

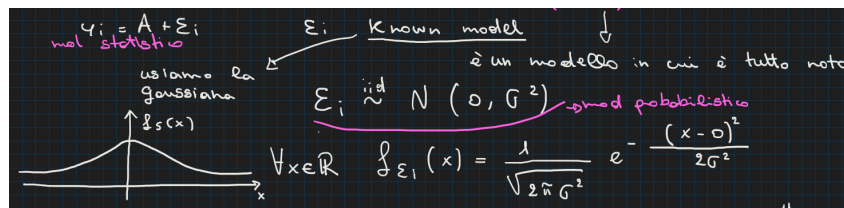
# Spiegazione MLE, Approccio Bayesiano, Dimostrazione test MSE - 17/09 + 19/09

## Maximum Likelihood Estimator (MLE) - 17/09

Ci serviranno strumenti per costruire un nostro stimatore a partire dai dati.

L'approccio classico richiede un po' di calcoli, si basa sul concetto di verosimiglianza e si chiama stimatore a massima

**verosimiglianza**. Maximum Likelihood estimator (MLE).



Dobbiamo conoscere il modello probabilistico di epsilon con i per usare MLE.

Questo è il classico approccio parametrico, ho tutto le forme note, il modello probabilistico eccetera.

Il nostro obiettivo è trovare A. ( $\theta$  è A, ha cambiato il come ma è a stessa cosa)

$$\text{LIKELIHOOD function } L(\theta) \propto P(D; \theta) \text{ given } D = \{y_i\}_{i=1}^n$$

↓  
proportionale

Avremo una funzione di Likelihood L che sarà funzione del parametro theta (che può essere un vettore).

Per prima cosa si osservano i dati, le mie y, possono essere un vettore ovviamente. I dati sono anche detti **variabile osservata**.

$$\underline{y} = (y_1, \dots, y_n) \text{ observed var.}$$
$$\epsilon_i \sim N(0, \sigma^2) \Rightarrow y_i \sim N(A, \sigma^2)$$

La somma di una costante non modifica la varianza, sposta solo la PDF.

Dobbiamo costruire il MLE (Maximum Likelihood estimator) per il parametro A. Il parametro A è la media delle  $y_i$ .

MLE per  $\mu$  cioè la media delle  $y$ :

$$f_y(y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}} \text{ pdf}$$

Conosciamo la PDF delle  $y$  che è uguale a quella delle epsilon solo che la media non sarà 0 ma  $\mu$ .



La verosimiglianza sarà funzione della  $\mu$  e dei dati raccolti.

Per scrivere la verosimiglianza bisogna ricordare come si scrive un modello probabilistico, la verosimiglianza è la PDF congiunta delle nostre  $y$ .

$$L(\underline{\theta}) = L(\mu; \underline{y}) = f_{y_1, \dots, y_n}(\mu; \underline{y})$$

produttoria  $\rightarrow \prod_{i=1}^n f_y(y_i)$

I dati sono iid, questo ci importa perché se le  $x$  sono indipendenti le  $y$  sono indipendenti. La congiunta di variabili aleatorie indipendenti è il prodotto delle marginali.

Visto che le marginali sono tutte uguali possiamo scrivere una produttoria.

La costante moltiplicativa sposta in basso o in alto il tutto ma non è importante, a noi interessa trovare il massimo della funzione ma questo massimo non è influenzato dalla costante  $c$  ( $c$  indica una costante generica, penso sia riferita al contributo di  $\mu$  che sposta la Gaussiana ma non ne cambia il massimo quindi è irrilevante).

Trovare  $\underline{\theta}$  che massimizza la prob di occorrenza di  $D$

$\downarrow$   $\mu$   $\rightarrow$  non la  $\mu$  massima ma il val di

$$A = \arg \max L(A) \text{ che massimizza}$$

$$= \arg \max \prod_{i=1}^n f_y(y_i) \Rightarrow \text{media campionaria}$$

Dobbiamo trovare il valore  $\mu$  per cui la funzione likelihood viene massimizzata.

