

Regressione logistica, funzioni surrogate e consistenza

Come abbiamo visto in precedenza, un approccio molto seguito per affrontare i problemi di classificazione binaria è quello di apprendere un modello predittivo $g : \mathcal{X} \rightarrow \mathbb{R}$ che rappresenta il classificatore $\text{sgn}(g(\mathbf{x})) \in \{-1, +1\}$ e poi utilizzare una funzione di perdita convessa $\ell(z)$ che maggiori la funzione indicatrice dell'errore di classificazione, $\ell(z) \geq \mathbb{I}\{z \leq 0\}$ dove $z = yg(\mathbf{x})$ e $y \in \{-1, +1\}$ è l'etichetta di \mathbf{x} . Funzioni convesse che maggiorano $\mathbb{I}\{z \leq 0\}$ vengono dette surrogate. Un esempio di funzione surrogata è la hinge loss $h(z) = [1 - z]_+$.

In classificazione lineare, $g(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$ e la hinge loss prende la forma $h(\mathbf{w}) = [1 - y\mathbf{w}^\top \mathbf{x}]_+$. In molte situazioni, però, si vuole utilizzare $g(\mathbf{x})$ per modellare $\eta(\mathbf{x}) = \mathbb{P}(Y = +1 \mid \mathbf{x})$. In altre parole, oltre a utilizzare $g(\mathbf{x})$ come classificatore $\text{sgn}(g(\mathbf{x}))$, lo si vuole utilizzare anche come stima della probabilità $\eta(\mathbf{x})$ che l'etichetta su un dato \mathbf{x} abbia valore +1. Questo evidentemente richiede che g abbia codominio $[0, 1]$ e di fatto rende innaturale l'utilizzo di modelli lineari $g(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$ che di norma possono assumere valori positivi e negativi. Una maniera elegante per aggirare questo ostacolo è porre

$$\ln \frac{\mathbb{P}(Y = +1 \mid \mathbf{x})}{\mathbb{P}(Y = -1 \mid \mathbf{x})} = \mathbf{w}^\top \mathbf{x}$$

dove la quantità a sinistra dell'uguale prende il nome di *log-odds ratio*. Risolvendo l'equazione per $\mathbb{P}(Y = +1 \mid \mathbf{x})$, ricordando che $\mathbb{P}(Y = -1 \mid \mathbf{x}) = 1 - \mathbb{P}(Y = +1 \mid \mathbf{x})$, otteniamo

$$\mathbb{P}(Y = +1 \mid \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}} .$$

La quantità a destra dell'uguale si chiama funzione logistica, ha codominio $[0, 1]$ e viene comunemente indicata con

$$\sigma(z) = \frac{1}{1 + e^{-z}} .$$

La regressione logistica apprende quindi un classificatore utilizzando un modello logistico $g(\mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x}) \in [0, 1]$ che rappresenta la probabilità $\mathbb{P}(Y = +1 \mid \mathbf{x})$. Dato che la predizione è probabilistica, una funzione di perdita appropriata è quella logaritmica,

$$\mathbb{I}\{y = +1\} \ln \frac{1}{\hat{y}} + \mathbb{I}\{y = -1\} \ln \frac{1}{1 - \hat{y}}$$

dove $\hat{y} = \sigma(\mathbf{w}^\top \mathbf{x})$. Osservando che $1 - \sigma(z) = \sigma(-z)$, possiamo infine scrivere la formula della funzione di perdita logistica $\ell(z) = \ln(1 + e^{-z})$. La funzione di perdita logistica è convessa. Per renderla una funzione surrogata, è sufficiente dividerla per $\ln 2$ in quanto è facile dimostrare che

$$\frac{\ln(1 + e^{-z})}{\ln 2} = \log_2(1 + e^{-z}) \geq \mathbb{I}\{z \leq 0\} .$$

Da un esame dei grafici delle due funzioni surrogate si nota come la funzione logistica $\ell_{\log}(z) =$

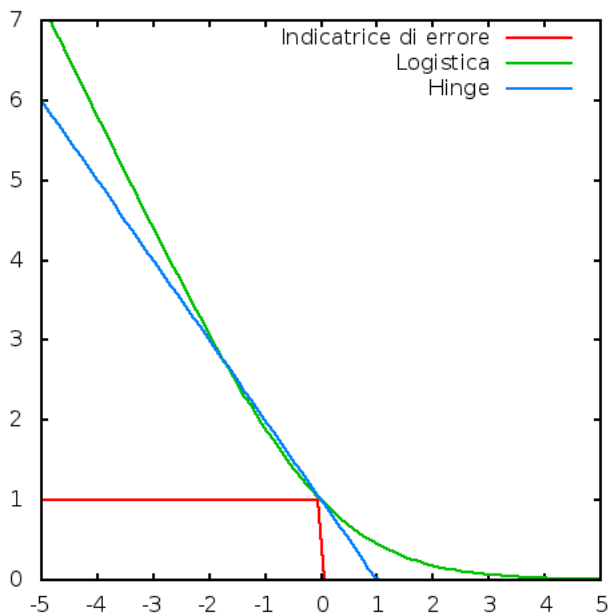


Figura 1: Confronto fra le due funzioni surrogate hinge $h(z) = [1 - z]_+$ e logistica $\ell_{\log}(z) = \log_2(1 + e^{-z})$.

