

Non-linear classification (continued)

- As discussed earlier the Adaline can readily be extended to training data that is not linearly separable.
- This was done by introducing a non-linear function

$$\underline{\Phi}: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

from n -dim feature space to an m -dim space for sufficiently large m so that the training data becomes linearly separable again.

- This can also be done with the SVM which in dual formulation reads

$$\max \tilde{L}(\alpha) \text{ for}$$

$$\tilde{L}(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} \langle x^{(i)}, x^{(j)} \rangle$$

$$\text{Subject to } 0 \leq \alpha_i \leq \mu, i=1, \dots, n, \sum_{j=1}^n \alpha_j y^{(j)} = 0$$

- The scalar product $\langle x^{(i)}, x^{(j)} \rangle$ must only be replaced by

$$S(x^{(i)}, y^{(j)}) = \langle \underline{\Phi}(x^{(i)}), \underline{\Phi}(x^{(j)}) \rangle_{\mathcal{X}}$$

where $\underline{\Phi}$ must be chosen such that

$$\underline{\Phi}: \mathbb{R}^n \rightarrow \mathcal{X} = \text{Relevant space with } \dim \mathcal{X} \gg n$$

- However, this is computationally very costly as computation of $S(x^{(i)}) \sim \mathcal{O}(\dim \mathcal{X})$.

- Idea: replace the scalar product $S(x^{(i)})$ by a function $k(x^{(i)})$ that mimics the properties of a scalar product.

This leads to the following definition:

DEF (pos. def. sym. kernels) A map

$K: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ is called positive definite symmetric kernel if

$\forall \{x_1, \dots, x_n\} \subset \mathbb{R}^n$: the matrix
 $(K(x_i, x_j))_{\substack{1 \leq i \leq n \\ 1 \leq j \leq n}}$

is symmetric and positive semi-definite

Recall: A matrix M is sym. & pos. semi-def.

iff it is symmetric and

$$\forall v \in \mathbb{R}^n: \langle v, Mv \rangle \geq 0$$

The idea is that for such kernel K there is a Φ such that $K(x, y) = \langle \Phi(x), \Phi(y) \rangle$

Example: Polynomial kernel

We define $\forall x, y \in \mathbb{R}^n$:

$$K(x, y) = (x \cdot y + c)^d$$

for tuning parameters $c \in \mathbb{R}$, $d \in \mathbb{N}$

For $n=2$, $d=2$ we find

$$\Phi: \mathbb{R}^2 \rightarrow \mathbb{R}^6, \quad \Phi(x) = \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2} x_1 x_2 \\ \sqrt{2c} x_1 \\ \sqrt{2c} x_2 \\ c \end{pmatrix}$$

because

$$\begin{aligned} (x \cdot y + c)^2 &= (x_1 y_1 + x_2 y_2 + c)^2 \\ &= (x_1 y_1 + x_2 y_2)^2 + 2(x_1 y_1 + x_2 y_2)c + c^2 \\ &= x_1^2 y_1^2 + 2x_1 y_1 x_2 y_2 + x_2^2 y_2^2 \\ &\quad + 2(x_1 y_1 + x_2 y_2)c + c^2 \end{aligned}$$

LEM (Cauchy-Schwarz): k pos. def. sym.,

then $\forall x, y \in \mathbb{R}^n$,

$$k(x, y)^2 \leq k(x, x) k(y, y)$$

Proof: Def. matrix

$$M := \begin{pmatrix} k(x, x) & k(x, y) \\ k(y, x) & k(y, y) \end{pmatrix}$$

Since k pos. def. sym. M

is symmetric and pos. semi-def.

$$\Rightarrow \det M \geq 0$$

$$= k(x, x)k(y, y) - k(x, y)^2 \quad \square$$

THM: Let $k: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ be a pos. def. sym. kernel, then

i) \exists a Hilbert space \mathcal{H} of real-valued functions and a map

$$\Phi: \mathbb{R}^n \rightarrow \mathcal{H} \text{ s.t.}$$

$$k(x, y) = \langle \Phi(x), \Phi(y) \rangle_{\mathcal{H}}$$

ii) $\forall f \in \mathcal{H}, x \in \mathbb{R}^n$

$$f(x) = \langle f, k(x, \cdot) \rangle_{\mathcal{H}}$$

Proof: For $x, y \in \mathbb{R}^n$ define

$$\Phi(x)(y) := k(x, y)$$

$$\text{Def. } \mathcal{H}^0 := \left\{ \sum_{i \in I} a_i \Phi_i(x_i) \mid a_i \in \mathbb{R}, x_i \in \mathbb{R}^n, |I| < \infty \right\}$$

and an inner product

$$\langle \cdot, \cdot \rangle: \mathcal{H}^0 \times \mathcal{H}^0 \rightarrow \mathbb{R} \text{ s.t.}$$

such that for $f, g \in \mathcal{X}^0$, i.e.,

$$f = \sum_{i \in I} a_i \Phi(x_i)$$

$$g = \sum_{j \in J} b_j \Phi(x_j)$$

$$\langle f, g \rangle := \sum_{\substack{i \in I \\ j \in J}} a_i b_j k(x_i, x_j)$$

$$= \sum_{i \in I} a_i g(x_i)$$

$$= \sum_{j \in J} b_j f(x_j)$$

Hence, i) $\langle f, g \rangle = \langle g, f \rangle$

ii) Representer independent

iii) $\langle f, f \rangle = \sum_{\substack{i \in I \\ j \in J}} a_i b_j k(x_i, x_j) \geq 0$

i.e. pos. semi-def.

