

Ottimismo

Data Mining

CLAMSES - University of Milano-Bicocca

Aldo Solari

Riferimenti bibliografici

- AS §3.4, §3.5.1, §3.5.2, §3.5.3
- HTF §7.4, §7.5

Table of Contents

Ottimismo

Criteri basati sull'informazione

Metodo della convalida incrociata

Visto che non conosciamo $f(x)$...

La conclusione della lezione precedente è stata che dobbiamo operare un compromesso tra le componenti di distorsione e di varianza. Operativamente però non possiamo utilizzare la conoscenza di $f(x)$, ovviamente ignota in pratica

Presenteremo due approcci per stimare l'errore di previsione e scegliere il modello:

- Il concetto di ottimismo e i criteri basati sull'informazione
- Il metodo della convalida incrociata

Stima dell'errore di previsione

Calcolare l'errore quadratico medio sugli stessi dati con i quali abbiamo stimato il modello non è sembrato molto utile

Si ricordi l'esempio della regressione polinomiale: MSE_{Tr} decresce al crescere di d , fino ad arrivare a 0 per $d = n - 1$

Tuttavia se risultasse

$$\mathbb{E}(\text{MSE}_{\text{Te}}) = \mathbb{E}(\text{MSE}_{\text{Tr}}) + \text{costante}$$

allora si potrebbe stimare l'errore (atteso) di previsione come

$$\begin{aligned}\mathbb{E}(\widehat{\text{MSE}_{\text{Te}}}) &= \mathbb{E}(\widehat{\text{MSE}_{\text{Tr}}}) + \text{costante} \\ &= \text{MSE}_{\text{Tr}} + \text{costante}\end{aligned}$$

a patto di conoscere il valore della costante

Ottimismo

Si consideri il setting Fixed-X.

Chiameremo *ottimismo* la differenza

$$\mathbb{E}(\text{MSE}_{\text{Te}}) - \mathbb{E}(\text{MSE}_{\text{Tr}}) = \frac{2}{n} \sum_{i=1}^n \text{Cov}(Y_i, \hat{f}(x_i)) = \text{OptF}$$

Maggiore è la correlazione tra Y_i e $\hat{f}(x_i)$, maggiore è l'ottimismo

Abbiamo

$$\begin{aligned}\mathbb{E}[(Y_i - \hat{f}(x_i))^2] &= \mathbb{V}\text{ar}(Y_i - \hat{f}(x_i)) + (\mathbb{E}[Y_i - \hat{f}(x_i)])^2 \\ &= \mathbb{V}\text{ar}(Y_i) + \mathbb{V}\text{ar}(\hat{f}(x_i)) - \\ &\quad - 2\text{Cov}(Y_i, \hat{f}(x_i)) + (\mathbb{E}[Y_i] - \mathbb{E}[\hat{f}(x_i)])^2\end{aligned}$$

e

$$\begin{aligned}\mathbb{E}[(Y_i^* - \hat{f}(x_i))^2] &= \mathbb{V}\text{ar}(Y_i^* - \hat{f}(x_i)) + (\mathbb{E}[Y_i^* - \hat{f}(x_i)])^2 \\ &= \mathbb{V}\text{ar}(Y_i^*) + \mathbb{V}\text{ar}(\hat{f}(x_i)) - \\ &\quad - 2\text{Cov}(Y_i^*, \hat{f}(x_i)) + (\mathbb{E}[Y_i^*] - \mathbb{E}[\hat{f}(x_i)])^2\end{aligned}$$

Si noti che Y_i^* è indipendente da Y_i , ma hanno la stessa distribuzione, quindi $\mathbb{E}(Y_i^*) = \mathbb{E}(Y_i)$, $\mathbb{V}\text{ar}(Y_i^*) = \mathbb{V}\text{ar}(Y_i)$ e $\mathbb{C}\text{ov}(Y_i^*, \hat{f}(x_i)) = 0$.

