

- Kernel Methods

- SVMs

$$\begin{array}{c} \theta^T \phi(x) \\ \uparrow \text{input} \\ \phi: \mathbb{R}^d \rightarrow \mathbb{R}^P \quad \text{feature map} \\ \text{attributes} \quad \text{features} \\ \phi(x): " \text{features}" \end{array}$$

p: very high $p > n$ or ∞

LMS using gradient descent

Loop {

$$\theta := \theta + \alpha \sum_{i=1}^n (y^{(i)} - \underbrace{\theta^T \phi(x^{(i)})}_{\mathbb{R}^P}) \underbrace{\phi(x^{(i)})}_{\mathbb{R}^P}$$

$O(np)$ time

Key observation:

$$\theta = \sum_{i=1}^n \beta_i \phi(x^{(i)}) \quad \text{for some } \underbrace{\beta_1, \dots, \beta_n}_{\beta \in \mathbb{R}^n} \in \mathbb{R}$$

New algo: update β

p parameters \rightarrow n parameters

$$\beta_i := \beta_i + \alpha \left(y^{(i)} - \sum_{j=1}^n \beta_j \underbrace{\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle}_{p} \right)$$

① Precompute $\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$

$$\langle a, b \rangle = \sum_{i=1}^p a_i b_i$$

② $\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$ can be computed faster than by explicitly computing $\phi(\cdot)$

Cubic polynomials

$$\Phi(x) = \begin{bmatrix} 1 \\ x_i \\ x_i x_j \\ x_i x_j x_k \end{bmatrix} \quad \begin{array}{l} 1 \\ \} d \\ \} d^2 \\ \} d^3 \end{array}$$

$$\langle \Phi(x), \Phi(z) \rangle = [1 \dots x_i \dots x_i x_j \dots x_i x_j x_k \dots] \begin{bmatrix} 1 \\ z_i \\ z_i z_j \\ z_i z_j z_k \end{bmatrix}$$

$$= 1 + \sum_{i=1}^d x_i z_i + \sum_{i,j=1}^d x_i x_j z_i z_j + \sum_{i,j,k=1}^d x_i x_j x_k z_i z_j z_k$$

$$\sum_{i,j=1}^d u_i w_j = \left(\sum_{i=1}^d u_i \right) \left(\sum_{j=1}^d w_j \right)$$

$$u_i \rightarrow x_i z_i \quad w_j \rightarrow x_j z_j$$

$$= \left(\sum_{i=1}^d x_i z_i \right) \left(\sum_{j=1}^d x_j z_j \right) = \langle x, z \rangle^2$$

$$= \left(\sum_{i=1}^d x_i z_i \right) \left(\sum_{j=1}^d x_j z_j \right) \left(\sum_{k=1}^d x_k z_k \right)$$

$$= \langle x, z \rangle^3 \quad O(d) \text{ time}$$

$$\langle \Phi(x), \Phi(z) \rangle = 1 + \langle x, z \rangle + \langle x, z \rangle^2 + \langle x, z \rangle^3$$

$$O(d) \text{ time} \quad p = 1 + d + d^2 + d^3$$

$$= \Theta(d^3)$$

$$K(x, z) = \langle \Phi(x), \Phi(z) \rangle$$

$K(\cdot, \cdot)$ is Kernel function

$$K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$$

Compute $K(x^{(i)}, x^{(j)}) \neq i, j$

n^2 entries $O(n^2 p)$ $O(n^2 d)$ time

$$\beta = 0$$

$$\text{Loop } \left\{ \begin{array}{l} \beta_i := \beta_i + \alpha \left(y^{(i)} - \sum_{j=1}^n \beta_j \langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle \right) \\ \quad = \beta_i + \alpha \left(y^{(i)} - \sum_{j=1}^n \beta_j K(x^{(i)}, x^{(j)}) \right) \end{array} \right.$$

}

$K \in \mathbb{R}^{n \times n}$ Kernel matrix

$$K_{ij} = K(x^{(i)}, x^{(j)})$$

$$\beta := \beta + \alpha (\vec{y} - K\beta) \quad O(n^2) \text{ time}$$

Test time: given x , predict $\theta^T \phi(x)$

$$\theta^T \phi(x) = \left(\sum_{i=1}^n \beta_i \phi(x^{(i)}) \right)^T \phi(x)$$

$$= \sum_{i=1}^n \beta_i \langle \phi(x^{(i)}), \phi(x) \rangle = \sum_{i=1}^n \beta_i \cdot K(x^{(i)}, x)$$

