

## Da rischio sequenziale a rischio statistico

Ci occupiamo ora mettere in relazione il rischio sequenziale di un algoritmo online (come OGD) con il rischio di un predittore da esso prodotto, assumendo che i dati sui quali l'algoritmo viene eseguito siano generati da una sorgente statistica. Questo ci permette di creare un ponte fra il modello di apprendimento sequenziale e quello statistico.

Fissiamo una funzione di perdita  $\ell$ . Dato un predittore lineare  $\mathbf{w} \in \mathbb{R}^d$  e un esempio  $(\mathbf{x}_t, y_t) \in \mathbb{R}^d \times \mathbb{R}$ , sia  $\ell_t(\mathbf{w}) = \ell(\mathbf{w}^\top \mathbf{x}_t, y_t)$  la perdita associata. Per esempio,  $\ell(\mathbf{w}^\top \mathbf{x}_t, y_t) = (\mathbf{w}^\top \mathbf{x}_t - y_t)^2$  in regressione o  $\ell(\mathbf{w}^\top \mathbf{x}_t, y_t) = [1 - y_t \mathbf{w}^\top \mathbf{x}_t]_+$  in classificazione con hinge loss. Nel seguito, assumeremo che la funzione  $\ell_t$  sia convessa.

Consideriamo il caso dell'apprendimento statistico, in cui gli esempi  $(\mathbf{x}_t, y_t)$  sono realizzazioni indipendenti di variabili casuali  $(\mathbf{X}_t, Y_t)$  con distribuzione comune  $\mathcal{D}$  su  $\mathbb{R}^d \times \mathbb{R}$  fissata ma ignota. Il rischio statistico rispetto alla funzione di perdita  $\ell$  di un predittore lineare  $\mathbf{w}$  è definito da

$$\text{er}_{\mathcal{D}}(\mathbf{w}) = \mathbb{E} \left[ \ell(\mathbf{w}^\top \mathbf{X}, Y) \right]$$

dove l'esempio  $(\mathbf{X}, Y)$  è estratto dalla distribuzione congiunta  $\mathcal{D}$  su  $\mathbb{R}^d \times \mathbb{R}$ .

Consideriamo ora un training set  $S$  composto da realizzazioni  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$  di variabili casuali  $(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_m, Y_m)$  estratte da  $\mathcal{D}$ . Eseguiamo un algoritmo online come OGD sul training set  $S$  con la sequenza di funzioni di perdita  $\ell_1, \dots, \ell_m$  definite da  $\ell_t(\mathbf{w}) = \ell(\mathbf{w}^\top \mathbf{x}_t, y_t)$ . Per definizione di algoritmo online, otteniamo una sequenza  $\mathbf{w}_1, \dots, \mathbf{w}_m$  di predittori lineari (quelli prodotti da OGD durante il run). Vogliamo ora stabilire un maggiorante sul rischio statistico di un predittore lineare ottenuto in modo naturale dalla sequenza, ovvero il *predittore medio*

$$\bar{\mathbf{w}} = \frac{1}{m} \sum_{t=1}^m \mathbf{w}_t .$$

Anzitutto, dato che  $\ell$  è convessa, la disuguaglianza di Jensen ci dà che

$$\text{er}_{\mathcal{D}}(\bar{\mathbf{w}}) = \mathbb{E} \left[ \ell(\bar{\mathbf{w}}^\top \mathbf{X}, Y) \right] \leq \mathbb{E} \left[ \frac{1}{m} \sum_{t=1}^m \ell(\mathbf{w}_t^\top \mathbf{X}, Y) \right] = \frac{1}{m} \sum_{t=1}^m \text{er}_{\mathcal{D}}(\mathbf{w}_t)$$

dove l'ultima eguaglianza vale per linearità dell'aspettazione. Quindi il rischio del predittore medio è maggiorato dal rischio medio dei predittori della sequenza  $\mathbf{w}_1, \dots, \mathbf{w}_m$ .

