6 *Gli autovalori di R sono dati dalla relazione*

$$\lambda_k = \tilde{\lambda}_k + \sigma^2 \quad (k = 1, \dots, m)$$

dove $\{\tilde{\lambda}_k\}_{k=1}^m$ sono gli autovalori di $\Omega P \Omega^\top$ listati in ordine non crescente. L'insieme degli autovalori di R può così essere suddiviso in due sottoinsiemi:

$$\begin{cases} \lambda_k > \sigma^2 & k = 1, \dots, 2\nu \\ \lambda_k = \sigma^2 & k = 2\nu + 1, \dots, m \end{cases} \quad (138)$$

Dimostrazione. Per il teorema spettrale esiste una matrice ortogonale (di autovettori) $T \in \mathbb{R}^{m \times m}$ tale che

$$\Omega P \Omega^\top = T \tilde{\Lambda} T^{-1}$$

dove $\tilde{\Lambda} := \text{diag}\{\tilde{\lambda}_1, \tilde{\lambda}_2, \dots, \tilde{\lambda}_m\}$. Il risultato scende dall'osservazione che il termine $\sigma^2 I$ è diagonale nella stessa base, dato che si può scrivere $\sigma^2 I = T \sigma^2 I T^{-1}$. \square

Gli autovettori associati a ognuno di questi sottoinsiemi possiedono alcune proprietà, qui di seguito riportate, che vengono usate nell'algoritmo per la stima di frequenze. Siano

$$S = [s_1, \dots, s_{2\nu}] \quad (m \times n), \quad G = [g_1, \dots, g_{m-2\nu}] \quad (m \times (m - 2\nu)) \quad (139)$$

le matrici formate dagli autovettori ortonormali associati agli autovalori di R . Dalla definizione di R e da (138) si ottiene:

$$RG = G \begin{bmatrix} \lambda_{2\nu+1} & & 0 \\ & \dots & \\ 0 & & \lambda_m \end{bmatrix} = \sigma^2 G = \Omega P \Omega^\top G + \sigma^2 G \quad (140)$$

L'ultima uguaglianza implica che $\Omega P \Omega^\top G = 0$ ovvero, dato che la matrice AP ha rango di colonna pieno,

$$\Omega^\top G = 0 \quad (141)$$

cioé le colonne $\{g_k\}$ di G appartengono allo spazio nullo di Ω^\top .

L'equazione (141) sta all'abasse dei metodi MUSIC. Scritta esplicitamente

ha il seguente aspetto

$$\begin{bmatrix} 1 & \cos \omega_1 & \cos 2\omega_1 & \dots & \cos(m-1)\omega_1 \\ 0 & -\sin \omega_1 & \sin 2\omega_1 & \dots & \sin(m-1)\omega_1 \\ \dots & & & & \dots \\ \dots & & & & \dots \\ 1 & \cos \omega_n & \cos 2\omega_n & \dots & \cos(m-1)\omega_n \\ 0 & -\sin \omega_n & \sin 2\omega_n & \dots & \sin(m-1)\omega_n \end{bmatrix} \begin{bmatrix} g_1 & \dots & g_{m-2\nu} \end{bmatrix} = 0$$

Notiamo adesso che questo sistema consiste di n blocchi di coppie di equazioni reali della forma

$$\begin{bmatrix} 1 & \cos \omega_k & \cos 2\omega_k & \dots & \cos(m-1)\omega_k \\ 0 & -\sin \omega_k & \sin 2\omega_k & \dots & \sin(m-1)\omega_k \end{bmatrix} g_j = 0 \quad k = 1, \dots, n \quad j = 1, \dots, m-2\nu$$

che si possono scrivere in forma complessa moltiplicando a sinistra per $[1 \ i]$ nella forma,

$$\begin{bmatrix} 1 & e^{i\omega_k} & e^{i2\omega_k} & \dots & e^{i(m-1)\omega_k} \end{bmatrix} g_j = 0 \quad k = 1, \dots, n \quad j = 1, \dots, m-2\nu \quad (142)$$

che si può interpretare dicendo che

Proposizione 11 *Gli $m - 2\nu$ polinomi di grado $m - 1$*

$$a_j(z) := [1 \quad z \quad z^2 \quad \dots \quad z^{m-1}] g_j = 0 \quad j = 1, \dots, m - 2\nu$$

si annullano tutti nei punti $z = e^{\pm i\omega_k}$.

