



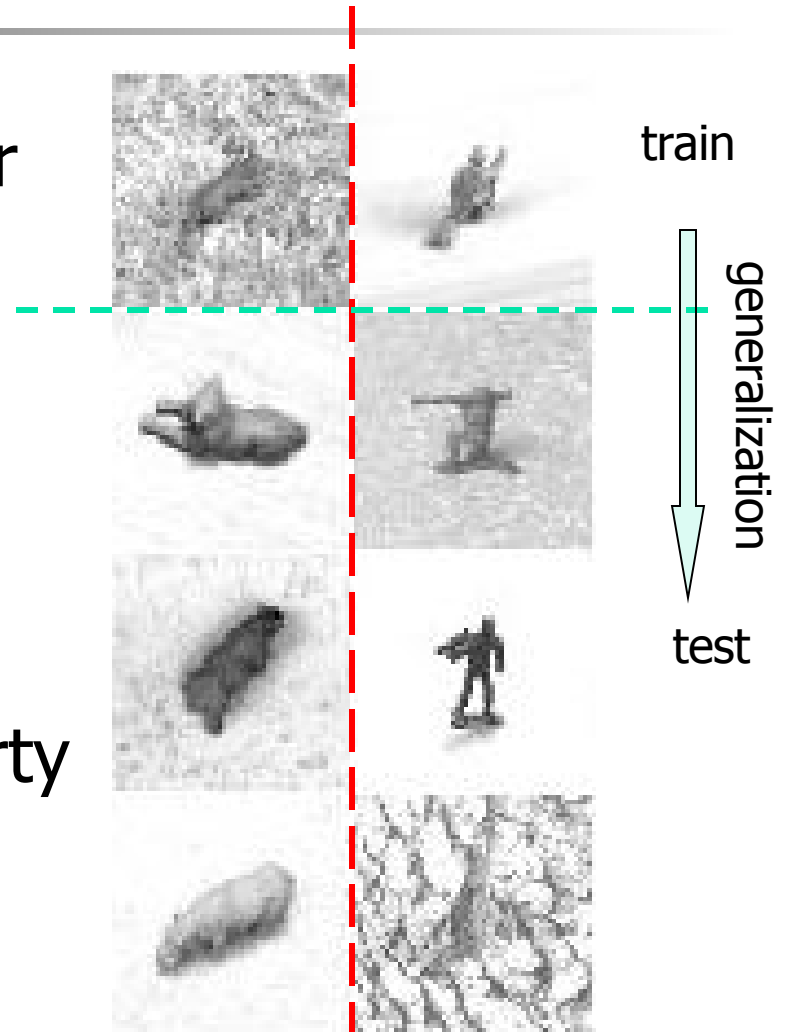
Support Vector Machines

Fu Jie Huang

Dec 5, 2006

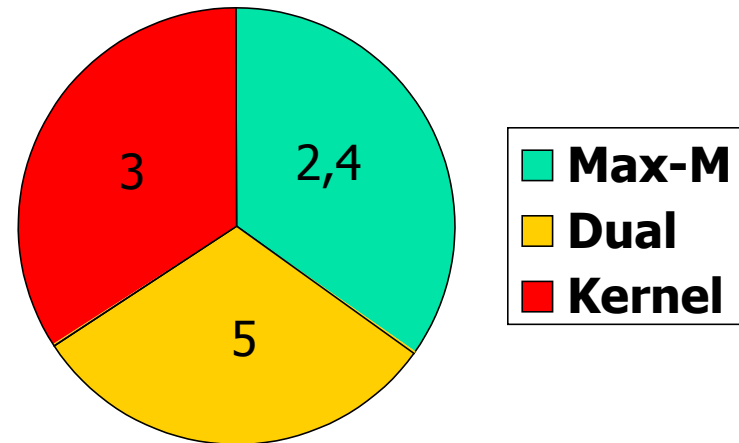
What's SVM

- A (mostly) binary classifier
- A linear classifier
- Supervised training
- Nice generalization property



