

I predittori lineari, che usano la mappa $\mathbf{x} \mapsto \mathbf{w}^\top \mathbf{x}$, possono soffrire di un errore di bias molto elevato. Ovvero, la distribuzione \mathcal{D} che genera i dati può essere tale che il rischio statistico $\ell_{\mathcal{D}}(h)$ è alto per ogni elemento h in una classe \mathcal{H} di predittori lineari. In questa situazione anche il miglior predittore in \mathcal{H} è destinato ad avere un test error tipicamente alto. D'altra parte, i predittori lineari sono semplici da implementare e da analizzare.

Utilizzando una mappa non lineare per trasformare le istanze di un problema di apprendimento, è però possibile utilizzare predittori lineari per rappresentare predittori non lineari. Consideriamo infatti una funzione nonlineare $\phi : \mathbb{R}^d \rightarrow \mathcal{V}$ che mappa le istanze $\mathbf{x} \in \mathbb{R}^d$ in un nuovo spazio \mathcal{V} con un numero molto maggiore (anche infinito) di dimensioni. Costruendo la funzione ϕ opportunamente, abbiamo che superfici nonlineari in \mathbb{R}^d vengono mappate in superfici lineari (iperpiani) in \mathcal{V} .

Per esempio, consideriamo la funzione $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^6$, tale che $\phi(x_1, x_2) = (1, x_1^2, x_2^2, x_1, x_2, x_1x_2)$. Un iperpiano omogeneo in \mathbb{R}^6 con coefficienti $\mathbf{w} = (w_1, \dots, w_6)$ ha forma

$$\{\mathbf{z} \in \mathbb{R}^6 : \mathbf{w}^\top \mathbf{z} = 0\} .$$

Questo iperpiano può essere utilizzato per descrivere una superficie nonlineare in \mathbb{R}^2 . Infatti, l'insieme dei $\mathbf{x} \in \mathbb{R}^2$ tali che $\mathbf{w}^\top \phi(\mathbf{x}) = 0$ per definizione di ϕ soddisfano l'equazione

$$w_1 + w_2x_1^2 + w_3x_2^2 + w_4x_1 + w_5x_2 + w_6x_1x_2 = 0$$

che è l'equazione di una conica nello spazio Euclideo bidimensionale (ovvero la famiglia a cui appartengono ellissi, parabole e iperboli).

Se ora consideriamo una generica mappa $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^N$, con $N \gg d$, un classificatore lineare in \mathbb{R}^N per le istanze mappate da \mathbb{R}^d avrà forma

$$h(\phi(\mathbf{x})) = \text{sgn}(\mathbf{w}^\top \phi(\mathbf{x})) \quad \text{dove} \quad \mathbf{w}^\top \phi(\mathbf{x}) = \sum_{i=1}^N w_i \phi(\mathbf{x})_i .$$

Se $N \gg d$ il calcolo di $\mathbf{w}^\top \phi(\mathbf{x})$, che richiede la somma di N prodotti, può però risultare computazionalmente eccessivo.

In realtà, questo problema computazionale può essere enormemente alleviato combinando certi tipi di mappe col modo con cui gli algoritmi di apprendimento lineare tipicamente rappresentano i loro modelli. Infatti, ricordando per esempio la regola di aggiornamento del Perceptrone $\mathbf{w} \leftarrow \mathbf{w} + y\mathbf{x}$, possiamo scrivere un generico classificatore lineare h generato dal Perceptrone come

$$h(\mathbf{x}) = \text{sgn} \left(\sum_{j=1}^M y_{t_j} \mathbf{x}_{t_j}^\top \mathbf{x} \right)$$

dove $(\mathbf{x}_{t_1}, y_{t_1}), \dots, (\mathbf{x}_{t_M}, y_{t_M})$ sono gli esempi sui quali il Perceptrone ha effettuato un aggiornamento. Se quindi facciamo girare il Perceptrone in \mathcal{V} , il classificatore lineare h diventa

$$\tilde{h}(\mathbf{x}) = \text{sgn} \left(\sum_{j=1}^M y_{t_j} \phi(\mathbf{x}_{t_j})^\top \phi(\mathbf{x}) \right).$$