Il passo cruciale sta ora nel legare il rischio statistico medio con il rischio sequenziale. Questo viene fatto osservando che, sotto l'ipotesi che  $S$  sia un campione casuale estratto da  $\mathcal{D}$ ,  $\mathbf{w}_t$  è determinato dai primi  $t-1$  esempi  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{t-1}, y_{t-1})$  estratti. Quindi, applicando la definizione di rischio al valore atteso della perdita di  $\mathbf{w}_t$  sul  $t$ -esimo esempio  $(\mathbf{x}_t, y_t)$ , possiamo scrivere

$$\mathbb{E} \left[ \text{er}_{\mathcal{D}}(\mathbf{w}_t) - \ell(\mathbf{w}_t^\top \mathbf{X}_t, Y_t) \mid (\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_{t-1}, Y_{t-1}) \right] = 0 . \quad (1)$$

La relazione qui sopra dice la seguente cosa: se condizioniamo sui primi  $t-1$  esempi, il valore atteso di  $\ell_t(\mathbf{w}_t)$  rispetto all'estrazione del  $t$ -esimo esempio è semplicemente (per definizione) il rischio di  $\mathbf{w}_t$ .

Per semplicità, denotiamo con  $\mathbb{E}_{t-1}$  il valore atteso condizionato come sopra. Se sommiamo entrambi i membri di (1) per  $t = 1, \dots, m$  e dividiamo per  $m$  otteniamo

$$\frac{1}{m} \sum_{t=1}^m \mathbb{E}_{t-1} \left[ \text{er}_{\mathcal{D}}(\mathbf{w}_t) - \ell(\mathbf{w}_t^\top \mathbf{X}_t, Y_t) \right] = 0 .$$

Per ogni  $t = 1, \dots, m$  sia  $Z_t$  la variabile casuale  $\text{er}_{\mathcal{D}}(\mathbf{w}_t) - \ell(\mathbf{w}_t^\top \mathbf{X}_t, Y_t)$ . Le variabili casuali  $Z_1, \dots, Z_m$  sono tutte funzioni del medesimo campione  $S$  e sono tali che

$$\frac{1}{m} \sum_{t=1}^m \mathbb{E}_{t-1}[Z_t] = 0 .$$

Assumiamo che  $\ell_t \in [0, M]$ , allora  $|Z_t| \leq M$ . Questo tipo di variabili casuali, o più precisamente di processo, viene definito sequenza di differenze di martingale con incrementi limitati da  $2M$ . Si noti che le  $Z_t$  non sono indipendenti. Da un punto di vista statistico, però, queste variabili si comportano come se lo fossero (almeno per certi aspetti). In particolare, vale una legge dei grandi numeri del tipo seguente

$$\frac{1}{m} \sum_{t=1}^m Z_t \leq 2M \sqrt{\frac{1}{2m} \ln \frac{1}{\delta}}$$

con probabilità almeno  $1 - \delta$  rispetto all'estrazione di  $S$ . Questo significa che

$$\frac{1}{m} \sum_{t=1}^m \text{er}_{\mathcal{D}}(\mathbf{w}_t) \leq \frac{1}{m} \sum_{t=1}^m \ell(\mathbf{w}_t^\top \mathbf{X}_t, Y_t) + M \sqrt{\frac{2}{m} \ln \frac{2}{\delta}} \quad (2)$$

con probabilità almeno  $1 - \delta/2$  rispetto all'estrazione di  $S$ . Ritornando al predittore medio  $\bar{\mathbf{w}}$ , il risultato che otteniamo può essere sintetizzato come

$$\text{er}_{\mathcal{D}}(\bar{\mathbf{w}}) \leq \frac{1}{m} \sum_{t=1}^m \ell(\mathbf{w}_t^\top \mathbf{x}_t, y_t) + \mathcal{O}\left(\frac{1}{\sqrt{m}}\right) \quad \text{con alta probabilità.}$$

In altre parole, il rischio del predittore medio è limitato in probabilità dal rischio sequenziale sul training set.