(theta è in generale, nel nostro caso è  $\mu$ )

$\arg \max$  sarebbe il valore di  $\mu$  che massimizza la funzione.

Piccolo spoiler: il valore che ne uscirà sarà la media campionaria.

**19/09**

Riprendendo il discorso sul massimizzare  $L(A)$ , dobbiamo trovare il massimo di una funzione, il che solleva la domanda di come si trova il massimo di una funzione, ovviamente basta porre la derivata uguale a 0.

A questo punto calcolo la derivata seconda e ci interessa che sia negativa. Questo nel caso monodimensionale.

$$\frac{d}{dA} L(A) = 0 \quad \text{con} \quad L''(A) < 0$$

Esempio multidimensionale:

$$\begin{aligned} &\text{e più olim } \underline{\sigma} \in \mathbb{R}^d \\ &\nabla L(\underline{\sigma}) = 0 \text{ e} \\ &H(\underline{\sigma}) = \{h_{ij}(\underline{\sigma})\}_{i,j=1,d} \\ &h_{ij} = \frac{\partial^2}{\partial \sigma_i \partial \sigma_j} L(\underline{\sigma}) \\ &\text{definita negativa (eigen} < 0) \end{aligned}$$

Per theta vettore a più dimensioni il gradiente (equivalente a più dimensioni della derivata prima, sono tipo tutte le derivate) di L di theta deve essere uguale al vettore nullo e la matrice Hessiana (matrice delle derivate rispetto a theta i e theta j, è l'equivalente a più dimensioni della derivata seconda) deve essere definita negativa (visto che è quadrata pari simmetrica basta che gli autovalori siano tutti dello stesso segno e negativi).

Invece della funzione verosimiglianza usiamo il log della funzione verosimiglianza, si può fare perché il log non cambia il massimo.

▼ Perché si usa la log likelihood?

- **Semplificazione dei calcoli:**

- La likelihood è spesso il prodotto di molte probabilità (o densità), che possono diventare numeri molto piccoli, portando a problemi di precisione numerica. La log-likelihood trasforma questo prodotto in una somma, che è più facile da gestire matematicamente e computazionalmente.
- Le derivate della log-likelihood (necessarie per ottimizzare o stimare parametri) sono più semplici da calcolare rispetto a quelle della likelihood.

- **Equivalenza nell'ottimizzazione:**

- La log-likelihood conserva l'ordine di grandezza della likelihood. Massimizzare la log-likelihood è equivalente a massimizzare la likelihood, poiché il logaritmo è una funzione monotona crescente. Quindi, entrambe portano agli stessi parametri stimati.

- **Stabilità numerica:**

- La likelihood può diventare estremamente piccola (vicina a zero), soprattutto con dataset grandi, il che può causare underflow numerico. Prendere il logaritmo di valori piccoli evita questi problemi.

- **Additività:**

- Il logaritmo trasforma il prodotto delle probabilità in una somma. Questo rende i modelli più facili da analizzare e permette di applicare tecniche di ottimizzazione standard (come il gradiente o il metodo di Newton).

$l(\theta) = \log L(\theta)$  *log likelihood*  
 $\sigma^2$  nota  
 $D = \{y_i\}_{i=1}^n$  (data size  $n$ )  
 $y_i \sim N(\mu, \sigma^2)$   $y_i = \mu + \epsilon_i$

Vogliamo stimare  $\mu$ .

**Supponiamo che sigma quadro è noto**, il caso incognito è quando sia  $\mu$  che sigma quadro sono da individuare, qui invece supponiamo di averlo per semplificare ma è poco realistico averlo.

Vettore  $y$  raccoglie tutti i dati da uno ad  $n$ , è il data set.

**Dimostrazione che la stima MLE è uguale alla media campionaria per il caso del calcolo della media di una gaussiana**

$$L(\mu; y) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_i - \mu)^2}{2\sigma^2}} = \frac{1}{\sigma^n (2\pi)^{n/2}} \cdot e^{-\frac{\sum_{i=1}^n (y_i - \mu)^2}{2\sigma^2}}$$

Possiamo cacciare fuori dalla produttoria la parte costante, ora abbiamo la produttoria di esponenziali che sarebbe l'esponenziale con esponente la somma degli esponenti.

