

2. REGRESSIONE LOGISTICA

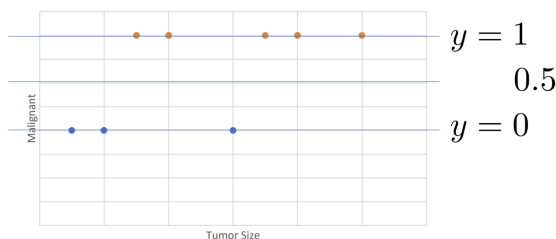
martedì 22 ottobre 2024 18:45

La regressione logistica è un modello statistico utilizzato per predire la probabilità che un'osservazione appartenga a una di due categorie possibili (binaria) sulla base di una o più variabili indipendenti. È particolarmente utile quando la variabile dipendente è dicotomica, ad esempio, "successo" o "fallimento", "sì" o "no", "malato" o "sano", ecc.

La regressione logistica è molto usata in vari campi, come:

- Medicina: Predire la probabilità di una malattia in base ai sintomi e ai fattori di rischio.
- Marketing: Stimare la probabilità che un cliente effettui un acquisto in base al suo comportamento passato.
- Finanza: Valutare il rischio di insolvenza di un cliente in base a dati finanziari.

Classifier threshold



ESEMPLARE quando si parla di Regressione Logistica

L'OUTPUT deve essere per forza $0 \leq h(x) \leq 1$

IPOTESI RAPPRESENTAZIONE

FUNZIONE SIGMOIDE

La funzione ipotetica per la regressione logistica è indicata con

$h_{\theta}(x) = \theta^T x$ che rappresenta la previsione del modello basata sui parametri θ su D e T x

Ricordiamo che $\theta^T x$ è una combinazione lineare.

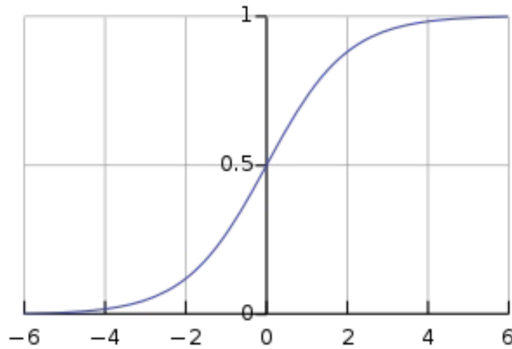
FUNZIONE SIGMOIDE

ORA che serve noi sappiamo $h_{\theta}(x) \in (-\infty, +\infty)$

$$Def: \frac{1}{1 + e^{-z}}$$

$$z = \theta^T x$$

$$f = \frac{1}{1 + e^{-\theta^T x}} \quad \xrightarrow{\text{def}} \quad h_{\theta}(x) = g(\theta^T x)$$



ORA FACCIAMO UN ESEMPIO:

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2) = g(-3 + x_1 + x_2)$$

$$\theta_0 = -3, \quad \theta_1 = 1, \quad \theta_2 = +1$$

$$h_{\theta}(x) \geq 0.5$$

Lo posso rappresentare con la sigmoide

$$h_{\theta}(x) = \frac{1}{1 + e^{-g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}} \geq 0.5$$

Affinche io possa rappresentare o assegnare 1 o 0

$$\text{Se } -3 + x_1 + x_2 \geq 0 \quad \text{ALLORA ASSEGNA } 1 \text{ e } y$$

$$f = y = 1$$

INTERPRETAZIONE PROBABILISTICA SIGMOIDE

Rappresentiamo la stima probabilistica :

$$h_{\theta}(x) = p(y = 1 | x; \theta)$$

Facendo l'esempio dei tumori come nelle slide

$$h_{\theta}(x) = 0.7 \quad \text{tumore benigno}$$

$$\text{Altra } h_{\theta}(x) = 0.3 \quad \text{tumore maligno}$$

$$\text{Generalmente } P(y=0 | x; \theta) = 1 - P(y=1 | x; \theta)$$

La regressione logistica non solo prevede una classe, ma restituisce anche una probabilità associata a quella previsione, rendendo il modello particolarmente utile per problemi in cui l'incertezza è importante, come nella diagnosi medica.

COST FUNCTION | FUNZIONE DI COSTO

L'obiettivo della regressione logistica è trovare i parametri θ che meglio spiegano i dati di addestramento. Questo viene fatto minimizzando una funzione di costo, che misura quanto le previsioni del modello differiscono dai valori effettivi.

