

- Kernel Methods

- SVMs

$$\theta^T x$$

$\uparrow$   
input

$$\phi : \mathbb{R}^d \rightarrow \mathbb{R}^p$$

attributes      features

$$\phi(x) : \text{"features"}$$

p: very high    p > n    p  $\infty$

LMS using gradient descent

Loop {

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

}

each iteration:  $O(n p)$  time

Key observation:

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

$$\theta \in \mathbb{R}^p$$

New algo: update  $\beta$

$$\theta = \sum_{i=1}^n (\underbrace{\beta_i + \alpha (y^{(i)} - \theta^T \phi(x^{(i)}))}_{\text{new } \beta_i}) \phi(x^{(i)})$$

p parameters  $\rightarrow$  n parameters

$$\begin{aligned}
 \beta_i &:= \beta_i + \alpha (y^{(i)} - \theta^T \phi(x^{(i)})) \\
 &= \beta_i + \alpha (y^{(i)} - (\sum_{j=1}^n \beta_j \phi(x^{(j)}))^T \phi(x^{(i)})) \\
 &= \beta_i + \alpha (y^{(i)} - \underbrace{\sum_{j=1}^n \beta_j}_{\text{n}} \underbrace{\langle \phi(x^{(j)}), \phi(x^{(i)}) \rangle}_P)
 \end{aligned}$$

Precompute

$$\begin{aligned}
 ① \quad &\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle \\
 &\langle a, b \rangle = \sum_{i=1}^p a_i b_i
 \end{aligned}
 \qquad \qquad \qquad
 \begin{aligned}
 a &= (a_1, \dots, a_p) \\
 b &= (b_1, \dots, b_p)
 \end{aligned}$$

②  $\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$  can often be computed much faster without explicitly computing  $\phi(\cdot)$

e.g. Cubic polynomials

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

$$\begin{aligned}
 \langle \phi(x), \phi(z) \rangle &= [1 \dots x_i \dots x_i x_j \dots x_i x_j x_k \dots] \begin{bmatrix} 1 \\ z_i \\ z_i z_j \\ z_i z_j z_k \end{bmatrix} \\
 &= 1 + \sum_{i=1}^d x_i z_i + \sum_{i,j=1}^d x_i x_j z_i z_j + \sum_{i,j,k=1}^d x_i x_j x_k z_i z_j z_k
 \end{aligned}$$

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

$$\begin{aligned}
 u_i &\rightarrow x_i z_i, \quad w_j \rightarrow x_j z_j \\
 &= \left( \sum_{i=1}^d x_i z_i \right) \left( \sum_{j=1}^d x_j z_j \right) = \langle x, z \rangle^2
 \end{aligned}$$

$\mathcal{O}(d)$  time

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

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

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

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

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

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

Mercer Kernels

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

$$\text{Compute } K(x^{(i)}, x^{(j)}) \quad \forall i, j$$

$$n^2 \text{ entries} \quad O(n^2 p) \quad O(n^2 d) \text{ time}$$

$$\beta = 0$$

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

$$= \beta_i + \alpha \left( y^{(i)} - \sum_{j=1}^n \beta_j \cdot K(x^{(i)}, x^{(j)}) \right) \}$$

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

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

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

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

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

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

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 you need is  $K(\cdot, \cdot)$  function  $K(\cdot, \cdot)$  is valid (Kernel fn)  
 $\nexists \phi$  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 also be "Kernelized"

perceptron, 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 fns:

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

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

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

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

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

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

$\phi$   $\infty$  dimensional

Valid Kernel?

Necessary cond'

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

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

Claim: Kernel matrix is positive semidefinite

$$K \succeq 0$$

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

Also sufficient

Theorem (Mercer, 1909)

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

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, ...