Questa espressione ci fa capire perché passiamo ai logaritmi.

Inoltre la costante visto che dobbiamo derivare per  $\mu$  si perderebbe non ci serve. **Il fatto che ignoriamo la costante si dice che scriviamo un termine di proporzionalità ad  $\mu$ .**

$$l(\mu; y) = \log L(\mu; y) = \log \left( \frac{1}{\sigma^n (2\pi)^{n/2}} \cdot e^{-\frac{\sum_{i=1}^n (y_i - \mu)^2}{2\sigma^2}} \right)$$

Ora calcoliamo il massimo facendo la derivata della log likelihood.

Cerchiamo  $\hat{\mu}$  che non sarà altro che l'arg max della log likelihood cioè il punto tale che la derivata di  $L$  rispetto ad  $\mu$  sia uguale a 0.

Cerchiamo A tale che questa espressione è uguale a 0, questa sarà la A che massimizza la funzione (dopo aver verificato le condizioni sulla derivata seconda).

$$\hat{A} = \arg \max \ell(A; \underline{y}) \text{ s.t. } \frac{d\ell(A; \underline{y})}{dA} = 0$$

s.t. = such that

$$\begin{aligned} \frac{d\ell(A; \underline{y})}{dA} &= + \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - A) = 0 \\ \Rightarrow \sum_{i=1}^n y_i &= nA \Rightarrow \hat{A} = \frac{\sum_{i=1}^n y_i}{n} = \bar{y}_n \quad \text{sample mean} \\ &\quad \text{solution} \end{aligned}$$

Abbiamo portato la parte di sommatoria relativa ad A dall'altro lato (nA va a destra).

La stima della media (sappiamo per dimostrazioni precedenti che A è la media) è la media campionaria (sample mean) intuitivamente, tutti questi calcoli lo confermano.

La A che risolve l'equazione è la A cappello.

In questo caso A, sigma quadro e i modelli sono noti, io devo arrivare alla regressione, in futuri casi avremo tutto incognito e aleatorio e in quei casi la verosimiglianza non vale.

Da ora A sarà una variabile aleatoria non più una costante  
monodimensionale

Ci servirà l'approccio supervisionato che tramite le y e le x ci permette di ottenere la f(x) che per noi ora è aleatoria.

Ci serve un approccio diverso che è l'approccio Bayesiano.

## Bayesian approach

$$\begin{aligned} x_i &= A + \varepsilon_i \\ A &\sim \sigma. \quad \text{Variabile aleatoria} \rightarrow \text{random variable} \\ f_A(\cdot) &\quad \text{p.d.f. to ext. from data} \end{aligned}$$

Nel mio modello semplice semplice,  $x_i = A + \varepsilon_i$ , ora anche A è una variabile aleatoria.

Così posso inglobare il fatto che può esserci aleatorietà anche in quello che devo stimare. La A avrà quindi una sua PDF che volendo va stimata a partire dai dati.

Cerchiamo sempre lo stimatore, lo stimatore Bayesiano in questo caso.

Lo stimatore sarà sempre funzione dei dati (che possiamo scrivere come x o come y, è solo una lettera).

$$\hat{A} = \hat{A}(x_1, \dots, x_n) \quad \text{stimatore} \quad \hat{A} = \hat{A}(\underline{x})$$

$\hat{A}$  visto che va ricavata dai dati sarà quindi funzione del vettore dei dati  $\underline{x}$ .

Anche ora cerchiamo lo stimatore che minimizza l'errore. Continuiamo a usare l'errore quadratico ma ce ne sarebbero altri.

$$\text{BMSE}(\hat{A}) \triangleq \mathbb{E} \{ (A - \hat{A})^2 \}$$

↳ Bayssiano

Talvolta si segna BMSE per indicare che è da trattare con l'approccio Bayesiano.

