

MSE, Stima del parametro deterministico

A - 17/09

Problema supervisionato

Il nostro scopo è la stima di f , quindi ottenere \hat{f} , tale che il MSE sia minimo.

Goal: estimate $f(x) \rightarrow \hat{f}(x)$

$y = f(x) + \varepsilon$ model

$\{(x_i, y_i)\}_{i=1}^n = D_{\text{tr}}$ training set

Optimality criterion $\hat{f}(x)$ minimizing Mean Squared Error (MSE)

Stimare f è proprio effettuare il training, in questo caso vediamo sia x_i che y_i nel training set, quindi si tratta di un problema supervisionato.

Ci sono varie metriche per la misura dell'errore, quella classica è la media quadratica visto che l'errore può essere sia positivo che negativo.

Stimare i parametri dato il modello

(qua si inizia un altro discorso da zero)

Il problema più semplice è: noto il modello stimare i parametri.

Stima di un parametro θ dato il modello
CASO SEMPLICE [θ 1-dim] quindi è un solo parametro

Non conoscendo μ e σ (media e varianza) il problema è bidimensionale, noi iniziamo da un caso monodimensionale dove dobbiamo trovare θ e ci interessa stimare la media della Gaussiana.

dati $\{x_i\}_{i=1}^n$ x_i appartengono tutti allo stesso modello
 $x_i \sim F(x; \theta)$
sono distribuiti

Abbiamo dei dati del modello di nostro interesse che chiamiamo x_i , questi dati sono distribuiti (tilde \sim) su una certa CDF (ci potrebbe dare la PDF, o la PMF se è un caso discreto).

Mi serve uno **stimatore di θ** .

Nell'approccio classico θ è un parametro deterministico che però noi non conosciamo.

$\hat{\theta}$ lo stimiamo con il modello e i dati raccolti. L'estimatore è funzione T dei dati raccolti.

STIMATORE $\hat{\theta} = T(x_1, \dots, x_n)$
statistica
La stima è in funzione dei dati

Lo stimatore se considero le maiuscole (variabili aleatorie) abbiamo una funzione di variabili aleatorie e quindi $\hat{\theta}$ sarà una variabile aleatoria.

Se invece considero proprio i valori numerici avremo una funzione che ci restituisce una variabile deterministica.



Lo **stimatore** è la funzione che, essenzialmente, se ci mettiamo i dati otteniamo dei risultati cioè la stima.



La **stima** è un numero prodotto dallo stimatore quando si usano i dati osservati

Errore quadratico medio - MSE (mean squared error)



Lo **stimatore ottimo** è quello che minimizza l'MSE per il parametro θ . θ è monodimensionale cioè un numero.

σ 1-dim $MSE(\hat{\theta}) \triangleq E[(\theta - \hat{\theta})^2]$ \rightarrow valore atteso (media)
 \downarrow
si deve trovare il $\hat{\theta}$ che minimizza l'MSE \rightarrow errore quadratico medio
 $\sigma \in \mathbb{R}^d$ $MSE(\hat{\theta}) = E[\|\sigma - \hat{\theta}\|^2]$ *inciso bidimensionale*
 \downarrow
 σ vettore $\|y\| = \sqrt{\sum_{i=1}^d y_i^2}$ *norma di \mathbb{R}^d*

definizione di $MSE(\hat{\theta})$ in caso monodimensionale e multidimensionale

MSE in funzione di varianza e bias

Segue la dimostrazione

$$\begin{aligned} \Rightarrow \text{MSE}(\hat{\theta}) &= E\left\{(\hat{\theta} - \underbrace{E[\hat{\theta}]}_{\text{sommo e sottratto la media di } \hat{\theta}} + (\underbrace{E[\hat{\theta}] - \theta}_{\text{cost}}))^2\right\} \\ &= \underbrace{E\{(\hat{\theta} - E[\hat{\theta}])^2\}}_{\text{VARIANZA DI } \hat{\theta}} + E\{(\underbrace{E[\hat{\theta}] - \theta}_{\text{cost}})^2\} \\ &\quad + 2E\{(\hat{\theta} - E[\hat{\theta}])(\underbrace{E[\hat{\theta}] - \theta}_{\text{cost}})\} \end{aligned}$$

Sommo e sottraggo lo stesso valore ($E[\text{theta cappello}]$), poi associo la somma e sottrazione per comodità.

Essendo la media quadratica lineare dividiamo le diverse E.

