

Imad domande

Silviu Filote

July 2023

Contents

1	Lezione 4: Stima a massima verosimiglianza	1
2	Lezione 5: Regressione logistica	4
3	Lezione 6: Fondamenti di machine learning	7
4	Lezione 7: Fondamenti di stima Bayesiana	16
5	Lezione 8: Processi stocastici	20
6	Lezione 9: Famiglie di modelli stocastici	23
7	Lezione 10: Predizione	24
8	Lezione 11: Identificazione: concetti fondamentali	26
9	Lezione 12: Identificazione analisi e complementi	29
10	Lezione 13: Identificazione - valutazione del modello	34

1 Lezione 4: Stima a massima verosimiglianza

Cos'è la massima verosimiglianza e calcolarla

Il metodo della massima verosimiglianza (MLE – Maximum Likelihood Estimation) è una procedura di stima che, **dato un modello probabilistico**, stima i suoi parametri in modo tale che siano più coerenti con i dati osservati.

Supponiamo di avere a disposizione N osservazioni $y(i) \sim \mathcal{N}(\mu, \sigma^2)$ iid e $Y = [y(1), y(2), \dots, y(N)]^T$.

La pdf congiunta dei dati o la probability density function:

Indica la probabilità che si realizzi il vettore di dati osservato

$$\begin{aligned} f_Y(y(1), y(2), \dots, y(N) | \mu, \sigma^2) &= f_Y(Y | \mu, \sigma^2) = \prod_{i=1}^N \mathcal{N}(y(i) | \mu, \sigma^2) \\ &= f_y(y(1) | \mu, \sigma^2) \cdot \dots \cdot f_y(y(N) | \mu, \sigma^2) \end{aligned}$$

Quando la **pdf multivariabile** $f_Y(Y | \mu, \sigma^2)$ è vista in funzione dei parametri μ e σ è chiamata **funzione di likelihood** $\mathcal{L}_Y(\mu, \sigma^2 | Y)$, cambia solo interpretazione, ma sono lo stesso oggetto matematico, dunque $\theta_{ML} = [\mu, \sigma^2]$.

$$f(Y | \mu, \sigma^2) = \mathcal{L}(\mu, \sigma^2 | Y) = \prod_{i=1}^N \mathcal{N}(y(i) | \mu, \sigma^2)$$

La stima a massima verosimiglianza può essere espressa come:

$$\hat{\theta}_{ML} = \begin{bmatrix} \hat{\mu} \\ \hat{\sigma}^2 \end{bmatrix} = \arg \max_{\theta} \mathcal{L}(\theta | Y) = \arg \max_{\theta} \prod_{i=1}^N \mathcal{N}(y(i) | \mu, \sigma^2)$$

In generale posso attribuire ai dati qualsiasi distribuzione di probabilità $d(\cdot)$, sia continua che discreta:

$$\hat{\theta}_{ML} = \arg \max_{\theta} \mathcal{L}(\theta | Y)$$

Lo stimatore a massima verosimiglianza gode di buone proprietà:

- **Asintoticamente corretto:** $\lim_{N \rightarrow +\infty} \mathbb{E}[\hat{\theta}_{ML}] = \theta^0$
- **Consistente:** più N è grande, più la stima è precisa
- **Asintoticamente efficiente:** $\lim_{N \rightarrow +\infty} \text{Var}[\hat{\theta}_{ML}] = M^{-1}$
- **Asintoticamente normale:** $\hat{\theta}_{ML} \sim \mathcal{N}(\theta^0, M^{-1})$ per $N \rightarrow +\infty$

Note:

- Spesso, anziché massimizzare $\mathcal{L}(\theta | Y)$, si massimizza il suo logaritmo naturale $\ln \mathcal{L}(\theta | Y)$
- Dato che il logaritmo è una funzione monotona crescente, $\ln(\mathcal{L}(\theta | Y))$ ha lo stesso massimo di $\mathcal{L}(\theta | Y)$
- Usare il logaritmo è efficiente da un punto di vista implementativo, perchè evita possibili underflow dati dal prodotto di piccole probabilità (sostituendolo con la somma delle log-probabilità)
- A meno di casi particolari fortunati, l'ottimizzazione è effettuata con metodi iterativi

$$\hat{\theta}_{ML} = \arg \max_{\theta} \ln \mathcal{L}(\theta | Y)$$

NB: massimizzare la log-verosimiglianza equivale a **minimizzare la meno log-verosimiglianza:**

$$\hat{\theta}_{ML} = \arg \max_{\theta} \ln [\mathcal{L}(\theta | Y)] = \arg \min_{\theta} -\ln [\mathcal{L}(\theta | Y)]$$

Formulando il problema in questo modo abbiamo un problema di minimizzazione proprio come con lo **stimatore a minimi quadrati**: $\hat{\theta}_{LS} = \arg \min_{\theta} J(\theta)$

Stima a massima verosimiglianza di un modello lineare

Imponiamo un modello probabilistico alle osservazioni $y(i)$:

$$y(i) = \theta_0 + \theta_1 \varphi_1(i) + \dots + \theta_{d-1} \varphi_{d-1}(i) + \epsilon(i) = \boldsymbol{\varphi}^T(i) \cdot \boldsymbol{\theta} + \epsilon(i) \quad \text{modello lineare}$$

dove $\boldsymbol{\varphi}$ é il vettore delle features e $\boldsymbol{\theta}$ é il vettore dei parametri

$$\boldsymbol{\varphi}_{d \times 1} = \begin{bmatrix} 1 \\ \varphi_1 \\ \vdots \\ \varphi_{d-1} \end{bmatrix} \quad \boldsymbol{\theta}_{d \times 1} = \begin{bmatrix} 1 \\ \theta \\ \vdots \\ \theta_{d-1} \end{bmatrix}$$

In particolare, se assumiamo che $\epsilon \sim \mathcal{N}(0, \lambda^2)$ iid

$$\Rightarrow y(i) \sim \mathcal{N}(\boldsymbol{\varphi}^T(i) \cdot \boldsymbol{\theta}, \lambda^2) \quad \text{iid}$$
$$\mu(i) = \boldsymbol{\varphi}^T(i) \cdot \boldsymbol{\theta}$$

La media $\mu(i)$ di $y(i)$ é espressa come combinazione lineare dei regressori

La **distribuzione congiunta** dei dati è:

$$\begin{aligned} f_Y(y(1), y(2), \dots, y(N) | X, \boldsymbol{\theta}, \lambda^2) &= \prod_{i=1}^N f_y(y(i) | \boldsymbol{\varphi}(i), \boldsymbol{\theta}, \lambda^2) \\ &= \prod_{i=1}^N \mathcal{N}(y(i) | \boldsymbol{\varphi}(i), \boldsymbol{\theta}, \lambda^2) \\ &= \prod_{i=1}^N \frac{1}{\sqrt{2\pi\lambda^2}} \cdot \exp\left[-\frac{1}{2} \left(\frac{y(i) - \boldsymbol{\varphi}^T(i) \cdot \boldsymbol{\theta}}{\lambda}\right)^2\right] \\ &= \prod_{i=1}^N \frac{1}{\sqrt{2\pi\lambda^2}} \cdot \exp\left[-\frac{1}{2} \left(\frac{y(i) - \mu(i)}{\lambda}\right)^2\right] \\ &= \mathcal{L}(\boldsymbol{\theta} | Y, X, \lambda^2) \quad \text{supponiamo noto } \lambda^2 \end{aligned}$$

Calcolo la **log-verosimiglianza**:

$$\begin{aligned} \ln[\mathcal{L}(\boldsymbol{\theta} | Y, X, \lambda^2)] &= \ln\left[\prod_{i=1}^N \frac{1}{\sqrt{2\pi\lambda^2}} \cdot \exp\left[-\frac{1}{2} \left(\frac{y(i) - \boldsymbol{\varphi}^T(i) \cdot \boldsymbol{\theta}}{\lambda}\right)^2\right]\right] \\ &= \sum_{i=1}^N \ln\left[\frac{1}{\sqrt{2\pi\lambda^2}} \cdot \exp\left[-\frac{1}{2} \left(\frac{y(i) - \boldsymbol{\varphi}^T(i) \cdot \boldsymbol{\theta}}{\lambda}\right)^2\right]\right] \\ &= \sum_{i=1}^N \ln\left(\frac{1}{\sqrt{2\pi\lambda^2}}\right) + \sum_{i=1}^N \ln\left[\exp\left[-\frac{1}{2} \left(\frac{y(i) - \boldsymbol{\varphi}^T(i) \cdot \boldsymbol{\theta}}{\lambda}\right)^2\right]\right] \\ &= N \cdot \ln\left(\frac{1}{\sqrt{2\pi\lambda^2}}\right) + \sum_{i=1}^N -\frac{1}{2} \cdot \left(\frac{y(i) - \boldsymbol{\varphi}^T(i) \cdot \boldsymbol{\theta}}{\lambda}\right)^2 \\ &= N \cdot \ln(2\pi\lambda^2)^{-\frac{1}{2}} - \frac{1}{2\lambda^2} \sum_{i=1}^N (y(i) - \boldsymbol{\varphi}^T(i) \cdot \boldsymbol{\theta})^2 \\ &= -\frac{1}{2} N \cdot \ln(2\pi\lambda^2) - \frac{1}{2\lambda^2} \sum_{i=1}^N (y(i) - \boldsymbol{\varphi}^T(i) \cdot \boldsymbol{\theta})^2 \end{aligned}$$

Calcolare il massimo di $\ln[\mathcal{L}(\boldsymbol{\theta} | Y, X, \lambda^2)]$ è equivalente a calcolare il minimo di $-\ln[\mathcal{L}(\boldsymbol{\theta} | Y, X, \lambda^2)]$:

Siccome non dipende da $\boldsymbol{\theta}$, questo termine non contribuisce al calcolo del minimo

$$\begin{aligned}
-\ln[\mathcal{L}(\boldsymbol{\theta}|Y, X, \lambda^2)] &= \frac{1}{2}N \cdot \ln(2\pi\lambda^2) + \frac{1}{2\lambda^2} \sum_{i=1}^N (y(i) - \boldsymbol{\varphi}^T(i) \cdot \boldsymbol{\theta})^2 \\
&= \frac{1}{2\lambda^2} \sum_{i=1}^N (y(i) - \boldsymbol{\varphi}^T(i) \cdot \boldsymbol{\theta})^2 \\
\hat{\boldsymbol{\theta}}_{ML} &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2\lambda^2} \sum_{i=1}^N (y(i) - \boldsymbol{\varphi}^T(i) \cdot \boldsymbol{\theta})^2
\end{aligned}$$

La stima ML del modello lineare $y(i) = \boldsymbol{\varphi}^T(i) \cdot \boldsymbol{\theta} + \epsilon(i)$, con $\epsilon \sim \mathcal{N}(0, \lambda^2)$ iid, é **equivalente alla stima a minimi quadrati** (che non aveva assunzioni sulla pdf dei dati).

Sia dato il modello lineare seguente:

$$y(i) = \boldsymbol{\varphi}(i)^T \cdot \boldsymbol{\theta} + \epsilon(i)$$

dove:

- $y(i) \in \mathbb{R}^{1 \times 1}$ è la variabile risposta per l'osservazione i
- $\boldsymbol{\varphi}(i) \in \mathbb{R}^{d \times 1}$ il vettore delle features
- $\boldsymbol{\theta} \in \mathbb{R}^{d \times 1}$ è il vettore dei parametri ed $\epsilon(i) \approx \mathcal{N}(0, \lambda^2)$
- $\epsilon(i) \in \mathbb{R}^{1 \times 1}$ è un rumore gaussiano *iid* ed è indipendente da $\boldsymbol{\varphi}(i)$
- $y(i) \approx \mathcal{N}(\boldsymbol{\varphi}(i)^T \cdot \boldsymbol{\theta}, \lambda^2)$ dove $\boldsymbol{\varphi}(i)^T \cdot \boldsymbol{\theta} = \mu(i)$

Ricavare la stima a Massima Verosimiglianza del vettore dei parametri $\boldsymbol{\theta}$, supponendo di aver osservato N valori della variabile y .

Vedi domanda precedente.

2 Lezione 5: Regressione logistica

Descrivere il modello di regressione logistica indicando:

- la tipologia di dati che modella
- funzione di costo da minimizzare
- interpretazione della funzione di costo
- quali algoritmi si possono usare per la minimizzazione della funzione di costo

Tipologia di dati

Il modello di regressione logistica viene applicato quando la variabile di risposta risulta essere **qualitativa** e non quantitativa. I dati di tipo categorico descrivono una categoria di appartenenza, non hanno un ordinamento e non hanno alcuna metrica di distanza. Il processo di stima di output categorici, utilizzando un insieme di regressori φ , è chiamato **classificazione**. Tramite classificazione la categoria più probabile viene scelta come **classe** (categoria) per l'osservazione φ .

Supponiamo di avere a disposizione un dataset $\mathcal{D} = \{\{\varphi(1), y(1)\}, \dots, \{\varphi(N), y(N)\}\}$, dove $\varphi \in \mathbb{R}^{d \times 1}$ e $y(i) \in \{0, 1\}$, con $i = 1, \dots, N$ iid (y variabile categorica).

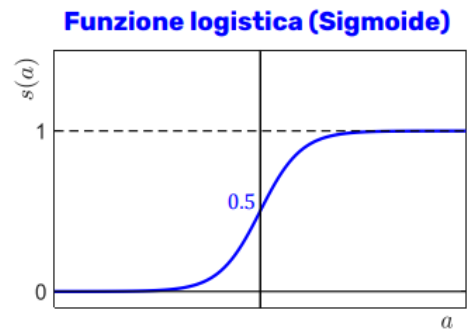
Obiettivo: Stimare la probabilità che le osservazioni $\varphi \in \mathbb{R}^{d \times 1}$ appartengano ad una di due classi $y \in \{0, 1\}$

Definiamo la combinazione lineare:

$$a = \sum_{j=0}^{d-1} \varphi_j \cdot \theta_j = \varphi^T \cdot \theta \in \mathbb{R}$$

La formula $s(a)$ è la **funzione logistica**:

$$s(a) = \frac{1}{1 + e^{-a}} = \frac{e^a}{a + e^a} = \begin{cases} a \gg 0 \Rightarrow s(a) \approx 1 \\ a \ll 0 \Rightarrow s(a) \approx 0 \end{cases}$$



Il modello di regressione logistica modella la probabilità che $y = 1$ tramite un modello lineare:

$$P(y = 1|\varphi) = s(a) = s(\varphi^T \cdot \theta) = \frac{1}{1 + e^{-\varphi^T \cdot \theta}} \equiv \pi$$

L'output di $s(a) = s(\varphi^T \theta)$ è interpretato come una probabilità:

$$\begin{cases} s(\varphi^T \theta) \gg 0.5 \Rightarrow P(y = 1|\varphi) \approx 1 & \varphi \text{ é classificato nella classe } 1 \\ s(\varphi^T \theta) \ll 0.5 \Rightarrow P(y = 1|\varphi) \approx 0 & \varphi \text{ é classificato nella classe } 0 \end{cases}$$

Interpretiamo i dati come realizzazioni di una distribuzione di Bernoulli:

$$y \sim \text{Bernoulli}(\pi) = \pi^y \cdot (1 - \pi)^{1-y}$$

Una volta che il modello stima la probabilità di una classe, possiamo classificare un punto φ in una particolare classe se la probabilità per quella classe è **superiore a una soglia** (di solito 0.5).

La funzione che stiamo stimando è:

$$f(\varphi) = P(y = 1|\varphi)$$

La regressione logistica tenta di modellare f con:

$$s(\varphi^T \theta) = \frac{1}{1 + e^{-\varphi^T \theta}}$$

Il punto φ può quindi essere classificato alla classe $y = 1$ se:

$$s(\varphi^T \cdot \theta) = \frac{1}{1 + e^{-\varphi^T \cdot \theta}} \geq 0.5$$

Funzione di costo da minimizzare

La probabilità che $\varphi(i) \in \mathbf{1}$, ossia $y(i) = 1$:

$$\pi(i) \equiv P(y(i) = 1 | \varphi(i)) = s(\varphi(i)^T \cdot \theta) = \frac{1}{1 + e^{-\varphi(i)^T \cdot \theta}}$$

Calcoliamo la **verosimiglianza**:

$$\mathcal{L}(\pi|Y) = \prod_{i=1}^N \pi(i)^{y(i)} \cdot (1 - \pi(i))^{1-y(i)}$$

Calcoliamo la **funzione di costo**:

$$\begin{aligned} -\ln[\mathcal{L}(\pi|Y)] &= -\ln \left[\prod_{i=1}^N \pi(i)^{y(i)} \cdot (1 - \pi(i))^{1-y(i)} \right] \\ &= -\sum_{i=1}^N \ln \left[\pi(i)^{y(i)} \cdot (1 - \pi(i))^{1-y(i)} \right] \\ &= -\sum_{i=1}^N \left(\ln \left[\pi(i)^{y(i)} \right] + \ln \left[(1 - \pi(i))^{1-y(i)} \right] \right) \\ &= -\sum_{i=1}^N \left(y(i) \cdot \ln[\pi(i)] + (1 - y(i)) \cdot \ln[1 - \pi(i)] \right) \equiv J(\theta) \end{aligned}$$

Interpretazione della funzione di costo

Interpretazione della **funzione di costo**:

Assumiamo ci sia un solo dato $\mathcal{D} = \{\{\varphi, y\}\}$

$$J(\theta) = \begin{cases} -\ln[\pi] & \text{se classifichiamo } y = 1 \\ -\ln[1 - \pi] & \text{se classifichiamo } y = 0 \end{cases}$$

Caso $\varphi \in y = 1$:

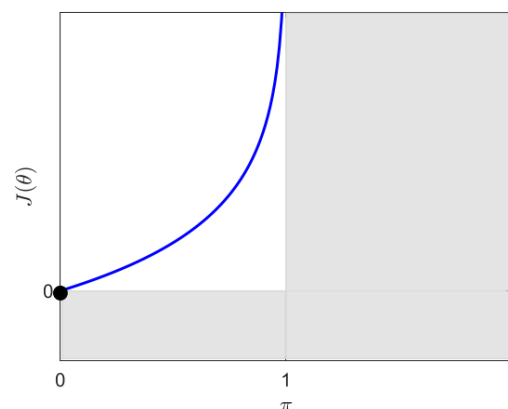
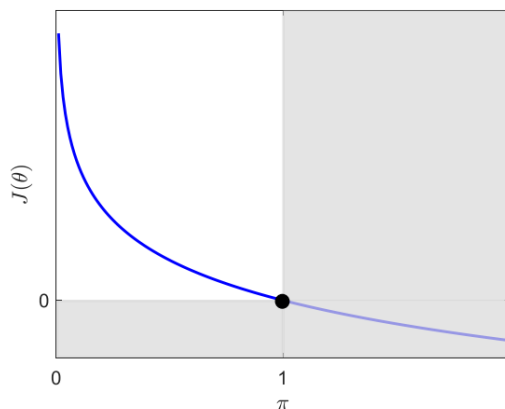
$$J(\theta) = -\ln[\pi]$$

$$\begin{cases} J(\theta) \approx 0, \pi \approx 1 \\ J(\theta) \approx +\infty, \pi \approx 0 \end{cases}$$

Caso $\varphi \in y = 0$:

$$J(\theta) = -\ln[1 - \pi]$$

$$\begin{cases} J(\theta) \approx 0, \pi \approx 0 \\ J(\theta) \approx +\infty, \pi \approx 1 \end{cases}$$



NB: π viene calcolato come la probabilità che l'osservazione appartenga alla classe 1 ($\varphi(i) \in y = 1$) infatti: $\pi = P(y(i) = 1 | \varphi(i))$. Allo stesso modo, la probabilità che l'osservazione appartenga alla classe 2 ($\varphi(i) \in y = 0$) è data da: $1 - \pi$.

