## PROBABILISTI INTERPRETATION OF LINEAR REGRESSION

sabato 19 ottobre 2024 18:22

L'interpretazione probabilistica di un modello di Machine Learning si riferisce al modo in cui possiamo comprendere e descrivere il comportamento di un modello statistico, come quello lineare, utilizzando i concetti della probabilità.

Quindi abbiamo detto che un modello di regressione lineare cerca di fare previsioni o meglio pronosticare i valori di output.

output.

QUIDI 43BMO X INPUT FEOTORES

MBIATIO 5 (i) OUTOP

Utilizzando la funzione ipotesi: h/l/=  $\Theta^{T}$  (i)

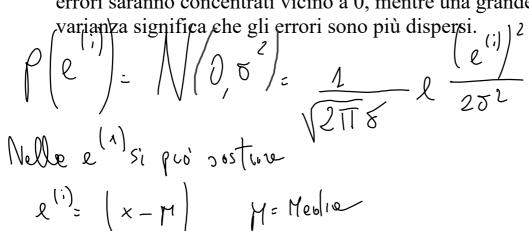
Ipotesi --> In quanto una fetta di problemi hanno un equazione di questo tipo ( Ma puo' anche variare)

(i) = GRORE ASSOCIATO KUA PREVISIONE

## PRIMA DI PARLARE DI INTERPRETAZIONE PROBABILISTICA DEGLI ERRORI DI UN MODELLO **BISOGNA ACCERTARSI CHE:**

- ERRORI INDIPENDENTI
- $\mathcal{L}^{\prime\prime}$  viene modellato come una variabile casuale stiamo considerando eventi che seguono una certa distribuzione probabilistica.
- Si assume che tutti gli errori seguano una distribuzione normale.
  - Media 0 : Questo implica che ci si aspetta che l'errore medio sia zero, il che ha senso in quanto, in media, il modello dovrebbe produrre previsioni corrette.
    - Varianza 🍎: La varianza rappresenta la dispersione degli errori attorno alla media. Se è piccola, gli

errori saranno concentrati vicino a 0, mentre una grande



Distribuzione di probabilità condizionale

Vediono ou coso é regressione lineare con rumore gaussiano,

$$p(y^{(i)}|\mathbf{x^{(i)}};\theta) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y^{(i)} - \theta^T\mathbf{x^{(i)}})^2}{2\sigma^2}} \qquad i = 1\dots m.$$

La verosimiglianza  $L(\theta)$  rappresenta la probabilità che, dato il modello  $\theta$ , i dati osservati (X,y) siano stati generati.



$$L(\theta) = P(y|X;\theta)$$

SOSTITUENDO L'ESPRESSIONE CON AL FORMA **GAUSSIANA** 

OTTENIATO:
$$L(0) = \prod_{\sigma \in \mathbb{Z}_{\overline{\Pi}}} \frac{1}{e^{-\left(\frac{g(i)}{2\sigma^2} + \frac{G(i)}{2\sigma^2}\right)^2}}$$

L'obiettivo finale sarà appunto quello di massimizzare la funzione di verosomiglianza per ottenere parametri ottimali. In altre parole, vogliamo trovare i parametri che rendono più probabile osservare i dati forniti.

Poiché la funzione di verosimiglianza può essere complessa da massimizzare direttamente, spesso si lavora con il logaritmo della verosimiglianza, dato che il logaritmo è una funzione monotona. Questo significa che massimizzare il logaritmo della verosimiglianza è equivalente a massimizzare la verosimiglianza stessa.

L'idea è trasformare il prodotto di probabilità condizionate, che rappresentano la distribuzione di probabilità dei dati, in una somma, utilizzando il logaritmo naturale.

ABBIAT70 
$$l(\theta) = log \prod_{i} p(g^{(i)} | x^{(i)}; \theta)$$

Il logaritmo trasforma la produttoria il una sommatoria. Si sfrutta la proprietà del logaritmo che permette di trasformare il logaritmo di un prodotto in una somma dei

Ora il log viene moltiplicato per A E B

Ora che succede possiamo semplificare tale equazione ? COME ?

