

Machine Learning per la fisica applicata e la fisica delle alte energie 1

Lezione 5: Modelli lineari - regressione

Come introdotto nelle lezioni precedenti, lo scopo del machine learning (supervised) consiste nel determinare il mapping:

$$y = f(\vec{\theta}/x)$$

Una classe di modelli per f molto importante è quella dei modelli lineari in cui f è una funzione lineare nei parametri $\vec{\theta}$:

$$f(\vec{\theta}/x) = \beta_0 + \sum_{j=1}^n x_j \beta_j \quad \vec{\theta} = \{\beta_0, \beta_1, \dots, \beta_p\} \quad j=1, \dots, p$$

N.B. Le x_i possono essere: input quantitativi; la loro trasformazione ($\log, \sqrt{}, \text{quadrato}$); un'espansione in una base p. es. di tipo polinomiale come $x_2 = x_1^2$; $x_3 = x_1^3$; "interazioni" tra variabili, p. es. $x_3 = x_1 \cdot x_2$.

Dato un data set di training $(x_1, y_1), \dots, (x_n, y_n)$ i parametri β vengono tipicamente stimati col metodo dei minimi quadrati

$$\begin{aligned} \text{RSS}(\beta) &= \sum_{i=1}^N (y_i - f(x_i))^2 \quad \text{RSS: residual square sum} \\ &= \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \end{aligned}$$

La RSS altro non è che la NLL per variabili casuali INDIPENDENTI distribuite in modo gaussiano. Riscriviamo la $\text{RSS}(\beta)$ come

$$\text{RSS}(\beta) = (\vec{y} - X\beta)^T (\vec{y} - X\beta)$$

e differenziamo

$$\begin{cases} \frac{\partial \text{RSS}}{\partial \beta} = -2 X^T (\vec{y} - X\beta) \\ \frac{\partial^2 \text{RSS}}{\partial \beta \partial \beta^T} = 2 X^T X \end{cases} \quad \begin{aligned} \frac{\partial \text{RSS}}{\partial \beta} = 0 &\Leftrightarrow X^T (\vec{y} - X\beta) = 0 \\ &\Leftrightarrow \hat{\beta} = (X^T X)^{-1} X^T \vec{y} \end{aligned}$$

Le predizioni per un generico input vector x_0 sono date da

$$\hat{f}(x_0) = (1, \dots, x_0)^T \hat{\beta} \Rightarrow \hat{\vec{y}} = X \hat{\beta} = X (X^T X)^{-1} X^T \vec{y} \text{ con } \hat{y}_i = \hat{f}(x_i)$$

Osservazione: può capitare che le colonne di X non siano linearmente indipendenti, quindi X non è full rank. Questo accade se due degli input sono perfettamente correlati ($x_2 = 3x_1$). In tal caso $X^T X$ è singolare e i coefficienti $\hat{\beta}$ non sono definiti. Occorre allora pre-trattare X , per esempio rimuovendo le colonne linearmente dipendenti. Questo è fatto automaticamente negli algoritmi comuni.

Supponiamo ora che le y_i siano variabili non correlate, che abbiano varianza costante e che le x_i siano fissate. La matrice di varianza è $\text{Var}(\hat{\beta}) = (X^T X)^{-1} \sigma^2$

Assumiamo ora che il modello lineare sia il modello corretto per il valore medio, cioè il valore di aspettazione condizionale per Y è lineare in X_1, \dots, X_p . Assumiamo inoltre che le deviazioni di Y intorno al suo valore di aspettazione siano gaussiane e additive. Allora

$$y_p = \beta_0 + \sum_{j=1}^p x_j \beta_j + \varepsilon \quad \text{con } \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

Segue che

$$\hat{\beta} \sim \mathcal{N}(\beta, (X^T X)^{-1} \sigma^2)$$

Supponiamo di testare l'ipotesi per cui un particolare coefficiente $\beta_j = 0$. Si definisce lo Z-score

$$z_j = \frac{\hat{\beta}_j}{\hat{\sigma} \sqrt{v_j}} \quad \text{dove } v_j \text{ è il } j\text{-esimo elemento diagonale di } (X^T X)^{-1}$$

Sotto l'ipotesi nulla che $\beta_j = 0$, allora z_j è distribuita come t_{N-p-1} e quindi un grande valore di z_j conduce a rigettare l'ipotesi. Se $\hat{\sigma}$ è sostituito da un valore noto di σ , allora z_j è distribuita gaussiane.

mente. La differenza nei quantili nella coda della distribuzione t o delle gaussiane diventa trascurabile per N grande. L'intervallo di credibilità per $\hat{\beta}$ è:

$$\left[\hat{\beta}_j - z^{(1-\alpha)} v_j^{1/2} \hat{\sigma}, \hat{\beta}_j + z^{(1-\alpha)} v_j^{1/2} \hat{\sigma} \right]$$

dove $z^{(1-\alpha)}$ è il percentile $(1-\alpha)$ della distribuzione gaussiana.