Quindi le frequenze incognite $\pm\omega_k$ si possono (in teoria) trovare calcolando gli zeri di modulo unitario di ciascun polinomio $a_j(z)$. Il fatto che si possa costruire un numero arbitrario ($m - 2\nu$) di polinomi $a_j(z)$ può essere usato per migliorare la stima delle frequenze incognite. Il metodo di Pisarenko [?] che è stato il primo metodo di questo tipo proposto in letteratura inizialmente costruiva solo un vettore g e quindi un solo polinomio $a(z)$.

In realtà, dato che

- la covarianza R dovrà essere stimata in base ai dati osservati e quindi sarà sempre affetta da rumore e quindi di rango pieno,

- le equazioni (142) sono quindi solo delle uguaglianze approssimate,
- il modello “vero” del segnale potrebbe essere più complesso di quello ipotizzato; in particolare potrebbe contenere rumore additivo arbitrario (non bianco) il che comporta autovalori del “rumore” $\{\lambda_{2\nu+1}, \dots, \lambda_m\}$ tra loro diversi,
- il modello “vero” potrebbe poi avere un numero più grande di frequenze nella componente p.d. e quindi potrebbero esserci più zeri (approssimativamente) a modulo unitario degli n del modello ipotizzato,

Queste difficoltà richiedono una trattazione statistica del problema. In particolare è necessaria una regola per decidere quali sono gli autovalori “piccoli” ($\{\lambda_{2\nu+1}, \dots, \lambda_m\}$) e approssimativamente uguali alla varianza σ^2 .

È in particolare necessaria una trattazione statistica del sistema di equazioni $a_j(z) = 0$, $j = 1, 2, \dots, m - 2\nu$. Questo problema è risolto in letteratura con

vari “trucchi”, i più noti dei quali, *pseudospectrum MUSIC* e *Root MUSIC*, vengono descritti qui sotto.

Il metodo dello Pseudospectrum determina le stime di frequenza considerando la posizione dei ν picchi più alti della funzione:

$$\frac{1}{a(z)^T g g^T a(z)}$$

sul cerchio unitario $\{|z| = 1\}$. Il metodo RootMusic determina invece le stime di frequenza come le posizioni angolari delle ν coppie di radici più vicine al cerchio di raggio unitario dell'equazione

$$a(z)^T g g^T a(z) = 0.$$

IL METODO ESPRIT

Questo metodo è basato su una proprietà della matrice di osservabilità di un sistema lineare che in letteratura è chiamata qualche volta *shift-invariance*. Ritroveremo questa idea più in dettaglio quando studieremo i metodi a sottospazi.

Sia

$$\Omega := \begin{bmatrix} C \\ CA \\ \vdots \\ CA^m \end{bmatrix}$$

la matrice di osservabilità “estesa” ($m \geq n$) di un sistema lineare di ordine n

$$\begin{cases} \mathbf{x}(t+1) = A\mathbf{x}(t) \\ \mathbf{y}(t) = C\mathbf{x}(t) \end{cases}$$

che assumeremo osservabile. Introducendo le matrici “traslate”

$$\uparrow\Omega := \begin{bmatrix} C \\ CA \\ \vdots \\ CA^{m-1} \end{bmatrix} \quad \downarrow\Omega := \begin{bmatrix} CA \\ CA^2 \\ \vdots \\ CA^m \end{bmatrix} \quad (143)$$

evidentemente si ha

$$(\downarrow\Omega) = (\uparrow\Omega)A \quad (144)$$

e, dato che $\uparrow\Omega$ ha le colonne linearmente indipendenti (per l’osservabilità), esiste un’unica matrice $A \in \mathbb{R}^{n \times n}$ che soddisfa l’equazione (144). In altri termini, la matrice A è *univocamente determinata dalla matrice di osservabilità estesa del sistema*. Su questa idea è basato il metodo cosiddetto ESPRIT per la stima delle N frequenze incognite del segnale.

Consideriamo il nostro segnale (??) con rumore additivo bianco, la cui dinamica è descritta dal sistema lineare osservabile (1) di ordine $2v$. Un modo numericamente affidabile per calcolare una base di vettori per lo

spazio immagine della matrice di osservabilità è di ricondursi all'espressione (113) della covarianza di \mathbf{y}^m .

