

I predittori lineari possono soffrire di un errore di bias molto elevato in quanto predittore il Bayesiano ottimo  $f^*$  è spesso molto lontano dall'essere lineare. Questo crea un potenziale rischio di underfitting tutte le volte che utilizziamo un algoritmo di apprendimento che genera predittori lineari.

Utilizzando una mappa non lineare per trasformare le istanze di un problema di apprendimento diventa possibile utilizzare predittori lineari per rappresentare predittori non lineari, alleviando il problema dell'underfitting. Consideriamo infatti una funzione non lineare  $\phi : \mathbb{R}^d \rightarrow \mathcal{V}$  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  $\phi$  opportunamente, abbiamo che superfici non lineari 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  $\mathbf{w} = (w_1, \dots, w_6)$  è descritto dall'insieme

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

Iperpiani di questa forma possono essere utilizzati per descrivere superfici non lineari in  $\mathbb{R}^2$ . Infatti, per definizione di  $\phi$ ,  $\mathbf{w}^\top \phi(\mathbf{x}) = w_1 + w_2x_1^2 + w_3x_2^2 + w_4x_1 + w_5x_2 + w_6x_1x_2$ . Ora notiamo che gli insiemi  $\left\{ \mathbf{x} \in \mathbb{R}^2 : w_1 + w_2x_1^2 + w_3x_2^2 + w_4x_1 + w_5x_2 + w_6x_1x_2 = 0 \right\}$  descrivono curve coniche nello spazio Euclideo bidimensionale (ovvero la famiglia a cui appartengono ellissi, parabole e iperboli).

Data una generica mappa  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^N$ , con  $N \gg d$ , consideriamo il classificatore  $h : \mathbb{R}^d \rightarrow \{-1, +1\}$  definito da

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

Questo classificatore, che non è lineare in  $\mathbb{R}^d$ , equivale ad un classificatore lineare nello spazio  $\phi(\mathbb{R}^d) \subseteq \mathbb{R}^N$  definito come  $\phi(\mathbb{R}^d) \equiv \{ \mathbf{z} \in \mathbb{R}^N : (\exists \mathbf{x} \in \mathbb{R}^d) \phi(\mathbf{x}) = \mathbf{z} \}$ . 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 utilizzando certi tipi di mappe  $\phi$ . Infatti, ricordando per esempio la regola di aggiornamento del Perceptrone  $\mathbf{w}_{t+1} = \mathbf{w}_t + y_t \mathbf{x}_t$ , possiamo scrivere un generico classificatore lineare  $h$  generato dal Perceptrone come

$$h(\mathbf{x}) = \text{sgn} \left( \sum_{s \in S} y_s \mathbf{x}_s^\top \mathbf{x} \right)$$

dove  $S$  è l'insieme degli indici  $s$  di esempi  $(\mathbf{x}_s, y_s)$  dove il Perceptrone ha eseguito un aggiornamento. Se quindi facciamo girare il Perceptrone in  $\phi(\mathbb{R}^d)$ , il classificatore lineare  $h$  diventa

$$h_\phi(\mathbf{x}) = \text{sgn} \left( \sum_{s \in S} y_s \phi(\mathbf{x}_s)^\top \phi(\mathbf{x}) \right) .$$

Per calcolare  $h_\phi(\mathbf{x})$  efficientemente basta poter calcolare ciascun termine  $\phi(\mathbf{x}_s)^\top \phi(\mathbf{x})$  velocemente. La tecnica delle funzioni kernel ci dice come scegliere  $\phi$  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

$$h_K(\mathbf{x}) = \text{sgn} \left( \sum_{s \in S} y_s K(\mathbf{x}_s, \mathbf{x}) \right) . \quad (2)$$

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  $\phi_n$  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

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

Quindi,

$$\phi_n(\mathbf{x}) = \left( \sqrt{\binom{n}{k}} \prod_{s=1}^k x_{v_s} \right)_{k=0, \dots, n, \mathbf{v} \in \{1, \dots, d\}^k} . \quad (4)$$

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.

