

Rischio statistico e sua analisi

Per analizzare un algoritmo di apprendimento dobbiamo costruire un modello formale di generazione dei dati etichettati, ovvero un modello di generazione degli esempi (\mathbf{x}, y) . Nel modello statistico di apprendimento assumiamo che ogni esempio (\mathbf{x}, y) sia ottenuto tramite un'estrazione indipendente da una distribuzione di probabilità su $\mathcal{X} \times \mathcal{Y}$ fissata ma ignota. Per evidenziare il fatto che \mathbf{x} e y sono variabili casuali spesso scriveremo (\mathbf{X}, Y) . È abbastanza plausibile modellare con una distribuzione di probabilità il fatto che non tutti i dati \mathbf{x} sono osservati con la stessa frequenza (pensate al caso in cui \mathbf{x} è un'immagine o un testo). Analogamente, dato che come abbiamo detto le etichette sono per varie ragioni affette da rumore, è pure plausibile modellare con una distribuzione di probabilità la varietà di etichette a cui un dato può essere associato.

Dato che nel modello statistico ogni esempio (\mathbf{X}, Y) è ottenuto tramite un'estrazione indipendente dalla stessa distribuzione di probabilità congiunta, ogni dataset (p.es., training set o test set) sarà un **campione casuale** nel senso statistico del termine. In generale, l'assunzione di indipendenza è comoda dal punto di vista della trattabilità analitica del problema, ma poco plausibile in realtà. Si pensi ad esempio al problema di categorizzare le notizie di prodotte da un'agenzia giornalistica. È ovvio che, queste non saranno statisticamente indipendenti, ma seguiranno i diversi temi che di volta in volta si alterneranno.

Le prestazioni di un predittore $h : \mathcal{X} \rightarrow \mathcal{Y}$ rispetto ad un modello statistico dato ed a una funzione di perdita $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ data vengono valutate con il **rischio statistico**, definito da

$$\text{er}(h) = \mathbb{E}[\ell(Y, h(\mathbf{X}))]$$

ovvero, il valore atteso della funzione di perdita rispetto ad un esempio (\mathbf{X}, Y) generato dal modello statistico. Ipotizzando di conoscere il modello statistico, possiamo costruire il **predittore Bayesiano ottimo** $f^* : \mathcal{X} \rightarrow \mathcal{Y}$. Esso è definito come

$$f^*(\mathbf{x}) = \underset{\hat{y} \in \mathcal{Y}}{\text{argmin}} \mathbb{E}[\ell(Y, \hat{y}) \mid \mathbf{X} = \mathbf{x}]$$

ovvero la predizione \hat{y} che minimizza il rischio condizionato, cioè la perdita attesa rispetto alla distribuzione di Y condizionata su $\mathbf{X} = \mathbf{x}$. Si noti che, per definizione di f^* , vale

$$\mathbb{E}[\ell(Y, f^*(\mathbf{x})) \mid \mathbf{X} = \mathbf{x}] \leq \mathbb{E}[\ell(Y, h(\mathbf{x})) \mid \mathbf{X} = \mathbf{x}]$$

per ogni classificatore $h : \mathcal{X} \rightarrow \mathcal{Y}$. Dato che la disuguaglianza sopra vale per ogni $\mathbf{x} \in \mathcal{X}$, vale anche in media rispetto all'estrazione di \mathbf{X} . Ma siccome la media del rischio condizionato coincide col rischio, abbiamo che $\text{er}(f^*) \leq \text{er}(h)$. Il rischio $\text{er}(f^*)$ del predittore Bayesiano ottimo è detto **Bayes error**. In generale il Bayes error è maggiore di zero in quanto i predittori sono deterministici mentre le etichette sono probabilistiche.

