

## Compression bounds

In questa lezione vediamo come limitare il rischio statistico del classificatore prodotto da un algoritmo di apprendimento in grado di rappresentare tale classificatore usando un piccolo sottoinsieme del training set.

Consideriamo una sequenza di esempi  $S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m))$  da utilizzare come training set per un problema di classificazione binaria e consideriamo un algoritmo di classificazione  $A$  che dato il training set  $S$  in input genera un classificatore  $A(S)$ . Chiamiamo sketch una qualsiasi sottosequenza  $S_0$  di  $S$  tale che  $A(S_0) = A(S)$ . Ovvero  $A$  con input  $S_0$  genera lo stesso classificatore di  $A$  con input  $S$  (si noti che questo è vero per qualsiasi algoritmo  $A$  dato che  $S_0$  può essere anche uguale ad  $S$ ). Vogliamo limitare il rischio  $\text{er}(\hat{h})$  in termini di  $\tilde{\text{er}}(\hat{h})$ , dove

$$\tilde{\text{er}}(\hat{h}) = \frac{1}{|S \setminus S_0|} \sum_{(\mathbf{x}_t, y_t) \in S \setminus S_0} \mathbb{I}\{\hat{h}(\mathbf{x}_t) \neq y_t\}$$

è la frazione di errori commessi da  $h$  sulla sottosequenza di training set  $S \setminus S_0$  che non include gli esempi dello sketch.

Indichiamo con  $\sigma(A, S)$  lo sketch  $S_0$  che l'algoritmo  $A$  produce con input  $S$ . Quindi, se  $\sigma(A, S) = S_0$  allora  $A(S_0) = A(S)$ . Inoltre, sia  $|S_0|$  la dimensione dello sketch  $S_0$ . Per ogni  $J \subseteq \{1, \dots, m\}$  sia  $S_J$  la sottosequenza del training set che contiene soltanto gli esempi indicizzati da  $J$ .

Procediamo ora a limitare il rischio di  $A(S)$ . Per comodità, nel resto dell'analisi consideriamo solo algoritmi  $A$  tali che  $|\sigma(A, S)| \leq \frac{m}{2}$ . Introducendo  $\varepsilon_k > 0$  da determinare in seguito notiamo che, per ogni training set  $S$  fissato,

$$\text{er}(A(S)) > \tilde{\text{er}}(A(S)) + \varepsilon_{|\sigma(A, S)|} \quad \text{implica} \quad \exists J, |J| \leq \frac{m}{2}, \text{er}(A(S_J)) > \tilde{\text{er}}(A(S_J)) + \varepsilon_{|\sigma(A, S)|}$$

dove  $S_J = \sigma(A, S)$ . Quindi possiamo scrivere

$$\begin{aligned} \mathbb{P}\left(\text{er}(A(S)) > \tilde{\text{er}}(A(S)) + \varepsilon_{|\sigma(A, S)|}\right) &\leq \mathbb{P}\left(\exists J, |J| \leq \frac{m}{2}, \text{er}(A(S_J)) > \tilde{\text{er}}(A(S_J)) + \varepsilon_{|\sigma(A, S)|}\right) \\ &\leq \sum_{k=0}^{m/2} \sum_{J: |J|=k} \mathbb{P}\left(\text{er}(A(S_J)) > \tilde{\text{er}}(A(S_J)) + \varepsilon_k\right) \end{aligned}$$

dove abbiamo usato la regola della somma  $\mathbb{P}(A \cup B) \leq \mathbb{P}(A) + \mathbb{P}(B)$ . Si noti che  $\tilde{\text{er}}(A(S_J))$  denota la frazione di errori di  $A(S_J)$  sugli esempi  $(\mathbf{x}_t, y_t)$  del training set tali che  $t \notin J$ . Ora, in ciascuna probabilità

$$\mathbb{P}\left(\text{er}(A(S_J)) > \tilde{\text{er}}(A(S_J)) + \varepsilon_k\right)$$

il classificatore  $A(S_J)$  per definizione è indipendente da tutti gli  $m - k$  esempi  $(\mathbf{x}_t, y_t)$  di training tali che  $t \notin J$ . Quindi  $\tilde{\text{er}}(A(S_J))$ , che è proprio determinato da questi  $m - k$  esempi, è una media

campionaria di un classificatore fissato ed ha valore atteso  $\text{er}(A(S_J))$ . Possiamo quindi applicare il maggiorante di Chernoff-Hoeffding ottenendo

$$\mathbb{P}\left(\text{er}(A(S)) > \tilde{\text{er}}(A(S)) + \varepsilon_{|\sigma(A,S)|}\right) \leq \sum_{k=0}^{m/2} \sum_{J:|J|=k} e^{-2(m-k)\varepsilon_k^2} .$$

Ora, scegliendo

$$\varepsilon_k = \sqrt{\frac{1}{m} \left( \ln \frac{1}{w_k} + \ln \frac{1}{\delta} \right)} \quad (1)$$

dove i pesi  $w_k \geq 0$ , che determineremo in seguito, soddisfano

$$\sum_{k=0}^{m/2} \sum_{J:|J|=k} w_k \leq 1 \quad (2)$$

otteniamo

$$\begin{aligned} \sum_{k=0}^{m/2} \sum_{J:|J|=k} e^{-2(m-k)\varepsilon_k^2} &\leq \sum_{k=0}^{m/2} \sum_{J:|J|=k} e^{-m\varepsilon_k^2} \quad \text{dato che } k \leq m/2 \\ &\leq \sum_{k=0}^{m/2} \sum_{J:|J|=k} \delta w_k \quad \text{usando (1)} \\ &\leq \delta \quad \text{usando (2)}. \end{aligned}$$

Quindi, con probabilità almeno  $1 - \delta$  rispetto all'estrazione del training set abbiamo

$$\text{er}(A(S)) \leq \tilde{\text{er}}(A(S)) + \sqrt{\frac{1}{m} \left( \ln \frac{1}{w_{|\sigma(A,S)|}} + \ln \frac{1}{\delta} \right)} .$$

Si noti ora che per soddisfare (2) è sufficiente definire

$$w_k = \frac{1}{m \binom{m}{k}} .$$

Utilizzando il maggiorante  $\binom{m}{k} \leq \left(\frac{em}{k}\right)^k$  vediamo che

$$\ln \frac{1}{w_k} = \ln m + \ln \binom{m}{k} \leq \ln m + k \ln \frac{em}{k} \leq k + (k+1) \ln m . \quad (3)$$

Quindi, con probabilità almeno  $1 - \delta$  rispetto all'estrazione del training set abbiamo infine

$$\text{er}(A(S)) \leq \tilde{\text{er}}(A(S)) + \sqrt{\frac{1}{m} \left( |\sigma(A,S)| + (|\sigma(A,S)| + 1) \ln m + \ln \frac{1}{\delta} \right)} .$$