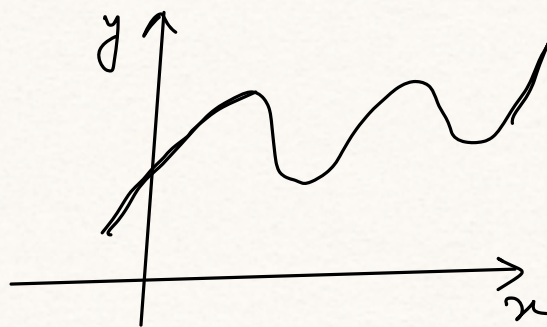
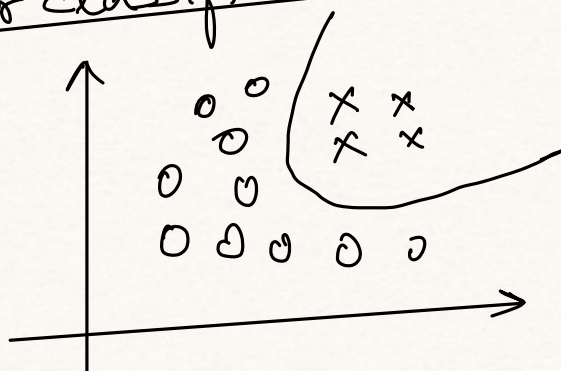


Kernel Methods

$$\underbrace{x}_{\text{attributes}} \rightarrow \underbrace{\phi(x)}_{\text{Features}} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$



For classification:



For L.R in G.D

$$\theta^{(t+1)} = \theta^{(t)} + \alpha \cdot \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)}) \cdot x^{(i)} \quad ; \theta \in \mathbb{R}^d$$

With a feature map,

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

$$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^p \quad ; \theta \in \mathbb{R}^p$$

Eg: $\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_1^2 \\ x_1 x_2 \\ x_1 x_3 \\ \vdots \\ x_1^3 \end{bmatrix}$

→ monomial terms
of order (≤ 3)
 $p \approx O(d^3)$

Obs :

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

$$\theta^{(1)} = \sum_{i=1}^n \underbrace{\alpha y^{(i)}}_{\beta_i^{(1)}} \cdot \phi(x^{(i)}) \rightarrow \text{since } \theta^{(0)} = 0.$$

$$\theta^{(t+1)} = \boxed{\theta^{(t)}} + \alpha \sum_{i=1}^n (y^{(i)} - \boxed{\theta^{(t)}}^T \phi(x^{(i)})) \cdot \phi(x^{(i)})$$

Induceⁿ

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

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

For $i = 1, 2, \dots, n$

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

α and $\alpha = 1/d$

$$x \in \mathcal{X} \quad x \in \mathbb{R}^d \quad \mathcal{X} = \mathbb{R}^d$$

$$\text{Kernel} \triangleq K: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$

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

$$= \phi(x)^T \phi(z)$$

$$\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

$$p \approx O(d^3)$$

$$K(x, z) = 1 + \overset{O(d)}{\langle x, z \rangle} + \overset{O(d)}{\langle x, z \rangle^2} + \overset{O(d)}{\langle x, z \rangle^3}$$

$$= \phi(x)^T \phi(z)$$

Linear Regression (Kernelized)

1. Precompute

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

2. Loop:

$$\forall i \in \{1, 2, \dots, n\}$$

$$\beta_i^{(t+1)} = \beta_i^{(t)} + \alpha \left(y^{(i)} - \sum_{j=1}^n \beta_j^{(t)} K_{ij} \right)$$

$$\beta^{(t+1)} = \beta^{(t)} + \alpha (\bar{y} - K \beta^{(t)})$$

Prediction

$$h_{\theta}(x) = \theta^T \phi(x) \\ = \sum_{i=1}^n \beta_i \phi(x^{(i)})^T \phi(x) \\ = \sum_{i=1}^n \beta_i K(x^{(i)}, x)$$

