

## 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  $z = y g(\mathbf{x})$  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 surrogate è la hinge loss  $h(z) = [1 - z]_+$ .

In molte situazioni, oltre a utilizzare  $g(\mathbf{x})$  per classificare, lo si vorrebbe utilizzare anche come stima della probabilità  $\eta(\mathbf{x}) = \mathbb{P}(Y = +1 \mid \mathbf{x})$ . Un modo elegante per utilizzare una qualunque  $g : \mathcal{X} \rightarrow \mathbb{R}$  per stimare  $\eta$  è attraverso la funzione logistica

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

dove  $0 < \sigma(z) < 1$  per ogni  $z \in \mathbb{R}$ . La stima di  $\eta(\mathbf{x})$  tramite  $1/(1 + e^{-g(\mathbf{x})})$  prende quindi il nome di regressione logistica.

Chiaramente, se vogliamo usare  $g$  per fare regressione logistica dobbiamo apprendere usando una funzione di perdita appropriata. Dato che vogliamo usare  $g$  per stimare una probabilità, usiamo la funzione di perdita logaritmica,

$$\ell(\hat{y}, y) = \mathbb{I}\{y = +1\} \ln \frac{1}{\hat{y}} + \mathbb{I}\{y = -1\} \ln \frac{1}{1 - \hat{y}}$$

dove  $\hat{y} = \sigma(g(\mathbf{x}))$ . Osservando che  $1 - \sigma(z) = \sigma(-z)$ , possiamo scrivere la perdita logaritmica in funzione di  $z = y g(\mathbf{x})$  ottenendo

$$\ell(\hat{y}, y) = \ln(1 + e^{-z})$$

che vale quando  $\hat{y} = \sigma(g(\mathbf{x}))$ . La parte destra dell'identità prende il nome di funzione di perdita logistica. Dato che è convessa, per renderla una funzione surrogate è sufficiente dividerla per  $\ln 2$ . Infatti, 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 in Figura 1 si nota come la funzione logistica  $\ell_{\log}(z) = \log_2(1 + e^{-z})$  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

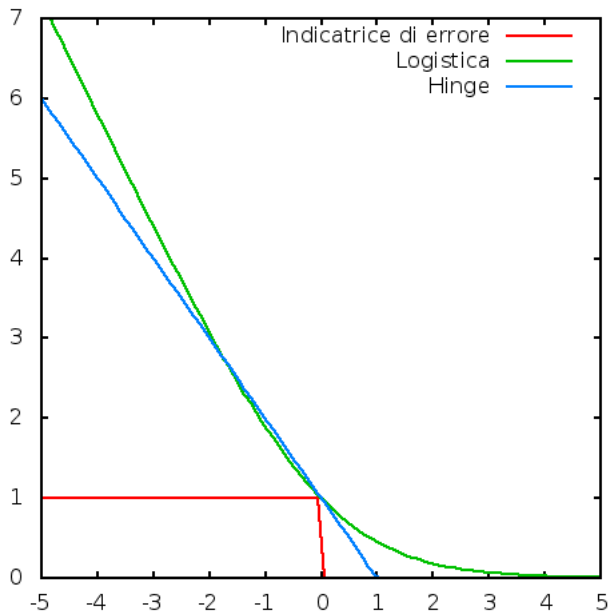


Figura 1: Confronto fra le due funzioni surrogate hinge  $h(z) = [1 - z]_+$  e logistica  $\ell_{\log}(z) = \log_2(1 + e^{-z})$ .

di altre. Una proprietà importante è quella di *consistenza*, che definiamo rispetto ad un modello statistico  $(D, \eta)$  —in realtà, come vedremo, la consistenza dipende solo da  $\eta$ .

Diciamo che una funzione di perdita surrogata  $\ell$  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 la consistenza della funzione di perdita logistica. Non è difficile verificare che

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

e quindi

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

La quantità  $y_{\mathbf{x}}^* = \ln \frac{\eta(\mathbf{x})}{1 - \eta(\mathbf{x})}$  si chiama *log-odds ratio*. Oltre a fornire il classificatore Bayesiano ottimo  $\text{sgn}(y_{\mathbf{x}}^*)$ ,  $y_{\mathbf{x}}^*$  è anche la migliore stima di  $\eta$  rispetto alla funzione di perdita logistica. Se calcoliamo il rischio condizionato di  $y_{\mathbf{x}}^*$  rispetto alla funzione logistica otteniamo

$$\mathbb{E} \left[ \log_2(1 + e^{-Y y_{\mathbf{x}}^*}) \mid \mathbf{X} = \mathbf{x} \right] = -\eta(\mathbf{x}) \log_2 \eta(\mathbf{x}) + (1 - \eta(\mathbf{x})) \log_2(1 - \eta(\mathbf{x})) .$$

La quantità a destra è l'entropia  $H(Y | \mathbf{X} = \mathbf{x})$  di  $Y$  dato  $\mathbf{x}$  e rappresenta il numero medio di bit che ricaviamo osservando  $Y$  quando  $\mathbf{x}$  è noto. Ciò implica che il rischio Bayesiano (cioè il rischio del Bayesiano ottimo  $g^*(\mathbf{x}) = y_{\mathbf{x}}^*$ ) per la perdita logistica è

$$\mathbb{E}\left[\log_2\left(1 + e^{-Y g^*(\mathbf{X})}\right)\right] = H(Y)$$

dove  $H(Y)$  è l'entropia della distribuzione Bernoulliana  $\{\mathbb{P}(Y = +1), \mathbb{P}(Y = -1)\}$ .

Verifichiamo ora che anche 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  $[-1, +1]$ .

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 descrivere l'algoritmo di Online Gradient Descent per regressione logistica dove  $g(\mathbf{x})$  è un modello lineare  $\mathbf{w}^\top \mathbf{x}$ . 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,

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

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

Se vogliamo fare regressione logistica per minimizzare il rischio empirico in un dato training set, possiamo aggiungere un termine di regolarizzazione per rendere la funzione risultante fortemente convessa e quindi convergere più velocemente tramite discesa del gradiente stocastico (e senza dover più utilizzare la proiezione sulla sfera),

$$\ell_{\text{reg}}(\mathbf{w}) = \log_2(1 + e^{-y\mathbf{w}^\top \mathbf{x}}) + \frac{\lambda}{2} \|\mathbf{w}\|^2 .$$

L'algoritmo risultante è simile a Pegasos dove la hinge loss è sostituita dalla perdita logistica.