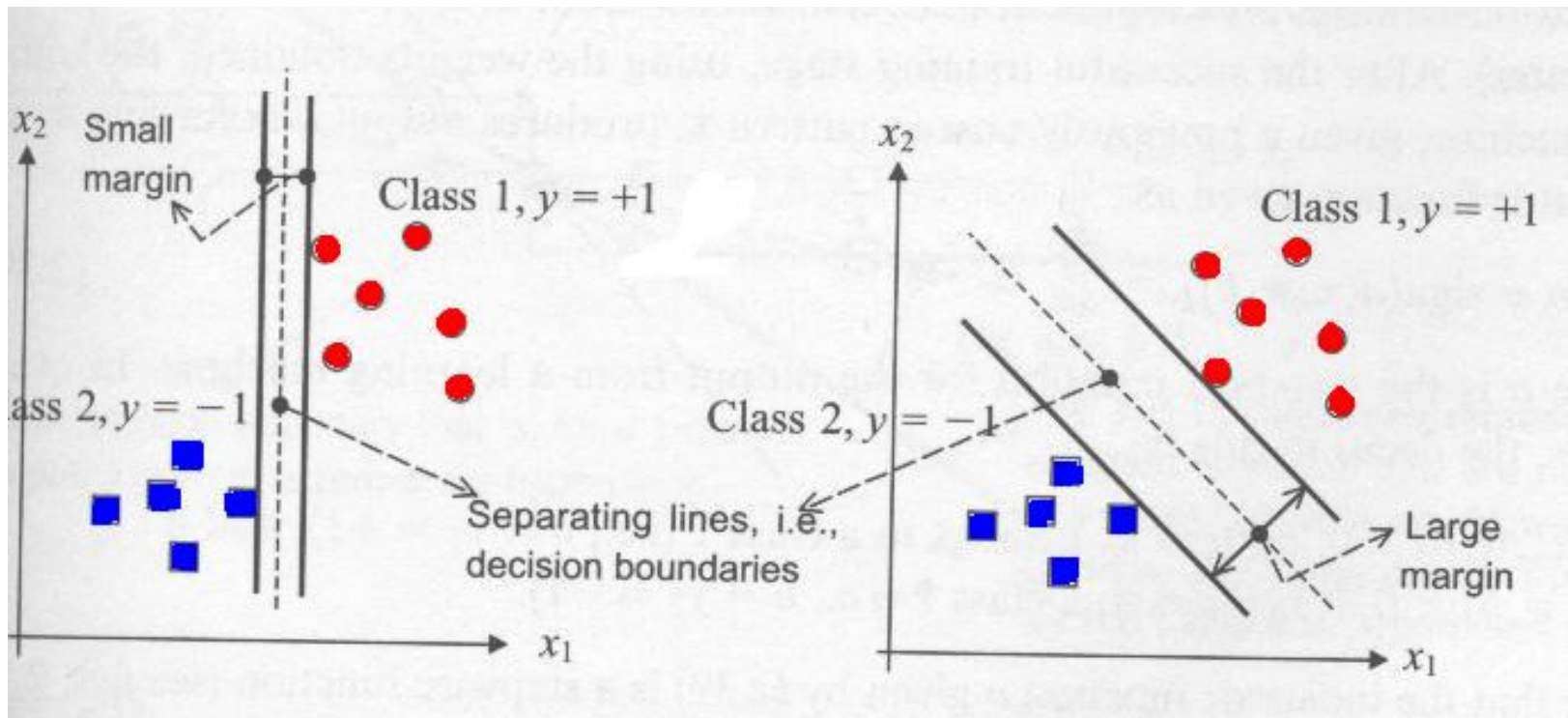
Main Ideas

- Maximal margin
- Dual space
- Kernel trick



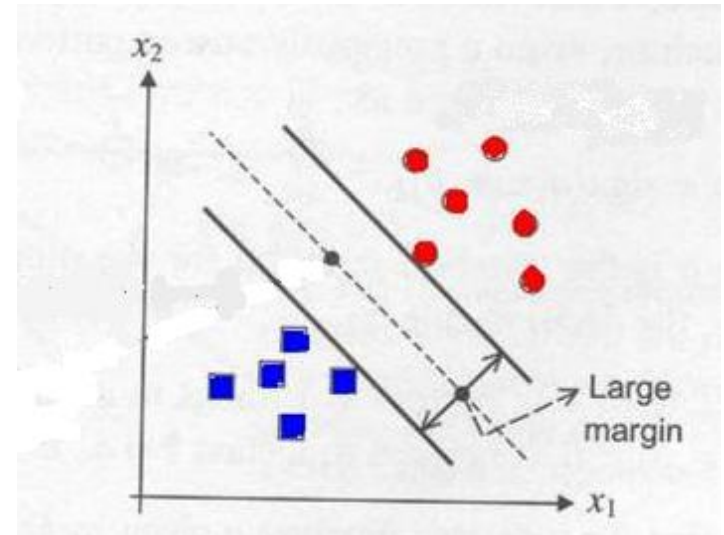
Maximal Margin

- Separating hyperplane is not unique
- Choose the one with... maximal margin



Motivation

- Generalization Error
 - Small
 - Predictable

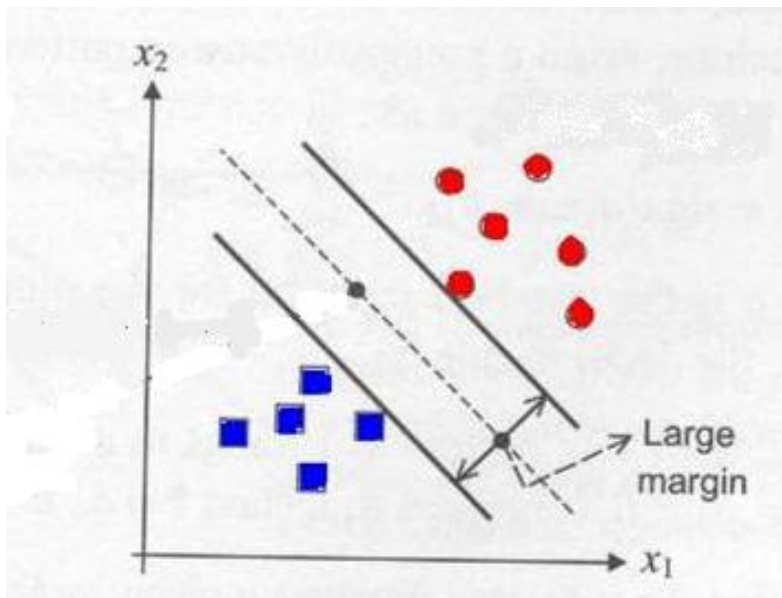


generalization error

$$E_{test} = E_{train} + E_{generalization}$$

Definition

- Maximize minimum distance
 - Data points to hyperplane
- Normalization needed



$$\max_{w, b} \min_i \frac{|\mathbf{w}^T \mathbf{x}_i + b|}{\|\mathbf{w}\|}$$



Canonical Form

- Normalize by scaling

$$f(w, b) = k \cdot w^T x + k \cdot b = 0$$

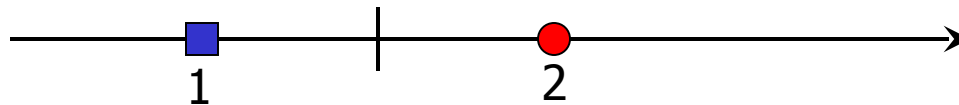
$$\min_i |w^T x_i + b| = 1$$

$$x - 1.5 = 0$$

$$2x - 3 = 0$$

