

# Regressione Lineare

---

L'obiettivo della regressione è quello di predire il valore di uno o più target continui  $t$ , dato il valore di un vettore  $x$   $D$ -dimensionale di variabili di input. La regressione polinomiale è un esempio specifico di **regressione lineare**: una classe di modelli di regressione che condividono la caratteristica di essere lineari nel rispetto dei parametri. Un sottogruppo di modelli lineari, ancora più semplici, prevede di essere lineare anche rispetto alle variabili di input, ma vedremo che può essere molto limitante. Conviene invece prendere combinazioni lineari di un set fissato di funzioni non lineari. Queste funzioni sono dette funzioni base (**basis functions**).

## Modelli di funzioni base lineari

Il modello lineare più semplice possibile per la regressione è quello formato da una combinazione lineare di variabili di input:

$$y(x, w) = \{ (w_0 + w_1 x_1 + \dots + w_D x_D) \}$$

Questo modello lineare è tale perché è lineare rispetto ai parametri  $w$ . Inoltre risulta lineare anche negli input  $x = (x_1, \dots, x_D)^T$ , ma questo impone delle limitazioni notevoli. È per questo motivo che estendiamo la classe di modelli di regressione lineare, includendo anche tutti quei modelli che sono combinazione lineare dei parametri  $w = (w_1, \dots, w_D)$  e combinazione lineare di un set fissato di funzioni non lineari di  $x$ . I modelli prendono la seguente forma:

$$y(X, W) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(X)$$

Dove  $\phi_j(X)$  sono le funzioni base. Questo modello è lineare rispetto ai parametri  $w$ , ma non è lineare rispetto ai parametri  $\phi$ . Questo modello è chiamato **modello di funzioni base lineari**. Il massimo valore di  $j$  è  $M-1$ , quindi il numero totale di parametri del modello sarà  $M$ .

Il parametro  $w_0$  è chiamato **bias** e permette qualsiasi offset fisso nei dati.

Analizziamo adesso il caso di modello di regressione lineare polinomiale: questo è un caso particolare di regressione lineare, in cui è presente una sola variabile di input  $x$ , e in cui le funzioni di base prendono la forma delle potenze di  $x$   $\phi_j(x) = x^j$ . Il più grande problema delle funzioni base polinomiali è che sono funzioni globali della variabile di input, ovvero che un cambiamento in una regione di  $x$  influisce in tutte le altre regioni. Questo può essere risolto dividendo lo spazio degli input, in diverse regioni e per ognuna usare un polinomio diverso, ottenendo così una spline, per esempio:

$$\phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\}$$

Questo tipo di funzioni base sono chiamate **Gaussiane**.

## Maximum likelihood e Least Squares

Assumiamo che la variabile di target  $t$  sia data da una funzione deterministica  $y(X, W)$  con l'aggiunta di un errore Gaussiano  $\epsilon$ :

$$t = y(X, W) + \epsilon$$

Dove  $\epsilon$  è una variabile casuale con distribuzione Gaussiana con media zero e varianza  $\beta^{-1}$ . Possiamo quindi scrivere:

$$p(t|X, W, \beta) = N(t|y(X, W), \beta^{-1})$$

Se consideriamo una squared loss function, allora la predizione ottimale per un nuovo valore di  $x$ , sarà data dalla media pesata della variabile di target:

$$E[t|X] = \int t p(t|X) dt = y(X, W)$$

Per via del rumore Gaussiano, la distribuzione condizionata di  $t$  data  $x$  è unimodale, ovvero presenta un solo picco di massimo, che può risultare limitante per alcune applicazioni.

Adesso consideriamo un set di input  $X = \{x_1, \dots, x_N\}$  e un set di target  $T = \{t_1, \dots, t_N\}$ . Raggruppiamo le variabili target in un vettore a colonna e lo chiamiamo  $t$ . Assumendo che questi dati sono indipendenti dalla distribuzione, otteniamo la seguente formula per la funzione di likelihood:

$$p(t|X, W, \beta) = \prod_{n=1}^N N(t_n|y(x_n, W), \beta^{-1})$$

Questa è una funzione dei parametri  $W$  e  $\beta$ . Applicando il logaritmo, ricaviamo:

$$\ln p(t|X, W, \beta) = \sum_{n=1}^N \ln N(t_n|y(x_n, W), \beta^{-1})$$

$$= -\frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \frac{1}{2\beta} E_D(W)$$