Iniziamo ora a calcolare il predittore Bayesiano ottimo per la funzione di perdita quadratica $\ell(y, \hat{y}) = (y - \hat{y})^2$ nel caso $\mathcal{Y} \equiv \mathbb{R}$,

$$\begin{aligned}
f^*(\mathbf{x}) &= \operatorname{argmin}_{\hat{y} \in \mathbb{R}} \mathbb{E}[(Y - \hat{y})^2 \mid \mathbf{X} = \mathbf{x}] \\
&= \operatorname{argmin}_{\hat{y} \in \mathbb{R}} \left(\mathbb{E}[Y^2 \mid \mathbf{X} = \mathbf{x}] + \hat{y}^2 - 2\hat{y}\mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}] \right) \\
&= \operatorname{argmin}_{\hat{y} \in \mathbb{R}} \left(\hat{y}^2 - 2\hat{y}\mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}] \right) \quad (\text{scartando il termine che non dipende da } \hat{y}) \\
&= \mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}] \quad (\text{minimizzando la funzione } F(\hat{y}) = \hat{y}^2 - 2\hat{y}\mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}].)
\end{aligned}$$

Ovvero, il predittore Bayesiano ottimo per la funzione di perdita quadratica è il valore atteso dell'etichetta condizionato sull'istanza.

Sostituendo nella formula del rischio condizionato $\mathbb{E}[(Y - f^*(\mathbf{X}))^2 \mid \mathbf{X} = \mathbf{x}]$ il predittore ottimo $f^*(\mathbf{x}) = \mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}]$ otteniamo

$$\mathbb{E}[(Y - f^*(\mathbf{X}))^2 \mid \mathbf{X} = \mathbf{x}] = \mathbb{E}[(Y - \mathbb{E}[Y \mid \mathbf{x}])^2 \mid \mathbf{X} = \mathbf{x}] = \operatorname{Var}[Y \mid \mathbf{X} = \mathbf{x}].$$

Ovvero il rischio condizionato del predittore Bayesiano ottimo per la perdita quadrata è la varianza dell'etichetta condizionata sull'istanza.

Focalizziamoci ora sul caso di classificazione binaria, dove $\mathcal{Y} = \{-1, +1\}$. In questo caso è conveniente specificare una distribuzione congiunta sugli esempi (\mathbf{X}, Y) con la coppia (D, η) , dove D è la marginale su \mathcal{X} e η rappresenta la distribuzione su $\{-1, +1\}$ condizionata su \mathcal{X} . Dato che η è una distribuzione concentrata su due valori, possiamo rappresentarla con la funzione $\eta : \mathcal{X} \rightarrow [0, 1]$ dove $\eta(\mathbf{x}) = \mathbb{P}(Y = +1 \mid \mathbf{X} = \mathbf{x})$ è la probabilità che \mathbf{x} assuma l'etichetta +1.

Sia $\mathbb{I}\{A\} \in \{0, 1\}$ la funzione indicatrice di un evento A : $\mathbb{I}\{A\} = 1$ se e solo se A è vero. Il rischio statistico rispetto alla funzione di perdita zero-uno $\ell(y, \hat{y}) = \mathbb{I}\{\hat{y} \neq y\}$ risulta quindi essere

$$\operatorname{er}(h) = \mathbb{E}[\ell(Y, h(\mathbf{X}))] = \mathbb{E}[\mathbb{I}\{h(\mathbf{X}) \neq Y\}] = \mathbb{P}(h(\mathbf{X}) \neq Y).$$

In altre parole, il rischio di h è la probabilità che h sbagli a classificare un \mathbf{X} estratto dalla distribuzione D su \mathcal{X} e avente etichetta Y estratta a caso dalla distribuzione $\{1 - \eta(\mathbf{X}), \eta(\mathbf{X})\}$ su $\{-1, +1\}$.