31/10/24, 14:31 OneNote

Praticamente le costanti le porto fuori dalla sommatoria fattorizzando il termine 1 con log

$$\begin{array}{c|c}
M & \frac{1}{\sqrt{2\pi}} - \sum_{i=1}^{M} \left( \frac{y^{(i)} - \theta^{T} x^{(i)}}{2\sigma^{2}} \right)^{2} \\
= M & \sqrt{\sqrt{2\pi}} - \frac{1}{2\sigma^{2}} \cdot \sum_{i=1}^{M} \left( y^{(i)} - \theta^{T} x^{(i)} \right)^{2}
\end{array}$$

TROVATA TALÉ EQ. ORA DOBBIATO MASSIPILENTE

Dobbiamo massimizzare la funzione di verosimiglianza perché l'obiettivo di molti algoritmi di apprendimento automatico, come la regressione lineare nel contesto di questa dimostrazione, è trovare i parametri ottimali  $\theta$   $\theta$  che rendano il nostro modello il più "verosimile" possibile rispetto ai dati osservati.

max 
$$l(g) = mox \left( m \log \frac{1}{l^{2}\pi \delta} - \frac{1}{2\delta^{2}} \sum_{i=1}^{m} (g^{ii} - g^{7}x^{ii})^{2} \right)$$

PONIMO Mlog 
$$\frac{1}{8\sqrt{2\pi}} = COSTANTE$$

$$\left| -\frac{1}{2\sqrt{2}} \cdot \sum_{i} \left( 5^{(i)} - \theta^T \chi^{(i)} \right)^2 \right|$$

STEDY (64) 
$$\sqrt{2}$$
 VARANZA  

$$= MOX \left( -\frac{1}{2} \ge \left( y^{(i)} - \theta^T x^{(i)} \right)^2 \right)$$

Per tale mot.

$$\int_{z}^{z} \min \left( 0.5 \cdot \sum_{i=1}^{m} \left( g^{(i)} - \Theta^{T} \chi^{(i)} \right)^{2} \right)$$

ORA:  

$$MSE$$
 Min  $\left(\frac{1}{2m}\sum_{i=1}^{m}\left(\frac{h}{x^{(i)}}\right)-\frac{y^{(i)}}{y^2}\right)^2$  MSE  
 $h\left(x^{(i)}\right) = PRONOSTICD$   
 $S\left(x^{(i)}\right) = VALORE REAGE OSJERATO$   
 $M = 6^{l}$  IL NUN TOINZE DI GMPIONI

## MINIMI QUADRATI

FORMUMZIONE INSUSSISTED VEROSONIUMA MAY

$$\min_{\theta} \left( \frac{1}{2} \sum_{i=1}^{m} \left( g^{(i)} - \theta^{T} \chi^{(i)} \right)^{2} \right)$$

Problema identico: Come detto nella slide, i due problemi sono identici a livello di ottimizzazione rispetto a. . La costante m ° Lom non influisce sul risultato finale, poiché non cambia la posizione del minimo della funzione da ottimizzare. L'unica differenza è che una delle formulazioni normalizza l'errore rispetto al numero di campioni.

Costante di regolarizzazione: La frase "These problems are identical but a regularization constant that has no influence in the minimization problem" si riferisce al fatto che questo termine costante (ad esempio il fattore m) non influisce sulla posizione del minimo della funzione. Influisce solo sul valore numerico della funzione obiettivo, ma non cambia il valore di  $\theta$  che minimizza l'errore.

31/10/24, 14:31 OneNote

IMPORTANTE, minimizzare il problema dei minimi quadrati nella regressione lineare (come visto a destra) è equivalente a minimizzare l'errore quadratico medio (MSE) (come visto a sinistra). La differenza nei fattori costanti non influisce sull'ottimizzazione rispetto a θ, perciò possiamo usare indifferentemente una o l'altra formulazione per trovare i migliori parametri del modello che riducono l'errore di previsione.