Per calcolare $\tilde{h}(\mathbf{x})$ efficientemente dobbiamo poter calcolare ciascun termine $\phi(\mathbf{x}_{t_j})^\top \phi(\mathbf{x})$ efficientemente. La tecnica delle funzioni kernel ci dice come scegliere ϕ in modo tale che esista una funzione $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ tale che

$$K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}') \quad \text{per ogni } \mathbf{x}, \mathbf{x}' \in \mathbb{R}^d. \quad (1)$$

Per esempio, la funzione kernel corrispondente a $\phi(x_1, x_2) = (1, x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2)$ è $K(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^\top \mathbf{x}')^2$, come è immediato verificare.

Fissata una funzione kernel K , possiamo allora scrivere il classificatore lineare a soglia generato dal Perceptrone come

$$h_K(\mathbf{x}) = \text{sgn} \left(\sum_{j=1}^M y_{t_j} K(\mathbf{x}_{t_j}, \mathbf{x}) \right). \quad (2)$$

Si veda qui sotto per una descrizione del Perceptrone con kernel.

Algoritmo: Perceptrone con kernel

Inizializza S con l'insieme vuoto

Per ogni $t = 1, 2, \dots$

1. Leggi il prossimo esempio (\mathbf{x}_t, y_t)

2. Calcola $\hat{y}_t = \text{sgn} \left(\sum_{s \in S} y_s K(\mathbf{x}_s, \mathbf{x}_t) \right)$

3. Se $\hat{y}_t \neq y_t$ aggiungi t a S

La generalizzazione ovvia del kernel per apprendere coniche in \mathbb{R}^2 è il kernel polinomiale $K_n(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^\top \mathbf{x}')^n$ con $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$. Questo kernel permette di apprendere classificatori di grado n in \mathbb{R}^d rappresentandoli come iperpiani in uno spazio \mathbb{R}^N . Possiamo esplicitare la funzione ϕ_n tale che $K_n(\mathbf{x}, \mathbf{x}') = \phi_n(\mathbf{x})^\top \phi_n(\mathbf{x}')$ usando il teorema del binomio di Newton,

$$(1 + \mathbf{x}^\top \mathbf{x}')^n = \sum_{k=0}^n \binom{n}{k} (\mathbf{x}^\top \mathbf{x}')^k. \quad (3)$$

Osserviamo poi che vale

$$(\mathbf{x}^\top \mathbf{x}')^k = \left(\sum_{i=1}^d x_i x'_i \right)^k = \sum_{v \in \{1, \dots, d\}^k} \left(\prod_{s=1}^k x_{v_s} x'_{v_s} \right).$$

Quindi,

$$\phi_n(\mathbf{x}) = \left(\sqrt{\binom{n}{k}} \prod_{s=1}^k x_{v_s} \right)_{k=0, \dots, n, v \in \{1, \dots, d\}^k} .$$

In altre parole, il mapping $\phi_n : \mathbb{R}^d \rightarrow \mathbb{R}^N$ associato al kernel polinomiale K_n è un vettore le cui componenti sono indicizzate da tutti i monomi $\prod_{s=1}^k x_{v_s}$ di grado al più n pesati dalle radici dei coefficienti binomiali.

Un ulteriore tipo di kernel è il kernel Gaussiano,

$$K_\gamma(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2\gamma} \|\mathbf{x} - \mathbf{x}'\|^2\right) \quad \gamma > 0 .$$

