

Rischio in Nearest Neighbour

In questa lezione vediamo come limitare il valore atteso del rischio statistico del classificatore prodotto da 1-NN rispetto all'estrazione del training set. Il risultato che vogliamo ottenere è del tipo

$$\mathbb{E}[\text{er}(\hat{h}_S)] \leq 2 \text{er}(f^*) + \varepsilon_m \quad (1)$$

dove \hat{h}_S indica il classificatore 1-NN prodotto sul training set S di taglia m , $\text{er}(f^*)$ è il rischio del classificatore Bayesiano ottimo e ε_m è una quantità che dipende da m . Questo tipo di risultato è diverso da quelli ottenuti in precedenza per il classificatore ERM (minimizzatore del training error), qui denotato con \tilde{h}_S , che erano del tipo

$$\mathbb{P}\left(\text{er}(\tilde{h}_S) > \min_{h \in \mathcal{H}} \text{er}(h) + \varepsilon\right) \leq \delta_{m,\varepsilon} \quad (2)$$

dove $\delta_{m,\varepsilon} > 0$ è una quantità che dipende sia da m che da ε . Le differenze sono di due tipi:

1. Il maggiorante (1) maggiora il rischio del classificatore 1-NN in termini assoluti (cioè si confronta direttamente col Bayes risk), mentre il maggiorante (2) lo maggiora relativamente al rischio del miglior classificatore in \mathcal{H} (che potrebbe essere arbitrariamente peggiore del Bayes risk).
2. Il maggiorante (2) descrive una proprietà della distribuzione del rischio dei classificatori ERM al variare di $\varepsilon > 0$, mentre il maggiorante (1) limita semplicemente il rischio di un "tipico" classificatore 1-NN

Dato un training set $S = (\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_m, Y_m)$ estratto da un modello statistico (D, η) definiamo la mappa $\pi_S : \mathbb{R}^d \rightarrow \{1, \dots, m\}$ come

$$\pi_S(\mathbf{x}) = \underset{t=1, \dots, m}{\text{argmin}} \|\mathbf{x} - \mathbf{X}_t\| \quad .$$

Il predittore 1-NN su input S è quindi definito da $\hat{h}_S(\mathbf{x}) = y_{\pi_S(\mathbf{x})}$.

Nel seguito, assumiamo che i dati \mathbf{x} generati dalla sorgente siano tali che $\max_i |x_i| \leq 1$ con probabilità uno. Ovvero, le componenti x_i delle istanze $\mathbf{x} = (x_1, \dots, x_d)$ generate dalla sorgente hanno sempre valori compresi fra -1 e $+1$. Sia $\mathcal{X} \subset \mathbb{R}^d$ il sottoinsieme delle istanze \mathbf{x} con questa proprietà. Esprimeremo il maggiorante sul valore atteso del rischio di 1-NN in termini di una quantità che caratterizza il modello statistico (D, η) , ovvero il più piccolo $c > 0$ tale che

$$|\eta(\mathbf{x}) - \eta(\mathbf{x}')| \leq c \|\mathbf{x} - \mathbf{x}'\| \quad \text{per ogni } \mathbf{x}, \mathbf{x}' \in \mathcal{X} \quad .$$

Si noti che $c < \infty$ implica che η è una funzione continua. Possiamo quindi scrivere

$$\eta(\mathbf{x}') \leq \eta(\mathbf{x}) + c \|\mathbf{x} - \mathbf{x}'\| \quad (3)$$

$$1 - \eta(\mathbf{x}') \leq 1 - \eta(\mathbf{x}) + c \|\mathbf{x} - \mathbf{x}'\| \quad (4)$$

Siccome i dati sono estratti in modo indipendente, per ogni (\mathbf{x}, y) e (\mathbf{x}', y') vale che

$$\mathbb{P}(Y = y, Y' = y' \mid \mathbf{X} = \mathbf{x}, \mathbf{X}' = \mathbf{x}') = \mathbb{P}(Y = y \mid \mathbf{X} = \mathbf{x})\mathbb{P}(Y' = y' \mid \mathbf{X}' = \mathbf{x}') . \quad (5)$$