$\log_2(1 + e^{-z})$ penalizza (seppure di poco) anche le classificazioni corrette, ovvero $\ell_{\log}(z) > 0$ per ogni $z \in \mathbb{R}$, mentre la funzione hinge h penalizza le classificazioni corrette soltanto quando hanno margine piccolo, ovvero $h(z) > 0$ per $z < 1$. Inoltre, la logistica penalizza molto più della hinge le classificazioni scorrette con margine grande, ovvero $\ell_{\log}(z) \gg h(z)$ quando $z \ll 0$.

Chiaramente, hinge e logistica non sono le uniche due possibili funzioni surrogate. Ma allora ci chiediamo quale sia un criterio che ci permetta di dire che alcune funzioni surrogate sono migliori di altre. Una prima proprietà è quella di *consistenza* che definiamo rispetto ad un modello statistico (D, η) —in realtà, come vedremo, la consistenza dipende solo da η .

Diciamo che una funzione di perdita surrogate ℓ è consistente quando, per ogni \mathbf{x} ,

$$\text{sgn}(y_{\mathbf{x}}^*) = f^*(\mathbf{x}) \quad \text{per} \quad y_{\mathbf{x}}^* = \underset{\hat{y} \in \mathbb{R}}{\text{argmin}} \mathbb{E}[\ell(Y\hat{y}) \mid \mathbf{X} = \mathbf{x}] . \quad (1)$$

In altre parole, il segno della predizione che minimizza il valore atteso della perdita su \mathbf{x} dev'essere uguale alla classificazione del Bayesiano ottimo su \mathbf{x} .

Dato che $\mathbb{P}(Y = +1 \mid \mathbf{x}) = \eta(\mathbf{x})$, possiamo scrivere $\mathbb{E}[\ell(Y\hat{y}) \mid \mathbf{X} = \mathbf{x}] = \eta(\mathbf{x})\ell(\hat{y}) + (1 - \eta(\mathbf{x}))\ell(-\hat{y})$.

Verifichiamo ora che la hinge loss è consistente,

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

Nella seconda uguaglianza, abbiamo potuto sostituire $\hat{y} \in \mathbb{R}$ con $\hat{y} \in [-1, +1]$ in quanto entrambe le funzioni $[1 - \hat{y}]_+$ e $[1 + \hat{y}]_+$ non aumentano se tronchiamo l'argomento \hat{y} nell'intervallo $[-1, +1]$.

Nel caso della funzione logistica è possibile dimostrare che

$$y_{\mathbf{x}}^* = \operatorname{argmin}_{\hat{y} \in \mathbb{R}} \left(\eta(\mathbf{x}) \log_2 (1 + e^{-\hat{y}}) + (1 - \eta(\mathbf{x})) \log_2 (1 + e^{\hat{y}}) \right) = \ln \frac{\eta(\mathbf{x})}{1 - \eta(\mathbf{x})}$$

e quindi

$$\operatorname{sgn}(y_{\mathbf{x}}^*) = \operatorname{sgn} \left(\ln \frac{\eta(\mathbf{x})}{1 - \eta(\mathbf{x})} \right) = \operatorname{sgn} \left(\eta(\mathbf{x}) - \frac{1}{2} \right) = f^*(\mathbf{x}) .$$

Ovvero, anche la funzione logistica è consistente.

In generale, è possibile dimostrare il seguente risultato.

Teorema 1 *Se una funzione di perdita surrogata $\ell : \mathbb{R} \rightarrow \mathbb{R}$ è convessa, differenziabile sullo zero e tale che $\ell'(0) < 0$ allora è anche consistente nel senso di (1).*

Questo implica immediatamente che anche le seguenti funzioni di perdita surrogate sono consistenti.

- Boosting loss: $\ell(z) = e^{-z}$.
- Square loss: $\ell(z) = (1 - z)^2$.
- Hinge loss quadratica: $\ell(z) = ([1 - z]_+)^2$.

Passiamo ora a esplicitare l'algoritmo di Online Gradient Descent con la funzione di perdita logistica. Ovvero, data una sequenza di esempi $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots \in \mathbb{R}^d \times \{-1, +1\}$ definiamo $\ell_t(\mathbf{w}) = \log_2 (1 + e^{-y_t \mathbf{w}^\top \mathbf{x}_t})$ e calcoliamo $\nabla \ell_t(\mathbf{w})$. Prima di tutto, osserviamo che

$$\frac{d}{dz} \ell_{\log}(yz) = \frac{d}{dz} \log_2 (1 + e^{-yz}) = \frac{1}{\ln 2} \times \frac{-ye^{-yz}}{1 + e^{-yz}} = \frac{1}{\ln 2} \times \frac{-y}{1 + e^{yz}} = \frac{-y\sigma(-yz)}{\ln 2} .$$

Quindi,

$$\nabla \ell_t(\mathbf{w}) = \frac{d}{dz} \ell_{\log}(y_t z) \Big|_{z=\mathbf{w}^\top \mathbf{x}_t} \times \mathbf{x}_t = \frac{-y_t \sigma(-\mathbf{w}^\top \mathbf{x}_t)}{\ln 2} y_t \mathbf{x}_t .$$

Il passo del gradiente si può scrivere allora come

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \eta_t \sigma(-y_t \mathbf{w}^\top \mathbf{x}_t) y_t \mathbf{x}_t$$

dove il fattore $\ln 2$ è stato incorporato nel coefficiente η_t . Dato che la funzione logistica non è fortemente convessa, dovremo usare $\eta_t = \eta/\sqrt{t}$ con $\eta > 0$ e far seguire ogni passo di discesa del gradiente da una proiezione sulla sfera di raggio $U > 0$.