

Consistenza e algoritmi nonparametrici

Fissiamo un modello statistico (D, η) per classificazione binaria. Per ogni algoritmo di apprendimento A , denotiamo con $A(S_m)$ il classificatore prodotto da A quando il training set S_m contenente m esempi è fornito in input. Denotiamo con $\text{er}(A(S_m))$ il rischio di questo classificatore rispetto al modello (D, η) . Quindi, se $A(S_m) = h$ allora $\text{er}(A(S_m)) = \mathbb{P}(h(\mathbf{X}) \neq Y)$. Infine, indichiamo con $\text{er}(f^*)$ il rischio del classificatore Bayesiano ottimo f^* , sempre rispetto a (D, η) .

L'algoritmo A è detto consistente rispetto a (D, η) quando

$$\lim_{m \rightarrow \infty} \mathbb{E}[\text{er}(A(S_m))] = \text{er}(f^*)$$

dove il valore atteso è rispetto all'estrazione del training set S_m da (D, η) . Quindi, la consistenza è una proprietà asintotica che certifica la capacità dell'algoritmo di raggiungere in media le prestazioni del Bayesiano ottimo al crescere del training set.

Se indichiamo con k -NN l'algoritmo k -Nearest Neighbor, allora per ogni (D, η) vale

$$\lim_{m \rightarrow \infty} \mathbb{E}[\text{er}(k\text{-NN}(S_m))] \leq \text{er}(f^*) + 2\sqrt{\frac{\text{er}(f^*)}{k}}. \quad (1)$$

Perciò per ogni k fissato k -NN non è consistente. Possiamo confrontare (1) con il risultato

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

enunciato in precedenza sotto l'assunzione

$$|\eta(\mathbf{x}) - \eta(\mathbf{x}')| \leq c \|\mathbf{x} - \mathbf{x}'\| \quad \text{per ogni } \mathbf{x}, \mathbf{x}' \in \mathcal{X}. \quad (2)$$

Se in (1) rendiamo k dipendente da m in modo che cresca non troppo velocemente, cioè se $k = k_m$ tale che $k_m \rightarrow \infty$ e $k_m = o(m)$, allora è possibile dimostrare che k -NN diventa consistente.

Sotto determinate assunzioni su (D, η) , è possibile dimostrare che anche l'algoritmo greedy di costruzione di classificatori ad albero diventa consistente quando ogni foglia è finale per un numero sufficiente di punti del training set (dettagli omessi).

Il motivo per cui un algoritmo consistente può risultare in pratica peggiore di uno non consistente è legato alla velocità di convergenza. Infatti, la consistenza per ogni modello (D, η) può essere pagata con una velocità di convergenza *arbitrariamente lenta* al Bayes error, come stabilito dal seguente risultato.

Teorema 1 (No Free Lunch) *Sia a_1, a_2, \dots una successione convergente a zero di numeri positivi tali che $\frac{1}{16} \geq a_1 \geq a_2 \geq \dots$. Per ogni algoritmo di apprendimento A esiste un modello statistico (D, η) tale che $\text{er}(f^*) = 0$ e simultaneamente $\mathbb{E}[\text{er}(A(S_m))] \geq a_m$ per ogni $m \geq 1$.*

Si noti che se $er(f^*) = 0$ allora dev'essere $\eta(\mathbf{x}) \in \{0, 1\}$ per ogni \mathbf{x} . Ovvero η è discontinua e la condizione (2) non può essere soddisfatta. Si noti anche che il Teorema 1 non impedisce a un algoritmo consistente di convergere rapidamente per specifiche distribuzioni (D, η) . Analogamente, un algoritmo non consistente può convergere rapidamente a specifici valori di rischio a patto che siano superiori al Bayes error. Il Teorema 1 dice solo che se A arriva al Bayes error, allora non è possibile che ci arrivi in fretta per tutte le distribuzioni

Gli algoritmi consistenti, come k -NN e gli algoritmi per i predittori ad albero, vengono anche detti algoritmi **nonparametrici**. Questa terminologia si riferisce al fatto che i classificatori prodotti non sono rappresentabili con un numero predeterminato di parametri (come ad esempio succede per i classificatori lineari). Infatti, la struttura di un classificatore nonparametrico non è fissata, ma viene determinata dai dati di training (si pensi alla costruzione di un predittore ad albero guidata dal training set).

In generale,

- Se il training set è piccolo conviene utilizzare un algoritmo parametrico.
- Se d è grande conviene usare un algoritmo parametrico perché anche quando η soddisfa (2) la convergenza nonparametrica è comunque dell'ordine $\mathcal{O}(m^{-1/(d+1)})$ mentre quella parametrica è dell'ordine $\mathcal{O}(m^{-1/2})$.
- Negli altri casi (dataset grandi con pochi attributi) possiamo provare ad usare un algoritmo nonparametrico.