Nel caso del kernel polinomiale di grado  $n$ , la dimensione  $N$  del codominio di  $\phi$  è

$$\sum_{k=0}^n |\{1, \dots, d\}^k| = \sum_{k=0}^n d^k = \frac{d^{n+1} - 1}{d - 1}$$

ovvero  $N = \Theta(d^n)$ , esponenziale nel grado del polinomio.

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  $\phi_\gamma$  associata a  $K_\gamma$  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 una combinazione lineare di infiniti kernel polinomiali di grado crescente, ciascuno pesato dal reciproco del fattoriale del suo grado. Il parametro  $\gamma$  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  $\mathbf{x} \in \mathbb{R}^d$ .

Notiamo che predittori della forma (2) che utilizzano kernel Gaussiani classificano un punto  $\mathbf{x}$  calcolando una combinazione lineare, con coefficienti  $y_s$ , di funzioni Gaussianhe  $e^{-z^2/(2\gamma)}$  delle distanze  $z = \|\mathbf{x}_s - \mathbf{x}\|$  per ogni  $s \in S$ . Questi classificatori sono quindi confrontabili con classificatori  $k$ -NN, dato che anche questi ultimi utilizzano la distanza dai punti di training per determinare la classificazione di un punto di test.

Il kernel Gaussiano è *universale*: per ogni  $\gamma > 0$ , per ogni funzione continua  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , e per ogni  $\varepsilon > 0$ , esiste  $g \in \mathcal{H}_\gamma$ , con

$$\mathcal{H}_\gamma = \left\{ \sum_{i=1}^N \alpha_i K_\gamma(\mathbf{x}_i, \cdot) : \mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d, \alpha_1, \dots, \alpha_N \in \mathbb{R}, N \in \mathbb{N} \right\} \quad (5)$$

tale che  $h$  approssima  $f$  con errore al più  $\varepsilon$ . Questo significa che se  $(D, \eta)$  è un modello statistico tale che  $\eta$  sia 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  $g \in \mathcal{H}_\gamma$ .

Un algoritmo che genera classificatori della forma  $h(\mathbf{x}) = \text{sgn}(g(\mathbf{x}))$ , con  $g \in \mathcal{H}_\gamma$ , è non parametrico. Infatti, pur essendo  $h$  un classificatore lineare, l'iperpiano  $g \in \mathcal{H}_\gamma$  non è definito tramite un numero predeterminato di parametri. Mentre nel caso di kernel polinomiali di grado  $n$ , gli iperpiani sono

sempre esprimibili con un numero limitato (ancorché esponenziale in  $n$ ) di coefficienti, nel caso di kernel Gaussiani gli iperpiani sono definiti tramite le combinazioni lineari

$$\sum_{i=1}^N \alpha_i K_\gamma(\mathbf{x}_i, \cdot)$$

dove  $N$  è sempre finito, ma arbitrariamente grande.

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  $\phi_K$  come  $\phi_K(\mathbf{x}) = K(\mathbf{x}, \cdot)$  e  $\mathcal{H}_K$  come lo spazio

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

equivalente a quello definito in (5) per i kernel Gaussiani.

Nel caso di kernel polinomiali  $K_n$ , sappiamo che  $\mathcal{H}_K$  ha una rappresentazione equivalente come  $\mathbb{R}^N$  (con  $N = \Theta(d^n)$ ) e sappiamo scrivere  $\phi_n$  usando la formula (4). Quindi, nel caso dei kernel polinomiali, riusciamo a scrivere i classificatori lineari esplicitamente come  $h(\mathbf{x}) = \text{sgn}(\mathbf{w}^\top \phi_n(\mathbf{x}))$ , con  $\mathbf{w} \in \mathbb{R}^N$ . Per fare la stessa cosa con kernel arbitrari, per i quali non sappiamo scrivere  $\phi_K$  esplicitamente, usiamo il fatto che  $K(\mathbf{x}, \mathbf{x}')$  implementa il prodotto interno fra  $\phi_K(\mathbf{x})$  e  $\phi_K(\mathbf{x}')$ . Usando  $\langle \cdot, \cdot \rangle_K$  per denotare il prodotto interno in  $\mathcal{H}_K$ , possiamo quindi definire

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

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

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

dove abbiamo usato le proprietà di linearità che valgono per qualsiasi prodotto interno. Da ciò si vede che  $f$  calcola il prodotto interno fra se stessa e  $\phi_K$ . Questa proprietà si chiama *reproducing* e per questa ragione  $\mathcal{H}_K$  è anche chiamato un *reproducing kernel Hilbert space* (RKHS).