Quali algoritmi si possono usare per la minimizzazione della funzione di costo

Per minimizzare la funzione di costo si può ad esempio utilizzare il metodo iterativo del calcolo del gradiente.

Supponendo di avere a disposizione N osservazioni *iid* distribuite come una distribuzione di Bernoulli:

$$p(y_i|\pi) \sim \pi^{y_i} \cdot (1 - \pi)^{1-y_i}$$

- calcolare lo stimatore a massima verosimiglianza per π
- dimostrare che lo stimatore è corretto $\mathbb{E}[y_i] = \pi \quad \forall i$
- dare un'interpretazione della formula ottenuta

Calcolare lo stimatore a massima verosimiglianza

Guardare passaggi completi sopra.

$$-\ln[\mathcal{L}(\pi|Y)] = -\sum_{i=1}^N \left(y(i) \cdot \ln[\pi(i)] + (1 - y(i)) \cdot \ln[1 - \pi(i)] \right) \equiv J(\theta)$$

Dimostrazione correttezza

$$\text{Pongo } Y = \sum_{i=1}^N y(i) \quad e \quad N - Y = \sum_{i=1}^N (1 - y(i))$$

$$\begin{aligned} J(\theta) &\equiv -\sum_{i=1}^N \left(y(i) \cdot \ln[\pi(i)] + (1 - y(i)) \cdot \ln[1 - \pi(i)] \right) = \\ &= -Y \cdot \ln(\pi) - \ln(1 - \pi)(N - Y) \end{aligned}$$

$$\nabla_{\pi} J(\theta) = 0$$

$$\nabla_{\pi} (-Y \cdot \ln(\pi) - (N - Y) \cdot \ln(1 - \pi)) = 0$$

$$-\nabla_{\pi} (Y \cdot \ln(\pi) + (N - Y) \cdot \ln(1 - \pi)) = 0$$

$$\nabla_{\pi} (Y \cdot \ln(\pi) + (N - Y) \cdot \ln(1 - \pi)) = 0$$

$$\frac{Y}{\pi} - \frac{N - Y}{1 - \pi} = 0$$

$$\frac{Y(1 - \pi) - \pi(N - Y)}{\pi(1 - \pi)} = 0$$

$$Y - \pi Y - \pi N + \pi Y = 0$$

$$Y - \pi N = 0$$

Lo stimatore risulta:

$$\hat{\pi} = \frac{Y}{N} = \frac{1}{N} \cdot \sum_{i=1}^N y(i)$$

Lo stimatore risulta essere corretto, poichè:

$$\mathbb{E}[\hat{\pi}] = \mathbb{E}\left[\frac{1}{N} \cdot \sum_{i=1}^N y(i)\right]$$

Interpretazione della funzione di costo

Guardare domande precedenti.

3 Lezione 6: Fondamenti di machine learning

Spiegare i concetti di:

- bias e varianza
- bias-variance tradeoff
- come capire se il modello soffre di uno o l'altro (learning curves)

bias e varianza

Insieme delle features (input):	$\varphi \in \mathbb{R}^{d \times 1}$
Response variable (output categorico/scalare):	$y \in \mathbb{R}, y \in \{cat1, cat2, \dots\}$
Funzione target (da stimare/ignota):	$f : \mathbb{R}^{d \times 1} \rightarrow \mathcal{Y}$
Ipotesi scelta (approssimazione di f):	$g : \mathbb{R}^{d \times 1} \rightarrow \mathcal{Y}$

Osservazioni:

- \mathcal{M} è chiamato spazio delle ipotesi (o set dei modelli)
- $g \in \mathcal{M}$, è un'approssimazione di f
- Obiettivo è stimare f usando i dati \mathcal{D}
- La funzione f viene cercata, dall'algoritmo di learning, nello spazio delle ipotesi \mathcal{M}
- Vogliamo trovare una funzione $h \in \mathcal{M}$ che approssima bene f , non solo sui dati \mathcal{D} a disposizione, ma sull'intero dominio $\mathbb{R}^{d \times 1}$ di f .
- $h \approx f \Rightarrow$ Dobbiamo definire una misura di errore o di costo. **Misure di errore puntuali**, definite su un singolo punto φ . **Misure di errore globale**, definite considerando tutte le N osservazioni.

Misure di errore puntuali:

Errore quadratico:	$l(f(\varphi), h(\varphi; \theta)) = (f(\varphi) - h(\varphi; \theta))^2$	usata per regressione
Errore binario:	$l(f(\varphi), h(\varphi; \theta)) = \mathbb{I}\{f(\varphi) \neq h(\varphi; \theta)\}$	usata per classificazione

Misure di errore globali:

Errore in-sample (errore di train):	$E_{in}(h(\theta)) \equiv J(\theta) = \frac{1}{N} \cdot \sum_{i=1}^N l(f(\varphi), h(\varphi; \theta))$
Errore out-of-sample (errore di validazione):	$E_{out}(h(\theta)) = \mathbb{E}_{\varphi}[l(f(\varphi), h(\varphi; \theta))]$

- Non è possibile conoscere con certezza come sarà il comportamento della funzione f su punti che non ho osservato (problema **dell'induzione di Hume**).
- Nel caso del learning di modelli, ciò che ci interessa veramente stimare è E_{out} , non E_{in} , in quanto E_{in} non è un buon indicatore della bontà del modello.
- L'algoritmo di learning è usato per scandagliare lo spazio delle ipotesi \mathcal{M} , al fine di trovare la miglior $h \in \mathcal{M}$ che approssima bene i dati osservati \rightarrow chiamiamo questa ipotesi g . Con tante ipotesi in \mathcal{M} , c'è un rischio maggiore di trovare una funzione g che spiega benissimo i dati misurati ma fa malissimo su dati nuovi
- Esiste quindi tradeoff tra **approssimazione** e **generalizzazione**. Si vuole: avere un buon modello sui dati misurati (training set) e avere un buon modello su dati non visti (e quindi non usati per la stima del modello)

Errore di generalizzazione:	$E_{out}(g) - E_{in}(g)$
Spazio delle ipotesi \mathcal{M} più complesso:	Migliori possibilità di approssimare f in-sample
Spazio delle ipotesi meno \mathcal{M} complesso:	Migliori possibilità di generalizzare f out-of-sample

Un modo per studiare questo trade-off è valutare i **concetti di bias e varianza** di un modello di learning. L'approccio bias-varianza decompone E_{out} in:

1. Quanto bene \mathcal{M} può approssimare $f \rightarrow$ **bias**
2. Quanto bene riusciamo a scegliere una buona $h \in \mathcal{M}$, usando i dati \rightarrow **varianza**

Sia $g = h \in \mathcal{M}$ la miglior ipotesi che approssima bene sui dati osservati e supponiamo di osservare i dati **senza rumore**, cioè che $y = f(\varphi)$. L'errore out-of-sample può essere espresso come (rendendo esplicita la dipendenza di g da \mathcal{D})

$$E_{out}(g^{\mathcal{D}}) = \mathbb{E}_{\varphi} \left[\left(g^{\mathcal{D}}(\varphi) - f(\varphi) \right)^2 \right]$$

L'errore out-of-sample atteso del modello è indipendente dalla particolare realizzazione dei dati utilizzati per stimare $g^{\mathcal{D}}$:

$$\mathbb{E}_{\mathcal{D}} \left[E_{out}(g^{\mathcal{D}}) \right] = \mathbb{E}_{\mathcal{D}} \left[\mathbb{E}_{\varphi} \left[\left(g^{\mathcal{D}}(\varphi) - f(\varphi) \right)^2 \right] \right] = \mathbb{E}_{\varphi} \left[\mathbb{E}_{\mathcal{D}} \left[\left(g^{\mathcal{D}}(\varphi) - f(\varphi) \right)^2 \right] \right]$$

Definiamo **ipotesi media**: $\bar{g} = \mathbb{E}_{\mathcal{D}}[g^{\mathcal{D}}(\varphi)]$

Ipotesi che deriva dall'usare K dataset

$\mathcal{D}_1, \dots, \mathcal{D}_K$ e costruendola come

$$\bar{g}(\varphi) \approx \frac{1}{K} \sum_{k=1}^K g^{\mathcal{D}_k}(\varphi)$$

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} \left[\left(g^{\mathcal{D}}(\varphi) - f(\varphi) \right)^2 \right] &= \mathbb{E}_{\mathcal{D}} \left[\left(g^{\mathcal{D}}(\varphi) - \bar{g}(\varphi) + \bar{g}(\varphi) - f(\varphi) \right)^2 \right] \\ &= \mathbb{E}_{\mathcal{D}} \left[\left(g^{\mathcal{D}}(\varphi) - \bar{g}(\varphi) \right)^2 + \left(\bar{g}(\varphi) - f(\varphi) \right)^2 + 2 \cdot \left(g^{\mathcal{D}}(\varphi) - \bar{g}(\varphi) \right) \cdot \left(\bar{g}(\varphi) - f(\varphi) \right) \right] \\ &= \underbrace{\mathbb{E}_{\mathcal{D}} \left[\left(g^{\mathcal{D}}(\varphi) - \bar{g}(\varphi) \right)^2 \right]}_{Var(\varphi)} + \underbrace{\left(\bar{g}(\varphi) - f(\varphi) \right)^2}_{bias^2(\varphi)} \end{aligned}$$

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} \left[E_{out}(g^{\mathcal{D}}) \right] &= \mathbb{E}_{\varphi} \left[\mathbb{E}_{\mathcal{D}} \left[\left(g^{\mathcal{D}}(\varphi) - f(\varphi) \right)^2 \right] \right] \\ &= \mathbb{E}_{\varphi} \left[bias^2(\varphi) + var(\varphi) \right] \\ &= bias^2 + var \end{aligned}$$

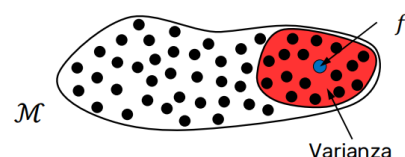
Interpretazione:

- Il termine $bias^2(\varphi) = (\bar{g}(\varphi) - f(\varphi))^2$ misura quanto il nostro modello (cioè la nostra funzione stimata \bar{g}) è lontano dalla funzione target f . Infatti \bar{g} ha il vantaggio di apprendere da un numero illimitato di datasets. Quindi \bar{g} , nella capacità di approssimare f , è limitata solo dai limiti di \mathcal{M} .
- $var(\varphi) = \mathbb{E}_{\mathcal{D}} \left[\left(g^{\mathcal{D}}(\varphi) - \bar{g}(\varphi) \right)^2 \right]$ misura quanto $g^{\mathcal{D}}$ si disperde da \bar{g} e si può vedere quanto l'ipotesi finale $g^{\mathcal{D}}$ differisca dall'ipotesi migliore ovvero quella media \bar{g} .

$$bias^2 = (\bar{g}(\varphi) - f(\varphi))^2$$



$$varianza = \mathbb{E}_{\mathcal{D}} \left[\left(g^{(\mathcal{D})}(\varphi) - \bar{g}(\varphi) \right)^2 \right]$$



Aniché osservare $y = f(\varphi)$, osserviamo $y = f(\varphi) + \eta(\varphi)$ dove η é un rumore stocastico con media zero e varianza σ^2

$$\mathbb{E}_{\mathcal{D}, \varphi, \eta} \left[\left(g^{\mathcal{D}}(\varphi) - (f(\varphi) + \eta(\varphi)) \right)^2 \right] = bias^2 + var + \sigma^2$$

L'errore stocastico σ^2 non può essere portato a zero e contribuisce alla varianza dell'ipotesi scelta, causando overfitting → **errore irriducibile**.

Bias-variance tradeoff

Il trade-off bias-varianza indica che ridurre il bias di un modello aumenta la sua varianza e viceversa. In altre parole, esiste una relazione inversa tra la capacità del modello di adattarsi ai dati di addestramento (approssimazione) e la sua capacità di generalizzare su nuovi dati. L'obiettivo è trovare il giusto equilibrio tra bias e varianza per ottenere un modello con una buona capacità di generalizzazione.

I principali fattori che influenzano questo aspetto sono:

- **Dimensione del set di addestramento:** un set di addestramento più grande risulta essere molto informativo e tende a ridurre la varianza. Al contrario, con un set di addestramento più piccolo, il modello può soffrire di una maggiore varianza, poiché è meno probabile che rappresenti l'intera distribuzione dei dati.
- **Complessità del modello:** un modello più complesso ha una maggiore capacità di adattarsi ai dati di addestramento, riducendo il bias. Tuttavia, un modello più complesso tende anche ad avere una maggiore varianza, poiché può adattarsi anche al rumore o alle peculiarità casuali dei dati.
- **Rumore nei dati:** la presenza di rumore nei dati può aumentare la varianza del modello. Se i dati contengono rumore o errori casuali, il modello può adattarsi anche a queste irregolarità, introducendo un'elevata varianza. Ridurre il rumore nei dati può ridurre la varianza del modello.
- **Selezione delle feature:** l'utilizzo di feature rilevanti e informative può aiutare a ridurre il bias, ma l'uso di feature non informative o irrilevanti può aumentare la varianza.
- **Tecniche di validazione:** aiutano a valutare il trade-off bias-varianza.

Un modello con:

- **alto bias e bassa varianza:** il modello ha una capacità limitata di adattarsi ai dati di addestramento ed è troppo semplice per catturare la complessità del problema → non é in grado di catturare i pattern nei dati
- **basso bias e alta varianza:** il modello ha una capacità elevata di adattarsi ai dati di addestramento ed è molto complesso, possibile overfitting con cattiva generalizzazione.
- **bias adeguato e varianza adeguata:** il modello ha una buona capacità di adattarsi ai dati di addestramento (approssimazione) e di generalizzare su nuovi dati.

modelli con basso bias sono più complessi → tendono ad overfittare e a inglobare errore
modelli con bassa varianza sono meno complessi → potrebbero non cogliere la relazione dei dati

Modello semplice:

Bias ↑
Varianza ↓

Approssimazione e generalizzazione cattive

All'aumentare di N rimane uguale

Modello complesso:

Bias ↓
Varianza ↑

Approssimazione positiva generalizzazione negativa

All'aumentare di N migliora

Il **trade-off tra approssimazione e generalizzazione** si traduce nel valutare i concetti di bias e varianza di un modello di learning. L'approccio bias-varianza decompone E_{out} in:

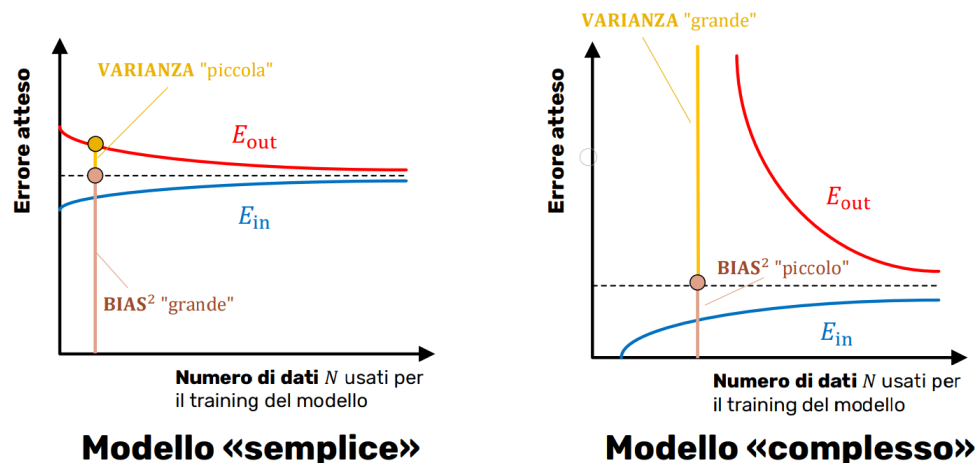
1. Quanto bene \mathcal{M} (spazio delle ipotesi) può approssimare f (funzione target) → **bias**
2. Quanto bene riusciamo a scegliere una buona $h \in \mathcal{M}$, usando i dati → **varianza**

Learning curves

Le learning curves sono uno strumento grafico per capire se un modello di learning soffre di **problemi di bias o varianza**. L'idea è di rappresentare, al variare del numero di dati N usati per stimare il modello:

- l'errore out-of-sample atteso $\rightarrow \mathbb{E}_{\mathcal{D}}[E_{out}(g^{\mathcal{D}})]$
- l'errore in-sample atteso $\rightarrow \mathbb{E}_{\mathcal{D}}[E_{in}(g^{\mathcal{D}})]$

Le curve vengono calcolate usando un solo dataset, oppure dividendolo in più parti e prendendo la **curva media** risultante dai vari sub-datasets.



