

Stability bounds e controllo del rischio in SVM

Come sappiamo, un modo per controllare l'overfitting è quello di impedire all'algoritmo di apprendimento di produrre dei predittori che abbiano un training error arbitrariamente basso su qualunque training set. Questo comportamento può essere ottenuto imponendo una condizione di "stabilità" sull'algoritmo stesso. Ovvero, richiedendo che il training error del predittore prodotto dall'algoritmo non cambi significativamente quando il training set viene leggermente perturbato.

Fissiamo un training set $S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m))$ e usiamo \mathbf{z}_t per denotare il t -esimo esempio $(\mathbf{x}_t, y_t) \in \mathbb{R}^d \times \mathbb{R}$. Data una funzione di perdita ℓ e un predittore h , indichiamo con $\ell(h, \mathbf{z}_t)$ la perdita di h sull'esempio \mathbf{z}_t e usiamo

$$\ell_S(h) = \frac{1}{m} \sum_{t=1}^m \ell(h, \mathbf{z}_t)$$

per denotare il training error di h . Fissato inoltre un modello statistico \mathcal{D} , denotiamo con $\ell_{\mathcal{D}}(h)$ il rischio statistico di h . Nel seguito, assumiamo che S sia un campione statistico estratto da \mathcal{D} e usiamo la notazione $S^{(t)}$ per indicare il training set S in cui in t -esimo esempio (\mathbf{x}_t, y_t) è stato sostituito con un altro esempio $\mathbf{z}'_t = (\mathbf{x}'_t, y'_t)$ estratto da \mathcal{D} in modo indipendente da S .

Fissato un algoritmo di apprendimento A , indichiamo con h_S e $h_{S^{(t)}}$ i predittori generati da A con input, rispettivamente, S e $S^{(t)}$. Diciamo che A è ε -stabile se, per ogni $t = 1, \dots, m$

$$\mathbb{E}[\ell(h_S, \mathbf{z}'_t) - \ell(h_{S^{(t)}}, \mathbf{z}'_t)] \leq \varepsilon$$

dove il valore atteso è calcolato rispetto alle estrazioni di S e $\mathbf{z}'_t = (\mathbf{x}'_t, y'_t)$. In altri termini, un algoritmo è stabile se, in media rispetto all'estrazione del training set, l'algoritmo produce un predittore il cui errore su un qualunque elemento del training set è vicino a quello commesso sullo stesso elemento da un predittore addestrato su un training set dove quell'elemento è sostituito da un altro.

I due prossimi risultati mostrano che un algoritmo stabile non soffre di overfitting.

Teorema 1 *Se A è ε -stabile allora $\mathbb{E}[\ell_{\mathcal{D}}(h_S) - \ell_S(h_S)] \leq \varepsilon$.*

DIMOSTRAZIONE. Consideriamo $S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m))$ e $S' = ((\mathbf{x}'_1, y'_1), \dots, (\mathbf{x}'_m, y'_m))$ estratti in modo indipendente da \mathcal{D} . Allora

$$\mathbb{E}[\ell_S(h_S)] = \mathbb{E}\left[\frac{1}{m} \sum_{t=1}^m \ell(h_S, \mathbf{z}_t)\right] = \frac{1}{m} \sum_{t=1}^m \mathbb{E}[\ell(h_S, \mathbf{z}_t)] = \frac{1}{m} \sum_{t=1}^m \mathbb{E}[\ell(h_{S^{(t)}}, \mathbf{z}'_t)] .$$

Inoltre,

$$\ell_{\mathcal{D}}(h_S) = \mathbb{E}[\ell(h_S, \mathbf{z}'_t) | S] = \frac{1}{m} \sum_{t=1}^m \mathbb{E}[\ell(h_S, \mathbf{z}'_t) | S]$$

che implica, facendo la media rispetto all'estrazione di S ,

$$\mathbb{E}[\ell_{\mathcal{D}}(h_S)] = \frac{1}{m} \sum_{t=1}^m \mathbb{E}[\ell(h_S, \mathbf{z}'_t)] .$$