Dove  $E_D(W)$  è la somma degli errori al quadrato:

$$E_D(W) = \frac{1}{2} \sum_{n=1}^N \{t_n - W^T \phi(X_n)\}^2$$

Dato che abbiamo scritto la funzione di likelihood, possiamo massimizzarla per determinare  $w$  e  $\beta$ . Massimizziamola nel rispetto di  $w$ . Grazie alla forma logaritmica della funzione di likelihood è evidente che massimizzare nel rispetto di  $w$ , significa massimizzare la somma dei quadrati degli errori  $E_D(W)$ :

$$\nabla \ln p(t|W, \beta) = \sum_{n=1}^N \{t_n - W^T \phi(X_n)\} \phi(X_n)$$

Imponiamo il gradiente a 0 risolvendo nel rispetto di  $w$ , otteniamo la seguente equazione:

$$W_{ML} = (\Phi^T \Phi)^{-1} \Phi^T t$$

Detta **equazione normale**. Dove  $\Phi$  è la matrice di dimensioni  $N \times M$  e i cui elementi sono dati da  $\Phi_{nj} = \phi_j(x_n)$ .

$$\Phi = \begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \dots & \phi_{M-1}(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x_N) & \phi_1(x_N) & \dots & \phi_{M-1}(x_N) \end{bmatrix}$$

**Bias**

Riprendiamo la formula ricavata per  $E_D(W)$  e esplicitiamo il bias ( $w_0$ ):

$$E_D(W) = \frac{1}{2} \sum_{n=1}^N \{t_n - w_0 - \sum_{j=1}^{M-1} w_j \phi_j(X_n)\}^2$$

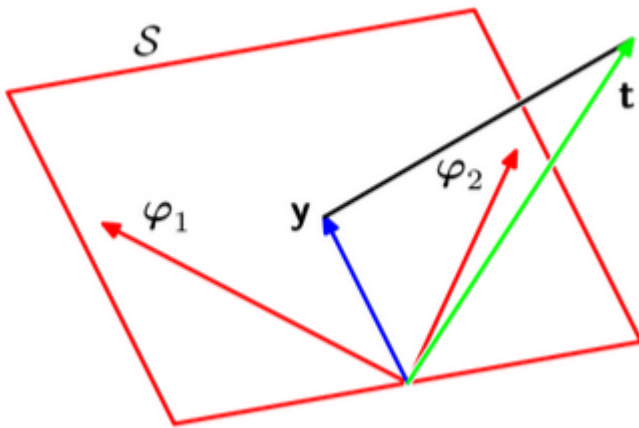
Derivando rispetto a  $w_0$  e uguagliando a zero, otteniamo il seguente risultato:

$$w_0 = \frac{1}{N} \sum_{n=1}^N t_n - \frac{\sum_{j=1}^{M-1} w_j \frac{1}{N} \sum_{n=1}^N \phi_j(X_n)}{\sum_{j=1}^{M-1} w_j \bar{\phi_j}}$$

Quindi il bias  $w_0$  **compensa la differenza tra la medie dei valori target e la somma pesata delle medie delle funzioni di base.**

## Intepretazione geometrica dei minimi quadrati

Consideriamo uno spazio N-dimensionale i cui assi sono dati da  $t_n$ , così che il vettore target  $t$  è un vettore in tale spazio. Ogni funzione base  $\phi_j(x)$  valorizzata negli N punti, può anch'essa essere vista come un vettore nello stesso spazio, chiamato  $\alpha_j$ .



$\alpha_j$  corrisponde alla  $j$ -esima colonna della matrice  $\Phi$ . Se  $M < N$ , allora i vettori  $\phi_j(X_n)$ , formeranno un sottospazio lineare  $S$  di dimensione  $M$  all'interno dello spazio  $N$ -dimensionale. Dato che  $y$  è una combinazione lineare arbitraria dei vettori  $\alpha_j$ , il vettore  $y$  giacerà nello spazio  $S$ . Il vettore target  $t$  giacerà in uno spazio  $N$ -dimensionale

Ripendiamo adesso la formula della somma dei quadrati degli errori:

$$E_D(W) = \frac{1}{2} \sum_{n=1}^N \{ t_n - W^T \phi(X_n) \}^2$$

Questa formula diventa uguale a:

$$E_D(W) = \frac{1}{2} \| t - \Phi W \|^2$$

ovvero la distanza Euclidea tra il vettore  $y$  e  $t$  (a meno di  $\frac{1}{2}$ ). Dato che  $y$  giace su  $S$ , la soluzione di minimo quadrato è data dalla proiezione ortogonale di  $t$  su  $S$ . Questa cosa è verificabile, notando che la soluzione per  $y$  è data da  $\Phi W_{ML}$ , dove  $W_{ML}$  è il vettore dei pesi che ottimizza i minimi quadrati.

Potrebbero sorgere delle difficoltà numeriche: se la matrice  $\Phi^T \Phi$  è quasi singolare (determinante vicino a zero). Questo accade quando due vettori o più, tra i vettori base  $\alpha_j$  sono quasi collineari. Questo può portare a parametri risultanti molto molto grandi.