- Il **bias** può essere presente quando l'errore atteso è piuttosto elevato e E_{in} in è simile a E_{out} , è improbabile che ottenere più dati aiuti.

\Rightarrow Aggiungere features (per esempio combinazioni di features originarie) o boosting

- La **varianza** può essere presente quando c'è un tanto divario tra E_{in} e E_{out} , è probabile che ottenere più dati sia d'aiuto

\Rightarrow Usare meno features, acquisire più dati, usare regolarizzazione, bagging

Discutere il metodo delle learning curves specificando:

- il motivo del loro utilizzo
- la tipologia di problemi diagnosticabili con esse
- dei possibili rimedi per questi problemi

Guardare la domanda precedente.

Overfitting, da cosa è causato, come si può risolvere e tipologie di regolarizzazione

Overfitting e cause di overfitting

L'overfitting è un problema comune nel machine learning ed è causato da un modello che si adatta troppo ai dati di addestramento, perdendo la capacità di generalizzare correttamente su nuovi dati.

- Un modello che fa overfitting presenta:

$$Bias = basso \quad E_{in} = basso$$

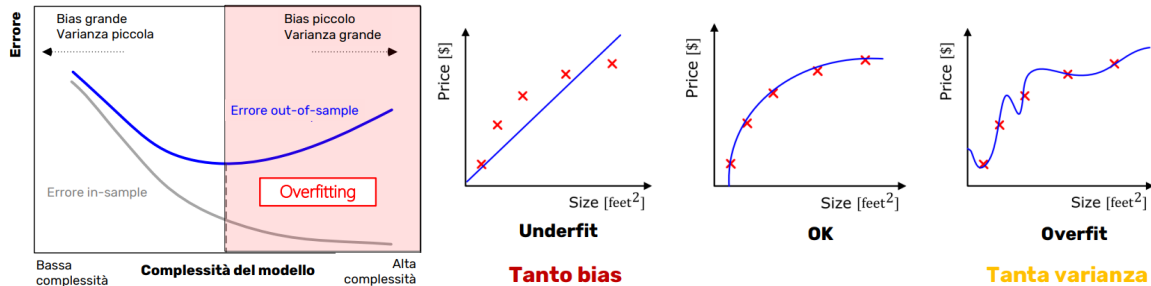
$$Var = alta \quad E_{out} = alto$$

Solitamente questo comportamento è tipico dei modelli complessi che hanno un potere espressivo molto alto e tendono ad adattarsi al rumore. Nel caso di pochi dati è necessario usare modelli più semplici indipendentemente dalla complessità della funzione target.

- Una seconda causa di overfitting è il rumore stocastico η (supponiamo media 0 e var σ^2) che affligge le misure. Al posto di osservare $y = f(\varphi)$ osservo $y = f(\varphi) + \eta(\varphi)$, dunque:

$$\mathbb{E}_{\mathcal{D}, \varphi, \eta} \left[\left(g^{\mathcal{D}}(\varphi) - (f(\varphi) + \eta(\varphi)) \right)^2 \right] = bias^2 + var + \sigma^2$$

L'errore stocastico σ^2 non può essere portata a zero e contribuisce alla varianza dell'ipotesi scelta, causando overfitting → **errore irriducibile**



Regularizzazione

I modelli più complessi sono più inclini all'overfitting, questo perché sono molto espressivi e quindi possono adattarsi anche al rumore. I modelli semplici mostrano meno varianza a causa della loro espressività limitata. La riduzione della varianza del modello è spesso maggiore dell'aumento del suo bias, per cui, nel complesso, errore atteso complessivo diminuisce ($bias^2 + var + \sigma^2$).

Tuttavia, se ci atteniamo solo a modelli semplici, potremmo non ottenere un'approssimazione soddisfacente della funzione target f .

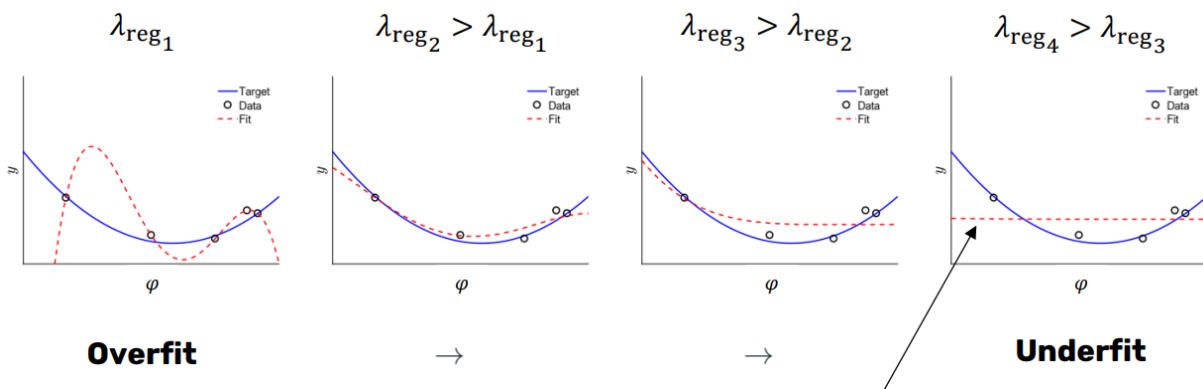
La regularizzazione oltre che minimizzare la funzione di costo $E_{in} \equiv J(\theta)$ minimizza anche la complessità del modello. Al posto di E_{in} minimizziamo dunque un **errore aumentato** $E_{aug}(\theta)$

$$E_{aug}(\theta) = \underbrace{\frac{1}{N} \cdot \sum_{i=1}^N (y(i) - h(\varphi(i); \theta))^2}_{\text{Quanto male il modello fitta i dati}} + \lambda_{reg} \cdot \Omega(\theta)$$

$\Omega(\theta) \rightarrow$ Regularizzatore: quanto il modello è complesso

$h(\cdot)$ è qualche funzione che rappresenta il nostro modello

Il termine λ_{reg} (iper-parametro) **pesa l'importanza** di minimizzare $E_{in} \equiv J(\theta)$ rispetto a minimizzare $\Omega(\theta)$, ossia la complessità del modello.



Se regularizzo troppo, imparerò la funzione più semplice possibile, ovvero una retta orizzontale (costante) con intercetta θ_0 .

Minimizzare E_{aug} rispetto ad E_{in} conduce ad un modello migliore (ovvero un modello con miglior capacità di generalizzare e quindi con E_{out} minore).

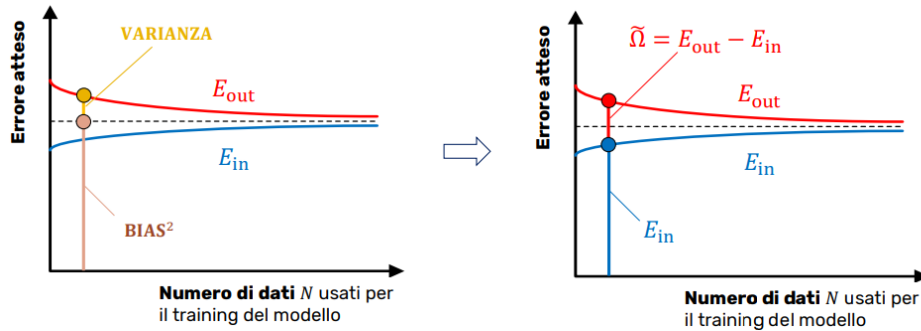
Possiamo interpretare E_{out} come la somma di due contributi:

$$E_{out}(\theta) = E_{in}(\theta) + \tilde{\Omega}(\theta)$$

Ricordando la definizione di E_{aug} abbiamo:

$$E_{aug}(\theta) = E_{in}(\theta) + \lambda_{reg} \cdot \Omega(\theta)$$

$\Rightarrow E_{aug}$ è migliore rispetto ad E_{in} come proxy per E_{out}



La regolarizzazione aiuta nello stimare la quantità $\Omega(\theta)$, che sommata ad E_{in} fornisce E_{aug} , il quale è una stima di E_{out} . Esistono diversi tipi di regolarizzazione. I più usati sono:

Regolarizzazione $L_1 \rightarrow$ Lasso

Tende a portare più coefficienti esattamente a zero

$$\Omega(\theta) = \sum_{j=0}^{d-1} |\theta_j|$$

Regolarizzazione $L_2 \rightarrow$ Ridge

Tende a ridurre tutti i coefficienti a un valore inferiore

$$\Omega(\theta) = \sum_{j=0}^{d-1} (\theta_j)^2$$

Regolarizzazione elastic-net:

Combinazione di L_1 e L_2

$$\Omega(\theta) = \beta \sum_{j=0}^{d-1} (\theta_j)^2 + (1 - \beta) \sum_{j=0}^{d-1} |\theta_j|$$

Gli effetti della regolarizzazione possono essere osservati nei termini di bias e varianza:

- La regolarizzazione **aumenta di poco il bias** (perché ottengo un modello più semplice) al fine di **ridurre considerevolmente la varianza** del modello di learning
- La regolarizzazione porta ad avere **ipotesi più smooth**, regolari, riducendo il rischio di overfitting
- L'iperparametro di regolarizzazione λ_{reg} deve essere scelto in modo specifico per ogni tipo di regolarizzatore. Solitamente si usa una procedura come la validazione o la cross-validazione

Cross-validazione

Validazione

$$E_{out}(\theta) = E_{in}(\theta) + \text{penalità per la complessità del modello}$$

La **REGOLARIZZAZIONE** stima questa quantità

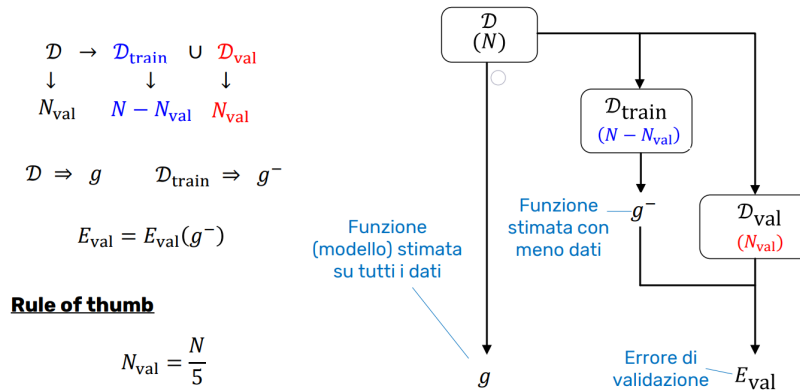
La **VALIDAZIONE** stima questa quantità

L'idea delle procedure di validazione è quella di stimare E_{out} utilizzando un dataset diverso (validation set) rispetto a quello usato per la stima del modello (training\identification set). La regolarizzazione e la validazione sono due tecniche che possono (e devono) essere usate insieme:

- la regolarizzazione aiuta a stimare un modello che può generalizzare meglio
- la validazione fornisce una stima dell'errore out-of-sample del modello

Una procedura comune che si segue è:

1. **Rimuovo** un subset di dati dai dati totali $\rightarrow \mathcal{D}_{val}$
2. **Stimo** il modello sulla parte di dati rimanente $\rightarrow \mathcal{D}_{train}$
3. **Valuto** le performance del modello sul subset di dati che ho rimosso al punto 1
4. **Ri-allo** il modello su tutti i dati



Le procedure di validazione possono essere utilizzate per due scopi:

- Valutare le performance del modello stimato \rightarrow stimare E_{out}
- Scegliere il modello migliore da un insieme di diversi modelli

Problema: se uso il dataset di validazione \mathcal{D}_{val} tante volte per compiere delle scelte, allora il dataset di validazione \mathcal{D}_{val} non fornisce più una buona stima dell'errore out-of-sample E_{out} .

Intuizione: usare \mathcal{D}_{val} per compiere delle scelte su quale modello usare fa sì che tali scelte siano dipendenti dai particolari valori dei dati contenuti in \mathcal{D}_{val} . Chi mi garantisce che con dati diversi avrei compiuto le medesime scelte? *Stiamo dunque overfittando il validation set.*

Soluzione: c'è bisogno di un terzo dataset. Il dataset di test, sul quale calcoleremo l'errore di test E_{test} .

Contaminazione: si riferisce all'effetto negativo che può verificarsi quando i dati in un determinato set (ad esempio, il training set o il validation set) influenzano eccessivamente le decisioni prese durante lo sviluppo di un modello. La **contaminazione** nei tre set di dati (training, validation e test) si riferisce al livello di presenza di dati sovrapposti o influenze reciproche tra i set. Quando dici che un set è "contaminato", significa che contiene dati che sono stati usati in modo improprio in altre fasi dello sviluppo, potenzialmente portando a una stima ottimistica delle prestazioni del modello.

- **Training set (60% dei dati):** totalmente contaminato
- **Validation set (20% dei dati):** un pò contaminato
- **Test set (20% dei dati):** totalmente pulito

Cross-validazione

La divisione del dataset in tre parti (train, validation, test) è fattibile se i dati a disposizione sono molti.

in teoria vorremmo che:

$$\underbrace{E_{out}(g) \approx E_{out}(\bar{g})}_{N_{val} \text{ piccolo}} \approx E_{val}(\bar{g})$$

$$E_{out}(g) \approx \underbrace{E_{out}(\bar{g}) \approx E_{val}(\bar{g})}_{N_{val} \text{ grande}}$$

- $E_{val}(\bar{g})$ è l'unico che posso calcolare
- N_{val} **grande:** in questo caso $E_{val}(\bar{g}) \approx E_{out}(\bar{g})$ poichè uso tanti dati N_{val} per la validazione. Ricordiamoci che l'obiettivo di E_{val} è proprio quello di stimare E_{out}
- N_{val} **piccolo:** in questo caso $E_{out}(\bar{g}) \approx E_{out}(g)$, poichè uso tanti dati $N - N_{val}$ per il train di \bar{g} . Questo è il valore che mi interessa ma che non posso calcolare direttamente

La cross-validazione permette di avere N_{val} sia grande che piccolo.

Leave-one-out cross-validation

- Usiamo $N - 1$ per il training e $N_{val} = 1$ dato per la validazione, dove \mathcal{D}_i è il dataset di training senza il dato i -esimo:

$$\mathcal{D}_i = \{(\varphi(1), y(1)), \dots, (\varphi(i), y(i)), \dots, (\varphi(N), y(N))\}$$

La fusione imparata usando \mathcal{D}_i è $\rightarrow g_i^-$

- L'errore di validazione sul **singolo punto rimosso** $\varphi(i)$ (usando l'errore puntuale) risulta essere:

$$l(i) = E_{val}(g_i^-) = l\left(y(i), g_i^-(\varphi(i))\right)$$

- E' possibile definire l'**errore di cross-validazione** E_{CV} come:

$$E_{CV} = \frac{1}{N} \sum_{i=1}^N l(i)$$

- Stimo N modelli usando $N - 1$ dati, e li valido usando N stime dell'errore di validazione (con $N_{val} = 1$). La cross-validazione $N_{val} = 1$ (**leave-one-out cross-validation**) **risulta dunque essere computazionalmente costosa**. Nel caso volessimo usarla per scegliere tra M modelli, richiederebbe un totale di N sessioni di training per ciascuno degli M modelli
- La stima dell'errore E_{CV} ha una **varianza elevata**, poiché si basa su un solo dato

Cross-validazione K-fold

- E' possibile riservare più punti per la validazione suddividendo il training set in **folds**. Per esempio, se $k = 10$ avremmo $\mathcal{D}_1, \dots, \mathcal{D}_{10}$ datasets, di cui ad ogni iterazione utilizzo 9 datasets per il training e 1 per la validazione.
- La K-fold cross-validation richiede $\frac{N}{N_{val}}$ sessioni di training, ognuna con $N - N_{val}$ dati. Un buon compromesso è usare $K = 10$ dove $N_{val} = \frac{N}{10}$
- Attenzione a non ridurre troppo il training set (guardare le learning curves)

Formule di complessità ottima

- Queste formule permettono di stimare l'errore out-of-sample E_{out} **utilizzando solo il dataset di train**. Per questo motivo, **si usano quando ho pochi dati** per poter usare validazione o cross-validazione.
- L'idea è simile alla regolarizzazione: modificare la funzione di costo dell'errore in-sample E_{in} , **aggiungendo un termine additivo che penalizza la complessità del modello**.
- Indichiamo la **stima dei parametri**, ottenuta con N dati, con $\hat{\theta}_N \in \mathbb{R}^{d \times 1}$.
- La stima è ottenuta minimizzando la funzione di costo $J(\theta; d)$, dove esplicitiamo la dipendenza del costo dal numero di parametri $d \Rightarrow AIC, BIC \equiv J(\theta)_{aug}$ minori sono i migliori

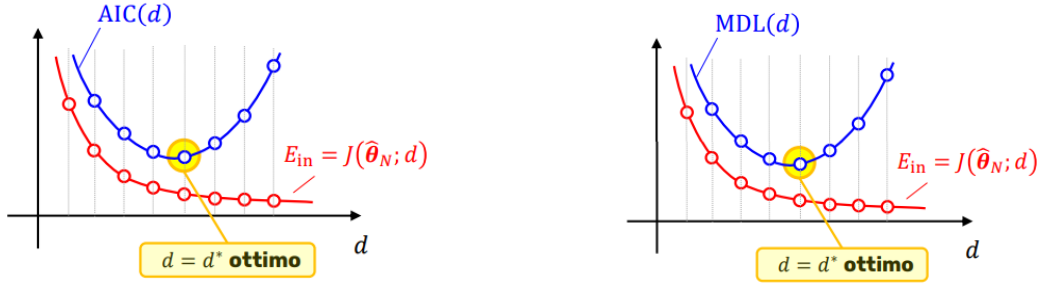
Akaike Information Criterion (AIC) \equiv FPE

$$AIC(d) = 2 \cdot \frac{d}{N} + \ln[J(\hat{\theta}_N; d)]$$

Minimum Description Length (MDL) *derivante dal BIC*

$$MDL(d) = \ln[N] \cdot \frac{d}{N} + \ln[J(\hat{\theta}_N; d)]$$

Se $\ln[N] > 2$ (ovvero se abbiamo più di 8 dati), MDL suggerisce di usare modelli più parsimoniosi.



$$AIC(d) = 2 \cdot \frac{d}{N} + \ln[J(\hat{\theta}_N; d)] \iff MDL(d) = \ln[N] \cdot \frac{d}{N} + \ln[J(\hat{\theta}_N; d)]$$

NB: Sotto l'assunzione che il meccanismo di generazione dei dati appartenga alla classe di modelli scelta, **FPE** e **AIC** hanno una probabilità non nulla di *soprastimare l'ordine del modello*, mentre **MDL** porta ad una *stima asintoticamente corretta dell'ordine*. Dato che raramente l'assunzione è verificata, si preferisce usare **AIC** o **FPE**, soprastimando leggermente d .