Quindi,

$$\mathbb{E}[\ell_{\mathcal{D}}(h_S) - \ell_S(h_S)] = \frac{1}{m} \sum_{t=1}^m \mathbb{E}[\ell(h_S, \mathbf{z}'_t) - \ell(h_{S^{(t)}}, \mathbf{z}'_t)] \leq \varepsilon$$

per l'assunzione di stabilità. □

Dato che la stabilità e la minimizzazione del training error sono obiettivi antagonisti, l'algoritmo ERM (che minimizza il rischio empirico) non è necessariamente stabile. D'altra parte, se un algoritmo di apprendimento è stabile e simultaneamente in grado di produrre predittori con basso rischio empirico (non necessariamente minimo) allora il suo errore di varianza è basso.

Teorema 2 *Se A è ε -stabile e inoltre minimizza in modo approssimato il rischio empirico in una classe \mathcal{H} di predittori, ovvero*

$$\ell_S(h_S) \leq \inf_{h \in \mathcal{H}} \ell_S(h) + \gamma$$

per un qualche $\gamma > 0$, allora

$$\mathbb{E}[\ell_{\mathcal{D}}(h_S)] \leq \inf_{h \in \mathcal{H}} \ell_{\mathcal{D}}(h) + \varepsilon + \gamma .$$

DIMOSTRAZIONE. Sia h^* il predittore a rischio minimo in \mathcal{H} . Allora

$$\begin{aligned} \mathbb{E}[\ell_{\mathcal{D}}(h_S)] &= \mathbb{E} \left[\underbrace{\ell_{\mathcal{D}}(h_S) - \ell_S(h_S)}_{\leq \varepsilon \text{ (stabilità)}} \right] + \mathbb{E} \left[\underbrace{\ell_S(h_S) - \ell_S(h^*)}_{\leq \gamma \text{ (ERM approssimato)}} \right] + \mathbb{E}[\ell_S(h^*)] \\ &\leq \varepsilon + \gamma + \mathbb{E}[\ell_S(h^*)] . \end{aligned}$$

La dimostrazione è conclusa osservando che $\mathbb{E}[\ell_S(h^*)] = \ell_{\mathcal{D}}(h^*)$, dato che il valore atteso del rischio empirico è il rischio. □

Nel caso di predittori parametrizzati da un vettore $\mathbf{w} \in \mathbb{R}^d$ (come quelli lineari), il minimizzatore del rischio empirico, rispetto ad una funzione di perdita data ℓ , può essere reso stabile aggiungendo a ℓ un termine cosiddetto di regolarizzazione. In realtà abbiamo bisogno anche di un'altra condizione, ovvero che la funzione $\ell(\cdot, \mathbf{z})$, dove $\ell(\mathbf{w}, \mathbf{z})$ è l'errore di \mathbf{w} sull'esempio \mathbf{z} , sia convessa e Lipschitz. Ricordiamo che Lipschitz significa che esiste una costante L tale che $|\ell(\mathbf{w}, \mathbf{z}) - \ell(\mathbf{w}', \mathbf{z})| \leq L \|\mathbf{w} - \mathbf{w}'\|$ per ogni $\mathbf{w}, \mathbf{w}' \in \mathbb{R}^d$ e per ogni $\mathbf{z} = (\mathbf{x}, y)$. Non sono richieste altre assunzioni su ℓ .

Teorema 3 *Sia $\ell(\mathbf{w}, \mathbf{z})$ una funzione di perdita tale che $\ell(\cdot, \mathbf{z})$ è convessa, differenziabile e Lipschitz con costante $L > 0$. Allora l'algoritmo di apprendimento che, su input S , produce*

$$\mathbf{w}_S = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \left(\ell_S(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \right)$$

è $(2L)^2/(\lambda m)$ -stabile per ogni $\lambda > 0$.