BMSE resta la media statistica di  $A$  meno  $A$  cappello al quadrato.

Obiettivo: trovare  $A$  cappello che minimizza l'MSE.

$$\begin{aligned} \text{Goal: } \hat{A} \text{ minimizza } \text{MSE} &= \mathbb{E}_{\underline{x}, A} \{ (A - \hat{A})^2 \} \\ \text{MSE}(\hat{A}) &= \iint (A - \hat{A})^2 \underbrace{f(\underline{x}, A)}_{\text{pdf congiunta}} d\underline{x} dA \end{aligned}$$

Prima il valore atteso in questione agiva solo sui dati  $\underline{x}$  ma ora che anche  $A$  è random variable questo agisce anche su di  $A$  e quindi c'è un doppio integrale.

tale che

$$\hat{A} \stackrel{!}{\Rightarrow} \text{MSE}(\hat{A}) \min \Leftrightarrow \hat{A} = \arg \min_{\hat{A}} \text{MSE}(\hat{A})$$

Arg min e arg max sono il minimizzatore e il massimizzatore.

Consideriamo la pdf congiunta visto che abbiamo detto che il valore atteso è su due variabili.

(Congiunta = joint distribution)

$$\begin{aligned} f(\underline{x}, A) &= f(A | \underline{x}) f(\underline{x}) \\ \Rightarrow \text{MSE}(\hat{A}) &= \int \left[ \int (A - \hat{A})^2 f(A | \underline{x}) dA \right] f(\underline{x}) d\underline{x} \end{aligned}$$

↳ congiunta

La congiunta si può sempre scrivere come probabilità di eventi condizionati.

Con integrali doppi può essere comodo dividerli facendo prima all'interno un'espressione e poi il risultato nell'integrale esterno.

Una proprietà di una pdf è che non può essere negativa. Questo vuol dire che se voglio minimizzare l'MSE mi basta minimizzare la funzione interna.

$$f(x) \geq 0 \text{ (perché pdf)} \forall x \Rightarrow \text{minimizzare } \int (A - \hat{A})^2 f(A|x) dA$$

Per trovare il minimo di nuovo devo fare la derivata e porla uguale a 0.

$$\begin{aligned} \theta &= \frac{\partial}{\partial \hat{A}} \int (A - \hat{A})^2 f(A|x) dA = \\ &= \int \frac{\partial}{\partial \hat{A}} (A - \hat{A})^2 \underbrace{f(A|x)}_{\substack{\text{dato} \\ \hookrightarrow \text{non dipende da } A}} dA = \end{aligned}$$

La PDF condizionata è indipendente da A cappello quindi la derivata non ha effetto su di lei.

$$\begin{aligned} &= -2 \int (A - \hat{A}) f(A|x) dA = \\ &= -2 \int A f(A|x) dA + 2\hat{A} \int f(A|x) dA \end{aligned}$$

1 (integrale di una pdf)

il -2 viene dalla derivazione

L'integrale della pdf su tutto il suo dominio fa 1.

Quindi rimane solo A cappello.

A questo punto ponendo tutto questo uguale a 0 trovo subito A cappello.

$$\begin{aligned} \Rightarrow \hat{A} &= \int A f(A|x) dA = E[A|x] \\ &\quad \hookrightarrow \text{valore atteso} \\ &= \hat{A}_{\text{MMSE}} \quad \hookrightarrow \text{Minimum MSE} \\ &\quad \hookrightarrow \text{stimatore bayesiano} \end{aligned}$$

$\Rightarrow$  valore atteso della variabile aleatoria dato tutto quello che abbiamo raccolto

Quello che abbiamo trovato è il valore atteso di A dato vettore x.



Lo stimatore Bayesiano è chiamato anche in ingegneria stimatore a minimo MSE.

$f(A|x)$  o la conosco o in qualche modo devo stimarla.

Supponiamo di doverla stimare.

$$f(A|x) \text{ to be known}$$

Questa è la pdf di A dopo aver raccolto i dati, quella prima si chiamava conoscenza a priori ( $f(A)$ ), ora è conoscenza a posteriori.