Data la funzione di costo definita con termini vettoriali e matriciali e regolarizzata con un Ridge Regularizer dimostrare che la minimizzazione della funzione è quella che si ottiene usando i minimi quadrati ma con il regolarizzatore in più da considerare

Formule:

dove x è un vettore $1 \times d$, A/B è una matrice $d \times d$ e I è identità

$$\begin{array}{ll} 1 & \|A\|_2^2 = A^T A \\ 2 & (AB)^T = B^T A^T \\ 3 & \nabla_x (x^T A x) = (A + A^T) \cdot x \end{array}$$

Consideriamo il modello di regressione lineare con un termine di regolarizzazione L_2 :

$$\begin{aligned} E_{aug}(\theta) \equiv J(\theta) &= \frac{1}{N} \cdot \sum_{i=1}^N (y(i) - \varphi^T(i)\theta)^2 + \lambda_{reg} \sum_{j=0}^{d-1} (\theta_j)^2 \\ &= \frac{1}{N} \cdot \|Y - X \cdot \theta\|_2^2 + \lambda_{reg} \cdot \|\theta\|_2^2 \\ &= \frac{1}{N} \cdot (Y - X\theta)^T (Y - X\theta) + \lambda_{reg} \theta^T \theta \\ &= \frac{1}{N} \cdot (Y^T Y - Y^T X \theta - \theta^T X^T Y + \theta^T X^T X \theta) + \lambda_{reg} \theta^T \theta \\ &= \frac{1}{N} \cdot (Y^T Y - 2\theta^T X^T Y + \theta^T X^T X \theta) + \lambda_{reg} \theta^T \theta \end{aligned}$$

Il minimo risulta essere:

$$\begin{aligned} \nabla_{\theta} J(\theta) = 0 &\Rightarrow \frac{1}{N} \cdot (-2X^T Y + 2X^T X \theta) + 2\lambda_{reg} \theta = 0 \\ \frac{1}{N} \cdot (-X^T Y + X^T X \theta) + \lambda_{reg} \theta &= 0 \\ -X^T Y + X^T X \theta + N\lambda_{reg} I_d \theta &= 0 \\ (X^T X + \lambda_{reg} I_d) \theta &= X^T Y \end{aligned}$$

Risulta dunque:

$$\hat{\theta}_{reg} = \begin{pmatrix} X^T X + \lambda_{reg} I_d \end{pmatrix}_{d \times d}^{-1} X^T Y_{d \times 1}$$

4 Lezione 7: Fondamenti di stima Bayesiana

Probabilità congiunta, marginale e condizionata. Stima bayesiana. Relazione stimatore lineare ottimo e caso in cui la distribuzione congiunta tra due variabili è gaussiana.

Supponiamo di avere due variabili casuali discrete e binarie a e b . Definiamo:

Distribuzione di probabilità congiunta:

$P(a, b)$ probabilità che sia a che b assumino un valore specifico.

La sommatoria di tutte le combinazioni specifiche di valori deve essere uguale a 1

$$\sum_{a=0}^1 \sum_{b=0}^1 p(a, b) = 1$$

$$P(a, b) = P(b, a)$$

Specifici valori:

$$P(a = 0, b = 0)$$

$$P(a = 1, b = 0)$$

$$P(a = 0, b = 1)$$

$$P(a = 1, b = 1)$$

Distribuzione di probabilità marginale:

La distribuzione marginale è la distribuzione di probabilità di un **sottoinsieme di variabili casuali**. Siccome abbiamo 2 variabili casuale (a e b), avremo due marginali $P(a)$ e $P(b)$. Se avessimo 3 v.c discrete a, b, c avremmo le marginali $P(a)$, $P(b)$, $P(c)$, $P(a, b)$, $P(a, c)$, $P(b, c)$.

$$n \text{ var, totale} = n!$$

La distribuzione marginale è ottenuta marginando (**sommando**) rispetto alle **variabili che non sono di interesse**. Nel caso di **v.c. continue**, si deve **integrare** anziché sommare.

$$P(b = 0) = P(a = 0, b = 0) + P(a = 1, b = 0)$$

$$P(b = 1) = P(a = 0, b = 1) + P(a = 1, b = 1)$$

Distribuzione di probabilità condizionata:

La distribuzione condizionata indica come la probabilità si **ridistribuisce** dato che si restringe la popolazione ad un particolare sottoinsieme.

$$P(A|B) = \frac{P(A, B)}{P(B)} \Rightarrow P(A, B) = P(A|B) \cdot P(B)$$

$P(A, B) = P(A) \cdot P(B)$ se e solo se $P(A|B) = P(A)$, ossia A e B sono eventi **indipendenti**, ovvero il verificarsi di B non modifica le probabilità di verificarsi di A .

Temora di Bayes

$$P(A, B) = P(B, A)$$

$$P(A|B) \cdot P(B) = P(B|A) \cdot P(A)$$

Allora

Teorema di bayes

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Il teorema di Bayes permette di **ridistribuire la probabilità**: prima conoscevamo $P(A)$, adesso conosco $P(A|B)$. La probabilità di A è cambiata in seguito all'informazione portata da B .

Stima Bayesiana

- Abbiamo finora considerato il vettore di parametri ignoto $\theta \in \mathbb{R}^{d \times 1}$ come una **variabile deterministica**. Spesso però, ancora prima di collezionare i dati, abbiamo delle informazioni (o supposizioni) sui possibili valori che potrebbe assumere θ .
- Ha quindi senso considerare θ come una **variabile casuale**: in questo modo, posso specificare una distribuzione di probabilità per θ , **per descriverne i valori** che io credo che possa assumere e la probabilità che θ li assuma. Assegno **maggior probabilità** ai valori che **io credo** siano più probabili che θ possa assumere e minor probabilità ai valori che io credo non possa assumere.
- La distribuzione a-priori $f_{\theta}(\theta)$ sui possibili valori di θ , ha dominio $[0, 1]$ poiché θ modellando una probabilità deve stare tra 0 e 1. Data $f_{\theta}(\theta)$, abbiamo già una stima del valore di θ **ancora prima di aver osservato i dati (STIMA A PRIORI)**. L'incertezza sulla stima sarà allora la varianza di 0 (**INCERTEZZA A PRIORI**).

Abbiamo quindi due elementi che portano informazione:

- La distribuzione a-priori $f_{\theta}(\theta)$ sui possibili valori di θ
- L'informazione che portano i dati sui possibili valori di θ , ovvero la likelihood $f_{Y|\theta}(Y|\theta)$

NB: Quello che veramente ci interessa è sapere **quanto può valere θ dato che ho osservato i dati**, ovvero la distribuzione $f_{\theta|Y}(\theta|Y)$

Usando il teorema di Bayes possiamo unire i due elementi di informazione:

$$f_{\theta|Y}(\theta|Y) = \frac{\overset{\text{LIKELIHOOD}}{f_{Y|\theta}(Y|\theta)} \cdot \overset{\text{PRIOR}}{f_{\theta}(\theta)}}{\underset{\text{POSTERIOR}}{f_Y(Y)} \quad \underset{\text{MARGINAL LIKELIHOOD}}{f_Y(Y)}}$$

$f_{\theta}(\theta)$	Distribuzione a-priori sui possibili valori di θ
$f_Y(Y)$	Distribuzione dei dati Y
$f_{Y \theta}(Y \theta)$	Distribuzione dei dati Y dato che ho supposto a priori θ
$f_{\theta Y}(\theta Y)$	Distribuzione di θ dato che ho osservato i dati Y

Osservazioni:

- $f_{\theta|Y}(\theta|Y)$ è una **distribuzione a-posteriori di possibili valori di θ** . Le probabilità di questi valori, rispetto a $f_{\theta}(\theta)$ sono state riallocate dall'aver osservato i dati Y (ecco perché a posteriori)
- Nel caso in cui $f_{Y|\theta}(Y|\theta)$ e $f_{\theta}(\theta)$ sono pdf continue allora:

$$f_Y(Y) = \int_{-\infty}^{+\infty} f_{Y|\theta}(Y|\theta) \cdot f_{\theta}(\theta) d\theta = \int_{-\infty}^{+\infty} f(Y, \theta) d\theta \quad \text{altrimenti sommatoria}$$

- Un altro problema è che $f_Y(Y)$ nel caso di dati intesi come v.c. continue, è un integrale che potremmo non sapere come risolvere possiamo utilizzare **Markov Chain Monte Carlo (MCMC)**
- In generale **non posso dire nulla sulla posterior**, solo in casi fortunati ha un'espressione analitica nota. Caso fortunato se $f_{\theta}(\theta) = \text{gaussiana}$ e anche $f_{Y|\theta}(Y|\theta) = \text{gaussiana}$ allora anche $f_{\theta|Y}(\theta|Y) = \text{gaussiana}$
- Un modo (computazionalmente oneroso ma semplice) per calcolare la posterior $f_{\theta|Y}(\theta|Y)$ è quello di **discretizzare** il range di valori del parametro θ tramite una griglia di valori.

Stima ottima

Indichiamo uno stimatore come una funzione $T(\cdot)$ dei dati \mathcal{D} :

$$\hat{\theta} = T(\mathcal{D})$$

Consideriamo il caso θ scalare. Vorremmo che la variabile casuale $\hat{\theta}$ fosse vicina alla variabile casuale θ per quantificare questa distanza, usiamo il concetto di **Mean Squared Error (MSE)**

$$MSE \equiv \mathbb{E}[(\hat{\theta} - \theta)^2] = \mathbb{E}[(T(\mathcal{D}) - \theta)^2]$$

Lo **stimatore ottimo di Bayes** è quella funzione $T^{opt}(\cdot)$ che **minimizza il MSE**

$$\mathbb{E}[(T^{opt}(\mathcal{D}) - \theta)^2] < \mathbb{E}[(T(\mathcal{D}) - \theta)^2], \quad \forall T(\mathcal{D})$$

Si dimostra che **lo stimatore che minimizza il MSE è il valore atteso condizionato** (al fatto che i dati \mathcal{D} abbiano assunto i valori in Y)

$$\hat{\theta} = T^{opt}(Y) = \mathbb{E}[\theta | \mathcal{D} = Y]$$

Nel caso in cui θ sia un **vettore di parametri**, il calcolo del MSE si modifica come segue:

$$MSE \equiv \text{tr} \left\{ \mathbb{E} \left[\underbrace{(\hat{\theta} - \theta)}_{d \times 1} \underbrace{(\hat{\theta} - \theta)^T}_{1 \times d} \right] \right\} = \mathbb{E} \left[\underbrace{(\hat{\theta} - \theta)^T}_{1 \times d} \underbrace{(\hat{\theta} - \theta)}_{d \times 1} \right] = \mathbb{E} \left[\underbrace{\|(\hat{\theta} - \theta)\|_2^2}_{1 \times 1} \right]$$

Stima ottima : il caso Gaussiano

Supponiamo ora di avere un dato interpretato come realizzazione di una variabile casuale Gaussiana $y \sim \mathcal{N}(0, \lambda_{yy}^2)$, e che anche il parametro ignoto (scalare per comodità) sia Gaussiano $\theta \sim \mathcal{N}(0, \lambda_{\theta\theta}^2)$

$$\underbrace{\begin{bmatrix} y \\ \theta \end{bmatrix}}_{\mathbf{z}} \sim \mathcal{N} \left(\underbrace{\begin{bmatrix} 0 \\ 0 \end{bmatrix}}_{\boldsymbol{\mu}}, \underbrace{\begin{bmatrix} \lambda_{yy}^2 & \lambda_{y\theta} \\ \lambda_{\theta y} & \lambda_{\theta\theta}^2 \end{bmatrix}}_{\Sigma} \right)$$

Congiunta è gaussiana: $f_{y\theta}(y, \theta)$

$$f_{y\theta}(y, \theta) = \frac{1}{\sqrt{(2\pi)^2 \det \Sigma}} \exp \left(-\frac{1}{2} (\mathbf{z} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{z} - \boldsymbol{\mu}) \right)$$

Pdf dei dati è gaussiana: $f_y(y)$

$$f_y(y) = \frac{1}{\sqrt{2\pi \cdot \lambda_{yy}^2}} \cdot \exp \left(-\frac{1}{2\lambda_{yy}^2} (y - 0)^2 \right)$$

Si dimostra che la **posterior** $f_{\theta|y}(\theta|y) = \frac{f_{y\theta}(y, \theta)}{f_y(y)}$ è ancora **Gaussiana** con:

Valore atteso: $\mu_{\theta|y} = \frac{\lambda_{\theta y}}{\lambda_{yy}^2} \cdot y$

Varianza: $\lambda_{\theta|y}^2 = \lambda_{\theta\theta}^2 - \frac{\lambda_{\theta y}^2}{\lambda_{yy}^2}$

Avendo osservato il valore $y(1)$ di y , la stima ottenuta dallo **stimatore ottimo Bayesiano nel caso Gaussiano** sarà:

$$\hat{\theta}_{opt} = \mathbb{E}[\theta | y = y(1)] = \frac{\lambda_{\theta y}}{\lambda_{yy}^2} \cdot y(1)$$

Stima ottima lineare

Vogliamo quindi trovare uno **stimatore che non faccia ipotesi sulla ddp congiunta** di y e θ . Supponiamo y e θ due variabili casuali scalari con valore atteso nullo e varianza λ_{yy}^2 e $\lambda_{\theta\theta}^2$ rispettivamente:

$$\mathbb{E}[y] = 0 \quad \mathbb{E}[y^2] = \lambda_{yy}^2 \quad \mathbb{E}[\theta] = 0 \quad \mathbb{E}[\theta^2] = \lambda_{\theta\theta}^2 \quad \mathbb{E}[\theta y] = \lambda_{\theta y}$$

Vogliamo stimare θ tramite uno **stimatore lineare** del tipo:

$$\hat{\theta}^{lin} = \alpha \cdot y + \beta \quad \alpha, \beta \in \mathbb{R}$$

Per trovare α e β , **minimizziamo la funzione di costo** data dal Mean square error:

$$MSE \equiv J(\alpha, \beta) = \mathbb{E}[(\hat{\theta} - \theta)^2] = \mathbb{E}[(\alpha \cdot y + \beta - \theta)^2]$$

Calcoliamo il gradiente e poniamolo uguale a zero (non verifichiamo sia un minimo):

$$\frac{\partial J(\alpha, \beta)}{\partial \alpha} = 0 \quad \Rightarrow \quad 2 \cdot \mathbb{E}[(\alpha y + \beta - \theta) \cdot y] = 0 \quad \Rightarrow \quad \alpha = \frac{\lambda_{\theta y}}{\lambda_{yy}^2}$$

$$\frac{\partial J(\alpha, \beta)}{\partial \beta} = 0 \quad \Rightarrow \quad 2 \cdot \mathbb{E}[(\alpha y + \beta - \theta) \cdot 1] = 0 \quad \Rightarrow \quad \beta = 0$$

Lo stimatore lineare ottimo **coincide con lo stimatore ottimo di Bayes per il caso Gaussiano**:

$$\hat{\theta}_{opt}^{lin} = \hat{\alpha} \cdot y + \hat{\beta} = \frac{\lambda_{\theta y}}{\lambda_{yy}^2} \cdot y$$

La varianza della stima si ricava essere uguale al caso Gaussiano:

$$Var[\hat{\theta}_{opt}^{lin} - \theta] = \lambda_{\theta\theta}^2 - \frac{\lambda_{\theta y}^2}{\lambda_{yy}^2}$$

5 Lezione 8: Processi stocastici

Proprietà della densità spettrale di potenza, stima dello spettro, proprietà statistiche, come migliorarne la stima

Densità spettrale di potenza e proprietà

Definizione: Dato un processo stocastico stazionario (sia in senso debole che in senso forte), si definisce **densità spettrale di potenza** $\Gamma_{vv}(\omega)$ come la DTFT di $\gamma_{vv}(\tau)$

$$\Gamma_{vv}(\omega) \equiv \mathcal{F}[\gamma_{vv}(\tau)] = \sum_{\tau=-\infty}^{+\infty} \gamma_{vv}(\tau) \cdot e^{-j\omega\tau}$$

$$\text{La trasformata } \mathcal{Z} \text{ di } \gamma_{vv}(\tau) \text{ è: } \Phi_{vv}(z) \equiv \mathcal{Z}[\gamma_{vv}(\tau)] = \sum_{\tau=-\infty}^{+\infty} \gamma_{vv}(\tau) \cdot z^{-\tau}$$

Data $\Phi_{vv}(z)$ si ha che $\Gamma_{vv}(\omega) = \Phi_{vv}(e^{j\omega})$ (stesso concetto tra \mathcal{Z} e \mathcal{F}).