Allora

$$\begin{aligned}\mathbb{E}[(Y_i^* - \hat{f}(x_i))^2] &= \mathbb{V}\text{ar}(Y_i) + \mathbb{V}\text{ar}(\hat{f}(x_i)) + (\mathbb{E}[Y_i] - \mathbb{E}[\hat{f}(x_i)])^2 \\ &= \mathbb{E}[(Y_i - \hat{f}(x_i))^2] + 2\mathbb{C}\text{ov}(Y_i, \hat{f}(x_i))\end{aligned}$$

Se facciamo la media su tutte le n osservazioni

$$\begin{aligned}\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^n (Y_i^* - \hat{f}(x_i))^2\right] &= \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^n (Y_i - \hat{f}(x_i))^2\right] + \frac{2}{n}\sum_{i=1}^n \mathbb{C}\text{ov}(Y_i, \hat{f}(x_i)) \\ \mathbb{E}(\text{MSE}_{\text{Te}}) &= \mathbb{E}(\text{MSE}_{\text{Tr}}) + \frac{2}{n}\sum_{i=1}^n \mathbb{C}\text{ov}(Y_i, \hat{f}(x_i))\end{aligned}$$

Ottimismo per il modello lineare

Si noti che

$$\frac{2}{n} \sum_{i=1}^n \mathbb{Cov}(Y_i, \hat{f}(x_i)) = \frac{2}{n} \text{tr}(\mathbb{Cov}(\mathbf{Y}, \hat{\mathbf{f}}))$$

dove $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ e $\hat{\mathbf{f}} = (\hat{f}(x_1), \dots, \hat{f}(x_n))^\top$

I valori previsti dal modello lineare si possono esprimere come

$$\hat{\mathbf{f}} = \mathbf{H}\mathbf{Y}$$

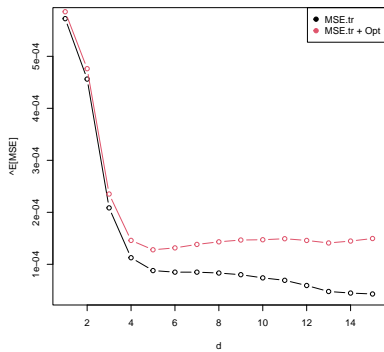
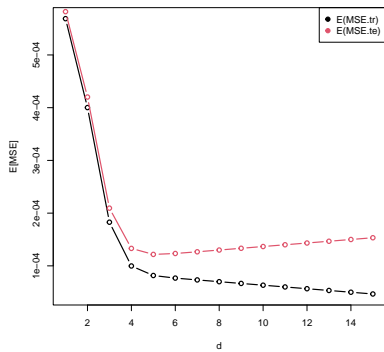
dove $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ è la matrice di proiezione (hat matrix),
simmetrica e idempotente

Abbiamo

$$\text{tr}\{\mathbb{Cov}(\mathbf{Y}, \mathbf{H}\mathbf{Y})\} = \text{tr}\{\mathbb{Cov}(\mathbf{Y}, \mathbf{Y})\mathbf{H}^\top\} = \text{tr}\{\sigma^2 \mathbf{I}_n \mathbf{H}\} = \sigma^2 \text{tr}\{\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top\}$$

e quindi

$$\text{OptF} = \frac{2\sigma^2 p}{n}$$



Cp di Mallows

$\text{OptF} = \frac{2\sigma^2 p}{n}$ richiede la conoscenza di σ^2 , che però è un valore incognito

Possiamo però sostituirlo con la sua stima

$$\hat{\sigma}^2 = \frac{\text{RSS}}{n - p} = \frac{\sum_{i=1}^n (y_i - \hat{f}(x_i))^2}{n - p}$$