In genere priori e posteriori dovrebbero somigliarsi.

Teorema di Bayes legava, del resto, conoscenze a priori e conoscenze a posteriori, il tramite è la verosimiglianza.

Approccio classico (verosimiglianza) vs Approccio Bayesiano (che ha anche le informazioni a priori).

## Errore di generalizzazione - dimostrazione

$$Y = f(x) + \varepsilon \rightarrow y_i = f(x_i) + \varepsilon_i$$

Modello della popolazione. L'espressione campionaria. Ora al posto di A c'è  $f(x)$  che è una funzione incognita.

Per approssimare bene  $f$  faccio il valore atteso della differenza tra la realtà e l'approssimazione.

$$\begin{aligned} & \mathbb{E} \left\{ (y - g(x))^2 \mid X=x \right\} \rightarrow \hat{g} \text{ minimizer} \\ & \quad \downarrow \quad \quad \quad \downarrow \\ & \text{tutte le stime di } g \text{ possibili} \quad \text{dopo aver raccolto i dati} \\ & \text{qualsiasi funzione tra le quali posso scegliere} \\ \Rightarrow \hat{g}(x) &= \mathbb{E}_y [Y \mid X=x] \quad \left. \vphantom{\mathbb{E}_y} \right\} \text{funzione di regressione} \\ & \quad \quad \quad \downarrow \\ & \quad \quad \quad \text{valore atteso su tutte le } y \end{aligned}$$

$F(x)$  resta una variabile aleatoria, per la variabile aleatoria abbiamo detto che lo stimatore migliore è quello che minimizza il BMSE.

$F$  è teorica e ha generato i dati, non la conoscerò mai al 100 per cento, conoscerò  $f$  cappello,  $g$  di  $x$  invece sono tutte le funzioni dello spazio delle funzioni di  $f$ .

$$\hat{g} \text{ stima di } f(x) \text{ dati } D$$



Consideriamo un test set costituito da un solo punto un test point. (vogliamo trovare il test mse).

test point  $(x_0, y_0)$  expected test MSE

$$\mathbb{E}[(y_0 - \hat{g}(x_0))^2] =$$

(+)

$$\mathbb{E}[(f(x_0) + \varepsilon - \hat{g}(x_0))^2] =$$

sostituiamo  $y_0$  con la sua definizione

Svolgiamo i quadrati.

$$= \mathbb{E}[(f(x_0) - \hat{f}(x_0) + \varepsilon)^2] =$$
$$\underbrace{\mathbb{E}[(f(x_0) - \hat{f}(x_0))^2]}_{\text{bias}} + \underbrace{\mathbb{E}[\varepsilon^2]}_{\text{var}} + 2 \underbrace{\mathbb{E}[\varepsilon(f(x_0) - \hat{f}(x_0))]}_{\text{per il principio di ortogonalità degli errori}} =$$
$$\text{var}[\hat{f}(x_0)] + b^2(\hat{f}(x_0)) + \text{var}(\varepsilon)$$

↙ bias                      ↘ perché ε è a media nulla

Si può dimostrare che il termine epsilon per la differenza tra  $f$  ed  $f$  cappello di  $x_0$  è uguale a 0, si chiama principio di ortogonalità degli errori per epsilon, la distanza tra  $f$  ed  $f$  cappello.

F cappello di  $x_0$  meno  $f(x_0)$  è un termine che riguarda il valore atteso di questa quantità.

Ricordiamo che Epsilon è a media nulla.

$$\Rightarrow \text{test MSE} \quad E[(y_0 - \hat{f}(x_0))^2] = \underbrace{\text{var}[\hat{f}(x_0)]}_{\text{rid.}} + \underbrace{b^2(f'(x_0))^2}_{\text{rid.}} + \text{var}(\varepsilon)$$

Oltre al classico termine di compromesso bias varianza c'è anche la varianza di epsilon, un errore irriducibile.

In statistica si sostituiscono alle medie le medie campionarie.