Possiamo spingerci ulteriormente per ottenere un maggiorante di rischio più esplicito che sfrutti direttamente i risultati dell'analisi online. Per esempio, per regressione con perdita quadratica, se facciamo girare OGD con proiezione nell'insieme  $\{\mathbf{u} \in \mathbb{R}^d : \|\mathbf{u}\| \leq U\}$ , e assumiamo che  $\|\mathbf{x}_t\| \leq X$  e  $|y_t| \leq UX$  per ogni  $t$ , otteniamo che, per ogni realizzazione  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$  di  $S$ ,

$$\frac{1}{m} \sum_{t=1}^m \ell(\mathbf{w}_t^\top \mathbf{x}_t, y_t) \leq \min_{\mathbf{u} \in \mathbb{R}^d : \|\mathbf{u}\| \leq U} \frac{1}{m} \sum_{t=1}^m \ell(\mathbf{u}^\top \mathbf{x}_t, y_t) + 8(UX)^2 \sqrt{\frac{2}{m}} .$$

Sostituendo il membro destro in (2) e notando che  $M = 4(UX)^2$  per la funzione di perdita quadratica, possiamo quindi scrivere che

$$\text{er}_{\mathcal{D}}(\bar{\mathbf{w}}) \leq \min_{\mathbf{u} \in \mathbb{R}^d : \|\mathbf{u}\| \leq U} \frac{1}{m} \sum_{t=1}^m \ell(\mathbf{u}^\top \mathbf{X}_t, Y_t) + 12(UX)^2 \sqrt{\frac{2}{m} \ln \frac{2}{\delta}}$$

con probabilità almeno  $1 - \delta/2$  rispetto all'estrazione di  $S$ .

Infine, possiamo notare che, detto

$$\mathbf{u}^* = \operatorname{argmin}_{\mathbf{u} \in \mathbb{R}^d: \|\mathbf{u}\| \leq U} \operatorname{er}_{\mathcal{D}}(\mathbf{u})$$

abbiamo, ovviamente,

$$\min_{\mathbf{u} \in \mathbb{R}^d: \|\mathbf{u}\| \leq U} \frac{1}{m} \sum_{t=1}^m \ell(\mathbf{u}^\top \mathbf{x}_t, y_t) \leq \frac{1}{m} \sum_{t=1}^m \ell(\mathbf{x}_t^\top \mathbf{u}^*, y_t).$$

Dato che, per ogni  $t = 1, \dots, m$  si ha  $\mathbb{E}[\ell(\mathbf{X}_t^\top \mathbf{u}^*, Y_t)] = \operatorname{er}_{\mathcal{D}}(\mathbf{u}^*)$ , possiamo applicare il maggiorante di Chernoff-Hoeffding e dedurre che

$$\frac{1}{m} \sum_{t=1}^m \ell(\mathbf{X}_t^\top \mathbf{u}^*, Y_t) \leq \operatorname{er}_{\mathcal{D}}(\mathbf{u}^*) + 4(UX)^2 \sqrt{\frac{1}{2m} \ln \frac{2}{\delta}} \quad \text{con probabilità almeno } 1 - \delta/2.$$

Abbiamo così ottenuto il seguente maggiorante esplicito sul rischio del predittore medio

$$\operatorname{er}_{\mathcal{D}}(\bar{\mathbf{w}}) \leq \operatorname{er}_{\mathcal{D}}(\mathbf{u}^*) + 14(UX)^2 \sqrt{\frac{2}{m} \ln \frac{2}{\delta}}$$

con probabilità almeno  $1 - \delta$  rispetto all'estrazione di  $S$ .

È possibile dimostrare che l'algoritmo che minimizza il rischio empirico,

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d: \|\mathbf{w}\| \leq U} \frac{1}{m} \sum_{t=1}^m \ell(\mathbf{w}^\top \mathbf{x}_t, y_t)$$

ha un errore di varianza dello stesso ordine di grandezza,

$$\operatorname{er}_{\mathcal{D}}(\hat{\mathbf{w}}) \leq \operatorname{er}_{\mathcal{D}}(\mathbf{u}^*) + \mathcal{O}\left(\frac{(UX)^2}{\sqrt{m}}\right)$$

con alta probabilità rispetto all'estrazione del training set.