

Kernel Methods

Terminology

$$\Phi: \mathbb{R}^d \rightarrow \mathbb{R}^p \quad \text{feature map}$$

\uparrow attributes $\quad \uparrow$ features
 x $\quad \quad \quad \Phi(x)$

What to do if p is very large?

$$\Phi(x) = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \\ x_1^2 \\ \vdots \\ x_d^2 \\ x_1 x_2 \\ \vdots \\ x_1 x_d \\ \vdots \\ x_d^3 \end{bmatrix} \quad \left. \begin{array}{l} \text{ } \\ \text{ } \\ \text{ } \end{array} \right\} \begin{array}{l} d \\ d^2 \\ d^3 \end{array}$$

$$\Theta^T \Phi(x) = -1 + x_1 + x_2 + x_1 x_2 + x_1 x_3 + x_2 x_3$$

LMS using kernel trick

$$\min_{\Theta} \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \Theta^T \Phi(x^{(i)}))^2$$

Gradient Descent Loop:

$$\Theta := \Theta + \alpha \sum_{i=1}^n (y^{(i)} - \Theta^T \Phi(x^{(i)})) \Phi(x^{(i)})$$

Problem: $\Phi(x)$ is high dimensional

$$p = 1 + d + d^2 + d^3 \quad O(d^3)$$

$$d \sim 10^3 \Rightarrow p \sim 10^9$$

Runtime for 1 iteration of GD is $O(np)$

Key observation:

If Θ initialized as 0, then at any time,

$$\Theta = \sum_{i=1}^n \beta_i \Phi(x^{(i)}) \quad \text{for } \beta_1, \dots, \beta_n \in \mathbb{R}$$

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

New algo: update β

$$\beta_1 := \beta_i + \alpha(y^{(i)} - \underbrace{\sum_{j=1}^n \beta_j}_{\hat{y}} \underbrace{\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle}_{\hat{K}_{ij}})$$

① Precompute $\langle \Phi(x^{(i)}), \Phi(x^{(j)}) \rangle$
 $\langle a, b \rangle = \sum_{i=1}^n a_i b_i$

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

e.g. cubic polynomials

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

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

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

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

$$q = \left(\sum_{i=1}^d x_i z_i \right) \left(\sum_{j=1}^d x_j z_j \right) \left(\sum_{k=1}^d x_k z_k \right) = \langle x, z \rangle^3 \quad O(d) \text{ time}$$

$$\langle \phi(x), \phi(z) \rangle = 1 + \langle x, z \rangle + \langle x, z \rangle^2 + \langle x, z \rangle^3 \quad O(d) \text{ time}$$

In general:

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

$k(\cdot, \cdot)$ is a kernel function

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

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

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

Loop {

$$\beta_i := \beta_i + \alpha(y^{(i)} - \sum_{j=1}^n \beta_j \langle \Phi(x^{(i)}), \Phi(x^{(j)}) \rangle)$$

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

}

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

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

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

Test time

Given x , predict $\Theta^T \Phi(x)$

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

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

$$= \sum_{i=1}^n \beta_i k(x^{(i)}, x)$$

linear in #exs
independent of p

Time complexity

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

Training $O(n^2) \times \# \text{ iters}$

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

Deeper observation

→ Only thing needed is $k(\cdot, \cdot)$

function k is valid kernel fn

if $\exists \Phi. k(x, z) = \langle \Phi(x), \Phi(z) \rangle$

Other algos can also be kernelized

e.g. perceptron, LR

→ algo for linear $\Theta^T x$

→ replace x by $\Phi(x)$

→ rewrite algo s.t. only depends on $\langle \Phi(x), \Phi(z) \rangle$

Kernel fns:

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

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

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

$$\Phi(x) = \begin{bmatrix} c \\ \sqrt{2c} x_1 \\ \vdots \\ x_1 x_i \end{bmatrix}$$

polynomial kernel

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

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

Φ : ∞ -dimensional

Kernel validity

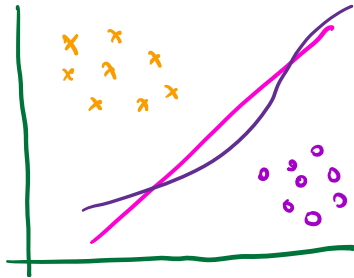
Necessary condition:

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

kernel matrix $k_{ij} = k(x^{(i)}, x^{(j)})$

kernel matrix is PSD ($k \succeq 0; z^T k z \geq 0 \quad \forall z \in \mathbb{R}^n$)

Support Vector Machines (SVM)



$\{x: w^T x + b = 0\}$ by way
kernel methods w/ feature map

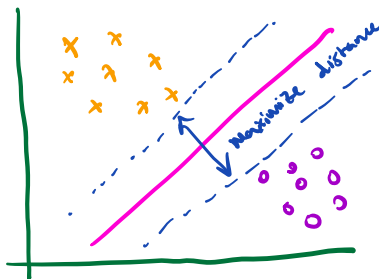
with SVM: let $y^{(i)} \in \{-1, 1\}$

Find w, b s.t.

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

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

Choose w, b that give most separation



amongst all (w, b) pairs that satisfy (1), (2) $(y^{(i)}(w^T x^{(i)} + b) > 0 \forall i)$

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

$$= \max_{w, b} \left[\min_i \frac{|w^T x^{(i)} + b|}{\|w\|_2} \right]$$

$$= \max_{w, b} \left[\min_{i \in \{1, \dots, n\}} \frac{y^{(i)}(w^T x^{(i)} + b)}{\|w\|_2} \right] \leftarrow \text{objective function}$$

Note: scaling invariant

Non-trivial facts (need KKT condition)

① optimal soln w^* , b^* satisfies

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

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

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

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

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

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

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