Per trovare la mappa ϕ_γ associata a K_γ procediamo come segue,

$$\begin{aligned} \exp\left(-\frac{1}{2\gamma} \|\mathbf{x} - \mathbf{x}'\|^2\right) &= \exp\left(-\frac{1}{2\gamma} (\|\mathbf{x}\|^2 + \|\mathbf{x}'\|^2)\right) \exp\left(\frac{1}{\gamma} (\mathbf{x}^\top \mathbf{x}')\right) \\ &= \exp\left(-\frac{\|\mathbf{x}\|^2}{2\gamma}\right) \exp\left(-\frac{\|\mathbf{x}'\|^2}{2\gamma}\right) \sum_{n=0}^{\infty} \frac{1}{n!} \frac{(\mathbf{x}^\top \mathbf{x}')^n}{\gamma^n} \end{aligned}$$

dove abbiamo usato l'espansione in serie di Taylor $e^x = 1 + x + \frac{x^2}{2!} + \dots$. Osservando la formula (3) ci rendiamo conto che il kernel Gaussiano è il limite di una successione di somme parziali in cui gli elementi della successione sommati sono gli elementi di un kernel polinomiale pesati dal reciproco del fattoriale. Il parametro γ svolge il ruolo di fattore di scala per i prodotti $\mathbf{x}^\top \mathbf{x}'$. Infine i fattori $e^{-\|\mathbf{x}\|^2/(2\gamma)} e^{-\|\mathbf{x}'\|^2/(2\gamma)}$ normalizzano rispetto a \mathbf{x} e \mathbf{x}' . Infatti, $K_\gamma(\mathbf{x}, \mathbf{x}) = 1$ per ogni $\mathbf{x} \in \mathbb{R}^d$.

Il kernel Gaussiano è *universale*, nel senso che un classificatore della forma (2), cioè tale che $h(\mathbf{x}) = \text{sgn}(g(\mathbf{x}))$, con $g : \mathbb{R}^d \rightarrow \mathbb{R}$ appartenente alla classe

$$\mathcal{H}_\gamma = \left\{ \sum_{i=1}^N \alpha_i K_\gamma(\mathbf{x}_i, \cdot) : \mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d, \alpha_1, \dots, \alpha_N \in \mathbb{R}, N \in \mathbb{N} \right\} \quad (4)$$

può approssimare con precisione arbitraria una qualunque funzione continua da \mathbb{R}^d a \mathbb{R} per ogni valore di $\gamma > 0$. Questo significa che se (D, η) è un modello statistico con η continuo, allora possiamo approssimare in modo arbitrariamente accurato il classificatore Bayesiano ottimo $f^*(x) = \text{sgn}(2\eta(x) - 1)$ usando funzioni della forma $\text{sgn}(g)$ con $g \in \mathcal{H}_\gamma$.

Dato un qualunque insieme \mathcal{X} a cui appartengono le istanze \mathbf{x} (non necessariamente \mathbb{R}^d), è possibile dimostrare che una qualunque funzione simmetrica $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ è un kernel se e solo se per ogni $m \in \mathbb{N}$ e per ogni $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathcal{X}$, la matrice \mathbf{K} di dimensioni $m \times m$ e tale che $\mathbf{K}_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j)$ è positiva semidefinita. Ovvero, $\mathbf{z}^\top \mathbf{K} \mathbf{z} \geq 0$ per ogni $\mathbf{z} \in \mathbb{R}^m$. La rappresentazione della mappa $\phi_K : \mathcal{X} \rightarrow \mathcal{H}_K$ associata ad un kernel K e la rappresentazione dello spazio \mathcal{H}_K non sono uniche. Però possiamo sempre rappresentare ϕ_K come $\phi_K(\mathbf{x}) = K(\mathbf{x}, \cdot)$ e \mathcal{H}_K come lo spazio

$$\mathcal{H}_K = \left\{ \sum_{i=1}^N \alpha_i K(\mathbf{x}_i, \cdot) : \mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d, \alpha_1, \dots, \alpha_N \in \mathbb{R}, N \in \mathbb{N} \right\}$$

equivalente a quello definito in (4) per i kernel Gaussiani.