Di conseguenza, mentre in  $\mathbb{R}^d$  il classificatore lineare era definito come  $\text{sgn}(\mathbf{w}^\top \mathbf{x})$ , in  $\mathcal{H}_K$  il classificatore diventa  $\text{sgn}(f(\mathbf{x}))$  dove  $f(\mathbf{x}) = \langle f, \phi_K(\mathbf{x}) \rangle_K$ .

Infine, vediamo come calcolare il prodotto interno  $\langle f, g \rangle_K$  fra due elementi arbitrari di  $\mathcal{H}_K$ ,

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

Abbiamo

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

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  $\|f\|_K^2 = \langle f, f \rangle_K$ .

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

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

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

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

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

È anche possibile dimostrare una versione kernel del maggiorante sul rischio sequenziale del Perceptrone nel caso non linearmente separabile.

Un predittore lineare in uno spazio kernel può subire overfitting. Nel caso di kernel polinomiali  $K_n$ , aumentando il grado  $n$  del kernel riduciamo il training error (in quanto separiamo i dati in  $\mathbb{R}^d$  con curve di grado  $n$  sempre più elevato). Se il grado  $n$  diventa troppo elevato rispetto alla dimensione del training set (e al livello di rumore nei dati) si crea quindi overfitting. Viceversa, quando  $n = 1$  ricadiamo nel caso particolare di classificatori lineari e siamo quindi a rischio di underfitting.

Nel caso di kernel Gaussiani  $K_\gamma$ , il parametro  $\gamma$  corrisponde all'ampiezza delle Gaussiane centrate sui punti  $\mathbf{x}_s$  —si veda (2). Se  $\gamma$  è molto piccolo rispetto alle distanze  $\|\mathbf{x}_s - \mathbf{x}_t\|^2$  fra i punti di training e test, allora la classificazione di un punto  $\mathbf{x}$  verrà essenzialmente determinata dal punto  $\mathbf{x}_s$  ad esso più vicino. Questo permette di ottenere training error vicini a zero in quanto, analogamente

a quanto succede per 1-NN, tutti i punti di training che sono inseriti in  $S$  vengono classificati correttamente dal predittore di forma (2). Anche in questo caso, quindi, si può creare overfitting scegliendo  $\gamma$  troppo piccolo. Viceversa, quando  $\gamma$  è grande rispetto alle distanze  $\|\mathbf{x}_s - \mathbf{x}_t\|^2$ , allora le Gaussiane centrate sui punti  $\mathbf{x}_s$  sono molto ampie. I classificatori diventano quindi simili a classificatori  $k$ -NN quando  $k$  è scelto vicino alla cardinalità del training set e siamo quindi a rischio di underfitting.

Abbiamo stabilito che una qualunque funzione simmetrica  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  è un kernel se e solo la matrice  $\mathbf{K}$  è positiva semidefinita. Questo risultato è importante perché non pone condizioni sullo spazio  $\mathcal{X}$  dei dati. Possiamo quindi definire kernel  $K$  su qualsiasi insieme  $\mathcal{X}$ : matrici, sequenze, alberi, grafi, eccetera. Il kernel può perciò essere visto come un modo di incapsulare lo spazio dei dati offrendo all'algoritmo di apprendimento un'interfaccia uniforme (ovvero, la funzione kernel). Per guidare l'intuizione nel disegno di una funzione kernel su un dato spazio  $\mathcal{X}$ , si può ricordare il fatto che  $K(x, x')$  implementa il prodotto interno fra i due elementi  $\phi_K(x)$  e  $\phi_K(x')$  di uno spazio lineare  $\mathcal{H}_K$ , e può quindi essere visto come una misura di similarità fra  $x$  e  $x'$  in  $\mathcal{X}$ . Possiamo allora partire definendo una misura di similarità  $K$  su  $\mathcal{X}$  e poi aggiustarne la definizione in modo che la matrice risultante rispetti la condizione di semidefinitezza positiva.