

Primal (Hard) SVM: When  $S = \{ (x^{(i)}, y^{(i)}) \}_{i=1}^n$  is linearly separable

$$\left[ \begin{array}{l} \theta^T x + \theta_0 \geq 0 \\ \theta^T x + \theta_0 < 0 \end{array} \right]$$

$$\min \frac{1}{2} \theta^T \theta$$

$$\text{s.t. } y^{(i)} (\theta^T x^{(i)} + \theta_0) \geq 1 \quad \forall i \in \{1, 2, \dots, n\}$$

Lagrangian:  $\mathcal{L}(\alpha, \theta, \theta_0) = \frac{1}{2} \theta^T \theta + \sum_{i=1}^n \alpha_i (1 - y^{(i)} (\theta^T x^{(i)} + \theta_0))$   
 $(\alpha_i \geq 0)$

$$\text{Primal} = \min_{\theta, \theta_0} \sup_{\alpha} \mathcal{L}(\alpha, \theta, \theta_0)$$

Since  $\mathcal{L}$  is convex

$$\sup_{\alpha} \min_{\theta, \theta_0} \mathcal{L}(\alpha, \theta, \theta_0) = \min_{\theta, \theta_0} \sup_{\alpha} \mathcal{L}(\alpha, \theta, \theta_0)$$

Achieved when  $\theta = \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)}$

Dual (Hard) SVM: Find  $\alpha$  such that

[Taking  $\theta_0 = 0$ ]

General case  $\theta_0 \neq 0$

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

$x^{(i)T} x^{(j)}$  [Inner prod / dot prod]

$$\text{s.t. } \alpha_i \geq 0 \quad \forall i \in \{1, 2, \dots, n\}$$

$$\sum \alpha_i y^{(i)} = 0$$

Prediction on new  $x$ : Compute  $\text{sign} \left( \sum_{i=1}^n \alpha_i y^{(i)} \langle x^{(i)}, x \rangle \right)$

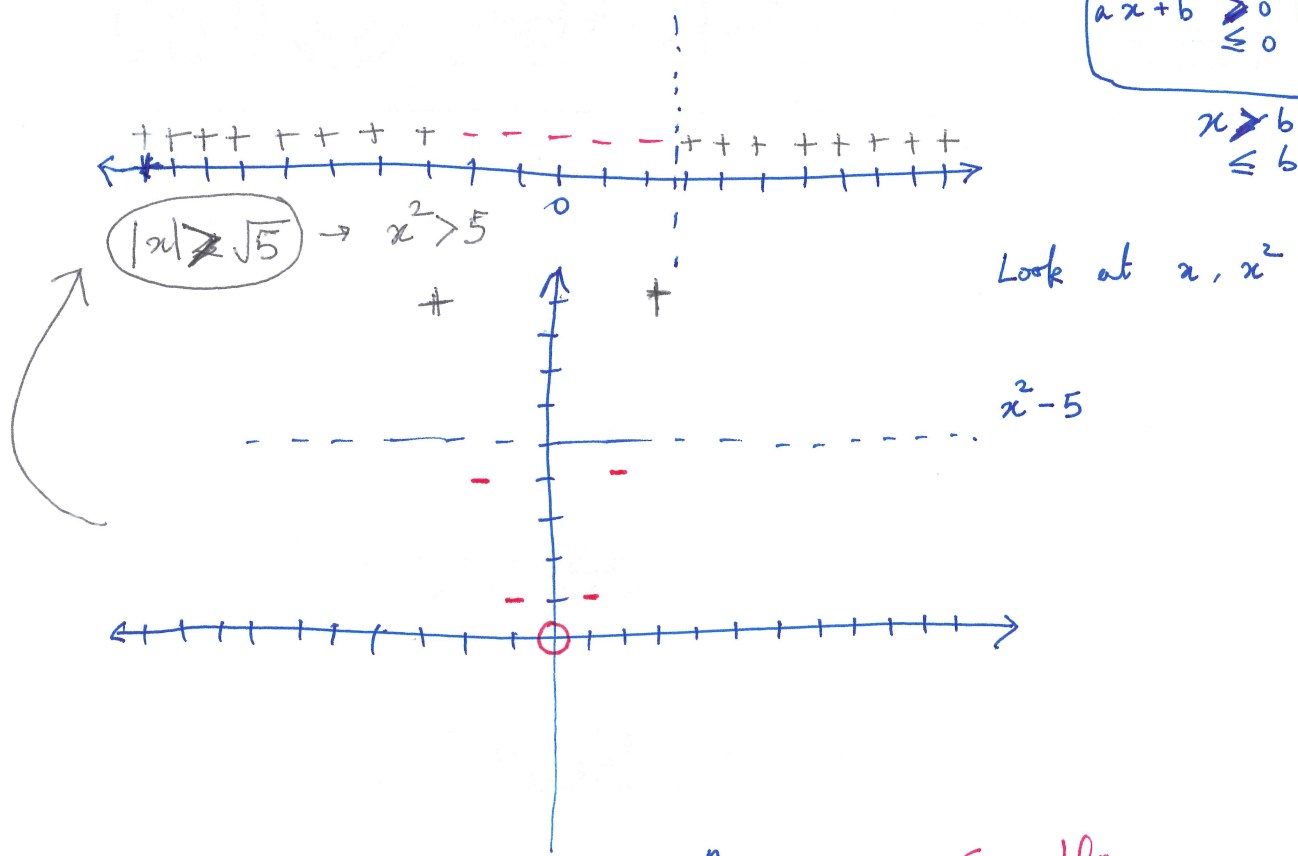
$$\text{sign}(a) = \begin{cases} +1 & \text{when } a > 0 \\ -1 & \text{when } a \leq 0 \end{cases}$$