5.1. Il teorema di Gauss - Markov

La stima ottenuta con il metodo dei minimi quadrati per i parametri β ha la varianza minore tra tutti gli stimatori lineari unbiased. Definiamo $\vartheta = a^T \beta$. Il metodo dei minimi quadrati su $a^T \hat{\beta}$

$$\hat{\vartheta} = a^T \hat{\beta} = a^T (X^T X)^{-1} X^T \bar{y}$$

Se consideriamo X una matrice fissa, dov'essere una funzione lineare $c^T \bar{y}$. Se il modello lineare è corretto, $a^T \hat{\beta}$ è unbiased poiché

$$\mathbb{E}(a^T \hat{\beta}) = \mathbb{E}(a^T (X^T X)^{-1} X^T \bar{y}) = a^T (X^T X)^{-1} X^T X \beta = a^T \beta$$

Il teorema di Gauss - Markov afferma che, se abbiamo qualunque altro stimatore $\tilde{\vartheta} = c^T \bar{y}$ che è unbiased per $a^T \beta$, cioè $\mathbb{E}(c^T \bar{y}) = a^T \beta$, allora $\text{Var}(a^T \hat{\beta}) \leq \text{Var}(c^T \bar{y})$. La dimostrazione segue dalla disuguaglianza triangolare.

Consideriamo l'errore quadratico medio di uno stimatore $\tilde{\vartheta}$ che stima ϑ

$$\text{MSE}(\tilde{\vartheta}) = \mathbb{E}(\tilde{\vartheta} - \vartheta)^2 = \text{Var}(\tilde{\vartheta}) + [\mathbb{E}(\tilde{\vartheta}) - \vartheta]^2$$

Il teorema di Gauss - Markov implica che lo stimatore ottenuto con il metodo dei minimi quadrati ha il più piccolo valore $\text{Var}(\hat{\vartheta})$ e bias nullo. Tuttavia potrebbe esistere uno stimatore con $\text{Var}(\tilde{\vartheta})$ più piccola e bias tale che $\text{MSE}(\tilde{\vartheta})$ sia complessivamente minore. Consideriamo ora di voler fare una predizione

$$y_0 = f(x_0) + \varepsilon_0$$

Allora l'errore nella predizione $\hat{f}(x_0) = x_0^T \tilde{\beta}$ è:

$$E(y_0 - \hat{f}(x_0))^2 = \sigma^2 + E(x_0^T \tilde{\beta} - f(x_0))^2 = \sigma^2 + MSE(\tilde{f}(x_0))$$

La differenza è la sola σ^2 , che è la varianza delle nuove osservabile

5.2 SUBSET SELECTION

Ci sono due ragioni per cui spesso non siamo soddisfatti del metodo dei minimi quadrati.

- 1) Il primo è l'accuratezza delle predizioni. Si può sacrificare il bias in favore di una varianza più piccola.
- 2) Il secondo è l'interpretabilità. Con un gran numero di predittori, vogliamo spesso determinare un sottoinsieme piccolo che esibisca gli effetti più rilevanti. In altri termini, siamo disposti a sacrificare alcuni dei dati più piccoli.

5.2.1 BEST SUBSET SELECTION

Si tratta di trovare $\forall k \in \{0, 1, \dots, p\}$ il sottoinsieme di dimensione k che minimizza $RSS(\tilde{\beta})$. La domanda è come scegliere k in modo da ottimizzare il trade-off tra bias e varianza. L'idea è di usare il cross-validation per stimare l'errore nella predizione e selezionare k .

- forward stepwise selection: si parte con l'intercetta e poi si aggiunge sistematicamente il predittore che migliora il fit. Si produce perciò una sequenza di modelli, con indice k , che va determinato. Vantaggi: 1) COMPUTAZIONALE: per grandi p non si può calcolare la sequenza migliore ($p \geq 40$). 2) STATISTICO: varianza minore, bias maggiore.

- backward stepwise selection: si parte con il modello completo e si rimuove in maniera sequenziale il predittore che ha l'impatto minore nel fit (valutato sulle base del me 2-score)

5.3 SHRINKAGE METHODS

L'idea è quella di rendere il processo (di selezione dei parametri) continuo.

5.3.1 RIDGE REGRESSION

Il metodo consiste nel ridurre i coefficienti di regressione imponendo una penalità sulla loro dimensione. I "ridge coefficient" minimizzano una RSS "penalizzata".

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}$$

λ è un parametro che controlla quanto "shrinkage" si applica ai parametri: maggiore è λ , più i parametri sono "schiacciati" a zero. In forma matriciale, otteniamo:

$$RSS(\lambda) = (\bar{y} - X\beta)^T (\bar{y} - X\beta) + \lambda \beta^T \beta \Rightarrow \hat{\beta}^{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T \bar{y}$$

5.3.2. LASSO

Il metodo è un metodo di shrinkage simile al precedente definito come

$$\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 \quad \text{con} \quad \sum_{j=1}^p |\beta_j| \leq t$$

$$\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2} \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}$$

La differenza sta nel fatto che la "ridge penalty" L_2 è stata sostituita con la lasso penalty L_1 . Questo implica che le soluzioni siano non-lineari in y . Ponendo t sufficientemente piccolo, il metodo realizza una sorta di selezione continua sui parametri del modello.

- RIDGE REGRESSION: shrinkage proporzionale
- LASSO: ogni coefficiente è traslato di λ

ridge regression e lasso possono essere generalizzati come

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j|^q \right\} \quad q \geq 0$$

Si può pensare a $|\beta_j|^q$ come la log-prior density per β_j o come i contorni della distribuzione a priori dei parametri.

$q=0$ subset selection; $q=1$ lasso; $q=2$ ridge regression
sono tutti Bayes estimates con un prior diverso.