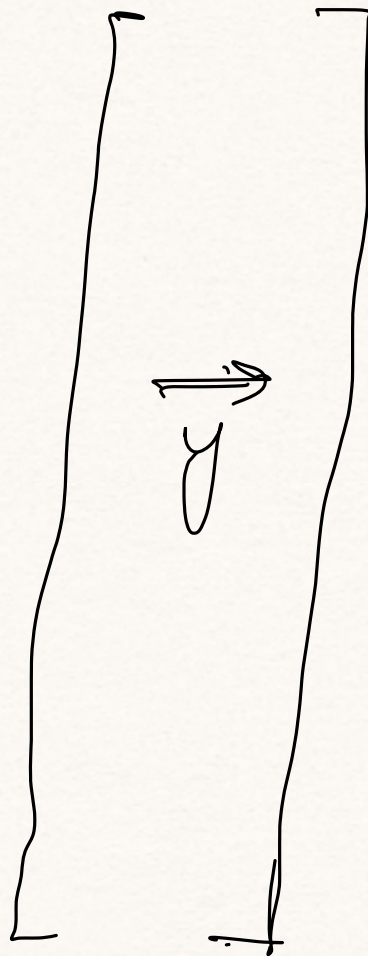
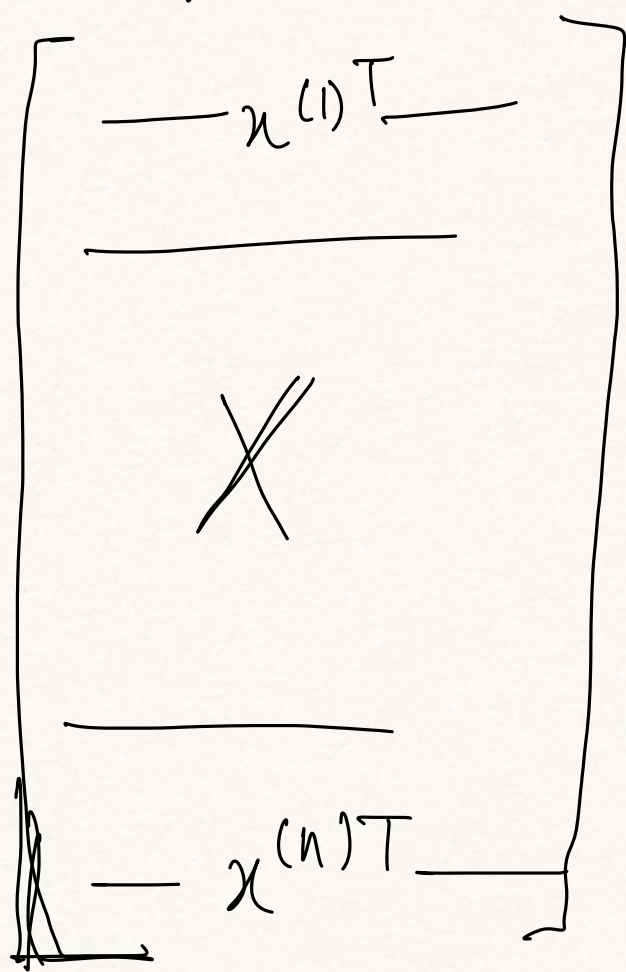
x -test example.