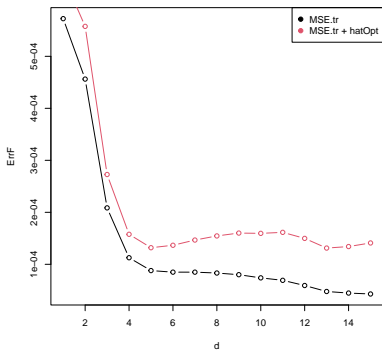
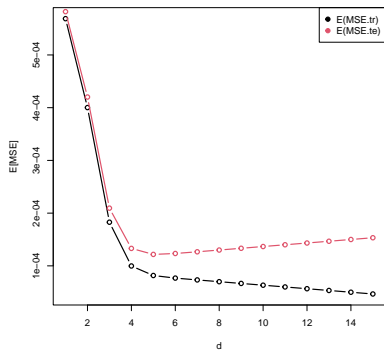
dove $\text{RSS} = n\text{MSE}_{\text{Tr}}$ è la somma dei quadrati dei residui (Residual Sum of Squares)

Quindi una stima per l'errore di previsione è data da

$$\mathbb{E}(\widehat{\text{MSE}}_{\text{Te}}) = \text{MSE}_{\text{Tr}} + \widehat{\text{OptF}}$$

Questo stimatore è noto come Cp di Mallows:

$$\text{Cp} = \text{MSE}_{\text{Tr}} + \frac{2\hat{\sigma}^2 p}{n}$$



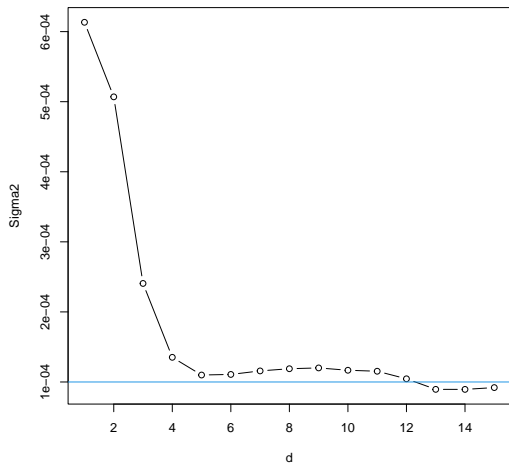


Table of Contents

Ottimismo

Criteri basati sull'informazione

Metodo della convalida incrociata

Criteri basati sull'informazione: AIC e BIC

AIC è definito come

$$\text{AIC} = -2 \cdot \text{loglikelihood}(\hat{\beta}, \hat{\sigma}^2) + 2p$$

dove per il modello lineare $-2 \cdot \text{loglikelihood}(\hat{\beta}, \hat{\sigma}^2) = n \log(\text{MSE}_{\text{Tr}})$

Per i modelli lineari, Cp e AIC sono proporzionali, e quindi il valore più piccolo per Cp corrisponde al valore più piccolo per AIC

BIC è definito come

$$\text{BIC} = -2 \cdot \text{loglikelihood}(\hat{\beta}, \hat{\sigma}^2) + \log(n)p$$

Dal momento che $\log(n) > 2$ per $n > 7$, BIC generalmente seleziona modelli meno complessi rispetto a quelli selezionati da AIC

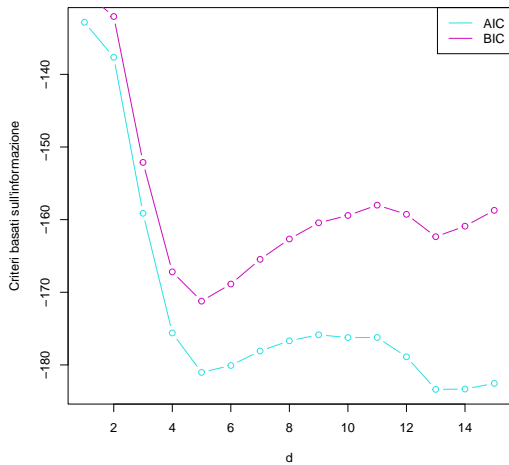


Table of Contents

Ottimismo

Criteri basati sull'informazione

Metodo della convalida incrociata

Metodo della convalida incrociata

Come notato in precedenza, stimare un modello e valutarne la performance sugli stessi dati produce un risultato troppo ottimistico