Hence, $\langle \cdot, \cdot \rangle$ is a pos. def. sym. kernel on \mathcal{X}^0 .

But $\forall f \in \mathcal{X}^0$

$$\langle f, \Phi(x) \rangle^2 \stackrel{(*)}{\leq} \langle f, f \rangle \langle \Phi(x), \Phi(x) \rangle$$

and for $f = \sum_{i \in I} a_i \Phi(x_i)$

$$f(x) = \sum_{i \in I} a_i \Phi(x_i)(x)$$

$$= \sum_{i \in I} a_i k(x_i, x)$$

$$= \langle f, \Phi(x) \rangle$$

Hence,

$$|f(x)|^2 = \langle f, \Phi(x) \rangle^2 \leq \langle f, f \rangle k(x, x)$$

$$\Rightarrow f = 0 \Leftrightarrow \langle f, f \rangle = 0$$

Thus $\langle \cdot, \cdot \rangle$ is an inner product on \mathcal{X}^0 .

$\mathcal{H} = \overline{\mathcal{H}^0}^{(\cdot, \cdot)}$, i.e., completion w.r.t. (\cdot, \cdot) , \Rightarrow activation function

\mathcal{H} is again a Hilbert space with inner product (\cdot, \cdot) .

As $f \mapsto \langle f, \Phi(x) \rangle$ is Lipschitz

continuous due to (*) and by construction

\mathcal{H}^0 is dense (***) holds on all of \mathcal{H} . \square

Conclusion: SVM for non-linear problems

1) choose pos. def. sym. kernel K

$$2) \max_{\alpha} \hat{L}(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} \underbrace{K(x^{(i)}, x^{(j)})}_{= \Phi(x^{(i)}) \cdot \Phi(x^{(j)})}$$

subject to $\mu = \alpha_i + \beta_i, \alpha_i \geq 0, \beta_i \geq 0$
 $\sum_{i=1}^n \alpha_i y^{(i)} = 0$ for $i=1 \dots n$

KKT(I) i) $\sum_{i=1}^n \alpha_i y^{(i)} = 0$

ii) $\underline{\omega} = \sum_{i=1}^n \alpha_i y^{(i)} \Phi(x^{(i)})$

iii) $\mu = \alpha_i + \beta_i$ for $i=1 \dots n$

KKT(II) iv) $\alpha_i (1 - y^{(i)} (\omega_0 + \underline{\omega} \cdot \Phi(x^{(i)})) - \xi_i) = 0$
 for $\alpha_i \geq 0$ $i=1 \dots n$

v) $\beta_j \xi_j = 0$

vi) for $\beta_j \geq 0$ $j=1 \dots n$

compute

$$\underline{\omega} = \sum_{i=1}^n \alpha_i y^{(i)} \underline{x}^{(i)}$$

for vectors fulfilling $\alpha_i \neq 0$

• these are $(y^{(i)}, x^{(i)})$ fulfilling

$$1 - y^{(i)} (\omega_0 + \underline{\omega} \cdot \Phi(x^{(i)})) - \xi_i = 0$$

• but must fulfill

$$\beta_i q_i = 0$$

$$\Rightarrow \begin{cases} q_i = 0 & \text{support vector } y^{(i)}(\omega_0 + \underline{\omega} \cdot \underline{\Phi}(x^{(i)})) = 1 \\ q_i \neq 0 & \text{outlier } y^{(i)}(\underline{\omega} \cdot x^{(i)}) = 1 - q_i \end{cases}$$

can be detected by
 $\beta_i = 0 \Rightarrow d_i = \mu$

So choose support vector $(x^{(i)}, y^{(i)})$ to compute

$$y^{(i)}(\omega_0 + \underline{\omega} \cdot \underline{\Phi}(x^{(i)})) = 1 \Rightarrow \omega_0 = y^{(i)} - \underline{\omega} \cdot \underline{\Phi}(x^{(i)})$$

This gives $\omega = (\omega_0, \underline{\omega})$ and defines the activation function:

$$\begin{aligned} h(x) &= \sigma(\omega_0 + \underline{\omega} \cdot \underline{\Phi}(x)) \\ &= \sigma(y^{(i)} - \underline{\omega} \cdot \underline{\Phi}(x^{(i)}) + \underline{\omega} \cdot \underline{\Phi}(x)) \\ &= \sigma\left(y^{(i)} - \sum_{j=1}^n \alpha_j y^{(j)} k(x^{(j)}, x^{(i)})\right) \\ &\quad + \sum_{j=1}^n \alpha_j y^{(j)} k(x^{(j)}, x) \end{aligned}$$

LEM: Pos. def. sym. kernel are closed under sum, product, tensor product, point-wise limit and composition with power series $\sum_{n=0}^{\infty} a_n x^n$ for $a_n \geq 0$.