Osserviamo il doppio prodotto.

$$\begin{aligned} &+ 2E\{(\hat{\theta} - E[\hat{\theta}])(\underbrace{E[\hat{\theta}] - \theta}_{\text{cost}})\} \\ &\quad \downarrow \\ &\quad \text{val att è un num e } \theta \text{ è un num quindi questo è un num e lo possiamo estrarre} \\ &\quad \downarrow \\ &2(E[\hat{\theta}] - \theta) \cdot \underbrace{E\{(\hat{\theta} - E[\hat{\theta}])\}}_{=0} = 0 \end{aligned}$$

Dopo aver tolto la costante di cui sopra usando la linearità dividiamo di nuovo il contenuto della stima, il risultato sarà 0 in quanto ne esce differenza di $E[\text{theta cappello}]$ e $E[E[\text{theta cappello}]]$, $E[\text{theta cappello}]$ è una costante e la stima di una costante è uguale alla costante stessa quindi i due termini sono uguali.

$$\begin{aligned} &\text{VARIANZA DI } \hat{\theta} \quad \text{sommo e sottratto la med} \\ &= E\{(\hat{\theta} - E[\hat{\theta}])^2\} + E\{(\underbrace{E[\hat{\theta}] - \theta}_{\text{cost}})^2\} \end{aligned}$$

(rigo dell'immagine di sopra)

L'MSE di theta cappello è pari al primo termine più il secondo termine, il primo termine è proprio la varianza di theta cappello, del secondo termine eliminiamo la E perché al suo interno c'è una costante ($E[\text{theta cappello}]$ e theta sono entrambe costanti), per la quale il valore atteso è irrilevante.

Il secondo termine prende il nome di bias al quadrato dello stimatore theta cappello.

$$\begin{aligned} \Rightarrow \text{MSE}(\hat{\theta}) &= \text{var}[\hat{\theta}] + \underbrace{(E[\hat{\theta}] - \theta)^2}_{\text{bias}} = \\ &= \text{var}[\hat{\theta}] + b^2(\hat{\theta}) \\ &\quad \downarrow \quad \downarrow \\ &\text{sempre positivo} \quad \text{positivo perché } b^2 \text{ meglio vicino a 0} \\ &\text{comp-omesso ottimo tra bias} \quad \text{(quanto è vicina la media al valore atteso)} \\ &\quad \text{E varianza} \end{aligned}$$

Notiamo che non compare la varianza di ϵ , perché il problema non ha errore, se ci fosse stato quello sarebbe stato il termine irriducibile, perché se con la magia si azzeravano varianza e bias quadro di $\hat{\theta}$ quello restava.



Il **bias** è la differenza tra la media dello stimatore e il valore vero che voglio stimare.

La maggior parte della teoria si studia con bias 0, in teoria possiamo parlare di stimatori unbiased.

L'MSE lo vogliamo 0, la varianza deve essere positiva, il bias è per forza positivo, quindi vogliamo ridurre il più possibile questi termini.

Mettere bias 0 non è per forza il caso migliore, magari con bias diverso da 0 otteniamo una varianza molto più piccola.

Si cerca il compromesso ottimo tra bias e varianza per ridurre l'MSE.

Stima del parametro deterministico A

$y_i = A + \epsilon_i$
 $\hat{A} = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y}_n$
 uno stimatore della media
 = la MEDIA CAMPIONARIA
 $E[\epsilon_i] = 0$
 ϵ_i iid \rightarrow indipendenti
 identicamente
 distribuite
 $i = 1, \dots, n$
 $\text{var}[\epsilon_i] = \sigma^2$
 $E[y_i] = E[A + \epsilon_i]$
 $= A + E[\epsilon_i] = A$

A destra dimostriamo che la media di y_i è uguale ad A. Vogliamo quindi stimare A per stimare la media di y_i .

Uno stimatore della media è la media campionaria.

A cappello è lo stimatore di A.

Come pedice i va da 1 alla nostra cardinalità n.

Assumiamo che la varianza per le nostre variabili aleatorie sarà sempre la stessa. (Omoschedasticità)

Non ci ha detto CDF o PDF o cose simili ma in questo caso ci ha fornito media e varianza.

Caratteristiche dello stimatore \hat{A}

BIAS

$E[\hat{A}] = \frac{1}{n} \sum_{i=1}^n E[y_i] = \frac{1}{n} \sum_{i=1}^n A = A$
 $b(\hat{A}) = E[\hat{A}] - A = 0 \Rightarrow \hat{A} \text{ unbiased}$

Proviamo a calcolare la media dello stimatore, se ci metto i numeri diventa un valore ma qui non mettiamo i numeri.

Lo stimatore \hat{A} ha come media proprio A , quindi il bias è 0, si dice che lo stimatore è non polarizzato o unbiased.