DIMOSTRAZIONE. Definiamo

$$F_S(\mathbf{w}) = \ell_S(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

e osserviamo che

$$\mathbf{w}_S = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} F_S(\mathbf{w}) \quad \text{e} \quad \mathbf{w}_{S^{(t)}} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} F_{S^{(t)}}(\mathbf{w}) .$$

Per dimostrare la stabilità dobbiamo maggiorare $\mathbb{E}[\ell(\mathbf{w}_S, \mathbf{z}'_t) - \ell(\mathbf{w}_{S^{(t)}}, \mathbf{z}'_t)]$. Nella dimostrazione otterremo un risultato più forte, andando a maggiorare la quantità $\ell(\mathbf{w}_S, \mathbf{z}'_t) - \ell(\mathbf{w}_{S^{(t)}}, \mathbf{z}'_t)$ per ogni S e \mathbf{z}'_t . Come primo passo, usiamo Lipschitz per scrivere

$$\ell(\mathbf{w}_S, \mathbf{z}'_t) - \ell(\mathbf{w}_{S^{(t)}}, \mathbf{z}'_t) \leq L \|\mathbf{w}_S - \mathbf{w}_{S^{(t)}}\| . \quad (1)$$

Quindi, procediamo a maggiorare $\|\mathbf{w}_S - \mathbf{w}_{S^{(t)}}\|$. Introduciamo le abbreviazioni $\mathbf{w} = \mathbf{w}_S$ e $\mathbf{w}' = \mathbf{w}_{S^{(t)}}$. Allora

$$\begin{aligned} F_S(\mathbf{w}') - F_S(\mathbf{w}) &= \ell_S(\mathbf{w}') - \ell_S(\mathbf{w}) + \frac{\lambda}{2} (\|\mathbf{w}'\|^2 - \|\mathbf{w}\|^2) \\ &= \ell_{S^{(t)}}(\mathbf{w}') - \ell_{S^{(t)}}(\mathbf{w}) + \frac{\ell(\mathbf{w}', \mathbf{z}_t) - \ell(\mathbf{w}, \mathbf{z}_t)}{m} - \frac{\ell(\mathbf{w}', \mathbf{z}'_t) - \ell(\mathbf{w}, \mathbf{z}'_t)}{m} + \frac{\lambda}{2} (\|\mathbf{w}'\|^2 - \|\mathbf{w}\|^2) \\ &= F_{S^{(t)}}(\mathbf{w}') - F_{S^{(t)}}(\mathbf{w}) + \frac{\ell(\mathbf{w}', \mathbf{z}_t) - \ell(\mathbf{w}, \mathbf{z}_t)}{m} - \frac{\ell(\mathbf{w}', \mathbf{z}'_t) - \ell(\mathbf{w}, \mathbf{z}'_t)}{m} \\ &\leq \frac{|\ell(\mathbf{w}', \mathbf{z}_t) - \ell(\mathbf{w}, \mathbf{z}_t)|}{m} + \frac{|\ell(\mathbf{w}', \mathbf{z}'_t) - \ell(\mathbf{w}, \mathbf{z}'_t)|}{m} \\ &\leq \frac{2L}{m} \|\mathbf{w} - \mathbf{w}'\| \end{aligned}$$

dove la prima disuguaglianza vale perché $\mathbf{w}' = \mathbf{w}_{S^{(t)}}$ minimizza $F_{S^{(t)}}$ e la seconda disuguaglianza vale per il fatto che $\ell(\cdot, \mathbf{z})$ è Lipschitz.

Proseguiamo notando che la funzione F_S è λ -fortemente convessa: infatti $\ell(\cdot, \mathbf{z})$ è convessa, $\frac{\lambda}{2} \|\mathbf{w}\|^2$ è λ -fortemente convessa e quindi la loro somma è λ -fortemente convessa. Allora, per definizione di funzione fortemente convessa,

$$F_S(\mathbf{w}') \geq F_S(\mathbf{w}) + \nabla F_S(\mathbf{w})^\top (\mathbf{w}' - \mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w} - \mathbf{w}'\|^2 .$$