Proof: sum) clear

product) give two pos. def. sym. kernels $K(x, y), K'(x, y)$

def. kernel matrix for $K(x, y)$

$$K_{ij} := K(x_i, x_j), 1 \leq i, j \leq n$$

$\exists R$ s.t. $K = R^* R$ since K pos. semi-def

Then $K(x, y) K'(x, y)$ has kernel matrix

$$\begin{aligned} & K_{ij} K'_{ij} \\ \forall v \in \mathbb{R}^n : & \sum_{i=1}^n \sum_{j=1}^n v_i K_{ij} K'_{ij} v_j \\ &= \sum_{i,j=1}^n \sum_{k=1}^n v_i R_{ik}^* K_{ij} R_{kj} v_j = \sum_{i=1}^n \omega \end{aligned}$$

$$= \sum_{k=1}^n \underbrace{\langle w_k, k' w_k \rangle}_{\geq 0 \text{ because } k' \text{ pos. def. sym.}} \text{ for } (w_k)_j = R_{kj} v_j$$

furthermore sym.

tensor product) k, \hat{k} pos. def. sym. kernels

$$\begin{aligned} \text{Say } k: (x_1, x_1', x_2, x_2') &\rightarrow k(x_1, x_2) \\ \hat{k}: (x_1, x_1', x_2, x_2') &\rightarrow \hat{k}(x_1', x_2') \end{aligned}$$

tensor product

$$k \otimes \hat{k}(x_1, x_1', x_2, x_2') = k(x_1, x_2) \hat{k}(x_1', x_2')$$

then argument as with product.

point-wise limit) Say $(k_n)_{n \in \mathbb{N}}$ pos. def. sym.

kernels and $k_n(x, y) \rightarrow k(x, y)$ then

for all $v \in \mathbb{R}^n$

$$\langle v, k_n v \rangle \geq 0 \Rightarrow \lim_{n \rightarrow \infty} \langle v, k_n v \rangle \geq 0$$

and symmetry likewise

power series) Say $\sum_{n=0}^{\infty} a_n x^n$ converges for

conv. radius $\rho > 0$ and suppose

$|k(x, y)| < \rho \forall x, y \in \mathbb{R}^n$, then

$\sum_{n=0}^N a_n k^n$ pos. def. sym. by products) (sum) and $a_n \geq 0$ and $\sum_{n=0}^{\infty} a_n k^n$ by limit). \square

Normalization of kernels

for any kernel k we def the

normalized kernel

$$\hat{k}(x, y) := \begin{cases} 0 & \text{if } k(x, x) = 0 \\ & \text{or } k(y, y) = 0 \\ \frac{k(x, y)}{\sqrt{k(x, x) k(y, y)}} & \text{otherwise} \end{cases}$$

LEM: Let k be pos. def. sym.

\Rightarrow the corresponding \hat{k} is pos. def. sym.

Proof: Let $\{x_1, \dots, x_n\} \subset \mathbb{R}^n, c \in \mathbb{R}^n$.

$$\text{If } k(x_i, x_i) = 0 \Rightarrow k(x_i, x_j) = 0 \\ \Rightarrow \hat{k}(x_i, x_j) = 0.$$

Hence, we only need to treat $k(x_i, x_i) > 0$.

Then

$$\sum_{i,j=1}^n c_i \frac{k(x_i, x_j)}{\sqrt{k(x_i, x_i) k(x_j, x_j)}} c_j \\ = \sum_{i,j=1}^n c_i \frac{\langle \phi(x_i), \phi(x_j) \rangle}{\|\phi(x_i)\| \|\phi(x_j)\|} c_j \\ = \left\| \sum_{i=1}^n c_i \frac{\phi(x_i)}{\|\phi(x_i)\|} \right\|^2 \geq 0$$

Symmetry follows from the definition. \square

Example (Gaussian kernel)

$$\hat{k}(x, y) = \exp\left(-\frac{|x-y|^2}{2\sigma^2}\right), \sigma \geq 0$$

Let us first consider

$$k(x, y) = \exp\left(+\frac{x \cdot y}{\sigma^2}\right) \\ = \lim_{N \rightarrow \infty} \sum_{n=0}^N \frac{(x \cdot y)^{2n}}{n! \sigma^{2n}}$$

is a pos. def. sym. kernel since $x \cdot y$ is and LEM product, sum, power series).

$$\hat{k}(x, y) = \frac{k(x, y)}{\sqrt{k(x, x) k(y, y)}} = \frac{e^{\frac{x \cdot y}{\sigma^2}}}{e^{-\frac{x^2}{2\sigma^2}} e^{-\frac{y^2}{2\sigma^2}}} \\ = \exp\left(-\frac{|x-y|^2}{2\sigma^2}\right)$$