

Lecture 7 Kernels

1. < continue > Optimal Margin Classifier

$$x^{(i)} \in \mathbb{R}$$

Suppose $w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$ (w can be represented as a linear combination of x)
Why reasonable? < representation theorem >

- Intuition 1: (refer to logistic regression)

$$\theta_0 := 0$$

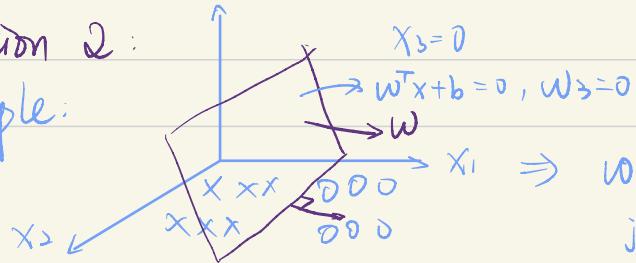
(stochastic) gradient descent: $\theta := \theta - \alpha (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$

(batch) gradient descent: $\theta := \theta - \alpha \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$

No matter iteration, still linear combination

- Intuition 2:

example:



w should be represented in the span of just feature x_1 and x_2 .

$$\begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \Rightarrow \begin{bmatrix} b \\ w_1 \\ \vdots \\ w_n \end{bmatrix}$$

$$\therefore w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$$

$$\min_{w,b} \frac{1}{2} \|w\|^2 \text{ s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1, i=1,2,\dots,m$$

$$\Downarrow \quad \Downarrow w^T w \quad \quad \quad \min \frac{1}{2} \left(\sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \right)^T \left(\sum_{j=1}^m \alpha_j y^{(j)} x^{(j)} \right)$$

$$y^{(i)} \left[\left(\sum_{j=1}^m \alpha_j y^{(j)} x^{(j)} \right)^T x^{(i)} + b \right] \geq 1$$

$$= \min \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y^{(i)} y^{(j)} \underbrace{x^{(i)T} x^{(j)}}_{\langle x, z \rangle = x^T z \text{ is the inner product}}$$

$$y^{(i)} \left(\sum_{j=1}^m \alpha_j y^{(j)} \underbrace{\langle x^{(j)}, x^{(i)} \rangle}_{\text{inner product}} + b \right) \geq 1$$

$\langle x, z \rangle = x^T z$ is the inner product $\langle x^{(i)}, x^{(j)} \rangle$ efficient

< a further simplification > : "Dual optimization Problem"

$$\max \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle \text{ s.t. } \alpha_i \geq 0, \sum_i y^{(i)} \alpha_i = 0$$

To predict :

1) Solve for α_i 's; b

2) Compute $h_w.b(x) = g(w^T x + b)$

$$= g\left((\sum_i \alpha_i y^{(i)} x^{(i)})^T x + b\right) = g(\sum_i \alpha_i y^{(i)} \langle x^{(i)}, x \rangle, x \rangle + b)$$

2. kernel trick:

- ① Write algorithm in terms of $\langle x^{(i)}, x^{(j)} \rangle$ (or $\langle x, z \rangle$)
- ② Let there be mapping from $\underset{2D}{x} \mapsto \phi(x)$ high-dimensional
- ③ Find way to compute $k(x, z) = \phi(x)^T \phi(z)$
- ④ Replace $\langle x, z \rangle$ in algorithm with $k(x, z)$

example:

$$x \in \mathbb{R}^n \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ \vdots \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix} \quad \phi(z) = \begin{bmatrix} z_1 z_1 \\ z_1 z_2 \\ z_1 z_3 \\ \vdots \\ z_3 z_1 \\ z_3 z_2 \\ z_3 z_3 \end{bmatrix}$$

↑ n^2 elements
↓ n^2 elements

Need $O(n^2)$ time to compute $\phi(x)$, or $\phi(x)^T \phi(z)$ explicitly

$k(x, z) = \phi(x)^T \phi(z) \quad \underline{\text{prone}} \quad (x^T z) \in \mathbb{R}^n \quad O(n)$ time

$$\Rightarrow \phi(x)^T \phi(z) = \left(\sum_{i=1}^n x_i z_i \right) \left(\sum_{j=1}^n x_j z_j \right)$$

$$= \sum_{i=1}^n \sum_{j=1}^n x_i z_i x_j z_j = \sum_{i=1}^n \sum_{j=1}^n (x_i x_j) (z_i z_j) = (x^T z)^2$$

3. Optimal Margin Classifier + kernel trick = SVM

4. How to make kernels?

| If x, z are "similar", $k(x, z) = \phi(x)^T \phi(z)$ is "large"
| If x, z are "dissimilar", $k(x, z)$ is "small"

= Does the kernel exists ϕ s.t. $k(x, z) = \phi(x)^T \phi(z)$

$$k(x, x) = \phi(x)^T \phi(x) \geq 0$$

Let $\{x^{(1)}, x^{(2)}, \dots, x^{(d)}\}$ be d points

Let $K \in \mathbb{R}^{d \times d}$, $k_{ij} = k(x^{(i)}, x^{(j)})$
↑ kernel matrix

Given any vector \mathbf{z} ,

$$\begin{aligned}\mathbf{z}^\top \mathbf{k} \mathbf{z} &= \sum_i \sum_j z_i k_{ij} z_j = \sum_i \sum_j z_i \phi(x^{(i)})^\top \phi(x^{(j)}) z_j \\ &= \sum_i \sum_j z_i \sum_k (\phi(x^{(i)}))_k (\phi(x^{(j)}))_k z_j \\ &= \sum_k \sum_i \sum_j z_i (\phi(x^{(i)}))_k (\phi(x^{(j)}))_k z_j \\ &= \sum_k \left(\sum_i z_i \phi(x^{(i)})_k \right)^2 \geq 0\end{aligned}$$

So $\mathbf{k} \geq 0$

Mercer's Theorem: \mathbf{k} is a valid kernel function (i.e. $\exists \phi$ s.t. $k(x, z) = \phi(x)^\top \phi(z)$) if and only if for any d points $\{x^{(1)}, \dots, x^{(d)}\}$, the corresponding kernel matrix $\mathbf{k} \geq 0$.

5. Other kernels:

Linear kernel: $k(x, z) = x^T z$; $\phi(x) = x$

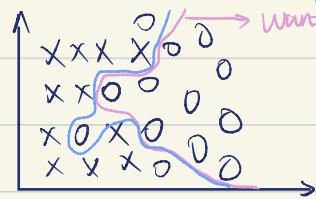
Gaussian kernel: $k(x, z) = \exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right)$; $\phi(x) \in \mathbb{R}^\infty$

Other Kernel: $k(x, z) = (x^T z + C)^d$ e.g.

Polynomial kernel: $k(x, z) = (x^T z + C)^d$

$$\begin{bmatrix} x_1 & x_3 \\ x_2 & x_3 \\ \vdots & \vdots \\ x_3 & x_3 \\ \sqrt{2} & x_1 \\ \vdots & \vdots \\ \sqrt{2} & x_3 \end{bmatrix}$$

b. L₁ norm soft margin SVM (when you don't want the algorithm to be so correct)



$$\text{Min } \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i$$

$$\text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1 - \xi_i \quad i=1, 2, \dots, m$$

$$\xi_i \geq 0$$

$$\Rightarrow \text{Max}_{\alpha} W(\alpha) = \sum_{i=1}^n \alpha_i - \sum_{i=1}^m \sum_{j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

$$\text{s.t. } \sum_{i=1}^m y^{(i)} \alpha_i = 0, \quad 0 \leq \alpha_i \leq C, \quad i=1, 2, \dots, m$$