Linear in # examples, independent of p

Training: Preprocessing: $O(n^2 d)$

training: $O(n^2) \times \# \text{ iterations}$

Test time: $O(nd)$ assuming $K(\cdot, \cdot)$ can be computed in $O(d)$ time

Deeper Observation

- the only thing needed is $K(\cdot, \cdot)$

function $K(\cdot, \cdot)$ is valid (Kernel fn.)

if $\exists \phi$ s.t. $K(x, z) = \langle \phi(x), \phi(z) \rangle$

Design some $K(\cdot, \cdot)$

Verify validity (by math)

run algo

Other algos can be "Kernelized"

perception, logistic regression

- algo for linear $\theta^T x$

- replace x by $\phi(x)$

- rewrite algo s.t. it only depends on $\langle \phi(x), \phi(z) \rangle$

Kernel fn's:

$$K(x, z) = 1 + x^T z + (x^T z)^2 + (x^T z)^3$$

$$K(x, z) = (x^T z)^2$$

$$K(x, z) = (x^T z + c)^2$$

$$\phi(x) = \begin{bmatrix} c \\ \sqrt{2c} \cdot x_i \\ \vdots \\ x_i \cdot x_j \end{bmatrix}$$

polynomial Kernel $K(x, z) = (x^T z + c)^k$ $\sim \binom{d+k}{k}$ monomials

$$K(x, z) = \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right) = \langle \phi(x), \phi(z) \rangle$$

ϕ : ∞ dimensional

Valid Kernel?

Necessary condn

n examples $x^{(1)}, \dots, x^{(n)}$

Kernel matrix $K_{ij} = K(x^{(i)}, x^{(j)})$

Claim: Kernel matrix is positive semidefinite

$$K \succeq 0$$

$$z^T K z \geq 0 \quad \forall z \in \mathbb{R}^n$$

$K(\cdot, \cdot)$ valid if $\exists \Phi$ s.t. $K_{ij} = \langle \Phi(x^{(i)}), \Phi(x^{(j)}) \rangle$

$$Z^T K Z = \sum_{i=1}^n \sum_{j=1}^n z_i K_{ij} z_j$$

$$= \sum_{i=1}^n \sum_{j=1}^n z_i \langle \Phi(x^{(i)}), \Phi(x^{(j)}) \rangle z_j$$

$$= \left\langle \sum_{i=1}^n z_i \Phi(x^{(i)}), \sum_{j=1}^n z_j \Phi(x^{(j)}) \right\rangle$$

$$\geq 0$$

also sufficient

Theorem (Mercer, 1909)

$K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a valid kernel fn

if for any $n < \infty$ and any $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^d$

the kernel matrix $K_{ij} := K(x^{(i)}, x^{(j)})$

is positive semidefinite



e.g. Protein sequence classification

20 amino acids

A, B, C, ...,

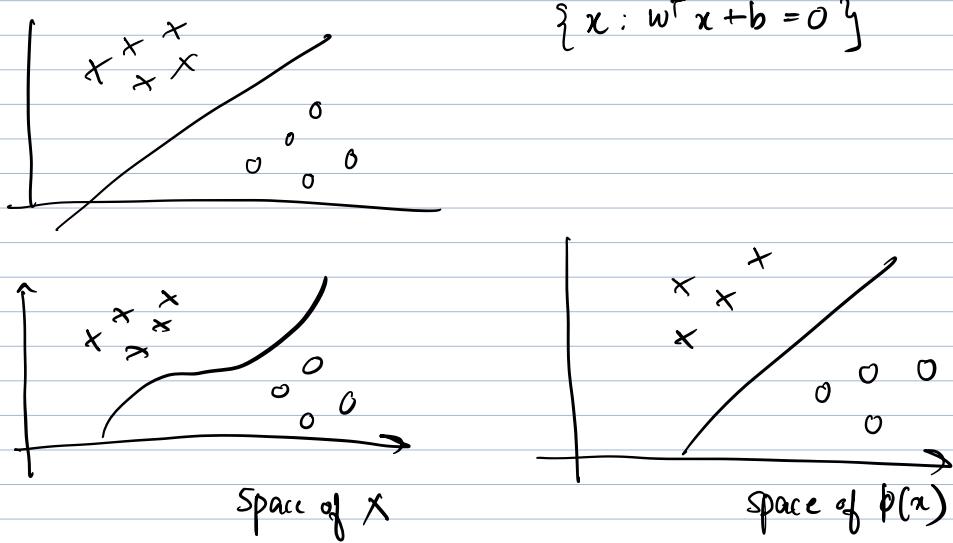
A B A C, B D A, ...