La matrice di autovettori S definita in (139) è una base per lo spazio immagine di Ω , in formule

$$\text{Im}\{S\} = \text{Im}\{\Omega\} \quad (145)$$

Dimostrazione. Abbiamo sostanzialmente già visto che $\Omega^\top G = 0$ e quindi le colonne della matrice G sono una base per lo spazio nullo di Ω^\top , ovvero $\text{Ker}\{\Omega^\top\} = \text{Im}\{G\}$ e quindi, dato che per una arbitraria matrice reale M vale la $\text{Ker}\{M^\top\} = \text{Im}\{M\}^\perp$, prendendo il complemento ortogonale in \mathbb{R}^m , si trova

$$\text{Im}\{\Omega\} = \left((\text{Im}\{\Omega\})^\perp \right)^\perp = (\text{Im}\{G\})^\perp = \text{Im}\{S\}$$

dato che $\text{Im}\{G\} = \text{Im}\{S\}^\perp$. □

Notiamo adesso che, data l'uguaglianza (145), deve esistere una matrice invertibile $T \in \mathbb{R}^{n \times n}$ tale che $S = \Omega T$. Ne viene che la relazione (144) vale

anche per le matrici traslate $\uparrow S$ e $\downarrow S$ definite in modo analogo alle $\downarrow \Omega$, $\uparrow \Omega$. Insomma,

$$(\downarrow S) = (\uparrow S)\hat{A} \quad (146)$$

dove $\hat{A} := T^{-1}AT$.

In conclusione, risolvendo il sistema (146)(che è in generale sovradeterminato) si può stimare una matrice simile alla A e quindi ricavare le frequenze incognite dagli autovalori, che in teoria dovrebbero stare tutti sul cerchio unitario. Un possibile metodo di soluzione è mediante i minimi quadrati. La soluzione calcolata con le equazioni normali (di solo uso concettuale) è

$$\hat{A} = \left[(\uparrow S)^\top (\uparrow S) \right]^{-1} (\uparrow S)^\top (\downarrow S).$$

Una volta stimata A , le frequenze di oscillazione si trovano prendendo gli n autovalori di \hat{A} più vicini al cerchio unitario.

Bontà statistica dei metodi a correlazione

In generale si può affermare che se la correlazione campionaria converge con probabilità uno a quella vera (il che accade in particolare con segnali ergodici del secondo ordine) le stime coi metodi basati sulla correlazione campionaria convergono ai parametri “veri”, quelli che soddisfano le relative relazioni limite che coinvolgono la covarianza vera. Sotto queste ipotesi i metodi basati sulla correlazione campionaria producono quindi in generale stimatori consistenti.

In realtà i segnali quasi periodici non sono ergodici (nemmeno del secondo ordine) e la correlazione limite per $N \rightarrow \infty$ dipende dall'ampiezza delle componenti periodiche. La linearità dei modelli fa sì che i parametri “veri” del modello in gioco, al limite soddisfino delle relazioni lineari che sono le stesse del caso di segnali ergodici. Quindi, anche se la covarianza limite è a stretto rigore aleatoria la relazione lineare è la stessa e il ragionamento continua a valere.

Viceversa è molto più difficile dare delle espressioni utili per la distribuzione e la varianza asintotica di questi stimatori. La letteratura (specialmente quella ingegneristica) sorvola su questi punti.

Limite di covarianze campionarie per segnali p.d. L'analisi dei metodi di stima che si usano in teoria dell'identificazione sono in genere basati sulla *ergodicità del secondo ordine* dei segnali in gioco. Questa proprietà si esprime dicendo che il limite delle covarianze campionarie

$$\hat{R}(\tau) := \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{t=1}^N y(t + \tau) y^\top(t)$$

calcolato per una singola traiettoria y del processo \mathbf{y} , è uguale all'aspettazione $R(\tau) = E\mathbf{y}(t + \tau)\mathbf{y}(t)^\top$, indipendentemente dalla traiettoria scelta. Sfortunatamente questa proprietà *non vale* per processi p.d.. Quindi il presupposto (comunemente dato per scontato in letteratura) che si possa sostituire la covarianza vera con il limite di quella campionaria anche quando sono in gioco segnali con componente p.d., non è affatto ovvio e va giustificato.