Interpretazione: la densità spettrale di potenza ci dice come, in media, le componenti in frequenza delle varie realizzazioni del processo stocastico $v(t, s)$ contribuiscono alla sua varianza. O in parole più povere **come l'energia del processo si distribuisce alle varie frequenze**.

Proprietà di $\Gamma_{vv}(\omega)$:

- **Reale:** dato che $\gamma_{vv}(\tau)$ è pari, i termini immaginari del tipo $\pm j \cdot \sin(\omega)$ si elidono
- **Positiva:** $\Gamma_{vv}(\omega) \geq 0, \forall \omega \in \mathbb{R}$
- **Pari:** $\Gamma_{vv}(\omega) = \Gamma_{vv}(-\omega), \forall \omega \in \mathbb{R}$
- **Periodica di periodo 2π :** $\Gamma_{vv}(\omega) = \Gamma_{vv}(\omega + k \cdot 2\pi), \forall \omega \in \mathbb{R}, \forall k \in \mathbb{Z}$
- Affinché $\Gamma_{vv}(\omega)$ **converga**, $\gamma_{vv}(\tau)$ deve tendere a zero in modo sufficientemente rapido.
- Ci basta valutare $\Gamma_{vv}(\omega)$ tra $[0, \pi]$
- È possibile risalire a $\gamma_{vv}(\tau)$ tramite l'antitrasformata

$$\gamma_{vv}(\tau) = \frac{1}{2\pi} \cdot \int_{-\pi}^{+\pi} \Gamma_{vv}(\omega) \cdot e^{j\omega\tau} d\omega$$

- è possibile esprimere la **varianza del processo** stazionario come l'area sottesa alla densità spettrale di potenza

$$\gamma_{vv}(0) = \frac{1}{2\pi} \int_{-\pi}^{+\pi} \Gamma_{vv}(\omega) d\omega$$

Stima dello spettro, proprietà statistiche, come migliorarne la stima

Vedere risposta successiva.

Descrivere come si può effettuare la stima dello spettro (densità spettrale di potenza) di un processo stocastico stazionario, partendo dai dati misurati. Discutere inoltre le proprietà dello stimatore. Cosa si può fare per migliorare la stima?

Definizione: Dato un processo stocastico stazionario (sia in senso debole che in senso forte), si definisce **densità spettrale di potenza** $\Gamma_{vv}(\omega)$ come la DTFT di $\gamma_{vv}(\tau)$

$$\Gamma_{vv}(\omega) \equiv \mathcal{F}[\gamma_{vv}(\tau)] = \sum_{\tau=-\infty}^{+\infty} \gamma_{vv}(\tau) \cdot e^{-j\omega\tau}$$

Autocovarianza (temporale) campionaria

$$\hat{\gamma}_{vv}(\tau) = \frac{1}{N-|\tau|} \sum_{t=0}^{N-|\tau|-1} v(t) \cdot v(t+|\tau|), \quad |\tau| < N$$

- $\mathbb{E}[\hat{\gamma}_{vv}(\tau)] = \gamma_{vv}(\tau)$, ovvero lo stimatore è **corretto**
- Per τ fissato, lo stimatore è **consistente**, sotto le ipotesi di ergodicità
- Per $\tau \approx N$, si ha che $\text{var}[\hat{\gamma}_{vv}(\tau)]$ è grande perché ci sono pochi addendi

Per risolvere quest'ultimo problema, possiamo pensare ad uno stimatore alternativo (seppur non corretto)

Autocovarianza (temporale) campionaria - versione alternativa

$$\hat{\gamma}'_{vv}(\tau) = \frac{1}{N} \sum_{t=0}^{N-|\tau|-1} v(t) \cdot v(t+|\tau|), \quad |\tau| < N$$

- $\mathbb{E}[\hat{\gamma}'_{vv}(\tau)] = \frac{N-|\tau|}{N} \cdot \gamma_{vv}(\tau)$, ovvero lo stimatore è **distorto**, ma **asintoticamente corretto**
- Per τ fissato, lo stimatore è **consistente**, sotto le ipotesi di ergodicità

Densità spettrale campionaria

Non conoscendo $\gamma_{vv}(\tau)$, uso $\hat{\gamma}_{vv}(\tau)$ oppure $\hat{\gamma}'_{vv}(\tau)$. Si definisce **periodogramma** il seguente **stimatore** della densità spettrale di potenza:

$$\text{Stimatore: periodogramma} \quad I_N(\omega) = \sum_{\tau=-(N-1)}^{N-1} \hat{\gamma}'_{vv}(\tau) \cdot e^{-j\omega\tau}$$

- A differenza di $\Gamma_{vv}(\omega)$, $I_N(\omega)$ è definito solo da $\tau = -(N-1)$ a $\tau = N-1$
- Essendo la DTFT di $\hat{\gamma}'_{vv}(\tau)$, $I_N(\omega)$ è una funzione **reale, continua, 2π periodica**

Proprietà dello stimatore

- Lo stimatore $I_N(\omega)$ **non è corretto**, ma è **asintoticamente corretto**. Notiamo che non lo sarebbe stato neanche se avessi usato $\hat{\gamma}_{vv}(\tau)$ al posto di $\hat{\gamma}'_{vv}(\tau)$, infatti:

$$\begin{aligned} \mathbb{E}[I_N(\omega)] &= \sum_{\tau=-(N-1)}^{N-1} \underbrace{\mathbb{E}[\hat{\gamma}'_{vv}(\tau)]}_{\frac{N-|\tau|}{N} \cdot \gamma_{vv}(\tau)} \cdot e^{-j\omega\tau} = \\ \mathbb{E}[I_N(\omega)] &= \sum_{\tau=-(N-1)}^{N-1} \underbrace{\frac{N-|\tau|}{N} \cdot \gamma_{vv}(\tau)}_{\neq \Gamma_{vv}(\tau)} \cdot e^{-j\omega\tau} \neq \Gamma_{vv}(\omega) \end{aligned}$$

- Si dimostra che come $\text{Var}[I_N(\omega)] \approx \Gamma_{vv}^2(\omega)$. Per cui, la varianza dello stimatore non decresce al crescere di N . Lo stimatore **non è consistente**
- Per $N \rightarrow +\infty$, $I_N(\omega_1)$ e $I_N(\omega_2)$ tendono a **diventare incorrelati**, $\forall \omega_1 \neq \omega_2$. Questo ci dà l'idea che il periodogramma sia una funzione **poco continua**, poiché la stima in una frequenza può non essere simile alla stima in una frequenza anche adiacente (una sorta di **rumore bianco in frequenza**).

Per migliorare la stima - Metodo di Bartlett

- Ipotizziamo di avere N dati a disposizione
- Dividiamo i dati in $K = \frac{N}{M}$ parti, dove M è la lunghezza di ogni porzione di dati
- Calcoliamo il periodogramma $I_{M,K}^{[i]}(\omega)$ per ciascuna parte $i = 1, 2, \dots, K$

- facciamo la media dei periodogrammi, ottenendo la stima

$$\bar{I}_{M,K}(\omega) = \frac{1}{K} \sum_{i=1}^K I_{M,K}^{[i]}(\omega)$$

Osservazioni:

- se $\gamma_{vv}(\tau) \rightarrow 0$ in modo sufficientemente rapido, i K periodogrammi sono **circa indipendenti**. In questo caso, si ha che

$$Var[\bar{I}_{M,K}(\omega)] = O\left(\frac{1}{K} \cdot \Gamma_{vv}^2(\omega)\right)$$

- Il $Bias[\bar{I}_{M,K}(\omega)]$ è maggiore rispetto a quello di $I_N(\omega)$. Questo comporta una maggior **perdita di informazione in frequenza**
 - Se so che $\Gamma_{vv}(\omega)$ ha **pichi molto stretti**, devo usare M **grande** in modo da avere abbastanza risoluzione in frequenza
-

6 Lezione 9: Famiglie di modelli stocastici

Teorema dei processi ARMA

Definizione: Un processo stocastico $y(t)$, generato a partire dal rumore bianco $e(t) \sim WN(\mu, \lambda^2)$, é detto di tipo **ARMA**(n_a, n_c), se:

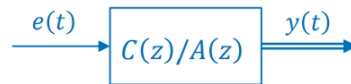
$$y(t) = a_1 \cdot y(t-1) + a_2 \cdot y(t-2) + \dots + a_{n_a} \cdot y(t-n_a) \quad \text{Parte AR}(n_a) \\ + e(t) + c_1 \cdot e(t-1) + c_2 \cdot e(t-2) + \dots + c_{n_c} \cdot e(t-n_c) \quad \text{Parte MA}(n_c)$$

$$a_1, \dots, a_{n_a} : \text{coefficienti del modello AR}(n_a) \quad n_a : \text{ordine del modello AR}(n_a) \\ c_0, c_1, \dots, c_{n_c} : \text{coefficienti del modello MA}(n_c) \quad n_c : \text{ordine del modello MA}(n_c)$$

La **funzione di trasferimento** di un $ARMA(n_a, n_c)$ risulta essere:

$$y(t)[1 - a_1 z^{-1} - a_2 z^{-2} - \dots - a_{n_a} z^{-n_a}] = [1 + c_1 z^{-1} + c_2 z^{-2} + \dots + c_{n_c} z^{-n_c}]e(t)$$

$$y(t) = \frac{1 + c_1 z^{-1} + c_2 z^{-2} + \dots + c_{n_c} z^{-n_c}}{1 - a_1 z^{-1} - a_2 z^{-2} - \dots - a_{n_a} z^{-n_a}} e(t) \quad \Rightarrow \quad y(t) = \frac{C(z)}{A(z)} e(t)$$



Il processo $y(t)$ é stazionario se e solo se $\frac{C(z)}{A(z)}$ é asintoticamente stabile.

Teorema: Dato un processo stocastico stazionario $ARMA(n_a, n_c)$, esso può essere descritto come un $MA(\infty)$.

Esempio: Supponiamo di avere un $AR(1)$ del tipo

$$y(t) = ay(t-1) + e(t) \quad e(t) \sim WN(0, \lambda^2)$$

$$y(t) = \frac{1}{1 - az^{-1}} e(t) \quad \text{può essere visto come il} \\ \text{limite di una serie} \\ \text{geometrica di ragione } az^{-1}$$

$$\Rightarrow \quad = \sum_{i=0}^{+\infty} (az^{-1})^i \cdot e(t) = \sum_{i=0}^{+\infty} a^i \cdot e(t-i) \quad \text{MA}(\infty)$$

7 Lezione 10: Predizione

Spiegare il concetto di filtro passa-tutto e le sue particolarità. Enunciare inoltre il teorema di fattorizzazione spettrale (forma canonica) e come il filtro passa tutto possa essere sfruttato per raggiungere la forma canonica.

Filtro passa-tutto

Definizione: il filtro passa-tutto è un filtro di ordine 1 definito come:

$$T(z) = \frac{1}{a} \cdot \frac{z + a}{z + \frac{1}{a}} \quad a \neq 0, a \in \mathbb{R}$$

Lo zero è reciproco del polo

$$\text{zero} : z = -a \quad \text{polo} : z = -\frac{1}{a}$$

Il fattore moltiplicativo è come il polo

Osservazioni:

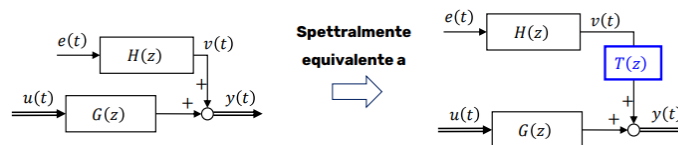
- La densità spettrale di potenza $\Gamma_{yy}(\omega)$ di un processo $y(t)$ in uscita dal passa tutto $T(z)$ alimentato da un generico processo stazionario in ingresso $v(t)$ è uguale a:

$$\Gamma_{yy}(\omega) = |T(e^{j\omega})|^2 \cdot \Gamma_{vv}(\omega)$$

Il filtro passa-tutto non modifica il modulo delle frequenze nella densità spettrale di potenza dell'ingresso. Quindi, si ha che:

$$\Gamma_{yy}(\omega) = \Gamma_{vv}(\omega)$$

- Il processo $y(t)$ in uscita al passa-tutto è **spetttralmente equivalente** al processo $v(t)$ in ingresso al passatutto, tuttavia i due processi $y(t)$ e $v(t)$ **non sono identici** poiché il passatutto introduce uno **sfasamento**



Filtro passa-tutto e forma canonica

Vogliamo risolvere il problema della predizione per **processi a spettro razionale**, ovvero processi $y(t)$ generati in uscita da un sistema dinamico lineare asintoticamente stabile con funzione di trasferimento $H(z)$ razionale fratta alimentato da $e(t) \sim WN(0, \lambda^2)$.

Il problema della fattorizzazione spettrale consiste nel trovare tutte le coppie $\{H(z), \lambda^2\}$ tali che:

$$\Phi_{yy}(z) = \lambda^2 \cdot H(z) \cdot H(z^{-1})$$

Per processi a spettro razionale, esistono **infiniti fattori spettrali** $\{H(z), \lambda^2\}$. Ai fini della predizione ottima ci servirà un fattore spettrale detto **canonico**.

Il filtro viene applicato alla funzione di trasferimento $H(z)$ in modo tale da eliminare eventuali radici che non sono interne al cerchio unitario e rendere dunque il processo in forma canonica.

Teorema della fattorizzazione spettrale

Dato un processo stocastico stazionario a spettro razionale, esiste **un solo fattore spettrale** $\{\tilde{H}(z), \tilde{\lambda}^2\}$, detto **fattore spettrale canonico**, dove $\tilde{H}(z) = \frac{C(z)}{A(z)}$ tale che:

- $C(z)$ e $A(z)$ hanno lo **stesso grado** (grado relativo nullo)
- $C(z)$ e $A(z)$ sono **coprime** (non ci son fattori in comune)
- $C(z)$ e $A(z)$ sono **monici** (il coefficiente del termine di grado massimo è 1)
- $C(z)$ e $A(z)$ hanno **radici interne al cerchio unitario**

Si consideri un processo ARMA dato in rappresentazione canonica. Spiegare perché la varianza dell'errore di predizione a un passo coincide con la varianza del rumore bianco che genera il processo.

Sia dato un processo $ARMA(n_a, n_c)$ in forma canonica:

$$y(t) = H(z) \cdot e(t) = \frac{C(z)}{A(z)} \cdot e(t) \quad e(t) \sim WN(0, \lambda^2)$$

$$C(z) = 1 + c_1 z^{-1} + c_2 z^{-2} + \dots + c_{n_c} z^{-n_c}$$

$$A(z) = 1 - a_1 z^{-1} - a_2 z^{-2} - \dots - c_{n_a} z^{-n_a}$$

Per scomporre la parte imprevedibile da quella predicibile viene impiegata la **lunga divisione**, in questo modo possiamo esprimere $C(z)/A(z)$ come:

$$\frac{C(z)}{A(z)} = E(z) + \frac{R(z)}{A(z)} = E(z) + \frac{z^{-k} \tilde{R}(z)}{A(z)} \quad \text{con } k = \text{passi di lunga divisione}$$

$$y(t) = \underbrace{E(z) \cdot e(t)}_{\text{parte impred.}} + \underbrace{\frac{\tilde{R}(z)}{A(z)} \cdot e(t-k)}_{\text{parte pred}}$$

Filtro sbiancante:

$$y(t) = \frac{C(z)}{A(z)} \cdot e(t) \quad \Rightarrow \quad e(t) = \frac{A(z)}{C(z)} \cdot y(t)$$

Il predittore ottimo dal rumore

$$\hat{y}(t|t-k) = \frac{\tilde{R}(z)}{A(z)} \cdot e(t-k)$$

Il predittore ottimo dai dati

$$\hat{y}(t|t-k) = \frac{\tilde{R}(z)}{A(z)} \cdot e(t-k) = \frac{\tilde{R}(z) \cdot z^{-k}}{A(z)} \cdot e(t) = \frac{\tilde{R}(z) \cdot z^{-k}}{A(z)} \cdot \frac{A(z)}{C(z)} \cdot y(t) = \frac{\tilde{R}(z)}{C(z)} \cdot y(t-k)$$

$$\hat{y}(t|t-k) = \frac{\tilde{R}(z)}{C(z)} \cdot y(t-k)$$

L'errore di predizione

$$\varepsilon_k(t) = y(t) - \hat{y}(t|t-k) = E(z) \cdot e(t)$$

Caso particolare: predizione ad un passo $k = 1$:

$$\bullet E(z) = 1 \quad \bullet R(z) = C(z) - A(z)$$

$$\hat{y}(t|t-1) = \frac{C(z) - A(z)}{C(z)} \cdot y(t)$$

$$\varepsilon_1(t) = E(z) \cdot e(t) = e(t)$$

Il predittore lineare ottimo dai dati é quello che minimizza il seguente criterio Mean Squared Error (MSE):

$$\text{var}[\varepsilon_k(t)] = \mathbb{E}[\varepsilon_k(t)^2] = \mathbb{E}[(y(t) - \hat{y}(t|t-k))^2]$$

Affiché il predittore sia ottimo occorre che:

- $\mathbb{E}[\varepsilon_k(t)] = \mathbb{E}[y(t) - \hat{y}(t|t-k)] = 0 \Rightarrow$ sia **Corretto**, ossia valore atteso nullo.
- $\mathbb{E}[\hat{y}(t|t-k) \cdot \varepsilon_k(t)] = 0 \Rightarrow$ il predittore e l'errore di predizione sono **incorrelati**
- $\text{Var}[\varepsilon_k(t)] = \mathbb{E}[\varepsilon_k(t)^2] = \mathbb{E}[(y(t) - \hat{y}(t|t-k))^2] = \mathbb{E}[(E(z) \cdot e(t))^2] = \mathbb{E}[e(t)^2]$ minima

8 Lezione 11: Identificazione: concetti fondamentali

Descrivere il metodo della massima verosimiglianza per la stima dei modelli ARMAX

I metodi di stima basati sulla minimizzazione dell'errore di predizione prendono il nome di **Prediction Error Methods (PEM)**.