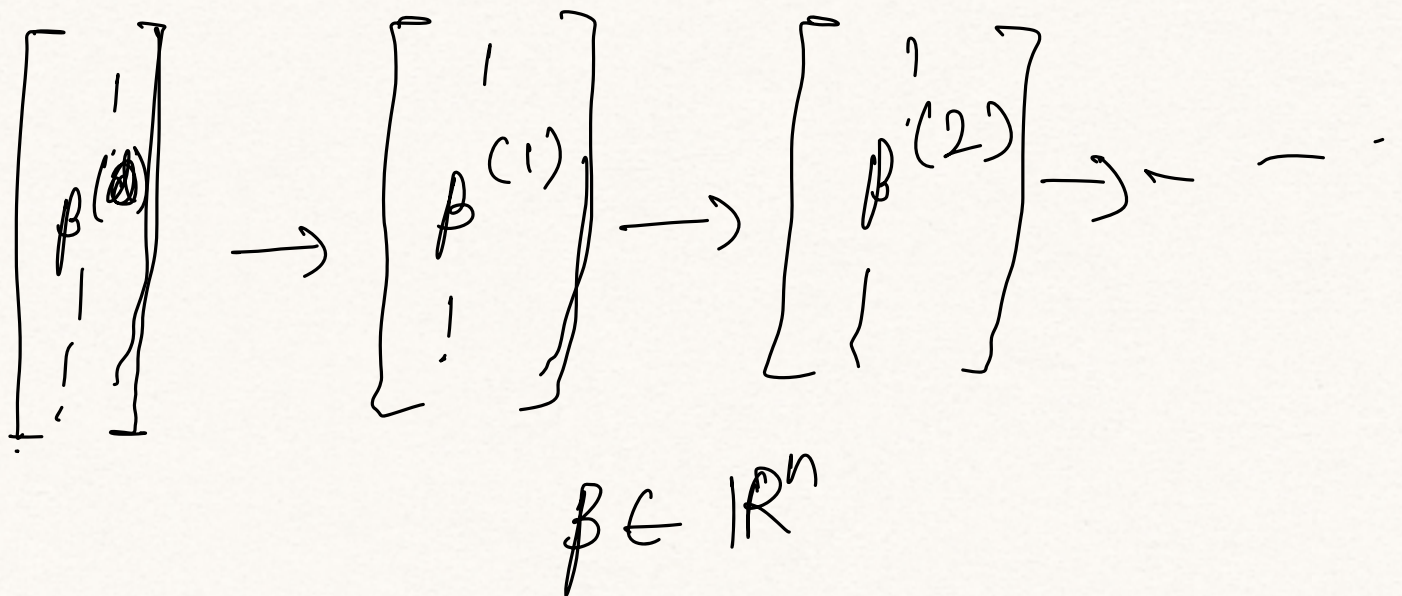
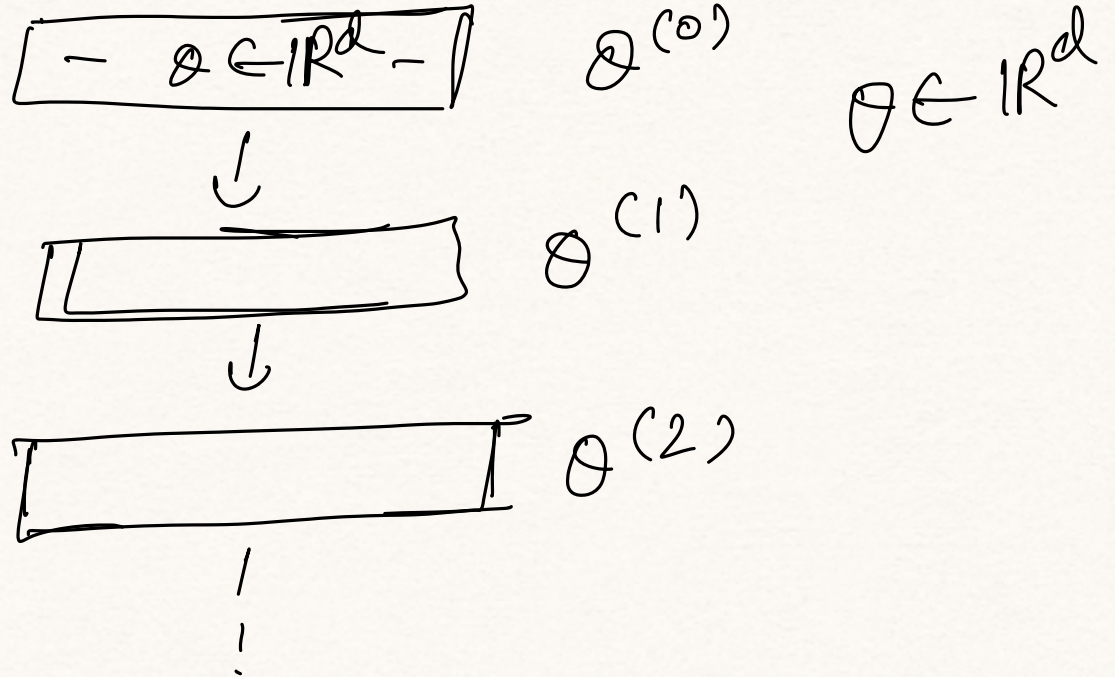
Obs:

$$\left. \begin{array}{l} \text{① Train: } \beta := \beta + \alpha (\vec{y} - K\beta) \\ \text{Test: } \hat{y} = \sum_{i=1}^n \beta_i K(x^{(i)}, x) \end{array} \right\} \phi(x) \text{ does not appear.}$$

② For prediction we need training examples to be stored in memory.

$$x \in \mathbb{R}^d$$





$$x \in \mathbb{R}^d \Rightarrow \phi(x) \in \mathbb{R}^p$$

Kernel Examples

Eg: A) $K(x, z) = \langle x, z \rangle^2$ $x \in \mathbb{R}^d$

$$\phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ \vdots \\ x_d x_d \end{bmatrix}$$

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

$$\phi(x) = \begin{bmatrix} x_1^2 \\ x_1 x_2 \\ \vdots \\ \sqrt{2}c \cdot x_1 \\ \sqrt{2}c \cdot x_2 \\ \vdots \end{bmatrix}$$

$K(x, z) \uparrow$ (for similar x, z)
 \downarrow (for not similar x, z)

$$K(x, z) = \phi(x)^T \phi(z)$$

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

↑ Gaussian Kernel

Necessary conditions for K to be Kernel

- K should be symmetric

$$K(x, z) = K(z, x)$$

$$K(x, z) \triangleq \phi(x)^T \phi(z)$$

- $\{x^{(1)}, \dots, x^{(m)}\}$

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

K is symmetric & P.S.D

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

$$= \dots = \sum_k \left(\sum_i z_i \phi_k(x^{(i)}) \right)^2 \geq 0$$

Mercer's Theorem

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

For K to be a kernel, it is necessary & sufficient for any $\{x^{(1)}, \dots, x^{(m)}\}$ the corresponding kernel matrix $K_{ij} = K(x^{(i)}, x^{(j)})$ is P.S.D & symm.