Ipotizzando di conoscere il modello statistico (D, η) , possiamo costruire il classificatore Bayesiano ottimo $f^* : \mathcal{X} \rightarrow \{-1, +1\}$. Allora

$$\begin{aligned}
f^*(\mathbf{x}) &= \operatorname{argmin}_{\hat{y} \in \{-1, +1\}} \mathbb{E}[\ell(Y, \hat{y}) \mid \mathbf{X} = \mathbf{x}] \\
&= \operatorname{argmin}_{\hat{y} \in \{-1, +1\}} \mathbb{E}[\mathbb{I}\{Y = +1\}\mathbb{I}\{\hat{y} = -1\} + \mathbb{I}\{Y = -1\}\mathbb{I}\{\hat{y} = +1\} \mid \mathbf{X} = \mathbf{x}] \\
&= \operatorname{argmin}_{\hat{y} \in \{-1, +1\}} \left(\mathbb{P}(Y = +1 \mid \mathbf{X} = \mathbf{x})\mathbb{I}\{\hat{y} = -1\} + \mathbb{P}(Y = -1 \mid \mathbf{X} = \mathbf{x})\mathbb{I}\{\hat{y} = +1\} \right) \\
&= \operatorname{argmin}_{\hat{y} \in \{-1, +1\}} \left(\eta(\mathbf{x})\mathbb{I}\{\hat{y} = -1\} + (1 - \eta(\mathbf{x}))\mathbb{I}\{\hat{y} = +1\} \right) \\
&= \begin{cases} -1 & \text{se } \eta(\mathbf{x}) < 1/2, \\ +1 & \text{se } \eta(\mathbf{x}) \geq 1/2. \end{cases}
\end{aligned}$$

Quindi il classificatore Bayesiano ottimo predice l'etichetta che massimizza la probabilità condizionata sull'istanza. Non è difficile verificare che il Bayes error è pari a $er(f^*) = \mathbb{E}[\min\{\eta(\mathbf{X}), 1 - \eta(\mathbf{X})\}]$.

Stima del rischio per classificazione binaria. Passiamo ora a capire come valutare il rischio di un algoritmo di apprendimento. Per interpretare il modello statistico, possiamo pensare che ogni esempio (\mathbf{x}, y) disponibile sia ottenuto attraverso un'estrazione indipendente di $\mathbf{x} \in \mathcal{X}$ secondo la distribuzione D seguito dall'attribuzione dell'etichetta $y \in \{-1, +1\}$ secondo la distribuzione $\{1 - q(\mathbf{x}), q(\mathbf{x})\}$.

Dovrebbe essere chiaro che, dato un qualunque classificatore h , non possiamo calcolare direttamente il suo rischio $er(h)$ rispetto ad un modello statistico (D, η) . Infatti, il calcolo di $er(h)$ richiede di conoscere esattamente D e η . Ma se conoscessimo η potremmo direttamente trovare il classificatore Bayesiano ottimo per il problema.

Consideriamo ora il problema di stimare il rischio di un dato classificatore h . Per calcolare questa stima si usa un cosiddetto **test set**, ovvero un insieme $(\mathbf{x}'_1, y'_1), \dots, (\mathbf{x}'_n, y'_n)$ di dati etichettati. Quindi, si stima $er_{D, \eta}(h)$ con il **test error**, ovvero la frazione degli esempi del test set scorrettamente classificati da h ,

$$\tilde{er}(h) = \frac{1}{n} \sum_{t=1}^n \ell(y'_t, h(\mathbf{x}'_t)) .$$

Sotto l'ipotesi che il test set sia stato generato mediante estrazioni indipendenti dal modello statistico (D, η) , il test error altro non è che la media campionaria del rischio in quanto, per ogni $t = 1, \dots, n$ abbiamo che (\mathbf{X}'_t, Y'_t) è un'estrazione indipendente da (D, η) . Quindi

$$\mathbb{E}[\ell(Y'_t, h(\mathbf{X}'_t))] = \mathbb{P}(h(\mathbf{X}'_t) \neq Y'_t) = er(h) .$$

Per valutare la precisione con la quale possiamo vedere il test error come stima del rischio per una data taglia di test set utilizziamo il risultato seguente legato alla legge dei grandi numeri.