Se ipotizzo che $S = \mathcal{M}(\theta^0)$ e $e(t) \sim WN$ Gaussiano, lo **stimatore PEM è circa uguale allo stimatore a massima verosimiglianza**. La differenza sta in come i due approcci trattano l'inizializzazione del predittore, ma se i dati sono molti, non c'è differenza.

Consideriamo un modello $ARMAX(n_a, n_c, n_b, k = 1)$, (dove k è il ritardo puro tra ingresso $u(t)$ ed uscita $y(t)$) e di avere a disposizione N dati $\{u(1), \dots, u(n)\}$ e $\{y(1), \dots, y(n)\}$

$$\begin{aligned}\mathcal{M}(\theta) : \quad y(t) &= \frac{B(z, \theta)}{A(z, \theta)} \cdot u(t-1) + \frac{C(z, \theta)}{A(z, \theta)} \cdot e(t) \quad e(t) \sim WN(0, \lambda^2) \\ B(z) &= b_0 + b_1 z^{-1} + \dots + b_{n_b} z^{-n_b} \\ A(z) &= 1 - a_1 z^{-1} - \dots - a_{n_a} z^{-n_a} \\ C(z) &= 1 + c_1 z^{-1} + \dots + c_{n_c} z^{-n_c}\end{aligned}$$

Calcoliamo l'espressione dell'errore di predizione ad un passo. In questo caso, si ha che $E(z) = 1$, e quindi $\varepsilon_1(t) = e(t)$. Di conseguenza, esprimendo $e(t)$ in funzione di $u(t)$ e $y(t)$:

$$\varepsilon_1(t; \theta) = e(t) = \frac{A(z, \theta)}{C(z, \theta)} y(t) - \frac{B(z, \theta)}{C(z, \theta)} u(t-1)$$

$$\text{oppure tramite:} \quad \varepsilon_1(t; \theta) = H^{-1}(z, \theta) \cdot [y(t) - G(z, \theta) \cdot u(t)]$$

Utilizziamo l'approccio predittivo:

$$\begin{aligned}J_N(\theta) &= \frac{1}{N} \sum_{t=1}^N \varepsilon_1(t; \theta)^2 = \frac{1}{N} \sum_{t=1}^N \left[\frac{A(z, \theta)}{C(z, \theta)} y(t) - \frac{B(z, \theta)}{C(z, \theta)} u(t-1) \right]^2 \\ \hat{\theta} &= \arg \min_{\theta} J_N(\theta)\end{aligned}$$

Osservazioni:

- Dato che ho $C(z, \theta)$ al denominatore, questa funzione di costo non è più convessa. In generale avrò dei **minimi locali**
- Per la risoluzione del problema di minimizzazione, devo utilizzare dei **metodi iterativi** per esempio **il metodo del gradiente**, oppure un'alternativa (più efficiente) al metodo del gradiente il metodo di ottimizzazione iterativo noto come **Metodo di Newton**. Questo metodo, oltre al gradiente, sfrutta anche l'informazione data dalla **matrice Hessiana**.

Rappresentare la cifra di merito e lo schema a blocchi del filtraggio dell'algoritmo di identificazione dei processi ARMAX. Si descriva l'approssimazione dell'Hessiana utilizzata, indicando le motivazioni che inducono a questa approssimazione

Cifra di merito o funzione di costo

$$J_N(\theta) = \frac{1}{N} \sum_{t=1}^N \varepsilon_1(t; \theta)^2 = \frac{1}{N} \sum_{t=1}^N \left[\frac{A(z, \theta)}{C(z, \theta)} y(t) - \frac{B(z, \theta)}{C(z, \theta)} u(t-1) \right]^2$$

Osservazioni:

- Dato che ho $C(z, \theta)$ al denominatore, questa funzione di costo non è più convessa. In generale avrò dei **minimi locali**
- Per la risoluzione del problema di minimizzazione, devo utilizzare dei **metodi iterativi** per esempio **il metodo del gradiente**, oppure un'alternativa (più efficiente) al metodo del gradiente il metodo di ottimizzazione iterativo noto come **Metodo di Newton**. Questo metodo, oltre al gradiente, sfrutta anche l'informazione data dalla **matrice Hessiana**.

Metodo di Newton

Sviluppo in serie di Taylor troncata al 2° ordine di $J_N(\theta)$, nell'intorno della stima all'iterazione i -esima $\hat{\theta}^{(i)}$ (approssimiamo la funzione di costo con il paraboloide).

$$J_N(\theta) \approx V(\theta)$$

La funzione $V(\theta)$ è un **paraboloide**
(è facile da calcolare il minimo)

$$V^{(i)}(\theta) = J_N(\hat{\theta}^{(i)}) + \underbrace{(\theta - \hat{\theta}^{(i)})^\top \cdot \frac{dJ_N(\theta)}{d\theta} \Big|_{\theta=\hat{\theta}^{(i)}}}_{\text{Gradiente}} + \frac{1}{2} \underbrace{(\theta - \hat{\theta}^{(i)})^\top \cdot \frac{d^2 J_N(\theta)}{d\theta^2} \Big|_{\theta=\hat{\theta}^{(i)}} \cdot (\theta - \hat{\theta}^{(i)})}_{\text{Matrice Hessiana}}$$

Una volta ottenuta l'approssimazione $V^{(i)}(\theta)$, si calcola $\hat{\theta}^{(i+1)}$ come minimo di $V^{(i)}(\theta)$. Troviamo un'espressione esplicita per $\hat{\theta}^{(i+1)}$ imponendo:

$$\frac{dV^{(i)}(\theta)}{d\theta} = \mathbf{0}_{d \times 1}$$

$$\frac{dV^{(i)}(\theta)}{d\theta} = \frac{dJ_N(\theta)}{d\theta} \Big|_{\theta=\hat{\theta}^{(i)}} + \frac{1}{2} \cdot 2 \cdot \frac{d^2 J_N(\theta)}{d\theta^2} \Big|_{\theta=\hat{\theta}^{(i)}} \cdot (\theta - \hat{\theta}^{(i)}) = \mathbf{0} \quad \Rightarrow \quad \text{Ricavo il minimo e lo chiamo } \hat{\theta}^{(i+1)}$$

Da seguente formula ricaviamo la regola di update per il metodo di Newton:

$$\hat{\theta}^{(i+1)} = \hat{\theta}^{(i)} - \left[\frac{d^2 J_N(\theta)}{d\theta^2} \Big|_{\theta=\hat{\theta}^{(i)}} \right]^{-1} \cdot \frac{dJ_N(\theta)}{d\theta} \Big|_{\theta=\hat{\theta}^{(i)}}$$

Gradiente di $J_N(\theta)$:

$$\frac{dJ_N(\theta)}{d\theta} = \frac{d}{d\theta} \cdot \frac{1}{N} \sum_{t=1}^N \varepsilon_1(t; \theta)^2 = \frac{1}{N} \sum_{t=1}^N \frac{d}{d\theta} \varepsilon_1(t; \theta)^2 = \frac{2}{N} \sum_{t=1}^N \varepsilon_1(t; \theta) \cdot \frac{d\varepsilon_1(t; \theta)}{d\theta}$$

Hessiano di $J_N(\theta)$:

$$\frac{d^2 J_N(\theta)}{d\theta^2} = \frac{2}{N} \sum_{t=1}^N \frac{d\varepsilon_1(t; \theta)}{d\theta} \cdot \frac{d\varepsilon_1(t; \theta)}{d\theta} + \frac{2}{N} \sum_{t=1}^N \varepsilon_1(t; \theta) \cdot \frac{d^2 \varepsilon_1(t; \theta)}{d\theta^2}$$

Ignoriamo il secondo termine, approssimando l'Hessiana, dato che:

- Se siamo vicini all'ottimo, $\varepsilon_1(t; \theta)$ è piccolo e il termine conta poco
- Possiamo evitare di calcolare

$$\frac{d^2 \varepsilon_1(t; \theta)}{d\theta^2}$$

- Ci assicuriamo una **Hessiana semi-definita positiva**. In questo modo, la direzione dell'algoritmo è sicuramente discendente (concetto simile ad avere learning rate ≥ 0)

Dopo aver introdotto l'approssimazione dell'Hessiana, la regola di update diventa:

$$\hat{\boldsymbol{\theta}}_{d \times 1}^{(i+1)} = \hat{\boldsymbol{\theta}}_{d \times 1}^{(i)} - \left[\frac{2}{N} \sum_{t=1}^N \frac{d\varepsilon_1(t; \boldsymbol{\theta})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}^{(i)}} \cdot \frac{d\varepsilon_1(t; \boldsymbol{\theta})^\top}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}^{(i)}} \right]^{-1} \cdot \left[\frac{2}{N} \sum_{t=1}^N \varepsilon_1(t; \hat{\boldsymbol{\theta}}^{(i)}) \cdot \frac{d\varepsilon_1(t; \boldsymbol{\theta})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}^{(i)}} \right]_{d \times 1}$$

$$\hat{\boldsymbol{\theta}}^{(i+1)} = \hat{\boldsymbol{\theta}}^{(i)} - [\text{hessiana}]^{-1} \cdot [\text{gradiente}]$$

Ricordando che:

$$\varepsilon_1(t; \boldsymbol{\theta}) = e(t) = \frac{A(z, \boldsymbol{\theta})}{C(z, \boldsymbol{\theta})} y(t) - \frac{B(z, \boldsymbol{\theta})}{C(z, \boldsymbol{\theta})} u(t-1)$$

$$\varepsilon_1(t; \boldsymbol{\theta}) = \frac{1 - a_1 z^{-1} - \dots - a_{n_a} z^{-n_a}}{1 + c_1 z^{-1} + \dots + c_{n_c} z^{-n_c}} y(t) - \frac{b_0 + b_1 z^{-1} + \dots + b_{n_b} z^{-n_b}}{1 + c_1 z^{-1} + \dots + c_{n_c} z^{-n_c}} u(t-1)$$

$$\boldsymbol{\theta} = \begin{bmatrix} a_1 \dots a_{n_a} & b_0 b_1 \dots b_{n_b} & c_1 \dots c_{n_c} \end{bmatrix}^T$$

$$d \times 1 = n_a + n_b + 1 + n_c \times 1$$

Derivate di $\varepsilon_1(t; \boldsymbol{\theta})$ rispetto a a_1, a_2, \dots, a_{n_a}

$$\frac{d\varepsilon_1(t)}{da_1} = -\frac{z^{-1}}{C(z)} y(t) = \alpha(t-1) \quad \alpha(t) \equiv -\frac{1}{C(z)} y(t)$$

Derivate di $\varepsilon_1(t; \boldsymbol{\theta})$ rispetto a b_0, b_1, \dots, b_{n_b}

$$\frac{d\varepsilon_1(t)}{db_0} = -\frac{1}{C(z)} u(t-1) = \beta(t-1) \quad \beta(t) \equiv -\frac{1}{C(z)} u(t)$$

Derivate di $\varepsilon_1(t; \boldsymbol{\theta})$ rispetto a c_1, c_2, \dots, c_{n_c}

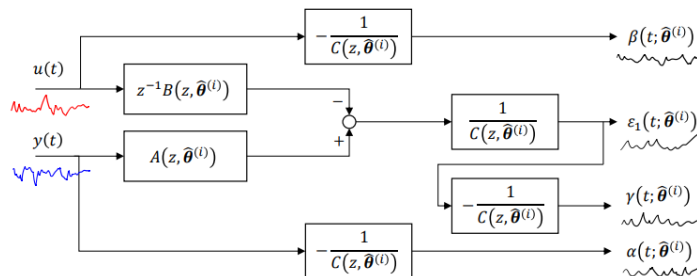
$$\frac{d[C(z) \cdot \varepsilon_1(t)]}{dc_1} = 0 \Rightarrow z^{-1} \varepsilon_1(t) + C(z) \frac{d\varepsilon_1(t)}{dc_1} = 0$$

$$\frac{d\varepsilon_1(t)}{dc_1} = -\frac{1}{C(z)} \varepsilon(t-1) = \gamma(t-1) \quad \gamma(t) \equiv -\frac{1}{C(z)} \varepsilon_1(t)$$

Riassumendo, il vettore gradiente è:

$$\frac{d\varepsilon_1(t)}{d\boldsymbol{\theta}}_{d \times 1} = \begin{bmatrix} \alpha(t-1) \\ \vdots \\ \alpha(t-n_a) \\ \beta(t-1) \\ \vdots \\ \beta(t-n_b-1) \\ \gamma(t-1) \\ \vdots \\ \gamma(t-n_c) \end{bmatrix} \quad t = 1, \dots, N$$

È possibile definire in modo elegante il calcolo del gradiente tramite una **serie di filtraggi dei segnali di ingresso e uscita**



9 Lezione 12: Identificazione analisi e complementi

Si dimostri che, se il sistema da identificare appartiene alla classe dei modelli e il minimo globale della cifra di merito è unico, tale minimo globale corrisponde al sistema che genera i dati.

Ipotesi di lavoro: Sia l'ingresso $u(t, s)$ e sia l'uscita $y(t, s)$ sono processi stocastici stazionari ed ergodici. Di conseguenza, i dati misurati saranno una realizzazione dei processi $u(t, s)$ e $y(t, s)$ in corrispondenza di un particolare esito \bar{s} .

La funzione di costo dipende anch'essa dall'esito \bar{s} poiché utilizza i dati misurati da cui ottengo la stima $\hat{\theta}_N(\bar{s})$, dunque la stima $\hat{\theta}_N(s)$ è una **variabile causale**.

$$J_N(\theta, \bar{s}) = \sum_{t=1}^N \varepsilon_1(t; \theta, \bar{s})^2$$

La funzione di costo $J_N(\theta, s)$ dovrebbe essere interpretata come un **insieme di curve** le quali dipendono da s e la stima $\hat{\theta}_N(s)$ come un **un insieme di punti**.

Osservazione: grazie all'**ipotesi di ergodicità** di $u(t, s)$ e $y(t, s)$ abbiamo che i momenti temporali convergono ai rispettivi momenti di insieme \Rightarrow cioè le curve $J_N(\theta, s)$ convergono ad **unica (deterministica) curva** $\bar{J}(\theta)$.

$$J_N(\theta, s) = \frac{1}{N} \sum_{t=1}^N \varepsilon_1(t; \theta, s)^2 \xrightarrow{N \rightarrow +\infty} \bar{J}(\theta) \equiv \mathbb{E}_s[\varepsilon_1(t; \theta)^2]$$

Definiamo inoltre l'**insieme dei punti di minimo globale** di $\bar{J}(\theta)$ come:

$$\Delta_\theta = \{\bar{\theta} \mid \bar{J}(\theta) \geq \bar{J}(\bar{\theta}), \quad \forall \theta\}$$

Da dimostrare: se $\mathcal{S} \in \mathcal{M}(\theta)$ e $\Delta_\theta = \bar{\theta}$, allora $\bar{\theta} = \theta^0$ dove $\mathcal{S} = \mathcal{M}(\theta^0)$.

Ipotesi di lavoro:

- Assumiamo che $\mathcal{S} \in \mathcal{M}(\theta)$, ovvero che esista $\theta^0 \in \Theta$ tale che $\mathcal{S} = \mathcal{M}(\theta^0)$.
- $\theta^0 \in \Delta_\theta$ di $\bar{J}(\theta)$? ossia la stima $\hat{\theta}_N$ tende asintoticamente a θ^0 ? I metodi PEM sono in grado di trovare la parametrizzazione vera del modello?