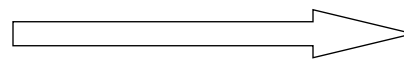
Canonical form ($w=2$)

$$|2 \cdot 2 - 3| = 1$$



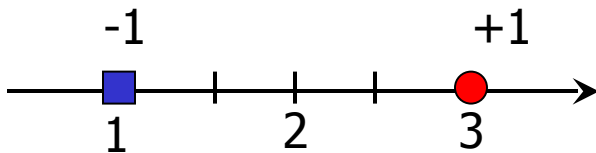
Quadratic Programming

$$\max_{w, b} \min_i \frac{|\mathbf{w}^T \mathbf{x}_i + b|}{\|\mathbf{w}\|}$$

$$\min_i |\mathbf{w}^T \mathbf{x}_i + b| = 1$$


$$\min_{w, b} \frac{1}{2} \langle \mathbf{w}^T \cdot \mathbf{w} \rangle$$

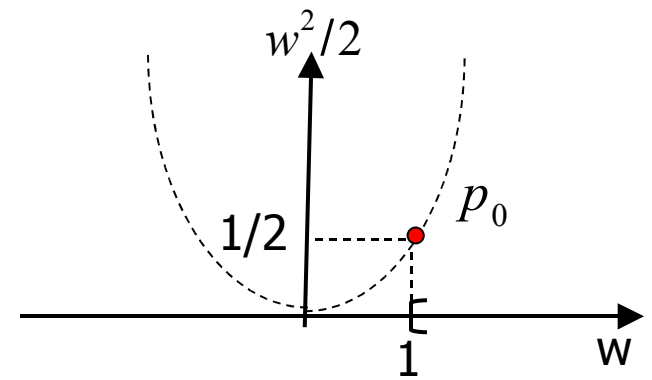
$$y_i (\langle \mathbf{w}^T \cdot \mathbf{x}_i \rangle + b) \geq 1$$



$$\min_w \frac{w^2}{2}$$

$$(+1)(w \cdot 3 + b) \geq 1$$

$$(-1)(w \cdot 1 + b) \geq 1$$





An Alternative Way

$$\begin{aligned} \min f_0(x) \\ f_i(x) \geq 0 \end{aligned}$$

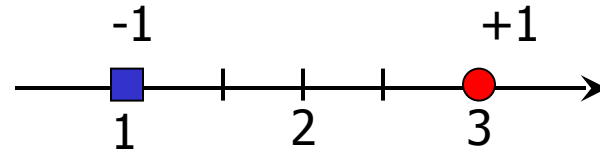
- Step 1-2-3:
 - Combine target and constraints
 - Minimize over primal
 - Maximize over dual