Ricordando che $\widehat{h}_S(\mathbf{x}) = y_{\pi_S(\mathbf{x})}$, notiamo che, per ogni coppia di istanze \mathbf{x}, \mathbf{x}' ,

$$\begin{aligned} \mathbb{P}(Y \neq Y' \mid \mathbf{X} = \mathbf{x}, \mathbf{X}' = \mathbf{x}') &= \mathbb{P}(Y = 1 \mid \mathbf{X} = \mathbf{x})\mathbb{P}(Y' = -1 \mid \mathbf{X}' = \mathbf{x}') \\ &\quad + \mathbb{P}(Y = -1 \mid \mathbf{X} = \mathbf{x})\mathbb{P}(Y' = 1 \mid \mathbf{X}' = \mathbf{x}') \\ &= \eta(\mathbf{x})(1 - \eta(\mathbf{x}')) + (1 - \eta(\mathbf{x}))\eta(\mathbf{x}') \end{aligned}$$

dove la probabilità è rispetto all'estrazione di S e (\mathbf{X}, Y) e abbiamo usato (5). Applicando le due disuguaglianze (3) e (4), possiamo allora scrivere

$$\begin{aligned} \mathbb{E}[\text{er}(\widehat{h}_S)] &= \mathbb{E}[\mathbb{I}\{\widehat{h}_S(\mathbf{X}) \neq Y\}] \quad \text{il secondo valore atteso è rispetto all'estrazione di } S \text{ e } (\mathbf{X}, Y) \\ &= \mathbb{P}(\widehat{h}_S(\mathbf{X}) \neq Y) \quad \text{ricordando che } \mathbb{P}(A) = \mathbb{E}[\mathbb{I}\{A\}] \text{ per ogni evento } A \\ &= \mathbb{P}(Y_{\pi_S(\mathbf{X})} \neq Y) \quad \text{per definizione di 1-NN} \\ &= \mathbb{E}[\mathbb{I}\{Y_{\pi_S(\mathbf{X})} \neq Y\}] \\ &= \mathbb{E}\left[\mathbb{E}[\mathbb{I}\{Y' \neq Y\} \mid \mathbf{X} = \mathbf{x}, \mathbf{X}_{\pi_S(\mathbf{x})} = \mathbf{x}']\right] \quad \text{usando } \mathbb{E}[Y] = \mathbb{E}[\mathbb{E}[Y \mid X]] \text{ che vale per ogni } X, Y \\ &= \mathbb{E}\left[\eta(\mathbf{X})(1 - \eta(\mathbf{X}')) + (1 - \eta(\mathbf{X}))\eta(\mathbf{X}')\right] \quad \text{dove } \mathbf{X}' = \mathbf{X}_{\pi_S(\mathbf{X})} \\ &\leq \mathbb{E}\left[\eta(\mathbf{X})(1 - \eta(\mathbf{X})) + \eta(\mathbf{X})c\|\mathbf{X} - \mathbf{X}'\| + (1 - \eta(\mathbf{X}))\eta(\mathbf{X}) + (1 - \eta(\mathbf{X}))c\|\mathbf{X} - \mathbf{X}'\|\right] \\ &\leq 2\mathbb{E}\left[\eta(\mathbf{X})(1 - \eta(\mathbf{X}))\right] + c\mathbb{E}\left[\|\mathbf{X} - \mathbf{X}_{\pi_S(\mathbf{X})}\|\right] \end{aligned}$$

dove i valori attesi e le probabilità sono rispetto alle estrazioni indipendenti di S e di (\mathbf{X}, Y) .

Ora ricordiamo che il rischio del classificatore Bayesiano ottimo f^* soddisfa

$$\text{er}(f^*) = \mathbb{E}\left[\min\{\eta(\mathbf{X}), 1 - \eta(\mathbf{X})\}\right] \geq \mathbb{E}\left[\eta(\mathbf{X})(1 - \eta(\mathbf{X}))\right] .$$

Quindi abbiamo che

$$\mathbb{E}[\text{er}(\widehat{h}_S)] \leq 2\text{er}(f^*) + c\mathbb{E}\left[\|\mathbf{X} - \mathbf{X}_{\pi_S(\mathbf{X})}\|\right] .$$

Per gestire il termine contenente il valore atteso di $\|\mathbf{X} - \mathbf{X}_{\pi_S(\mathbf{X})}\|$ suddividiamo lo spazio delle istanze \mathcal{X} in ipercubetti d -dimensionali di lato $\varepsilon > 0$ —si veda la Figura 1 per un esempio bidimensionale. Siano C_1, \dots, C_r gli ipercubetti. Ora possiamo limitare $\|\mathbf{X} - \mathbf{X}_{\pi_S(\mathbf{X})}\|$ distinguendo il caso in cui \mathbf{X} appartiene un C_i in cui c'è almeno un punto \mathbf{X}_t di S e il caso in cui \mathbf{X} appartiene ad un C_i dove non c'è neanche un punto di S . Nel primo caso $\|\mathbf{X} - \mathbf{X}_{\pi_S(\mathbf{X})}\|$ è al più la lunghezza della diagonale dell'ipercubetto, cioè $\varepsilon\sqrt{d}$ —si veda la parte sinistra della Figura 1. Nel secondo caso $\|\mathbf{X} - \mathbf{X}_{\pi_S(\mathbf{X})}\|$ è maggiorata dalla lunghezza della diagonale di \mathcal{X} , cioè $2\sqrt{d}$ —si veda la parte destra della Figura 1. Quindi possiamo scrivere