Dimostrazione: dimostriamo che, sotto le ipotesi fatte θ^0 appartiene sempre a Δ_θ

Supponiamo che i dati siano generati dal sistema \mathcal{S} , tale che:

$$y(t) = \hat{y}(t|t-1; \theta^0) + e(t), \quad e(t) \sim WN(0, \lambda^2)$$

*Consideriamo un generico modello $\mathcal{M}(\theta)$,
senza assumere che $\varepsilon_1(t; \theta)$ sia bianco*

$$y(t) = \hat{y}(t|t-1; \theta) + \varepsilon_1(t; \theta)$$

L'errore di predizione ad un passo commesso dal modello $\mathcal{M}(\theta)$ è dunque:

$$\varepsilon_1(t; \theta) = y(t) - \hat{y}(t|t-1; \theta)$$

Aggiungiamo e togliamo $\hat{y}(t|t-1; \theta^0)$, ovvero il predittore del sistema \mathcal{S} che genera i dati

$$\varepsilon_1(t; \theta) = y(t) - \hat{y}(t|t-1; \theta^0) + \hat{y}(t|t-1; \theta^0) - \hat{y}(t|t-1; \theta)$$

$$\varepsilon_1(t; \theta) = \underbrace{y(t) - \hat{y}(t|t-1; \theta^0)}_{\substack{\text{Err. di predizione ottimo} \\ \varepsilon_1(t; \theta^0) = e(t)}} + \hat{y}(t|t-1; \theta^0) - \hat{y}(t|t-1; \theta)$$

$$\text{Errore di predizione ottimo} = y(t) - \hat{y}(t|t-1; \theta^0) = \varepsilon_1(t; \theta^0) = e(t)$$

$$\Rightarrow \varepsilon_1(t; \theta) = e(t) + \hat{y}(t|t-1; \theta^0) - \hat{y}(t|t-1; \theta)$$

Calcoliamo la varianza dell'errore di predizione (che è a media nulla poiché il predittore è corretto):

$$\begin{aligned}\mathbb{E}[\varepsilon_1(t; \boldsymbol{\theta})^2] &= \mathbb{E}\left[\left(e(t) + \hat{y}(t|t-1; \boldsymbol{\theta}^0) - \hat{y}(t|t-1; \boldsymbol{\theta})\right)^2\right] \\ \Rightarrow \bar{J}(\boldsymbol{\theta}) &= \mathbb{E}[e(t)^2] + \mathbb{E}\left[\left(\hat{y}(t|t-1; \boldsymbol{\theta}^0) - \hat{y}(t|t-1; \boldsymbol{\theta})\right)^2\right] + \underbrace{2\mathbb{E}\left[e(t) \cdot \left(\hat{y}(t|t-1; \boldsymbol{\theta}^0) - \hat{y}(t|t-1; \boldsymbol{\theta})\right)\right]}_{\text{incorrelati}} \\ &= \lambda^2 + \underbrace{\mathbb{E}\left[\left(\hat{y}(t|t-1; \boldsymbol{\theta}^0) - \hat{y}(t|t-1; \boldsymbol{\theta})\right)^2\right]}_{\text{varianza}}\end{aligned}$$

È una varianza, quindi una quantità ≥ 0 .
In particolare si annulla solo per $\boldsymbol{\theta} = \boldsymbol{\theta}^0$

$$\Rightarrow \bar{J}(\boldsymbol{\theta}) \geq \lambda^2 = \bar{J}(\boldsymbol{\theta}^0), \quad \forall \boldsymbol{\theta}$$

$$\bar{J}(\boldsymbol{\theta}) \geq \bar{J}(\boldsymbol{\theta}^0) \quad \forall \boldsymbol{\theta}$$

$\boldsymbol{\theta}^0$ è un minimo di $\bar{J}(\boldsymbol{\theta})$

Conclusione (fondamentale): Se $\mathcal{S} \in \mathcal{M}(\boldsymbol{\theta})$ e $u(t), y(t)$ sono pss ergodici, allora, per $N \rightarrow +\infty$, un metodo PEM garantisce che il modello stimato è quello vero $\mathcal{S} = \mathcal{M}(\boldsymbol{\theta}^0)$ o un insieme equivalente di modelli $\{\mathcal{M}(\boldsymbol{\theta}) \mid \boldsymbol{\theta} \in \Delta_{\boldsymbol{\theta}}\}$ con la stessa capacità nello spiegare i dati.

- Se $\mathcal{S} \in \mathcal{M}(\boldsymbol{\theta})$, allora in corrispondenza di $\boldsymbol{\theta}^0$ si ha che $\varepsilon_1(t; \boldsymbol{\theta}^0) = e(t) \sim WN$ basta fare test bianchezza sui residui $\varepsilon_1(t; \hat{\boldsymbol{\theta}}_N)$.

$$\mathcal{S} \in \mathcal{M}(\boldsymbol{\theta}) \text{ e } \Delta_{\boldsymbol{\theta}} = \bar{\boldsymbol{\theta}}, \text{ allora } N_{\rightarrow \infty} \hat{\boldsymbol{\theta}}_N \approx \bar{\boldsymbol{\theta}} = \boldsymbol{\theta}^0 \quad \text{modello perfetto}$$

$$\mathcal{S} \in \mathcal{M}(\boldsymbol{\theta}) \text{ e } \Delta_{\boldsymbol{\theta}} = \{\boldsymbol{\theta}_i, \dots\}, \text{ allora } N_{\rightarrow \infty} \hat{\boldsymbol{\theta}}_N \approx \bar{\boldsymbol{\theta}} \neq \boldsymbol{\theta}^0 \quad \mathcal{M}(\bar{\boldsymbol{\theta}}) \approx \mathcal{M}(\boldsymbol{\theta}^0)$$

- Se $\mathcal{S} \notin \mathcal{M}(\boldsymbol{\theta})$, allora i metodi PEM non garantiscono di stimare correttamente TUTTE le componenti del sistema \mathcal{S} , ma il modello trovato $\mathcal{M}(\bar{\boldsymbol{\theta}})$ risulta essere la migliore approssimazione di \mathcal{S} nella famiglia di modelli $\mathcal{M}(\boldsymbol{\theta})$

$$\mathcal{S} \notin \mathcal{M}(\boldsymbol{\theta}) \text{ e } \Delta_{\boldsymbol{\theta}} = \bar{\boldsymbol{\theta}}, \text{ allora } N_{\rightarrow \infty} \hat{\boldsymbol{\theta}}_N \approx \bar{\boldsymbol{\theta}} \neq \boldsymbol{\theta}^0 \quad \mathcal{M}(\bar{\boldsymbol{\theta}}) \sim \mathcal{M}(\boldsymbol{\theta}^0)$$

$$\mathcal{S} \notin \mathcal{M}(\boldsymbol{\theta}) \text{ e } \Delta_{\boldsymbol{\theta}} = \{\boldsymbol{\theta}_i, \dots\}, \text{ allora } N_{\rightarrow \infty} \hat{\boldsymbol{\theta}}_N \approx \bar{\boldsymbol{\theta}} \neq \boldsymbol{\theta}^0 \quad \mathcal{M}(\bar{\boldsymbol{\theta}}) \sim \mathcal{M}(\boldsymbol{\theta}^0)$$

Identificabilità sperimentale e segnale persistentemente eccitante

Vedere risposta successiva.

Spiegare cosa si intende per identificabilità sperimentale e strutturale precisando la nozione di persistente eccitazione di un segnale

Identificabilità sperimentale e strutturale

L'analisi asintotica ci dice che se $\mathcal{S} \in \mathcal{M}(\boldsymbol{\theta})$ e $u(t), y(t)$ sono pss ergodici, allora, per $N \rightarrow +\infty$, un metodo PEM garantisce che il modello stimato è quello vero $\mathcal{S} = \mathcal{M}(\boldsymbol{\theta}^0)$ o un insieme equivalente di modelli $\{\mathcal{M}(\boldsymbol{\theta}) \mid \boldsymbol{\theta} \in \Delta_{\boldsymbol{\theta}}\}$.

Nel secondo caso dobbiamo capire in quali condizioni il sistema \mathcal{S} può essere **identificato univocamente** dai dati. Affinché un modello sia univocamente identificabile è necessario avere:

- **Identificabilità strutturale:** il modello $\mathcal{M}(\boldsymbol{\theta})$ non deve essere *sovraparametrizzato* rispetto al sistema \mathcal{S} .
- **Identificabilità sperimentale:** i dati $\{u(t), y(t)\}_{t=1}^N$ devono contenere *sufficiente informazione*

Il problema di non identificabilità più critico è quello sperimentale: se non abbiamo sufficiente informazione nei dati, non possiamo fare nulla. La non identificabilità strutturale è, invece, facilmente risolvibile riducendo l'ordine del modello.

Persistente eccitazione di un segnale

Vedere domande successive.

Dire cosa si intende per identificabilità di un processo ARX e descrivere i principali problemi di identificabilità

Identificabilità dei modelli ARX

Dato un modello $ARX(n_a, n_b, 1)$, avendo N dati $\{u(1), \dots, u(N)\}$, $\{y(1), \dots, y(N)\}$ risulta essere:

$$\begin{aligned}\mathcal{M}(\boldsymbol{\theta}) : \quad y(t) &= \frac{B(z, \boldsymbol{\theta})}{A(z, \boldsymbol{\theta})} \cdot u(t-1) + \frac{1}{A(z, \boldsymbol{\theta})} \cdot e(t) \quad e(t) \sim WN(0, \lambda^2) \\ B(z) &= b_0 + b_1 z^{-1} + \dots + b_{n_b} z^{-n_b} \\ A(z) &= 1 + a_1 z^{-1} + \dots + a_{n_a} z^{-n_a}\end{aligned}$$

La stima tramite il metodo dei minimi quadrati risulta essere:

$$\hat{\boldsymbol{\theta}}_N = \left[\sum_{t=1}^N \boldsymbol{\varphi}(t) \cdot \boldsymbol{\varphi}^T(t) \right]^{-1} \cdot \left[\sum_{t=1}^N \boldsymbol{\varphi}(t) \cdot y(t) \right]$$

Problema di identificabilità: quando $\hat{\boldsymbol{\theta}}_N$ esiste ed è unico? \Leftrightarrow quando $\sum_{t=1}^N \boldsymbol{\varphi}(t) \cdot \boldsymbol{\varphi}^T(t)$ è invertibile?

Definiamo:

$$\begin{aligned}S(N) &= \sum_{t=1}^N \boldsymbol{\varphi}(t) \cdot \boldsymbol{\varphi}^T(t) \quad \Rightarrow \quad \hat{\boldsymbol{\theta}}_N = S(N)^{-1} \cdot \left[\sum_{t=1}^N \boldsymbol{\varphi}(t) \cdot y(t) \right] \\ R(N) &= \frac{1}{N} S(N) \quad \Rightarrow \quad \hat{\boldsymbol{\theta}}_N = R(N)^{-1} \cdot \left[\frac{1}{N} \sum_{t=1}^N \boldsymbol{\varphi}(t) \cdot y(t) \right]\end{aligned}$$

Le matrici $S(N)$ e $R(N)$ sono **semidefinite positive** in quanto prodotto di un vettore per sé stesso. Affinché $\hat{\boldsymbol{\theta}}_N$ **esista** ed sia **unico**, è necessario che $S(N) > 0$ oppure $R(N) > 0$, ossia che:

$$\det(R(N)) > 0$$

Grazie all'ipotesi di **ergodicità**, abbiamo che $R(N) \xrightarrow{N \rightarrow +\infty} \bar{R}$, dove la matrice \bar{R} è la **matrice di autocovarianze** del processo congiunto $\{y(t), u(t)\}$. In generale, per un generico modello $ARX(n_a, n_b, 1)$ abbiamo che la matrice \bar{R} può essere divisa in quattro sotto-matrici:

$$\bar{R} = \begin{bmatrix} \bar{R}_{yy} & \bar{R}_{yu} \\ \bar{R}_{uy} & \bar{R}_{uu} \end{bmatrix}$$

$(n_a + n_b + 1) \times (n_a + n_b + 1)$ $(n_b + 1) \times n_a$ $n_a \times (n_b + 1)$ $(n_b + 1) \times (n_b + 1)$

Lemma di Schur

Data una matrice M nella forma $M = \begin{bmatrix} F & K \\ K^T & H \end{bmatrix}$, con F e K simmetriche. Condizione **necessaria e sufficiente** per l'invertibilità di M è che valgano:

- $H > 0$
- $F - KH^{-1}K^T > 0$

Ricordando che

$$\bar{R} = \begin{bmatrix} \bar{R}_{yy} & \bar{R}_{yu} \\ \bar{R}_{uy} & \bar{R}_{uu} \end{bmatrix}$$



Condizione **necessaria** per l'invertibilità di \bar{R} è che $\bar{R}_{uu} > 0$

La condizione (solo necessaria) sulla matrice \bar{R}_{uu} riguarda quindi solo il segnale di ingresso $u(t)$ che progettiamo noi. Possiamo quindi tenere conto di questa condizione in fase di progettazione dell'esperimento, e scegliere il segnale di eccitazione più opportuno al fine di ottenere dati informativi

Identificabilità sperimentale e strutturale

Vedere risposte precedenti.

Discutere il concetto di identificabilità di un modello ARX, trattando approfonditamente la nozione di persistente eccitazione di un segnale nella stima di modelli dinamici e fare degli esempi di segnali persistentemente eccitanti

Identificabilità dei modelli ARX

Vedi risposte precedenti.

Persistente eccitazione nella stima di modelli dinamici

Definizione (Persistente eccitazione)

Definiamo la matrice $\bar{R}_{uu}^{(i)}$ di autocovarianza di $u(t)$ di ordine i come

$$\bar{R}_{uu}^{(i)} = \begin{bmatrix} \gamma_{uu}(0) & \gamma_{uu}(1) & \cdots & \gamma_{uu}(i-1) \\ \gamma_{uu}(1) & \gamma_{uu}(0) & \cdots & \gamma_{uu}(i-2) \\ \vdots & \vdots & \ddots & \vdots \\ \gamma_{uu}(i-1) & \gamma_{uu}(i-2) & \cdots & \gamma_{uu}(0) \end{bmatrix}_{i \times i}$$

Il segnale $u(t)$ è detto **persistente eccitante di ordine n** se:

- $\bar{R}_{uu}^{(1)} > 0, \bar{R}_{uu}^{(2)} > 0, \dots, \bar{R}_{uu}^{(n)} > 0$
- $\bar{R}_{uu}^{(n+1)} \geq 0, \bar{R}_{uu}^{(n+2)} \geq 0, \dots \geq 0$

Ovvero se n è il massimo ordine per cui $\bar{R}_{uu}^{(i)}$ è invertibile

Possiamo quindi dire che **condizione necessaria** per l'identificabilità sperimentale di un modello $ARX(n_a, n_b, 1)$, usato per produrre i dati, sia **persistente eccitante di ordine pari ad almeno $n_b + 1$** (Infatti \bar{R}_{uu} ha dimensione $(n_b + 1) \times (n_b + 1)$).

Conclusione generale: supponendo quindi $S \in \mathcal{M}(\theta)$, definiamo il numero di parametri di $G(z; \theta)$ come n_g . Allora la soluzione del problema di identificazione

$$\bar{\theta} = \arg \min_{\theta} \mathbb{E}[\varepsilon_1(t, \theta)^2]$$

ha un'unica soluzione

$$\bar{\theta} = \theta^0$$

se il segnale $u(t)$ che genera i dati è
persistente eccitante di ordine $\geq n_g$

\Rightarrow Avere dunque un segnale eccitante è importante in ogni caso si voglia indentificare un modello dinamico.

Osservazioni:

- Se un segnale $u(t)$ è persistentemente eccitante di ordine n , allora è anche persistentemente eccitante di ordine $n - 1$
- badiamo che la **condizione vista è solamente necessaria**: anche se $\bar{R}_{uu} > 0$, La \bar{R} **potrebbe comunque non essere invertibile** per ragioni di **non identificabilità strutturale**, per le quali il minimo di $\bar{J}(\theta)$ non è unico
- Il concetto di persistente eccitazione che abbiamo visto è stato esemplificato per la stima di modelli ARX, ma avere **un segnale eccitante è importante in ogni caso** si voglia identificare un modello dinamico

Esempi di segnali:

- **Rumore bianco:** $u(t) \sim WN(0, \lambda^2)$

$$\bar{R}_{uu}^{(i)} = \begin{bmatrix} \gamma_{uu}(0) & \gamma_{uu}(1) & \cdots & \gamma_{uu}(i-1) \\ \gamma_{uu}(1) & \gamma_{uu}(0) & \cdots & \gamma_{uu}(i-2) \\ \cdots & \cdots & \cdots & \cdots \\ \gamma_{uu}(i-1) & \gamma_{uu}(i-2) & \cdots & \gamma_{uu}(0) \end{bmatrix} = \begin{bmatrix} \lambda^2 & 0 & 0 & \cdots & 0 \\ 0 & \lambda^2 & 0 & \cdots & 0 \\ 0 & 0 & \lambda^2 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & 0 \\ 0 & 0 & 0 & 0 & \lambda^2 \end{bmatrix} = \lambda^2 \cdot I_i > 0$$

Un **white noise** è un segnale persistentemente eccitante di ordine ∞ . Se usiamo un WN per eccitare il sistema, i dati generati saranno molto informativi. Questo perché eccitiamo tutte le frequenze del sistema (con la stessa energia).

un rumore bianco è un ottimo segnale di eccitazione. Nella pratica, però, **è impossibile generare un rumore bianco perfetto:**

- le sequenze di numeri saranno pseudo-casuali, e non casuali
- **a causa di limiti dell'elettronica** (capacità parassite), il segnale generato e trasmesso agli attuatori sarà filtrato passa basso per cui non si avrà uno spettro perfettamente piatto
- Inoltre, l'ampiezza del rumore bianco non è limitata. Talvolta, è necessario garantire che l'attuatore non saturi l'ingresso.