Lemma 1 (Chernoff-Hoeffding) *Siano Z_1, \dots, Z_n variabili casuali indipendenti, identicamente distribuite con media μ , e tali che $0 \leq Z_t \leq 1$ per ogni $t = 1, \dots, n$. Allora, per ogni $\varepsilon > 0$ fissato,*

$$\mathbb{P}\left(\frac{1}{n} \sum_{t=1}^n Z_t > \mu + \varepsilon\right) \leq e^{-2\varepsilon^2 n} \quad e \quad \mathbb{P}\left(\frac{1}{n} \sum_{t=1}^n Z_t < \mu - \varepsilon\right) \leq e^{-2\varepsilon^2 n} .$$

Utilizzando il maggiorante di Chernoff-Hoeffding con $Z_t = \ell(y_t, h(x_t))$ dove ℓ è la funzione di perdita per classificazione binaria, possiamo quindi valutare la precisione della nostra stima come segue

$$\mathbb{P}\left(|er(h) - \tilde{er}(h)| > \varepsilon\right) = \mathbb{P}\left(er(h) - \tilde{er}(h) > \varepsilon\right) + \mathbb{P}\left(\tilde{er}(h) - er(h) > \varepsilon\right) \leq 2e^{-2\varepsilon^2 n} \quad (1)$$

dove la probabilità è calcolata rispetto all'estrazione del test set. Questa disuguaglianza ci dice che la misura (secondo D e η) dei test set che producono stime $\tilde{er}(h)$ che differiscono dal valor medio $er(h)$ per più di ε decresce rapidamente con l'aumentare di n , cioè del numero di esempi nel test set.

In particolare, ponendo la parte destra di (1) pari a δ , per un qualunque $\delta \in (0, 1)$ fissato, otteniamo che

$$|\text{er}(h) - \widehat{\text{er}}(h)| \leq \sqrt{\frac{1}{2n} \ln \frac{2}{\delta}}$$

vale con probabilità almeno $1 - \delta$ rispetto all'estrazione del test set.

La disuguaglianza (1) ci indica come stimare il rischio di un'ipotesi prodotta da un qualche algoritmo di apprendimento mediante un test set. Viceversa, la stessa disuguaglianza ci mostra come il test set, che è il modo con cui in pratica misuriamo le prestazioni di un classificatore su dati ignoti, sia ben correlato col rischio nel modello di apprendimento statistico.

Studiamo ora un algoritmo di apprendimento nel modello statistico che abbiamo appena delineato e verifichiamo se compare l'overfitting. Come di consueto, assumiamo che l'algoritmo riceva in ingresso un training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$, dove $(\mathbf{x}_t, y_t) \in \mathcal{X} \times \{-1, +1\}$, e generi un classificatore $h : \mathcal{X} \rightarrow \{-1, +1\}$ appartenente a un dato spazio di modelli \mathcal{H} .

Sotto queste ipotesi, la cosa migliore che può fare l'algoritmo è scegliere il miglior classificatore possibile $h \in \mathcal{H}$, ovvero il classificatore h^* tale che

$$\text{er}(h^*) = \min_{h \in \mathcal{H}} \text{er}(h) .$$

Grazie alla legge dei grandi numeri, sappiamo che il training error $\widehat{\text{er}}(h^*)$ è vicino a $\text{er}(h^*)$ con alta probabilità rispetto all'estrazione del training set su cui $\widehat{\text{er}}(h^*)$ viene calcolato.

Consideriamo l'algoritmo che sceglie $\widehat{h} \in \mathcal{H}$ in modo da minimizzare il training error, ovvero

$$\widehat{h} = \underset{h \in \mathcal{H}}{\text{argmin}} \widehat{\text{er}}(h) .$$

Putroppo non possiamo applicare direttamente a \widehat{h} il maggiorante di Chernoff-Hoeffding al fine di dimostrare che $\text{er}(\widehat{h})$ è vicino a $\widehat{\text{er}}(\widehat{h})$, per poi concludere che $\text{er}(\widehat{h})$ è vicino a $\text{er}(h^*)$. Il motivo è che \widehat{h} è una funzione del training set e quindi una variabile casuale. Chernoff-Hoeffding dice che $\widehat{\text{er}}(h)$ è vicino a $\text{er}(h)$ per ogni h fissato, mentre \widehat{h} non è fissato, ma dipende dal campione rispetto al quale calcoliamo la probabilità.