Il metodo della convalida incrociata (Cross-Validation, abbreviato CV) valuta le previsioni del modello su dati "nuovi" al fine di fornire una stima $\widehat{\text{Err}}$ dell'errore di previsione $\mathbb{E}(\text{MSE}_{\text{Te}})$

L'idea alla base del metodo è di dividere le osservazioni (data-split): una parte dei dati (training set) è utilizzata per addestrare il modello, e i dati rimanenti (test set) sono utilizzati per misurare la performance del modello

Una delle principali caratteristiche della CV è la sua universalità. CV è un metodo non parametrico che può essere applicato a qualsiasi algoritmo/modello. Questa universalità non è condivisa ad es. da C_p , che è specifico della regressione lineare

Validation set

Una semplice soluzione consiste nel dividere casualmente le n osservazioni in due parti: un insieme di addestramento e un insieme di verifica

Si stima il modello \hat{f}^{-V} sull'insieme delle osservazioni di addestramento $T \subset \{1, \dots, n\}$, e lo si utilizza per prevedere le osservazioni sull'insieme di verifica $V = \{1, \dots, n\} \setminus T$

Questo approccio fornisce una stima dell'errore di previsione (atteso)

$$\widehat{\text{Err}} = \frac{1}{\#V} \sum_{i \in V} (y_i - \hat{f}^{-V}(x_i))^2$$

Tuttavia questo procedimento riduce la numerosità delle osservazioni (ma questo non è un problema se n è veramente elevato)

Se n non è molto grande, tuttavia, questa stima può essere molto variabile

Convalida incrociata

Un modo per superare parzialmente questa arbitrarietà è dividere i dati in parti uguali V_1, \dots, V_K .

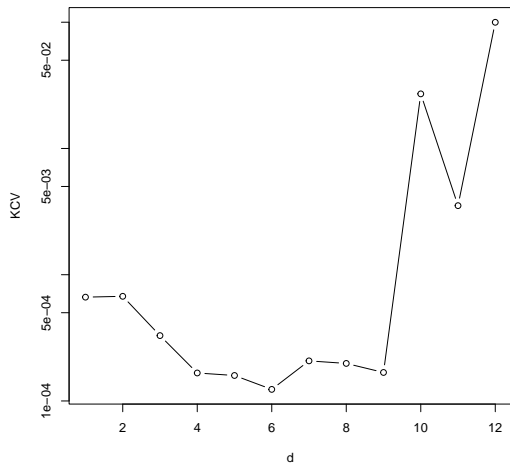
Si veda la figura Figure 5.5 del libro ISLR

Nel metodo della convalida incrociata con K porzioni (*K-fold cross-validation*) utilizziamo le osservazioni $i \notin V_k$ per addestrare il modello e le osservazioni $i \in V_k$ per valutarlo:

$$\frac{1}{\#V_k} \sum_{i \in V_k} (y_i - \hat{f}^{-V_k}(x_i))^2$$

e alla fine calcoliamo la media delle stime per stimare l'errore di previsione atteso:

$$\widehat{\text{Err}} = \frac{1}{K} \sum_{k=1}^K \left[\frac{1}{\#V_k} \sum_{i \in V_k} (y_i - \hat{f}^{-V_k}(x_i))^2 \right]$$



Leave-one-out cross validation

Nella leave-one-out cross validation (LOOCV), ciascuna osservazione viene esclusa a turno dall'insieme delle osservazioni per essere utilizzata per la verifica della previsione

Per $i = 1, \dots, n$:

- Escludere l' i -sima osservazione (x_i, y_i) - Utilizzare le rimanenti $n - 1$ osservazioni per stimare il modello \hat{f}^{-i} e valutarlo sull'osservazione esclusa (x_i, y_i) , calcolando $(y_i - \hat{f}^{-i}(x_i))^2$ - Infine, calcolare la media

$$\widehat{\text{Err}} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}^{-i}(x_i))^2$$

Si noti che LOOCV è un caso particolare di K -fold CV che corrisponde a $K = n$.