$$\begin{aligned} \mathbb{E}\left[\|\mathbf{X} - \mathbf{X}_{\pi_S(\mathbf{X})}\|\right] &\leq \mathbb{E}\left[\varepsilon\sqrt{d}\sum_{i=1}^r \mathbb{I}\{C_i \cap S \neq \emptyset\}\mathbb{I}\{\mathbf{X} \in C_i\} + 2\sqrt{d}\sum_{i=1}^r \mathbb{I}\{C_i \cap S = \emptyset\}\mathbb{I}\{\mathbf{X} \in C_i\}\right] \\ &= \varepsilon\sqrt{d}\mathbb{E}\left[\sum_{i=1}^r \mathbb{I}\{C_i \cap S \neq \emptyset\}\mathbb{I}\{\mathbf{X} \in C_i\}\right] + 2\sqrt{d}\sum_{i=1}^r \mathbb{E}\left[\mathbb{I}\{C_i \cap S = \emptyset\}\mathbb{I}\{\mathbf{X} \in C_i\}\right] \end{aligned}$$

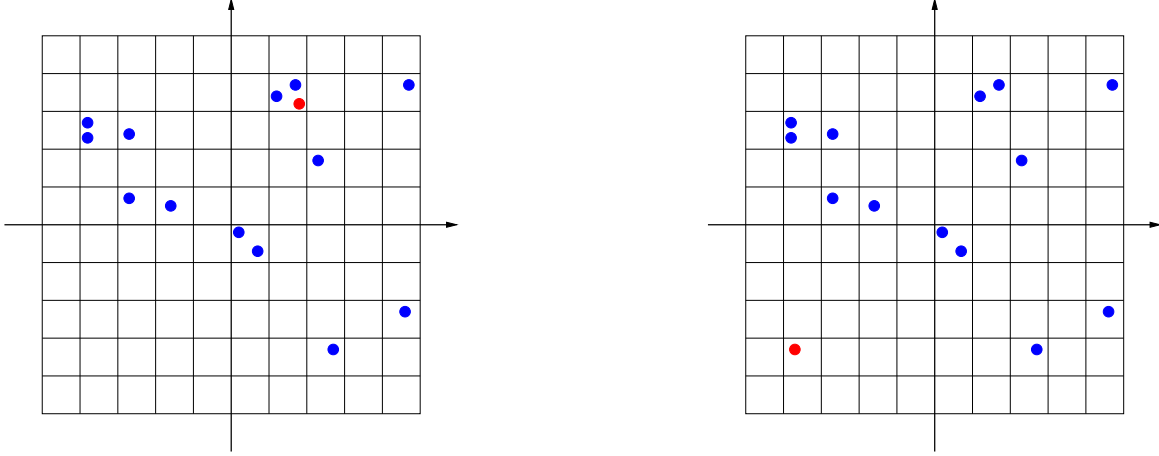


Figura 1: Esempio bidimensionale della costruzione utilizzata nell'analisi di 1-NN. Nella parte sinistra \mathbf{X} (il punto rosso) si trova nello stesso quadratino C_i rispetto al punto $\mathbf{X}_{\pi_S(\mathbf{X})}$ del training set S più vicino ad esso e quindi $\|\mathbf{X} - \mathbf{X}_{\pi_S(\mathbf{X})}\|$ è limitato dalla lunghezza della diagonale del quadratino. Nella parte destra, invece, non ci sono punti di training nel quadratino di \mathbf{X} e quindi $\|\mathbf{X} - \mathbf{X}_{\pi_S(\mathbf{X})}\|$ è limitato dalla lunghezza della diagonale del quadrato grande.

dove abbiamo usato la linearità del valore atteso nell'ultimo passo. Ora osserviamo che per ogni S e \mathbf{X} ,

$$\sum_{i=1}^r \mathbb{I}\{C_i \cap S \neq \emptyset\} \mathbb{I}\{\mathbf{X} \in C_i\} \in \{0, 1\}$$

dato che $\mathbf{X} \in C_i$ per un solo $i = 1, \dots, d$. Quindi,

$$\mathbb{E} \left[\sum_{i=1}^r \mathbb{I}\{C_i \cap S \neq \emptyset\} \mathbb{I}\{\mathbf{X} \in C_i\} \right] \leq 1 .$$