Per analizzare il rischio di \widehat{h} procediamo allora come segue:

$$\begin{aligned} \text{er}(\widehat{h}) &= \text{er}(\widehat{h}) - \text{er}(h^*) && \text{errore di "varianza"} \\ &+ \text{er}(h^*) - \text{er}(f^*) && \text{errore di "bias"} \\ &+ \text{er}(f^*) && \text{Bayes error} \end{aligned}$$

dove f^* è il classificatore Bayesiano ottimo per il modello statistico (D, η) soggiacente. Il Bayes error non è controllabile, dato che dipende unicamente dal modello (D, η) . L'errore di bias dipende dal fatto che \mathcal{H} può non contenere il classificatore Bayesiano ottimo. L'errore di varianza dipende dal fatto che $\text{er}(\widehat{h})$ è generalmente diverso dall'errore di \widehat{h} sul training set. Di conseguenza, scegliere \widehat{h} comporta un errore dovuto al fatto che la frazione di esempi classificati scorrettamente nel training set da un classificatore h è soltanto una stima (possibilmente imprecisa) del rischio $\text{er}(h)$.

Procediamo quindi a controllare l'errore di varianza. Per ogni training set fissato, abbiamo che

$$\begin{aligned}
\text{er}(\hat{h}) - \text{er}(h^*) &= \text{er}(\hat{h}) - \hat{\text{er}}(\hat{h}) + \hat{\text{er}}(\hat{h}) - \text{er}(h^*) \\
&\leq \text{er}(\hat{h}) - \hat{\text{er}}(\hat{h}) + \hat{\text{er}}(h^*) - \text{er}(h^*) \\
&\leq |\text{er}(\hat{h}) - \hat{\text{er}}(\hat{h})| + |\hat{\text{er}}(h^*) - \text{er}(h^*)| \\
&\leq 2 \max_{h \in \mathcal{H}} |\hat{\text{er}}(h) - \text{er}(h)|
\end{aligned}$$

dove abbiamo usato l'ipotesi che \hat{h} minimizza $\hat{\text{er}}(h)$ in \mathcal{H} . Quindi, per ogni $\varepsilon > 0$,

$$\text{er}(\hat{h}) - \text{er}(h^*) > \varepsilon \quad \Rightarrow \quad \max_{h \in \mathcal{H}} |\hat{\text{er}}(h) - \text{er}(h)| > \frac{\varepsilon}{2} \quad \Rightarrow \quad \exists h \in \mathcal{H} : |\hat{\text{er}}(h) - \text{er}(h)| > \frac{\varepsilon}{2}.$$

Dato che la catena di implicazioni qui sopra vale per qualsiasi realizzazione del training set, possiamo scrivere

$$\mathbb{P} \left(\text{er}(\hat{h}) - \text{er}(h^*) > \varepsilon \right) \leq \mathbb{P} \left(\exists h \in \mathcal{H} : |\hat{\text{er}}(h) - \text{er}(h)| > \frac{\varepsilon}{2} \right).$$

Studiamo il caso $|\mathcal{H}| < \infty$, ovvero lo spazio dei modelli contiene un numero finito di classificatori. Dato che l'evento

$$\exists h \in \mathcal{H} : |\hat{\text{er}}(h) - \text{er}(h)| > \frac{\varepsilon}{2}$$

