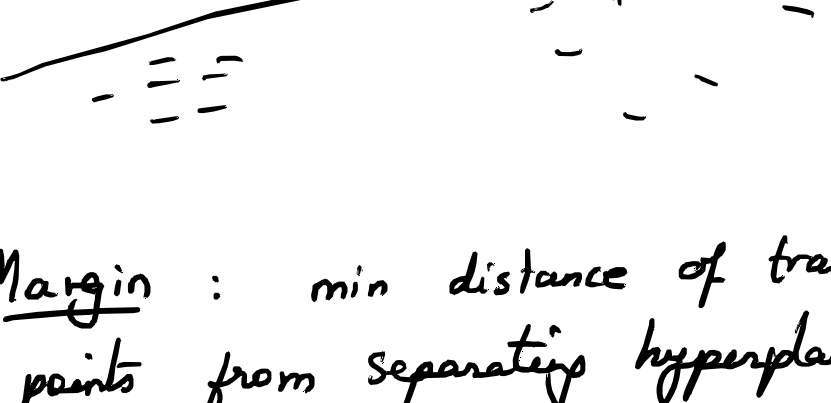


Perceptron

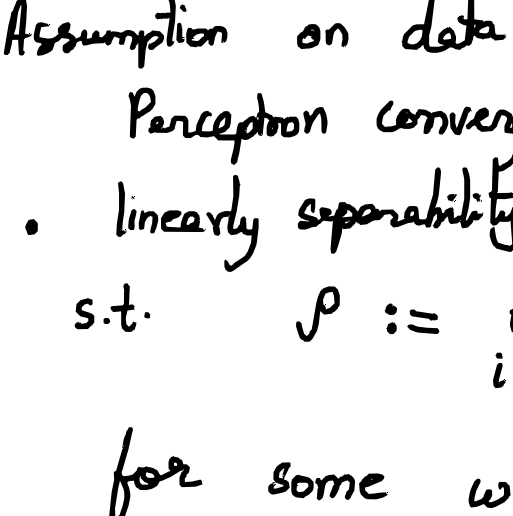
Binary classification under linear separability



Margin: min distance of training points from separating hyperplane

$$d(p, H) = \frac{|\omega \cdot p + b|}{\|\omega\|}$$

$$\rho_{\text{margin}} = \min_i \frac{y_i (\omega \cdot x_i + b)}{\|\omega\|}$$



$$[n] = \{1, \dots, n\} \quad x_i \in \mathbb{R}^d \quad y_i = \pm 1 \quad i=1, \dots, n$$

Assumption on data for theorem on Perceptron convergence:

• linearly separability: $\exists \rho > 0$

$$\text{s.t. } \rho := \min_{i \in [n]} \frac{y_i \omega \cdot x_i}{\|\omega\|}$$

for some $\omega \in \mathbb{R}^d$

• bounded data: $R > 0$ s.t.

$$\|x_i\| < R \quad \forall i.$$

Thm. Perceptron algorithm:

$$\omega_{i+1} = \begin{cases} \omega_i + \eta y_i x_i & \text{if } x_i \in I \text{ (incorrectly classified)} \\ \omega_i & \text{if } x_i \notin I \text{ (correct)} \end{cases}$$

makes at most $\frac{R^2}{\rho^2}$ updates.

Proof: Mohri et al 2018

Interlude:

Classifier returned by the Perceptron algorithm:

$$h(x) = \text{sgn}(\omega_T \cdot x) = \text{sgn}\left(\left(\sum_{i \in I} \eta y_i x_i\right) \cdot x\right)$$

WLOG: $\omega_0 = 0 \in \mathbb{R}^d$

linear comb. of incorrectly classified points

$$= \text{sgn}\left(\sum_{i \in I} \eta y_i \underbrace{(x_i \cdot x)}_{\text{dot product}}\right)$$

Dot product algorithms can be lifted to a higher-dimensional or infinite-dimensional spaces (kernel methods)

$$x \in \mathbb{R}^d \rightarrow \Phi(x) \in \mathbb{R}^D \quad D \gg d$$

$$\text{or } \Phi(x) \in \mathcal{H} \text{ (Hilbert space)}$$

features

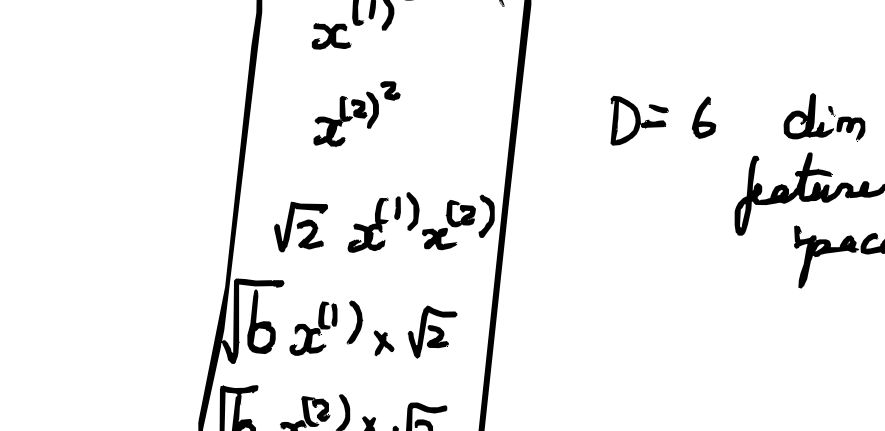
$$\text{if } \langle \Phi(x), \Phi(x') \rangle = \kappa(x, x') \text{ (kernel / similarity measure)}$$