Dato che \mathbf{w} è il minimo di F_S , $\nabla F_S(\mathbf{w}) = \mathbf{0}$ e quindi

$$F_S(\mathbf{w}') - F_S(\mathbf{w}) = \left(\ell_S(\mathbf{w}') + \frac{\lambda}{2} \|\mathbf{w}'\|^2 \right) - \left(\ell_S(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \right) \geq \frac{\lambda}{2} \|\mathbf{w} - \mathbf{w}'\|^2$$

Combinando assieme le due disuguaglianze otteniamo

$$\frac{\lambda}{2} \|\mathbf{w} - \mathbf{w}'\|^2 \leq \frac{2L}{m} \|\mathbf{w} - \mathbf{w}'\| \quad \text{ovvero} \quad \|\mathbf{w} - \mathbf{w}'\| \leq \frac{4L}{\lambda m}$$

che, combinato con (1) dimostra la stabilità di $\mathbf{w} = \mathbf{w}_S$. \square

Mostriamo ora come la nozione di stabilità può essere usata per ottenere un maggiorante sul rischio del predittore SVM. Prima di tutto, ricordiamo che la hinge loss $\ell(\mathbf{w}, (\mathbf{x}, y)) = [1 - y \mathbf{w}^\top \mathbf{x}]_+$ per $(\mathbf{x}, y) \in \mathbb{R}^d \times \{-1, +1\}$ è convessa e Lipschitz con costante $L = 1$.

Teorema 4 *Dato un training set S , la soluzione SVM*

$$\mathbf{w}_S = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \left(\ell_S(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \right)$$

soddisfa

$$\mathbb{E}[\ell_{\mathcal{D}}(\mathbf{w}_S)] \leq \min_{\mathbf{u} \in \mathbb{R}^d} \left(\ell_{\mathcal{D}}(\mathbf{u}) + \frac{\lambda}{2} \|\mathbf{u}\|^2 \right) + \frac{4}{\lambda m} .$$

DIMOSTRAZIONE. Chiaramente, per ogni $\mathbf{u} \in \mathbb{R}^d$ vale

$$\ell_S(\mathbf{w}_S) \leq \ell_S(\mathbf{u}) + \frac{\lambda}{2} \|\mathbf{w}_S\|^2 \leq \ell_S(\mathbf{u}) + \frac{\lambda}{2} \|\mathbf{u}\|^2 . \quad (2)$$

Quindi, dato che per il Teorema 3 la soluzione \mathbf{w}_S è $(4L^2)/(\lambda m)$ -stabile per $L = 1$,

$$\begin{aligned} \mathbb{E}[\ell_{\mathcal{D}}(\mathbf{w}_S)] &\leq \mathbb{E}[\ell_S(\mathbf{w}_S)] + \frac{2}{\lambda m} \quad \text{per il Teorema 1} \\ &\leq \mathbb{E} \left[\ell_S(\mathbf{u}) + \frac{\lambda}{2} \|\mathbf{u}\|^2 \right] + \frac{4}{\lambda m} \quad \text{usando (2)} \\ &= \ell_{\mathcal{D}}(\mathbf{u}) + \frac{\lambda}{2} \|\mathbf{u}\|^2 + \frac{4}{\lambda m} \end{aligned}$$

come volevamo dimostrare. □

Si noti che scegliendo $\lambda = \sqrt{8/m}$ e usando il fatto che il rischio di classificazione $\operatorname{er}_{\mathcal{D}}(\mathbf{w}_S) = \mathbb{P}(Y \mathbf{w}_S^\top \mathbf{X} \leq 0)$ è maggiorato dal rischio della hinge loss $\mathbb{E}[\ell_{\mathcal{D}}(\mathbf{w}_S)]$, otteniamo

$$\mathbb{E}[\operatorname{er}_{\mathcal{D}}(\mathbf{w}_S)] \leq \min_{\mathbf{u} \in \mathbb{R}^d} \left(\ell_{\mathcal{D}}(\mathbf{u}) + \sqrt{\frac{2}{m}} \|\mathbf{u}\|^2 \right) + \sqrt{\frac{2}{m}} .$$