$$\text{sign}(\theta^T x)$$

Example:  $S = \{(-60, +1), (-9, +1) \dots (-2, -1), (-1, +1), (0, -1) \dots (2, -1), (3, +1) \dots (10, +1)\}$

$$\begin{cases} ax + b \geq 0 \\ \leq 0 \end{cases}$$

$$\begin{cases} x \geq b \\ \leq b \end{cases}$$



Enhance Features:  $S = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$

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

Mapping  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{d_1}$  ( $d_1 \geq d$ )

Example  
 $\phi: \mathbb{R} \rightarrow \mathbb{R}^2$   
 $\phi(x) = (x, x^2)$

② Construct  $\phi(S) = \{(\underbrace{\phi(x^{(i)})}_{\in \mathbb{R}^{d_1}}, y^{(i)})\}_{i=1}^n$

③ Train a linear classifier on  $\phi(S)$ . Find  $\hat{\theta} \in \mathbb{R}^{d_1}, \theta_0 \in \mathbb{R}$ .

④ Prediction on  $x$ :  $\text{sign}(\hat{\theta}^T \phi(x) + \theta_0)$

Quadratic Classifier:

$$x^T = (x_1, x_2 \dots x_d)$$

$$\phi(x)^T = (1, x_1, x_2, \dots, x_d, x_1 x_1, x_1 x_2, x_1 x_3 \dots x_1 x_d, x_2 x_1 \dots x_j x_d)$$

$d^2 + d + 1$

~~Kernel Trick:~~  
~~Fixed~~

SVM in high dimensions: Suppose  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{d_1}$