then,

$$h(x) = \text{sgn}\left(\sum_{i \in I} \eta y_i \langle \Phi(x_i), \Phi(x) \rangle\right)$$

$$= \text{sgn}\left(\sum_{i \in I} \eta y_i \kappa(x_i, x)\right)$$

Example (XOR)



$$x^{(1)}(x_1) = 1 \quad x^{(2)}(x_1) = 1$$

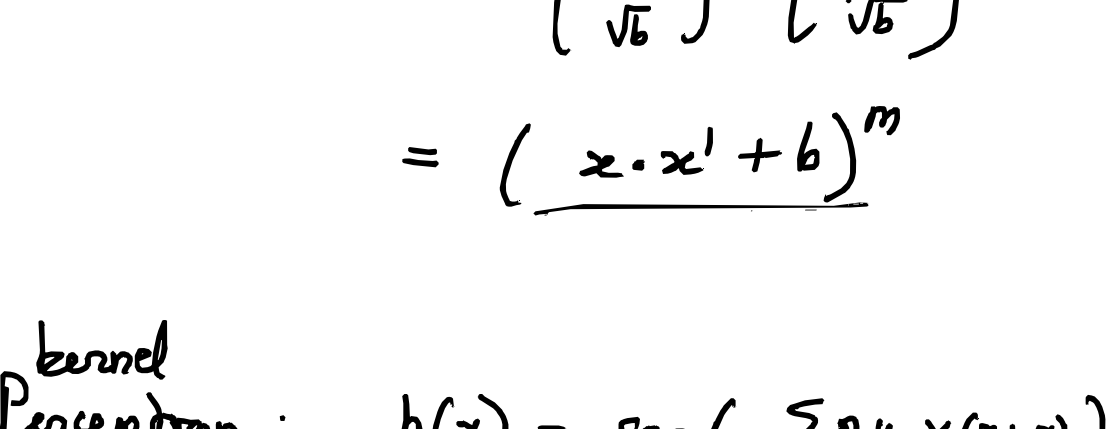
$$x^{(2)}(x_3) = -1$$

$$\begin{bmatrix} x^{(1)2}(x) \\ x^{(2)2}(x) \\ \sqrt{2} x^{(1)} x^{(2)} \end{bmatrix} \cdot \begin{bmatrix} x^{(1)2}(x') \\ x^{(2)2}(x') \\ \sqrt{2} x^{(1)} x^{(2)}(x') \end{bmatrix} = \Phi(x) \cdot \Phi(x')$$

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

$$x^{(1)2}(x) x^{(1)2}(x') + x^{(2)2}(x) x^{(2)2}(x') + 2 x^{(1)}(x) x^{(2)}(x) x^{(1)}(x') x^{(2)}(x')$$

$$= (x \cdot x')^2$$



Polynomial kernel: $x, x' \in \mathbb{R}^d$

$$\kappa(x, x') = (x \cdot x' + b)^m$$

Can be written as a dot product of features

$$b \neq 0, \quad d = 2, \quad D = 6, \quad m = 2 \quad (\text{dim of feature space})$$

$$(x \cdot x' + b)^2 = \left(\sum_{i=1}^d x^{(i)}(x) x^{(i)}(x') + b\right)^2$$

$$\begin{bmatrix} x^{(1)2}(x) \\ x^{(2)2}(x) \\ \sqrt{2} x^{(1)} x^{(2)} \\ \sqrt{b} x^{(1)} \times \sqrt{b} \\ \sqrt{b} x^{(2)} \times \sqrt{b} \\ \sqrt{b} \end{bmatrix}$$

D=6 dim feature space.

$$(x^{(1)}(x) x^{(1)}(x') + x^{(2)}(x) x^{(2)}(x') + b)^2$$

$$= \sum_{i=1}^2 x^{(i)2}(x) x^{(i)2}(x') + b^2 + 2 b x^{(1)}(x) x^{(1)}(x') + 2 b x^{(2)}(x) x^{(2)}(x')$$

$$= \Phi(x) \cdot \Phi(x').$$

To evaluate $\Phi(x) \cdot \Phi(x')$ (which consists of non-linear functions), we only need $\kappa(x, x')$

$$\kappa(x, x') = \Phi(x) \cdot \Phi(x')$$

$$= \begin{bmatrix} x^{(1)2}(x) \\ x^{(2)2}(x) \\ \sqrt{2} x^{(1)} x^{(2)} \\ \sqrt{b} x^{(1)} \\ \sqrt{b} x^{(2)} \\ \sqrt{b} \end{bmatrix} \cdot \begin{bmatrix} x^{(1)2}(x') \\ x^{(2)2}(x') \\ \sqrt{2} x^{(1)} x^{(2)}(x') \\ \sqrt{b} x^{(1)} \\ \sqrt{b} x^{(2)} \\ \sqrt{b} \end{bmatrix}$$

$$= (x \cdot x' + b)^m$$

kernel Perceptron: $h(x) = \text{sgn}\left(\sum \eta y_i \kappa(x_i, x)\right)$

Kernel evaluation \Leftrightarrow inner product on feature space

• does not require explicit knowledge of features

• computationally expensive

• have to know features

$$[\omega \quad b] \begin{bmatrix} x^{(1)}(x) \\ x^{(2)}(x) \\ 1 \end{bmatrix} = 0$$

$$\omega^{(1)} x^{(1)}(x) + \omega^{(2)} x^{(2)}(x) + b = 0$$

$$x^{(2)}(x) = -\frac{\omega^{(1)}}{\omega^{(2)}} x^{(1)}(x) - \frac{b}{\omega^{(2)}}$$