\mathcal{H}_K è uno spazio lineare di funzioni dotato di un'operazione di prodotto interno (l'analogo del prodotto scalare fra vettori) che soddisfa

$$\langle K(\mathbf{x}, \cdot), K(\mathbf{x}', \cdot) \rangle = K(\mathbf{x}, \mathbf{x}')$$

in modo che $\langle \phi_K(\mathbf{x}), \phi_K(\mathbf{x}') \rangle = K(\mathbf{x}, \mathbf{x}')$ —ovvero, il kernel implementa il prodotto interno su ϕ_K . In generale, se

$$f = \sum_{i=1}^N \alpha_i K(\mathbf{x}_i, \cdot) \quad \text{e} \quad g = \sum_{j=1}^M \alpha_j K(\mathbf{x}'_j, \cdot)$$

allora

$$\begin{aligned} \langle f, g \rangle &= \left\langle \sum_{i=1}^N \alpha_i K(\mathbf{x}_i, \cdot), \sum_{j=1}^M \beta_j K(\mathbf{x}'_j, \cdot) \right\rangle \\ &= \sum_{i=1}^N \alpha_i \left\langle K(\mathbf{x}_i, \cdot), \sum_{j=1}^M \beta_j K(\mathbf{x}'_j, \cdot) \right\rangle \\ &= \sum_{i=1}^N \sum_{j=1}^M \alpha_i \beta_j \langle K(\mathbf{x}_i, \cdot), K(\mathbf{x}'_j, \cdot) \rangle \\ &= \sum_{i=1}^N \sum_{j=1}^M \alpha_i \beta_j K(\mathbf{x}_i, \mathbf{x}'_j) \end{aligned}$$

dove abbiamo usato le proprietà di linearità che valgono per qualsiasi prodotto interno.

Dato che i predittori con kernel sono predittori lineari in \mathcal{H}_K , è possibile limitare il rischio (statistico o sequenziale) usando le tecniche utilizzate in precedenza per i predittori lineari. Il maggiorante sul rischio dipenderà dai parametri usuali, ovvero $\|\mathbf{w}\|$ e $\max_t \|\mathbf{x}_t\|$, dove però le norme sono calcolate in \mathcal{H}_K . In analogia con la definizione di norma di vettori $\|\mathbf{x}\|^2 = \mathbf{x}^\top \mathbf{x}$, la norma di $f \in \mathcal{H}_K$ è definita come $\|f\|^2 = \langle f, f \rangle$.

Questo ci dà subito il seguente risultato: $\|\phi_K(\mathbf{x})\| = \sqrt{\langle K(\mathbf{x}, \cdot), K(\mathbf{x}, \cdot) \rangle} = \sqrt{K(\mathbf{x}, \mathbf{x})}$.

Nel caso di un kernel K , per ogni $g \in \mathcal{H}_K$ abbiamo che

$$\|g\| = \left\| \sum_{i=1}^N \alpha_i K(\mathbf{x}_i, \cdot) \right\| = \sqrt{\left\langle \sum_{i=1}^N \alpha_i K(\mathbf{x}_i, \cdot), \sum_{j=1}^N \alpha_j K(\mathbf{x}_j, \cdot) \right\rangle} = \sqrt{\sum_{i,j=1}^N \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j)}.$$

Possiamo anche calcolare il prodotto interno fra $\phi_K(\mathbf{x}) = K(\mathbf{x}, \cdot)$ e un qualunque $f \in \mathcal{H}_K$,

$$\langle f, K(\mathbf{x}, \cdot) \rangle = \sum_{i=1}^N \alpha_i \langle K(\mathbf{x}_i, \cdot), K(\mathbf{x}, \cdot) \rangle = \sum_{i=1}^N \alpha_i K(\mathbf{x}_i, \mathbf{x}) = f(\mathbf{x}).$$

Da ciò si vede che f calcola il prodotto interno fra se stessa e ϕ_K . Questa proprietà si chiama *reproducing* e per questa ragione \mathcal{H}_K è anche chiamato un *reproducing kernel Hilbert space* (RKHS).

Consideriamo un training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$ linearmente separabile in \mathcal{H}_K per un qualche kernel K . Ovvero, esiste un $f \in \mathcal{H}_K$ tale che $\min_t y_t f(\mathbf{x}_t) \geq 1$. Il teorema di convergenza del Perceptrone (variante kernel) ci dà un maggiorante sul numero di aggiornamenti pari a

$$\|f\|^2 \max_t K(\mathbf{x}_t, \mathbf{x}_t) .$$