è l'unione su ogni $h \in \mathcal{H}$ degli eventi (non necessariamente disgiunti)

$$|\hat{\text{er}}(h) - \text{er}(h)| > \frac{\varepsilon}{2}$$

usando la regola della somma

$$\mathbb{P}(A_1 \cup \dots \cup A_n) \leq \sum_{i=1}^n \mathbb{P}(A_i)$$

che vale per qualsiasi collezione A_1, \dots, A_n di eventi, abbiamo che

$$\begin{aligned}
\mathbb{P} \left(\exists h \in \mathcal{H} : |\hat{\text{er}}(h) - \text{er}(h)| > \frac{\varepsilon}{2} \right) &= \mathbb{P} \left(\bigcup_{h \in \mathcal{H}} \left(|\hat{\text{er}}(h) - \text{er}(h)| > \frac{\varepsilon}{2} \right) \right) \\
&\leq \sum_{h \in \mathcal{H}} \mathbb{P} \left(|\hat{\text{er}}(h) - \text{er}(h)| > \frac{\varepsilon}{2} \right) \\
&\leq |\mathcal{H}| \max_{h \in \mathcal{H}} \mathbb{P} \left(|\hat{\text{er}}(h) - \text{er}(h)| > \frac{\varepsilon}{2} \right) \\
&\leq |\mathcal{H}| 2e^{-m\varepsilon^2/2}
\end{aligned} \tag{2}$$

dove nell'ultimo passo abbiamo usato il maggiorante di Chernoff-Hoeffding.

Quindi, in conclusione, otteniamo che

$$\mathbb{P} \left(\text{er}(\hat{h}) - \text{er}(h^*) > \varepsilon \right) \leq \mathbb{P} \left(\exists h \in \mathcal{H} : |\hat{\text{er}}(h) - \text{er}(h)| > \frac{\varepsilon}{2} \right) \leq 2|\mathcal{H}|e^{-m\varepsilon^2/2}. \tag{3}$$

Ponendo il membro destro di (3) uguale a δ e risolvendo rispetto a ε , otteniamo che

$$\text{er}(\hat{h}) \leq \text{er}(h^*) + \sqrt{\frac{2}{m} \ln \frac{2|\mathcal{H}|}{\delta}}$$

vale con probabilità almeno $1 - \delta$ rispetto all'estrazione casuale di un training set di cardinalità m .

Vediamo quindi che il rischio di \hat{h} si scompone in due contributi: $\text{er}(h^*)$ e $\sqrt{\frac{2}{m} \ln \frac{2|\mathcal{H}|}{\delta}}$. In mancanza di informazioni su (D, η) , e per una fissata cardinalità m del training set, per ridurre $\sqrt{\frac{2}{m} \ln \frac{2|\mathcal{H}|}{\delta}}$, ovvero il maggiorante sull'errore di varianza, dobbiamo ridurre \mathcal{H} . Ma questo fa potenzialmente aumentare $\text{er}(h^*)$ e quindi l'errore di bias. La necessità di bilanciare l'errore di varianza e l'errore di bias per controllare il rischio di \hat{h} è il materializzarsi nella teoria del fenomeno dell'overfitting.

Nella dimostrazione del maggiorante sull'errore di varianza abbiamo anche dimostrato in (2) che

$$\forall h \in \mathcal{H} \quad |\hat{\text{er}}(h) - \text{er}(h)| \leq \sqrt{\frac{1}{2m} \ln \frac{2|\mathcal{H}|}{\delta}}$$

con probabilità almeno $1 - \delta$ rispetto all'estrazione del training set. Questo significa che quando la cardinalità m del training set è abbastanza grande rispetto a $\ln |\mathcal{H}|$, allora il training error $\hat{\text{er}}(h)$ diventa una buona stima del rischio $\text{er}(h)$ *simultaneamente* per tutti i classificatori $h \in \mathcal{H}$. In queste condizioni, cioè quando la legge dei grandi numeri vale uniformemente rispetto alla scelta $h \in \mathcal{H}$, è chiaro che qualsiasi algoritmo che sceglie classificatori da \mathcal{H} è protetto dall'overfitting.