Per l'altro termine, usiamo il fatto che il training set S e \mathbf{X} sono indipendenti

$$\sum_{i=1}^r \mathbb{E} [\mathbb{I}\{C_i \cap S = \emptyset\} \mathbb{I}\{\mathbf{X} \in C_i\}] = \sum_{i=1}^r \mathbb{E} [\mathbb{I}\{C_i \cap S = \emptyset\}] \mathbb{E} [\mathbb{I}\{\mathbf{X} \in C_i\}] = \sum_{i=1}^r \mathbb{P}(C_i \cap S = \emptyset) \mathbb{P}(\mathbf{X} \in C_i) .$$

Dato che S è un campione statistico composto da m esempi indipendenti, possiamo scrivere

$$\mathbb{P}(C_i \cap S = \emptyset) = (1 - \mathbb{P}(\mathbf{X} \in C_i))^m \leq \exp(-m\mathbb{P}(\mathbf{X} \in C_i))$$

dove nell'ultimo passo abbiamo usato la maggiorazione elementare $(1-p)^m \leq e^{-pm}$. Ricapitolando,

$$\begin{aligned} \mathbb{E} [\|\mathbf{X} - \mathbf{X}_{\pi_S(\mathbf{X})}\|] &\leq \varepsilon\sqrt{d} + (2\sqrt{d}) \sum_{i=1}^r \exp(-m\mathbb{P}(\mathbf{X} \in C_i)) \mathbb{P}(\mathbf{X} \in C_i) \\ &\leq \varepsilon\sqrt{d} + (2\sqrt{d}) r \max_{i=1, \dots, r} \exp(-m\mathbb{P}(\mathbf{X} \in C_i)) \mathbb{P}(\mathbf{X} \in C_i) \\ &\leq \varepsilon\sqrt{d} + (2\sqrt{d}) r \max_{0 \leq p \leq 1} e^{-pm} p . \end{aligned}$$

Studiando la funzione $g(p) = e^{-pm}p$ troviamo il massimo per $p = \frac{1}{m}$. Quindi, sostituendo,

$$\mathbb{E}\left[\|\mathbf{X} - \mathbf{X}_{\pi_S(\mathbf{X})}\|\right] \leq \varepsilon\sqrt{d} + \left(2\sqrt{d}\right) \frac{r}{em} = \sqrt{d} \left(\varepsilon + \frac{2}{em} \left(\frac{2}{\varepsilon}\right)^d \right)$$

dove abbiamo usato il fatto che il numero r di ipercubetti è pari a $\left(\frac{2}{\varepsilon}\right)^d$. Mettendo tutto assieme troviamo che

$$\mathbb{E}[\text{er}(\widehat{h}_S)] \leq 2 \text{er}(f^*) + c\sqrt{d} \left(\varepsilon + \frac{2}{em} \left(\frac{2}{\varepsilon}\right)^d \right)$$

Dato che la scelta di $0 < \varepsilon < 1$ è libera nell'analisi, possiamo porre $\varepsilon = 2m^{-1/(d+1)}$. Questo ci dà

$$\varepsilon + \frac{2}{em} \left(\frac{2}{\varepsilon}\right)^d = 2m^{-1/(d+1)} + \frac{2^{d+1}2^{-d}m^{d/(d+1)}}{em} = 2m^{-1/(d+1)} \left(1 + \frac{1}{e}\right) \leq 4m^{-1/(d+1)} .$$

Sostituendo, otteniamo infine

$$\mathbb{E}[\text{er}(\widehat{h}_S)] \leq 2 \text{er}(f^*) + c4m^{-1/(d+1)}\sqrt{d} .$$

Da questa analisi possiamo trarre due conclusioni.

1. Per $m \rightarrow \infty$, $\text{er}(f^*) \leq \mathbb{E}[\text{er}(\widehat{h}_S)] \leq 2 \text{er}(f^*)$. Ovvero, il rischio di 1-NN è compreso fra il Bayes error e due volte il Bayes error.
2. Perché $\mathbb{E}[\text{er}(\widehat{h}_S)]$ sia al più $2 \text{er}(f^*) + \varepsilon$ il training set dev'essere almeno di taglia $m \geq \left(\frac{4c}{\varepsilon}\sqrt{d}\right)^{d+1}$, cioè esponenziale nel numero d di attributi. Questo mostra che 1-NN paga un prezzo alto in termini di overfitting per potersi confrontare direttamente con il Bayes risk.

La dimostrazione può essere generalizzata per studiare, sotto le medesime assunzioni, il rischio del classificatore \widehat{h}_S prodotto da k -NN, che risulta essere maggiorato come segue

$$\mathbb{E}[\text{er}(\widehat{h}_S)] \leq \left(1 + \sqrt{\frac{8}{k}}\right) \text{er}(f^*) + \mathcal{O}(km^{-1/(d+1)}) .$$