

Outline:

- Kernel Methods
- SVM

Feature map (Recap)

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

$$x \rightarrow \phi(x)$$

fitting $h_\theta(x) = \theta^T \phi(x)$ using gradient descent.

Issue: $\theta, \phi(x) \in \mathbb{R}^p$ can be very high dimensional
runtime per iteration $O(np)$

Goal: improve to $O(n^2)$ per iteration

Kernel Methods

Key observation:

θ can be represented as

$$\theta = \sum_{i=1}^n \beta_i \phi(x_i)$$

scalar (n variables instead of p)

Proof by induction

At iteration 0

$$\theta = 0 = \sum_{i=1}^n 0 \cdot \phi(x^{(i)})$$

$$\text{Suppose at iteration } t, \theta = \sum_{i=1}^n \beta_i \phi(x^{(i)})$$

Next iteration

$$\theta := \theta + \alpha \sum (y^{(i)} - \theta^T \phi(x^{(i)})) \phi(x^{(i)})$$

$$\theta := \sum_{i=1}^n (\beta_i + \alpha \sum_j (y^{(i)} - \theta^T \phi(x^{(i)})) \phi(x^{(i)}))$$

New β_i

Represent $\theta \in \mathbb{R}^p$ implicitly by $B \in \mathbb{R}^n$
works better when $p \gg n$

Update rule for β

$$\beta_i = \beta_i + \alpha (y^{(i)} - \theta^T \phi(x^{(i)}))$$

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

$$= \sum_j \beta_j \phi(x^{(j)})^T$$

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

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

Note: $a^T b = \langle a, b \rangle$

Observations

① $\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$ $\forall i, j$ can be precomputed

② Often $\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$ can be computed faster than $O(p)$ ex: $O(d)$

$$\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \\ x_1 x_1 \\ \vdots \\ x_d x_d \\ \vdots \\ x_d^3 \end{bmatrix}$$

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

$$\begin{aligned} \langle \phi(x), \phi(z) \rangle &= 1 + \sum_{i=1}^d x_i z_i + \sum_{i=1}^d \sum_{j=1}^d x_i x_j z_i z_j + \sum_i \sum_j \sum_k x_i x_j x_k z_i z_j z_k \\ &= 1 + \underbrace{\langle x, z \rangle}_{O(d)} + \langle x, z \rangle^2 + \langle x, z \rangle^3 \end{aligned}$$

Can compute $\langle \phi(x), \phi(z) \rangle$ in $O(d)$ time!

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

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

Algorithm for computing Kernel

Compute $K(x^{(i)}, x^{(j)}) = \langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$
for all $i, j \in \{1, \dots, n\}$

Set $\beta = 0$
 $\beta_i := \beta_i + \alpha (y^{(i)} - \sum_{j=1}^n \beta_j K(x^{(i)}, x^{(j)}))$

preprocessing: $O(n^2 d)$

Each iteration: $O(n^2)$ per iteration

Prediction: Given x , compute $\theta^T \phi(x)$

$$\begin{aligned} \theta^T \phi(x) &= \sum_{i=1}^n \beta_i \phi(x^{(i)})^T \phi(x) \\ &= \sum_{i=1}^n \beta_i K(x^{(i)}, x) \end{aligned}$$

$O(nd)$ time

Deeper Observations

The algo only depends on $K(\cdot, \cdot)$, no need to know ϕ

* design of features \rightarrow design of kernel functions

What Kernels are valid?

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

Necessary conditions:

n data points. $x^{(1)}, \dots, x^{(n)}$

Kernel matrix $K \in \mathbb{R}^{n \times n}$ $K_{ij} = K(x^{(i)}, x^{(j)})$

matrix K is positive semidefinite

$$K \succeq 0$$

$$\boxed{K_{ij}}$$

$$K \succeq 0 \Leftrightarrow z^T K z \geq 0$$

$$K_{ij} = \langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$$

$$\phi(x) = (\phi_1(x), \phi_2(x), \dots)$$

$$z^T K z = \sum_i \sum_j z_i K_{ij} z_j$$

$$= \sum_i \sum_j z_i \phi(x^{(i)})^T \phi(x^{(j)}) z_j$$

$$= \sum_i \sum_j z_i \sum_k \phi_k(x^{(i)}) \phi_k(x^{(j)}) z_j$$

$$= \sum_l \left(\sum_i z_i \phi_l(x^{(i)}) \right)^2 \geq 0$$

$$\Rightarrow K \succeq 0$$

Theorem (Mercer) K is a valid kernel fn (i.e $K(x, z) = \phi(x)^T \phi(z)$)
iff for any $n < \infty$ and $x^{(1)}, \dots, x^{(n)}$

the corresponding Kernel matrix K s.t. $K_{ij} = K(x^{(i)}, x^{(j)})$
is positive semidefinite