$$\Phi(x) = \begin{bmatrix} 0 \\ 2 \\ \vdots \\ \vdots \end{bmatrix} \begin{array}{c} \text{AAAA} \\ \text{AAAB} \\ \vdots \\ \vdots \end{array}$$

20^4 dimensional vector

$\langle \Phi(x), \Phi(z) \rangle$ can be computed via dynamic programming

SVM: support vector machines

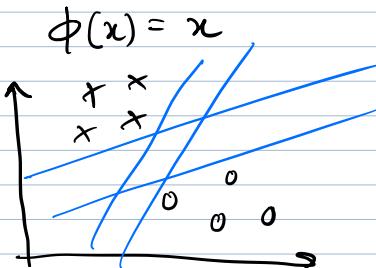


$$y^{(i)} \in \{-1, +1\}$$

$$\{x : w^T x + b = 0\}$$

$$\{x : w^T \phi(x) + b = 0\}$$

linear in kernel space

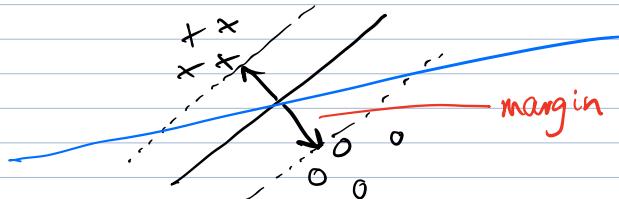


Find w, b s.t.

$$\text{If } y^{(i)} = 1, \quad w^T x^{(i)} + b > 0 \quad ①$$

$$y^{(i)} = -1, \quad w^T x^{(i)} + b < 0 \quad ②$$

choose (w, b) that give the most separation



amongst all (w, b) pairs that satisfy ①, ②

$$\max_{w, b} \left[\min_i \text{dist}(x^{(i)}, \text{decision boundary}) \right]$$

$$\text{①, ②} \Rightarrow y^{(i)}(w^T x^{(i)} + b) > 0 \quad \forall i$$

$$\begin{aligned} \text{Fact: } \text{dist}(x^{(i)}, \text{boundary}) &= \frac{|w^T x^{(i)} + b|}{\|w\|_2} \\ &= \frac{y^{(i)}(w^T x^{(i)} + b)}{\|w\|_2} \end{aligned}$$

$$\max_{w, b} \min_{i \in \{1, \dots, n\}} \frac{y^{(i)}(w^T x^{(i)} + b)}{\|w\|_2} \equiv \frac{\Theta^T x}{\|w\|_2}$$

scaling invariant $(w, b) \rightarrow (10w, 10b)$

$$\Leftrightarrow \min \frac{1}{2} \|w\|_2^2$$

$$\text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1 \quad \forall i$$

Facts (non trivial) (need KKT condn)

① Optimal soln w^*, b^* satisfies

$$w^* = \sum_{i=1}^n \alpha_i x^{(i)} y^{(i)} \quad \alpha_i \geq 0 \quad \alpha_i \in \mathbb{R}$$

② $\alpha = (\alpha_1, \dots, \alpha_n)$ is optimal soln of program

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

$$\text{s.t. } \alpha_i \geq 0$$

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

$$w^* = \sum_{i=1}^n \alpha_i \phi(x^{(i)}) y^{(i)}$$

$$\begin{aligned} \text{test time: } w^{*T} \phi(z) &= \sum_{i=1}^n \alpha_i \langle \phi(x^{(i)}), \phi(z) \rangle y^{(i)} \\ &= \sum_{i=1}^n \alpha_i K(x^{(i)}, z) y^{(i)} \end{aligned}$$

Random Features, Random Kitchen Sinks

Rahimi, Recht [NeurIPS 2007]

Fast Food [ICML 2013]

Le, Sarlos, Jordan