$\phi(x)$

AAAA  
AAA B

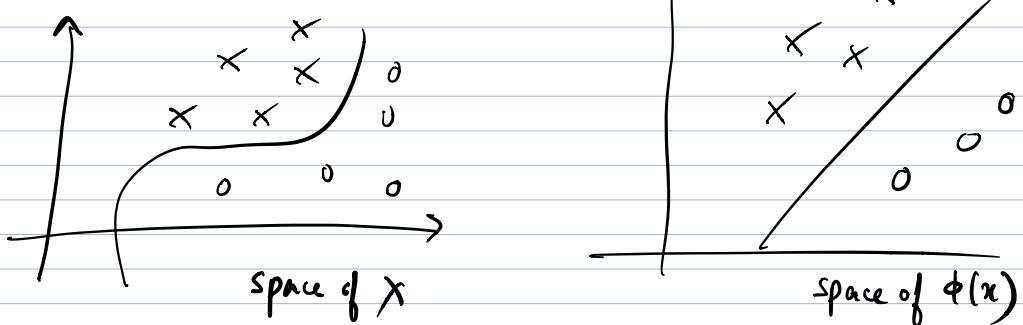
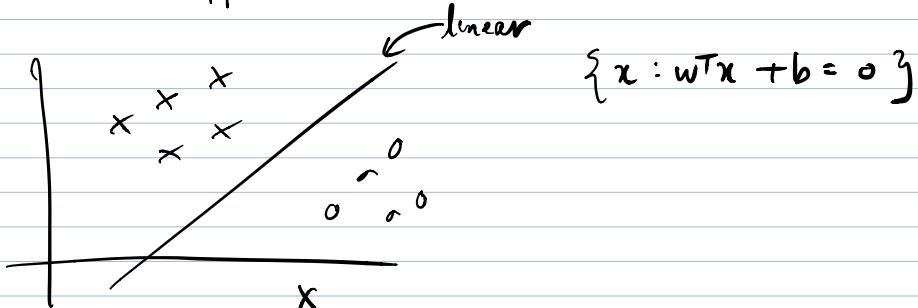
|

$$20^4 = 16,000$$

$$\text{AAAA} \begin{bmatrix} 1 \\ 2 \\ 0 \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

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

SVM: Support vector machines



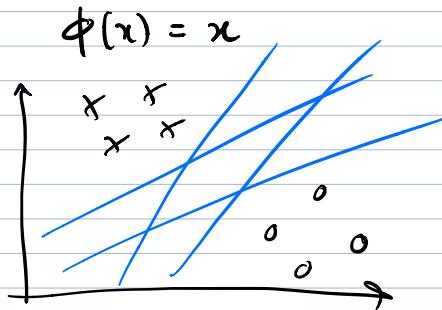
$$y^{(i)} \in \{-1, +1\}$$

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

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

linear in Kernel space

Fund  $w, b$



Find  $w, b$  s.t.

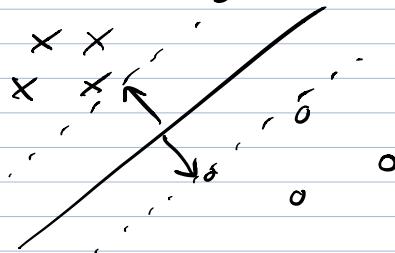
$$\text{if } y^{(i)} = 1, \quad w^T x^{(i)} + b > 0$$

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

①

②

choose the  $(w, b)$  that gives the most separation



Among all  $(w, b)$  satisfy ①, ②

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

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

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

$$= \frac{y^{(i)} (w^T x^{(i)} + b)}{\|w\|_2}$$

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

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 sol<sup>n</sup>  $w^*, b^*$  satisfies

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

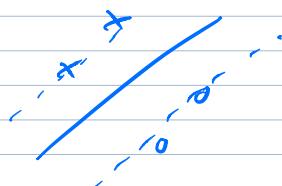
②  $\alpha$  is optimal sol<sup>n</sup> of program

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

$K(x^{(i)}, x^{(j)})$

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

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



$$w^* = \underbrace{\sum_{i=1}^n \alpha_i \phi(x^{(i)}) y^{(i)}}_{\text{test time:}}$$

test time:

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

$$= \sum_{i=1}^n \alpha_i K(x^{(i)}, x) y^{(i)}$$