LOOCV per il modello lineare

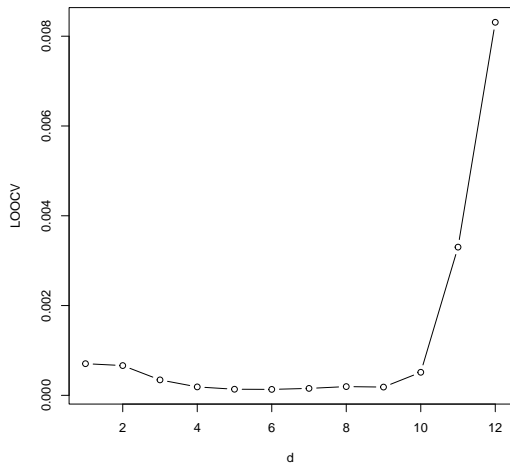
Per il modello lineare, c'è una scorciatoia per calcolare LOOCV:

$$\frac{1}{n} \sum_{i=1}^n \left(y_i - \hat{f}^{-i}(x_i) \right)^2 = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{f}(x_i)}{1 - h_{ii}} \right)^2$$

\mathbf{X} è la matrice del disegno
 $n \times p$

h_{ii} è l' i -simo elemento diagonale della matrice di proiezione

$$\mathbf{H}_{n \times n} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$



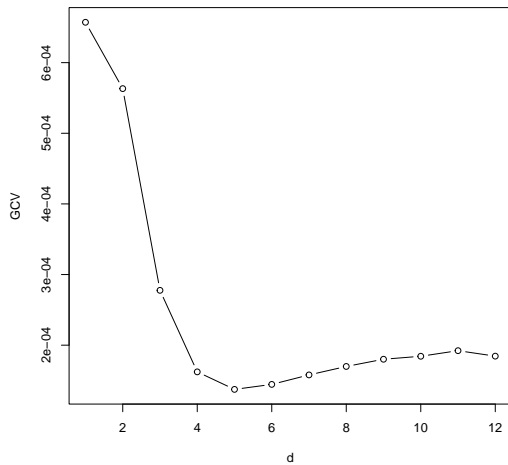
Metodo della convalida incrociata generalizzata

Nella convalida incrociata generalizzata calcoliamo

$$\widehat{\text{Err}} = \frac{\text{MSE}_{\text{Tr}}}{\left(1 - \frac{p}{n}\right)^2}$$

dove stiamo approssimando h_{ii} con la sua media

$$\frac{1}{n} \sum_{i=1}^n h_{ii} = \frac{p}{n}$$



La scelta di K

Una scelta comune per K oltre a $K = n$ è di scegliere $K = 5$ o $K = 10$

Tuttavia, trarre una conclusione generale sul CV è un compito quasi impossibile a causa della varietà delle situazioni che si possono incontrare

Compromesso distorsione-varianza per CV

Distorsione

K -fold CV con $K = 5$ o 10 fornisce una stima distorta (verso l'alto) dell'errore di previsione $\mathbb{E}(\text{MSE}_{\text{Te}})$ perchè utilizza meno osservazioni nella stima del modello ($4/5$ or $9/10$ delle osservazioni)

LOOCV ha distorsione molto bassa (utilizza $n - 1$ osservazioni)

Varianza

Solitamente LOOCV è fortemente variabile perchè è la media di n quantità estremamente correlate (perchè le stime \hat{f}^{-i} e \hat{f}^{-l} si basano su $n - 2$ osservazioni comuni), e K -fold CV con $K = 5$ o 10 a meno variabilità perchè è la media di quantità meno correlate

Si ricordi che la varianza della somma di quantità fortemente correlate è maggiore di quella con quantità mediamente correlate:

$$\text{Var}(A + B) = \text{Var}(A) + \text{Var}(B) + 2\text{Cov}(A, B)$$