Si desidera che la soluzione del problema di minimizzazione della funzione di costo corrisponda a risolvere il problema di massima verosimiglianza, cioè trovare i parametri θ che massimizzano la probabilità (verosimiglianza) dei dati osservati. Poiché l'output della regressione logistica può essere solo 0 o 1 (stiamo affrontando una distribuzione di Bernoulli), possiamo formulare la funzione di verosimiglianza per il modello. Sotto l'assunzione che i campioni siano indipendenti e seguano la stessa distribuzione, possiamo scrivere la funzione di verosimiglianza come segue:

$$L(\theta) = L(\theta; X, y) = P(y | X; \theta) = \prod_{i=1}^m P(y^{(i)} | x^{(i)}; \theta)$$

Altra:

$$\begin{cases} P(y^{(i)} = 1 | x^{(i)}; \theta) = h_{\theta}(x^{(i)}) \\ P(y^{(i)} = 0 | x^{(i)}; \theta) = 1 - h_{\theta}(x^{(i)}) \end{cases}$$

Nella forma compatta:

Sfrutta le proprietà delle potenze $e^0 = 1$

$$L(\theta) = \prod_{i=1}^m \left(h_{\theta}(x^{(i)})^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{1 - y^{(i)}} \right)$$

$$p(y^{(i)} | x^{(i)}; \theta) = h_{\theta}(x^{(i)})^{y^{(i)}} \cdot (1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}}$$

Quindi se $y^{(i)} = 1$

$$p(y^{(i)} | x^{(i)}; \theta) = h_{\theta}(x^{(i)})^{y^{(i)}} = h_{\theta}(x^{(i)})$$

E l'altro termine diventa 0

Se $y^{(i)} = 0$ Il termine

$$(1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}} = 1 - h_{\theta}(x^{(i)})$$

$$L(\theta) = \prod_{i=1}^m h_{\theta}(x^{(i)})^{y^{(i)}} \cdot (1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}}$$

COSTO FUNZIONE FORMA LOGARITMICA

RICORDIAMO $\ell(\theta) = \log L(\theta)$ *SOST.*

$$\ell(\theta) = \log \prod_{i=1}^m h_{\theta}(x^{(i)})^{y^{(i)}} \cdot (1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}}$$

$$= \sum_{i=1}^m \left(y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right)$$

Questa è chiamata log-verosimiglianza. La motivazione principale per utilizzare il logaritmo è la semplificazione: i prodotti complicati diventano somme, rendendo più semplice la derivazione della funzione di costo.

Quando si lavora con l'apprendimento automatico, in genere si preferisce minimizzare una funzione di costo piuttosto che massimizzare una verosimiglianza.

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right)$$

FUNZIONE DI ERRORE DI ENTROPIA INCROCIATA

$$\hat{\theta} = \arg \min_{\theta} J(\theta)$$

Dobbiamo trovare valori theta che minimizzino la nostra funzione Costo.

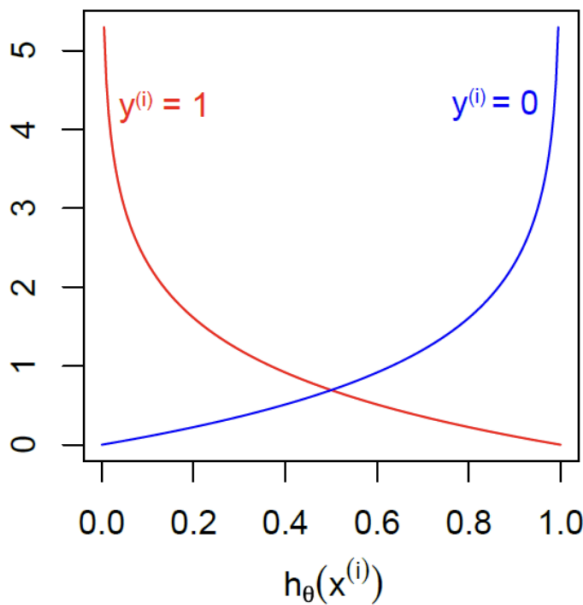
ERRORI | REGRESSIONE LOGISTICA

$$l^{(i)} = -y^{(i)} \log h_{\theta}(x^{(i)}) - (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)}))$$

Se $y^{(i)} = 1$

$$\hookrightarrow l^{(i)} = -\log h_{\theta}(x^{(i)})$$

SE H è vicino allo zero
allora ERRORE ALTO
SE H è alto quindi vicino
a 1 ERRORE BASSO



Se $y^{(i)} = 0$ $l^{(i)} = -\log (1 - h_{\theta}(x^{(i)}))$

H è basso cioè 0 allora l'errore è basso

H è alto cioè vicino a 1 l'errore tende a +inf