• To prove K is kernel

① construct ϕ st $K(\cdot) = \phi^T \phi$.

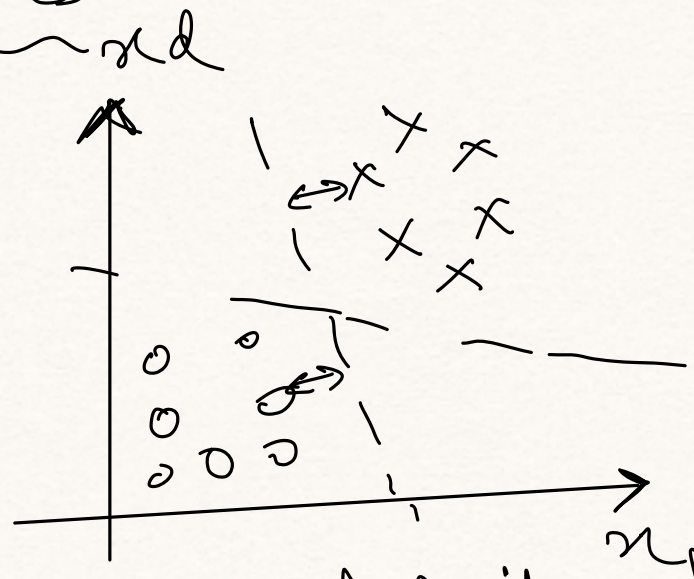
② Mercer's Thm.

$\{x^{(1)}, \dots, x^{(m)}\}$

$K_{ij} = K(x^{(i)}, x^{(j)})$ is P.S.D

Support Vector Machines

- Discriminative
- Classification



Notation

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

parameters: w, b

$$w \in \mathbb{R}^d \quad b \in \mathbb{R}$$

$$x \in \mathbb{R}^d$$

$$w^T x + b$$

$$\text{Functional Margin} = y^{(i)} (w^T x^{(i)} + b) > 0$$

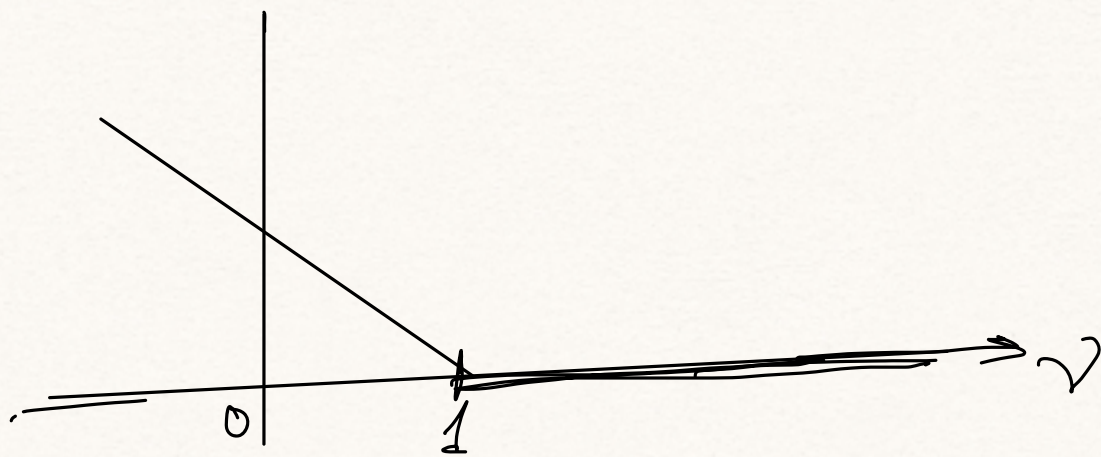
$$\begin{aligned} w^T x + b &> 0 \quad \text{for } y = +1 \\ &< 0 \quad \text{for } y = -1 \end{aligned}$$

Desire: Margin to be large.

→ Maximize the smallest margin.

$$\min_{w, b} \left(\sum_{i=1}^n \max[0, 1 - \underbrace{y^{(i)}(w^T x^{(i)} + b)}_{\text{margin} = \gamma^{(i)}}] \right) + \frac{\|w\|^2}{c}$$

Hinge loss / SVM loss



$$\min_{\xi, w, b} \quad \frac{1}{2} \|w\|^2 + c \sum_{i=1}^n \xi_i$$

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

$$\forall i \in \{1, 2, \dots, n\}$$

$$\xi_i \geq 0 ; i = 1, 2, \dots, n$$

Primal Convex Problem

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

$$\text{s.t.} \quad 0 \leq \alpha_i \leq C$$

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

Dual convex problem