Find  $\theta \in \mathbb{R}^{d_1}, \theta_0 \in \mathbb{R}$

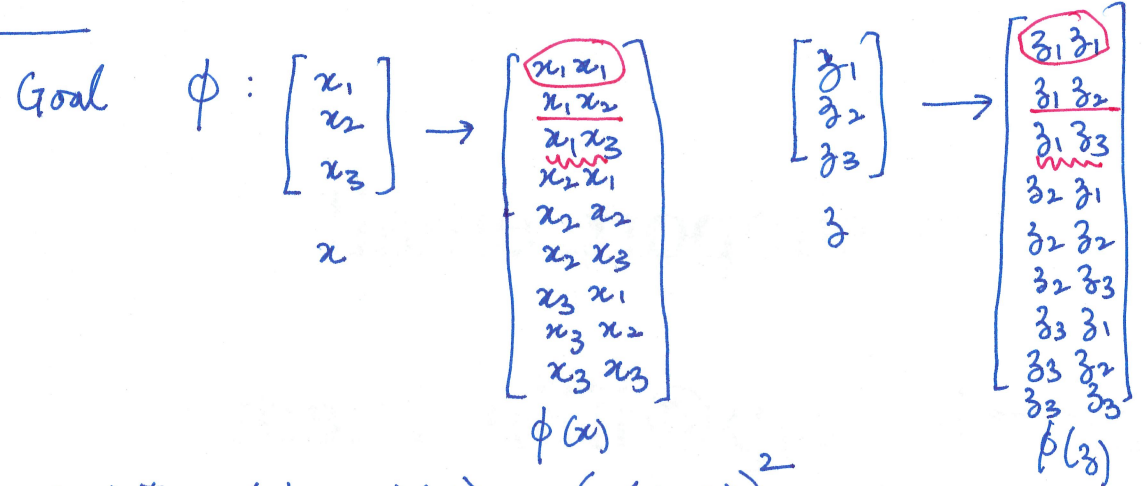
$$\min \frac{1}{2} \theta^T \theta$$
$$\text{s.t. } y^{(i)} (\theta^T \phi(x^{(i)}) + \theta_0) \geq 1$$

$$\text{Optimal } \theta = \sum_{i=1}^n \alpha_i y^{(i)} \phi(x^{(i)})$$

$$\sup_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} \langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$$
$$\text{s.t. } \alpha_i \geq 0$$
$$\sum \alpha_i y^{(i)} = 0$$

[Assume  $\theta_0 = 0$ ]

Prediction  $x$ :  $\text{sign} (\sum \alpha_i y^{(i)} \langle \phi(x^{(i)}), \phi(x) \rangle)$



Compute  $\langle \phi(x), \phi(z) \rangle = (\langle x, z \rangle)^2$

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

$$\phi_1: \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \rightarrow \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ \vdots \\ x_3 x_3 \\ x_1 \\ x_2 \\ \vdots \\ x_3 \\ 1 \end{bmatrix}$$

$$\langle \phi_1(x), \phi_1(z) \rangle = (1 + \langle x, z \rangle)^2 \quad (4)$$

Instead of quadratic features,  
we want all  $k$ -ary monomials  
 $\phi_k: \mathbb{R}^d \rightarrow \mathbb{R}^{d_k}$

$$\langle \phi_k(x), \phi_k(z) \rangle = (1 + \langle x, z \rangle)^k$$

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

Compute  $\sum x_i - \frac{1}{2} \sum \sum x_i x_j y^{(i)} y^{(j)} \left( \frac{x_i x_j}{x_i x_j} \right) K(x^{(i)}, x^{(j)})$   
 $\text{sign}(\sum x_i y^{(i)}) K(x^{(i)}, x_j)$

Gram matrix: Given  $S = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$ ,  $G$  is  $n \times n$  matrix

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

Kernel: Function that measures the similarity between two vectors

$$K(x, z) = \begin{cases} \text{large} & \text{when } x, z \text{ are similar} \\ \text{small} & \text{when } x, z \text{ are dissimilar} \end{cases}$$

RBF or Gaussian Kernel:

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

A function  $K$  is "kernel" if  $\exists \phi$  s.t.  $K(x, z) = \langle \phi(x), \phi(z) \rangle$   
 $\rightarrow$  is a kernel  $\phi$  maps  $d$  vectors to an  $\infty$ -dimensional space.

Kernels: "linear":  $\phi(x) = x$ .  $K(x, z) = x^T z$

"polynomial":  $\phi(x)$  - all monomials up to  $k$ .  $K(x, z) = (1 + x^T z)^k$

"RBF" or "gaussian":  $K(x, z) = e^{-\frac{\|x - z\|^2}{2\sigma^2}}$

Kernel Trick: ① Write your learning algorithm in terms of inner products.

② Replace inner products with kernels.