La linearità del valore atteso la usiamo in tutti i passaggi essenzialmente.

VARIANZA \hat{A} che è lo stimatore media campionaria

Calcoliamo anche la varianza dello stimatore usando la formula normale della varianza.

$$\begin{aligned} \text{Var} [\hat{A}] &= E \{ (\hat{A} - E[\hat{A}])^2 \} \\ &= E \left\{ \left[\frac{1}{n} \sum_{i=1}^n y_i - \frac{nA}{n} \right]^2 \right\} = \\ &= \frac{1}{n^2} E \left\{ \left(\sum_{i=1}^n (y_i - A) \right)^2 \right\} \end{aligned}$$

Moltiplico e divido per n la A che sta da sola al fine di poter fare i passaggi successivi, per prima cosa per prendere in evidenza $1/n$ che diventa fuori dal quadrato $1/n^2$.

Eseguiamo il quadrato:

$$\begin{aligned} &= \frac{1}{n^2} \sum_{i=1}^n E \{ (y_i - A)^2 \} + \frac{1}{n^2} \sum_{i,j} E \{ (y_i - A)(y_j - A) \} \\ &\quad \text{dove } E[\epsilon_i] = 0 \quad \text{dato che i.d.d. allora COVARIANZA} = 0 \end{aligned}$$

Abbiamo, non a caso, trovato la covarianza.

Se le variabili sono incorrelate la covarianza è 0.

La media del prodotto è uguale al prodotto delle medie (scritte in rosa) + la covarianza, visto che tutti questi termini valgono 0 la doppia sommatoria vale 0.

$y_i - A$ è uguale ad ϵ_i ma ϵ_i ha media nulla quindi abbiamo come primo termine la varianza di epsilon con i (secondo la formula normale per la varianza, solo che la media come detto è 0).

$$\begin{aligned}
 &= \frac{1}{n^2} \sum_{i=1}^n E\{(y_i - A)^2\} + \frac{1}{n^2} \sum_{i,j} E\{(y_j - A)(y_i - A)\} \\
 &\quad \text{dato che i.i.d allora COVARIANZA} = 0 \\
 &= \frac{1}{n^2} \sum_{i=1}^n \text{var}(\epsilon_i) = \frac{1}{n^2} n \sigma^2 = \frac{\sigma^2}{n} \\
 &\quad \text{se } n \rightarrow \infty \quad \text{var} \rightarrow 0 \\
 &\quad \text{quindi è consistente l'estimatore}
 \end{aligned}$$

$$E[(\epsilon_i - E[\epsilon_i])^2] = \text{Var}(\epsilon_i)$$

visto che la media di ϵ_i è 0:

$$E[(\epsilon_i)^2] = \text{Var}(\epsilon_i)$$



Se la varianza tende a 0 lo stimatore si dice **consistente**.



La **consistenza** consiste nella capacità dello stimatore di produrre risultati sempre più affidabili (varianza sempre più piccola) al crescere della quantità di dati.

$$\text{VAR}[\hat{A}] \xrightarrow[n \rightarrow \infty]{} 0 \Rightarrow \hat{A} \rightarrow A \quad \text{consistente est}$$

Quindi già sappiamo che era unbiased, poi più è grande il campione e più la stima è vicina al valore deterministico, cioè è consistente.

Da cui consegue che l'MSE tende a 0 se n tende all'infinito.

$$\begin{aligned}
 \text{MSE}(\hat{A}) &= \text{var}(\hat{A}) + b^2(A) = \text{var}(\hat{A}) \\
 &= \frac{\sigma^2}{n} \xrightarrow[n \rightarrow \infty]{} 0 \quad \text{on training set}
 \end{aligned}$$

Quindi in queste condizioni è facile portare l'MSE vicino a 0 nel training set, ma a noi interessa minimizzare l'**errore di generalizzazione**, che ha a che fare con il **test set**.

Errore di generalizzazione - formula

Supponiamo di avere un test set con un unico campione y_0 .

$$\text{test observation} \rightarrow \text{un unico valore di test } y_0$$
$$y_0 = A + \varepsilon_0$$

Epsilon con 0 ovviamente è diversa da prima ma è sempre in accordo con la famiglia delle epsilon di prima (quindi varianza e media come prima).

$$\text{test MSE}(\hat{A}) = \underset{\text{training}}{\text{MSE}(\hat{A})} + \text{var}(\varepsilon_0)$$

σ^2 (nel nostro caso)

parte irriducibile oltre la quale non potremo scendere



Ci sarà una componente dell'MSE di test, sigma quadro, che non possiamo azzerare. Si chiama **parte irriducibile del MSE di test**.