## Sequential Learning

Il problema della tecnica del maximum likelihood è che richiede di processare tutto il set di addestramento in una sola volta. Con un set molto grande questo inizia a essere oneroso. Se il set è sufficientemente grande può convenire addestrare il modello sequenzialmente, ovvero aggiungendo un dato alla volta (one-

line algoritmo), in modo da aggiornare i parametri del modello a ogni presentazione dal dataset. Questo approccio è l'ideale anche per applicazioni in real time, in cui i dati arrivano costantemente. Per applicare questo approccio, possiamo usare la tecnica conosciuta come sequenziale gradient descent. Se la funzione di errore, consiste nella somma degli errori dei singoli punti, allora dopo la presentazione del pattern  $n$  l'algoritmo del gradiente sequenziale, aggiornerà il vettore dei parametri  $W$  come segue:

$$W^{(\tau+1)} = W^{(\tau)} + \eta \nabla E_n(W)$$

dove  $\tau$  indica in numero dell'interazione,  $\eta$  è il learning rate e  $\nabla E_n(W)$  è il gradiente dell'errore della singola presentazione  $n$ . Per il caso della somma degli errori al quadrato, abbiamo:

$$w^{(\tau+1)} = w^{(\tau)} + \eta (t_n - w^{(\tau)T} \phi(X_n)) \phi(X_n)$$

Questa formula è nota col nome di least-mean square (LMS) algorithm.

## Regularized Least Squares

Andiamo a aggiungere un termine di regolarizzazione alla funzione di errore per evitare overfitting. La funzione di errore da minimizzare diventa:

$$E_D(W) + \lambda E_W(W)$$

dove  $\lambda$  è il coefficiente di regolarizzazione che controlla l'importanza dell'errore dipendente dai dati  $E_D(W)$  e il termine di regolarizzazione  $E_W(W)$ .

Il modo più comune di regolarizzare è quello di usare un regolarizzatore dato dalla somma dei quadrati dei pesi:

$$E_W(W) = \frac{1}{2} W^T W$$

Consideriamo la funzione di errore come somma dei quadrati degli errori  $E_D(W) = \frac{1}{2} \sum_{n=1}^N \{ t_n - W^T \phi(X_n) \}^2$ . Otteniamo:

$$\frac{1}{2} \sum_{n=1}^N \{ t_n - W^T \phi(X_n) \}^2 + \frac{1}{2} \lambda W^T W$$

Questa forma di regolarizzatore è detta **weight decay**, perché negli algoritmi di apprendimento sequenziale fa tendere i pesi dei valori a zero, a meno che non siano supportati dai dati.

Il vantaggio di questo tipo di regolarizzatore è che la funzione di errore rimane quadratica rispetto ai pesi, quindi la minimizzazione avviene nello stesso modo. Calcoliamo quindi il gradiente della funzione di errore regolarizzata:

$$\nabla E(W) = \sum_{n=1}^N \{ t_n - W^T \phi(X_n) \} \phi(X_n) + \lambda W$$

Uguagliando a zero il gradiente, otteniamo la seguente equazione:

$$W = (\lambda I + \Phi^T \Phi)^{-1} \Phi^T t$$

Ricordiamo a cosa è uguale  $\Phi$ :

$$\Phi = \begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \dots & \phi_{M-1}(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x_N) & \phi_1(x_N) & \dots & \phi_{M-1}(x_N) \end{bmatrix}$$

Possiamo generalizzare la formula del regolarizzatore ponendo l'esponente uguale a  $q$ :

$$\frac{1}{2} \sum_{n=1}^N \|t_n - W^T \phi(X_n)\|^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$

$q=2$  corrisponde al regolarizzatore visto fino ad ora.

## Output multipli

Fino ad ora abbiamo considerato il caso di un solo target. Possiamo estendere il modello in modo da predire contemporaneamente  $K > 1$  target diversi, raggruppati in un vettore target  $t$ . È possibile ottenere questo risultato introducendo un set di funzioni base differenti per ogni target, ottenendo problemi di regressione tra loro indipendenti. Un approccio più comune è quello di usare un singolo set di funzioni base per modellare tutti i termini del vettore target. Definiamo:

$$y(X, W) = W^T \phi(X)$$

dove  $y$  è  $K$ -dimensionale e  $W$  è  $M \times K$  (matrice di parametri), mentre  $\phi(X)$  è  $M$ -dimensionale.

Supponiamo di prendere una distribuzione condizionata Gaussiana per i target:

$$p(t|X, W, \beta) = N(t|W^T \phi(X), \beta^{-1} I)$$

Se abbiamo un set di osservazioni  $t_1, \dots, t_N$ , possiamo creare una matrice  $T$   $N \times K$  in modo che la riga  $n$ -esima contenga il target per l'osservazione  $n$ -esima  $t_n^T$ . In modo analogo creiamo la matrice  $X$ . Quindi la forma logaritmica della funzione di likelihood diventa:

$$\ln p(T|X, W, \beta) = \sum_{n=1}^N \ln N(t_n|W^T \phi(X_n), \beta^{-1} I)$$

che è uguale a:

$$-\frac{NK}{2} \ln \left( \frac{\beta}{2\pi} \right) - \frac{\beta}{2} \sum_{n=1}^N \|t_n - W^T \phi(X_n)\|^2$$

## Decomposizione bias-varianza