Other Kernels:

$$\bullet K(x, z) = (x^T z + c)^2 = \langle \phi(x), \phi(z) \rangle$$

$$\text{for } \phi(x) = \begin{bmatrix} \frac{c}{\sqrt{2c}} x_1 \\ \vdots \\ \frac{c}{\sqrt{2c}} x_d \\ x^2 \\ \vdots \\ x^{d^2} \end{bmatrix}$$

$$\bullet K(x, z) = (x^T z + c)^t$$

$$\text{Gaussian Kernel: } K(x, z) = \exp\left(\frac{-\|x - z\|^2}{2\sigma^2}\right)$$

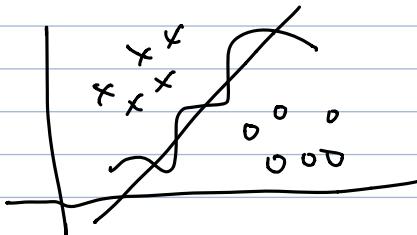
$$K(x, z) = \langle \phi(x), \phi(z) \rangle \Rightarrow \phi \text{ is actually infinite dimensional!}$$

Example: Protein sequence classification
seq of amino acids (A-T)

$$\begin{bmatrix} \text{AAAA} \\ \text{AAAB} \\ \vdots \\ \text{TTT} \end{bmatrix} \rightarrow 20^4 \text{ dim vector}$$

histograms: take min and sum up.

SVM for classification



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

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

SVM

$$\phi(x) = x \quad (\text{for now})$$

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

Warmups

$$\begin{aligned} & \text{find } w, b \text{ s.t.} \\ \text{if } y^{(i)} = 1 & \quad w^T x^{(i)} + b > 0 \quad (1) \\ y^{(i)} = -1 & \quad w^T x^{(i)} + b < 0 \quad (2) \end{aligned}$$

Many such w, b.

New goal: Among all w, b , satisfies ① & ②.

$$\text{Find } w, b \text{ s.t. } \max_{w, b} \min_{i \in \{1, \dots, n\}} \text{dist}(x^i, \text{boundary})$$

$$= \begin{cases} \frac{\mathbf{w}^T \mathbf{x}^{(i)} + b}{\|\mathbf{w}\|_2} & \text{if } \mathbf{x}^{(i)} \text{ is on the positive side} \\ -\frac{(\mathbf{w}^T \mathbf{x}^{(i)} + b)}{\|\mathbf{w}\|_2} & \text{if } \mathbf{w}^T \mathbf{x}^{(i)} + b < 0 \end{cases}$$

Scale invariant $(w, b) \rightarrow (100w, 100b)$

We want to find w, b

$$\text{s.t. } \min_{i \in \{1, n\}} y^{(i)} (w^T x^{(i)} + b) = 1 \quad (3)$$

$$\max \frac{1}{\|w\|_2} = \min \|w\|_2$$

$$\min \|w\|_2 \quad \text{equivalent}$$

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

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

(Convex!)

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

\nwarrow linear combination

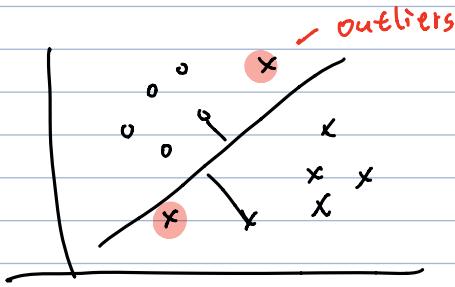
Facts: (i) optimal solution $w^* = \sum_{i=1}^h \alpha_i y^{(i)} x^{(i)}$ for some $\alpha_i \geq 0$ —④

(ii) The α in ④ is the optimizer of

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

s.t. $\alpha_i \geq 0, \sum_i \alpha_i y^{(i)} = 0$

Kernelize: replace $\langle x^{(i)}, x^{(j)} \rangle$ with $K(x^{(i)}, x^{(j)})$
n variables.



$$y^{(i)} (w^T x^{(i)} + b) \geq 1 - \epsilon_i \quad (\text{relax})$$

$$\min \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \epsilon_i$$

\nwarrow pay some penalty