- **Pseudo-Random Binary Signal (PRBS)**

- Il segnale di tipo PRBS è un segnale deterministico, periodico, a tempo discreto, che commuta tra due livelli. L'utente deve definire i due livelli $[-\bar{u}, +\bar{u}]$, il periodo e l'intervallo di clock.
- Di solito il periodo viene posto uguale al numero di dati N che si vuole collezionare, e l'intervallo di clock a un tempo di campionamento
- Quanto $T \rightarrow \infty$ il PRBS approssima un rumore bianco
- Il PRBS è persistentemente eccitante di ordine T

- **Multiseno**

- Il segnale multiseno è un segnale periodico, definito come una **media pesata di sinusoidi**, con frequenze multiple della risoluzione in frequenza della DFT $f_0 = \frac{f_s}{N}$

$$u(t) = \sum_{k=0}^F A_k \cdot \cos(2\pi \cdot k f_0 \cdot t + \phi_k)$$

- Il numero F di componenti in frequenza deve soddisfare il teorema del campionamento
- Gli sfasamenti ϕ_k sono in generale scelti in **modo casuale** e si possono anche ottimizzare per minimizzare il valore di picco del segnale
- Molto spesso le **ampiezze** A_k vengono scelte ad un **valore costante nella banda di frequenze di interesse e 0 altrove**.
- Quando l'ingresso $u(t)$ da un multiseno, anche il segnale di uscita $y(t)$ di un sistema LTI è un multiseno (dopo un transitorio), come conseguenza del principio di sovrapposizione degli effetti \Rightarrow di solito si scartano i primi periodi del segnale multiseno generato
- Quando si progetta un multiseno, si può **fissare la risoluzione in frequenza desiderata e la massima frequenza eccitata**, per calcolare automaticamente la lunghezza $N \cdot P$ del segnale

$$N = \text{numero dati per periodo} \\ P = \text{numero di periodi del multiseno}$$

10 Lezione 13: Identificazione - valutazione del modello

I metodi per scegliere la complessità del modello.

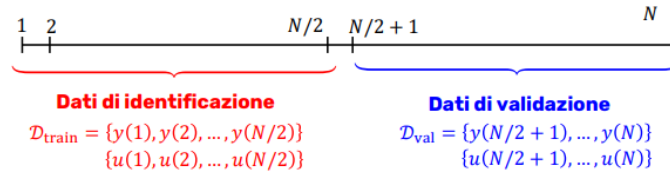
- Validazione e formule di complessità
- Analisi dei residui
- Analisi incertezza stima

Descrivere almeno due metodi alternativi per la selezione dell'ordine del modello dinamico nelle procedure di identificazione a minimizzazione dell'errore di predizione. Discutere quando conviene un metodo piuttosto che un altro.

Fissata la struttura di una famiglia di modelli $\mathcal{M}(\theta)$, dobbiamo scegliere la **complessità del modello** (numero dei parametri).

Validazione

Un metodo semplice ma efficace consiste nell'identificare un insieme di modelli di diversa complessità utilizzando un dataset di **identificazione**, e confrontarne la bontà (esempio calcolo $J(\hat{\theta}_N)$) su un dataset di **validazione**. Supponendo di avere N dati e li dividiamo in 2 sotto-sequenze:



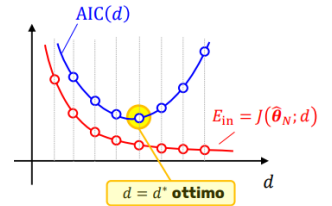
Per ogni ordine $m = 1, \dots, M$ si identifica un modello minimizzando $J(\theta, \mathcal{D}_{\text{train}})$ e si calcola $J(\hat{\theta}_{\frac{N}{2}}, \mathcal{D}_{\text{val}})$ sui dati di validazione infine si sceglie l'ordine m^* che minimizza $J(\hat{\theta}_{\frac{N}{2}}, \mathcal{D}_{\text{val}})$.

A differenza del caso statico, con i sistemi dinamici non è possibile estrarre i dati di identificazione e di validazione in modo casuale dal dataset completo, perché romperebbe la **causalità temporale dei dati**, inoltre per utilizzare questo approccio sono necessari molti dati.

Formule di complessità

In alternativa, se i dati sono pochi si possono usare le **formule di complessità** ottima: stimare d^*

- Akaike Information Criterion $AIC(d) = 2 \cdot \frac{d}{N} + \ln[J(\hat{\theta}_N; d)]$
- Final prediction error $FPE(d) = \frac{N + d}{N - d} \cdot J(\hat{\theta}_N; d)$
- Minimum Description Length $MDL(d) = \ln[N] \cdot \frac{d}{N} + \ln[J(\hat{\theta}_N; d)]$



Analisi residui identificazione modello dinamico per scelta della struttura e della complessità

Analisi dei residui

Date le seguenti ipotesi:

$$\begin{aligned} \mathcal{S}: y(t) &= G_0(z)u(t) + H_0(z)e(t) \\ \mathcal{M}(\theta): y(t) &= G(z, \theta)u(t) + H(z, \theta)e(t) \end{aligned}$$

$$\begin{aligned} \varepsilon_1(t; \theta) &= H^{-1}(z; \theta) \cdot [y(t) - G(z, \theta)u(t)] \\ &= H^{-1}(z; \theta) \cdot [G_0(z)u(t) + H_0(z)e(t) - G(z, \theta)u(t)] \\ &= H^{-1}(z; \theta) \cdot [(G_0(z) - G(z, \theta))u(t) + H_0(z)e(t)] \quad -e(t) + e(t) \\ &= H^{-1}(z; \theta) \cdot [(G_0(z) - G(z, \theta))u(t) + (H_0(z) - H(z, \theta))e(t)] \quad +e(t) \\ &= \frac{G_0(z) - G(z, \theta)}{H(z, \theta)}u(t) + \frac{H_0(z) - H(z, \theta)}{H(z, \theta)}e(t) + e(t) \end{aligned}$$

Se $\exists \theta^0$ t.c. $G(z, \theta^0) = G_0(z)$ e $H(z, \theta^0) = H_0(z)$, allora $\varepsilon_1(t; \theta) = e(t)$, dove $e(t) \sim WN(0, \lambda^2)$.

La stima asintotica può quindi essere ottenuta come:

$$\bar{\theta}_N = \arg \min_{\theta \in \Theta} \bar{J}(\theta) = \arg \min_{\theta \in \Theta} \mathbb{E}[\varepsilon_1(t; \theta)^2]$$

Dopo aver selezionato un modello $\mathcal{M}(\theta)$ e averne effettuato l'identificazione PEM, è possibile **validarne** (a-posteriori) la **struttura** e la **complessità** tramite analisi dei residui

Obiettivo: avendo la stima $\hat{\theta}_N$ e i dati $\{u(y), y(t)\}_{t=1}^N$, determinare se $\mathcal{M}(\theta)$ é tale che:

- $S \in \mathcal{M}(\theta)$
- $S \notin \mathcal{M}(\theta)$ con $G_0(z) \in \mathcal{G}(\theta)$
- $S \notin \mathcal{M}(\theta)$ con $G_0(z) \notin \mathcal{G}(\theta)$

Dove \mathcal{G} è l'insieme dei modelli che descrivono la relazione input-output del sistema

Consideriamo il caso asintotico $N \rightarrow +\infty$, in cui $\hat{\theta}_N \rightarrow \bar{\theta}$, abbiamo che

$$\varepsilon_1(t; \bar{\theta}) = H^{-1}(z; \bar{\theta})(y(t) - G(z, \bar{\theta})u(t)) = \frac{G_0(z) - G(z, \bar{\theta})}{H(z, \bar{\theta})}u(t) + \frac{H_0(z)}{H(z, \bar{\theta})}e(t) \quad e(t) \sim WN(0, \lambda^2)$$

La scelta della **struttura** e della **complessità** del modello $\mathcal{M}(\theta)$ può essere effettuata osservando:

- la funzione di **autocovarianza dei residui** $\rightarrow \gamma_{\varepsilon\varepsilon}(\tau)$
- la funzione di **cross-covarianza tra i residui** ed il segnale di **ingresso** $\rightarrow \gamma_{\varepsilon u}(\tau)$

Sapendo che $\delta(\tau)$ è un delta di Dirac centrata in τ , possiamo incorrere in 3 situazioni:

- **Situazione A:** $\gamma_{\varepsilon\varepsilon}(\tau) = \lambda^2 \cdot \delta(\tau)$ e $\gamma_{\varepsilon u}(\tau) = 0 \quad \forall \tau$

$$\gamma_{\varepsilon\varepsilon}(\tau) = \begin{cases} \lambda^2 & \tau = 0 \\ 0 & \tau \neq 0 \end{cases} \quad \gamma_{\varepsilon u} = 0 \quad \forall \tau$$

$$\varepsilon_1(t; \bar{\theta}) = \frac{G_0(z) - G(z, \bar{\theta})}{H(z, \bar{\theta})}u(t) + \frac{H_0(z)}{H(z, \bar{\theta})}e(t) = 0 \cdot u(t) + e(t)$$

*Overo se e solo se $G(z, \bar{\theta}) = G_0(z)$ e $H(z, \bar{\theta}) = H_0(z)$
Questo avviene se e solo se $S \in \mathcal{M}(\theta)$*

- **Situazione B:** $\gamma_{\varepsilon\varepsilon}(\tau) \neq \lambda^2 \cdot \delta(\tau)$ e $\gamma_{\varepsilon u}(\tau) = 0 \quad \forall \tau$

$$\varepsilon_1(t; \bar{\theta}) = \frac{G_0(z) - G(z, \bar{\theta})}{H(z, \bar{\theta})}u(t) + \frac{H_0(z)}{H(z, \bar{\theta})}e(t) = 0 \cdot u(t) + \underbrace{\frac{H_0(z)}{H(z, \bar{\theta})}}_{\neq 1} \cdot e(t)$$

*Overo se e solo se $G(z, \bar{\theta}) = G_0(z)$ e $H(z, \bar{\theta}) \neq H_0(z)$
Questo avviene se e solo se $S \notin \mathcal{M}(\theta)$ con $G_0(z) \in \mathcal{G}(\theta)$
per $\mathcal{M}(\theta)$ **OE, BJ, FIR***

- **Situazione C:** $\gamma_{\varepsilon\varepsilon}(\tau) \neq \lambda^2 \cdot \delta(\tau)$ e $\exists \tau$ t.c. $\gamma_{\varepsilon u}(\tau) \neq 0$

$$\varepsilon_1(t; \bar{\theta}) = \underbrace{\frac{G_0(z) - G(z, \bar{\theta})}{H(z, \bar{\theta})}u(t)}_{\neq 0} + \frac{H_0(z)}{H(z, \bar{\theta})}e(t)$$

Overo se e solo se $G(z, \bar{\theta}) \neq G_0(z)$

Questo avviene **se e solo se** $\begin{cases} S \notin \mathcal{M}(\theta) \text{ con } G_0(z) \in \mathcal{G}(\theta) \text{ per } \mathcal{M}(\theta) \text{ **ARX, ARMAX** } \\ S \notin \mathcal{M}(\theta) \text{ con } G_0(z) \notin \mathcal{G}(\theta) \end{cases}$

Considerando che $\mathcal{M}_1(\theta)$ è **OE, BJ o FIR** e $\mathcal{M}_2(\theta)$ è invece **ARX, ARMAX**.

- A:** $\gamma_{\varepsilon\varepsilon}(\tau) = \lambda^2 \cdot \delta(\tau)$ e $\gamma_{\varepsilon u} = 0 \quad \forall \tau$ $\mathcal{S} \in \mathcal{M}_1(\theta)$ $\mathcal{S} \in \mathcal{M}_2(\theta)$
- B:** $\gamma_{\varepsilon\varepsilon}(\tau) \neq \lambda^2 \cdot \delta(\tau)$ e $\gamma_{\varepsilon u} = 0 \quad \forall \tau$ $\mathcal{S} \notin \mathcal{M}_1(\theta), G_0(z) \in \mathcal{G}(\theta)$
- C:** $\gamma_{\varepsilon\varepsilon}(\tau) \neq \lambda^2 \cdot \delta(\tau)$ e $\exists \tau$ t.c. $\gamma_{\varepsilon u} \neq 0$ $\mathcal{S} \notin \mathcal{M}_1(\theta), G_0(z) \notin \mathcal{G}(\theta)$ $\mathcal{S} \notin \mathcal{M}_2(\theta), G_0(z) \notin \mathcal{G}(\theta)$

				$N \rightarrow +\infty$		N finito	
				$\gamma_{\varepsilon\varepsilon}(\tau)$	$\gamma_{\varepsilon u}(\tau)$	$\hat{\gamma}_{\varepsilon\varepsilon}(\tau)$	$\hat{\gamma}_{\varepsilon u}(\tau)$
• Situazione A:	$\gamma_{\varepsilon\varepsilon}(\tau) = \lambda^2 \cdot \delta(\tau)$	e	$\gamma_{\varepsilon u}(\tau) = 0 \quad \forall \tau$	$\mathcal{S} \in \mathcal{M}(\theta)$	0 $\forall \tau \neq 0$	0 $\forall \tau$	«piccola» ∈ intervallo di confidenza
• Situazione B:	$\gamma_{\varepsilon\varepsilon}(\tau) \neq \lambda^2 \cdot \delta(\tau)$	e	$\gamma_{\varepsilon u}(\tau) = 0 \quad \forall \tau$	$\mathcal{S} \notin \mathcal{M}(\theta)$ $G_0(z) \in \mathcal{G}(\theta)$	$\exists \tau \neq 0$ t.c. $\gamma_{\varepsilon\varepsilon}(\tau) \neq 0$	0 $\forall \tau$ OE, BJ, FIR	«piccola» ∈ intervallo di confidenza
• Situazione C:	$\gamma_{\varepsilon\varepsilon}(\tau) \neq \lambda^2 \cdot \delta(\tau)$	e	$\exists \tau$ t.c. $\gamma_{\varepsilon u}(\tau) \neq 0$	$\mathcal{S} \notin \mathcal{M}(\theta)$ $G_0(z) \notin \mathcal{G}(\theta)$	$\exists \tau \neq 0$ t.c. $\gamma_{\varepsilon\varepsilon}(\tau) \neq 0$	$\exists \tau$ t.c. $\gamma_{\varepsilon u}(\tau) \neq 0$	«grande» ∉ intervallo di confidenza

Una validazione della struttura che si conclude con un successo non **garantisce** che $G(z, \hat{\theta}_N)$ e $H(z, \hat{\theta}_N)$ siano buone stime di $G_0(z)$ e $H_0(z)$. È necessario controllare anche la varianza delle stime (sia dei parametri sia delle funzioni di trasferimento).

Descrivere tre modi per valutare la bontà della stima (identificazione) su un modello dinamico.

La bontà di un modello può essere valutata:

- Analizzando i residui (cioè gli errori di predizione a un passo), meglio con dati di validazione
- Analisi incertezza stima: parametri, funzione di trasferimento e posizione dei poli e degli zeri
- Confrontando l'uscita simulata o predetta con l'uscita misurata, su dati di validazione, e calcolandone un indicatore di **FIT**
- Confronto con stima nonparametrica

Analisi dei residui

Vedere precedentemente, nel caso diminuire la roba da scrivere.

Analisi incertezza stima

Con l'analisi dell'incertezza della stima intendiamo:

- **Incetezza sulla stima dei parametri:**
Si verifica la significatività statistica di un parametro, dove per significatività intendiamo quanto è probabile che il parametro vero sia effettivamente diverso da 0.
- **Incetezza sulla stima delle funzioni di trasferimento**
Nel caso $\mathcal{S} \in \mathcal{M}(\theta)$, l'espressione della varianza sulla stima del valore della funzione di trasferimento $G(z, \hat{\theta}_N)$ ad ogni frequenza $z = e^{j\omega}$ e assumendo che $u(t) \perp e(t)$, può essere approssimata come:

$$\text{Var}[G(e^{j\omega}, \hat{\theta}_N)] \approx \frac{n}{N} \cdot \frac{\Gamma_{vv}(\omega)}{\Gamma_{uu}(\omega)}$$

- n : n variabili di stato di G (ordine)
- $\Gamma_{vv}(\omega)$: $v(t) = H_0(z)e(t)$ disturbo

Se $\mathcal{S} \notin \mathcal{M}(\theta)$, esistono altre espressioni più complesse.

Quando $\sqrt{\text{Var}[G(e^{j\omega}, \hat{\theta}_N)]}$ può essere considerata piccola? La risposta dipende dall'uso del modello. Se è considerata troppo grande, non possiamo garantire che $G(e^{j\omega}, \hat{\theta}_N)$ sia una buona stima di $G_0(z)$.

- **Incertezza sulla posizione dei poli e degli zeri**

Una volta stimato il modello e analizzato i residui per capire se $\mathcal{S} \in \mathcal{M}(\theta)$, è utile rappresentare i poli e gli zeri del modello stimato, per verificare la possibilità di **cancellazioni** polo\zero (**e quindi evitare sovrapparametrizzazioni**).

Si verificano se gli ellissoidi di confidenza di poli e zeri si sovrappongono o sono molto vicini. In questo caso, è probabile che tali poli\zeri possano essere rimossi dal modello.

Simulazione e predizione del modello

Si simula o predice l'output del modello per confrontarlo con l'output misurato, a fronte del medesimo input. L'idea è che se la simulazione o predizione sono simili all'output misurato, allora il modello è buono.

L'errore di simulazione è $e_{sim}(t)$. Un modello buono non rende necessariamente $e_{sim}(t)$ piccolo, poiché la simulazione non considera il modello del rumore.

$$e_{sim} = y(t) - \hat{y}_{sim}(t) = y(t) - G(z, \hat{\theta}_N)u(t)$$

La simulazione può solo fornire informazioni sull'accuratezza di $G(z, \hat{\theta}_N)$. Inoltre, la simulazione è una rappresentazione più realistica di come il modello si comporta a fronte di un ingresso noto $u(t)$.

La simulazione permette anche di stimare il rumore $v(t) = H_0(z) \cdot e(t)$ tramite:

$$\hat{v}(t) = y(t) - \hat{y}_{sim}(t)$$

Confronto con stima nonparametrica

La stima nonparametrica lascia parlare i dati. È quindi importante valutare se le funzioni di trasferimento del modello aderiscono alle loro stime nonparametriche. Stime nonparametriche delle funzioni $G_0(z)$ e $H_0(z)$ si possono ottenere con la ETFE o con la stima spettrale.