$$L(x, \lambda) = f_0(x) - \sum \lambda_i f_i(x)$$

$$Q(\lambda) = \min_x L(x, \lambda)$$

$$\max_{\lambda} Q(\lambda), \lambda > 0$$

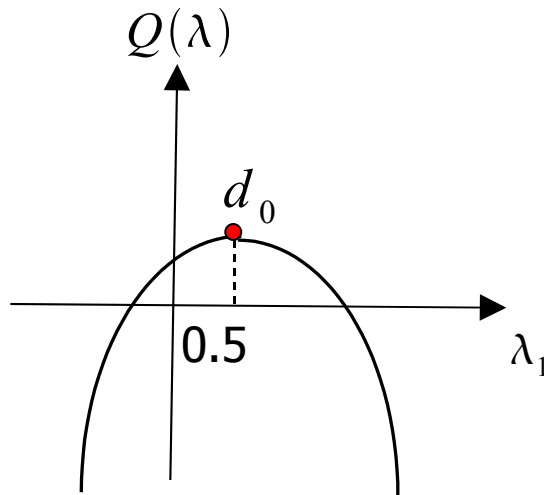
Dual Space



$$\min_w \frac{w^2}{2}$$

$$(+1)(w \cdot 3 + b) \geq 1$$

$$(-1)(w \cdot 1 + b) \geq 1$$



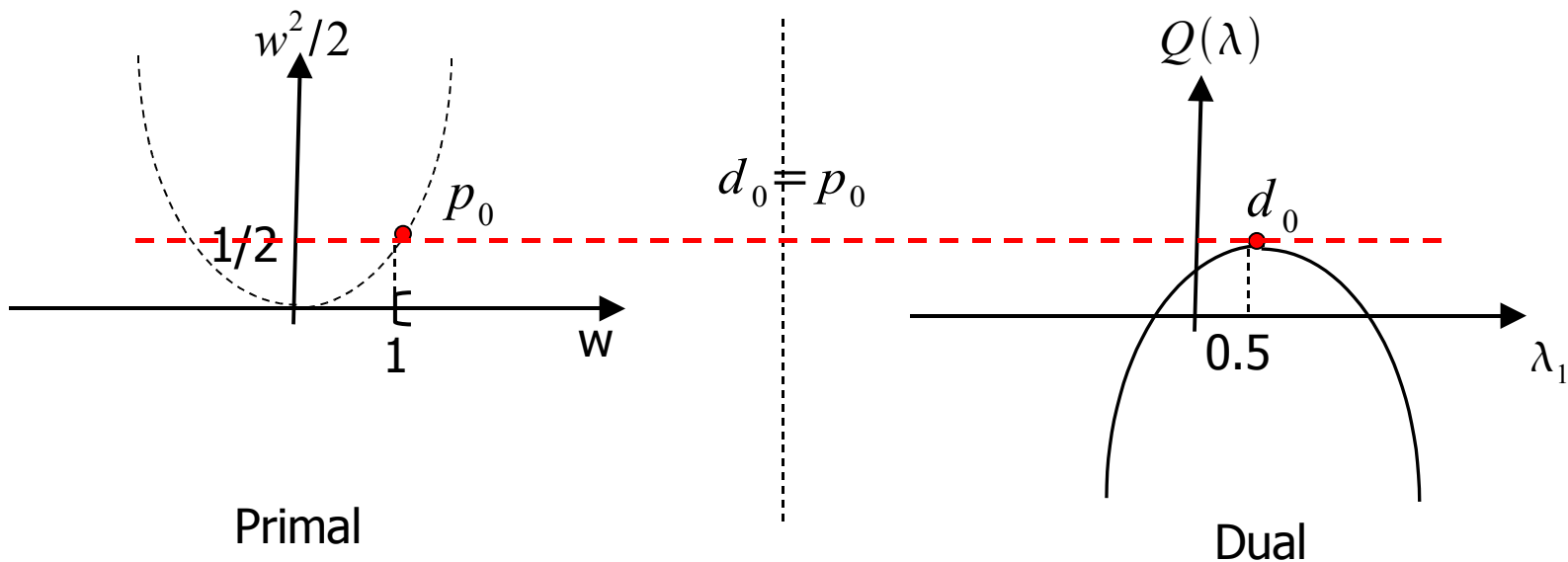
$$L(w, b, \lambda) = w^2/2 - \lambda_1(3w + b - 1) - \lambda_2(-w - b - 1)$$

$$\min_{w, b} L(w, b, \lambda) \Rightarrow \begin{cases} \lambda_1 = \lambda_2 \\ w = 3\lambda_1 - \lambda_2 = 2\lambda_1 \\ Q(\lambda) = Q(\lambda_1) = -2\lambda_1^2 + 2\lambda_1 \end{cases}$$

$$\max_{\lambda} Q(\lambda) \Rightarrow \lambda_1 = \lambda_2 = 1/2, w = 1, b = 2$$

Strong Duality

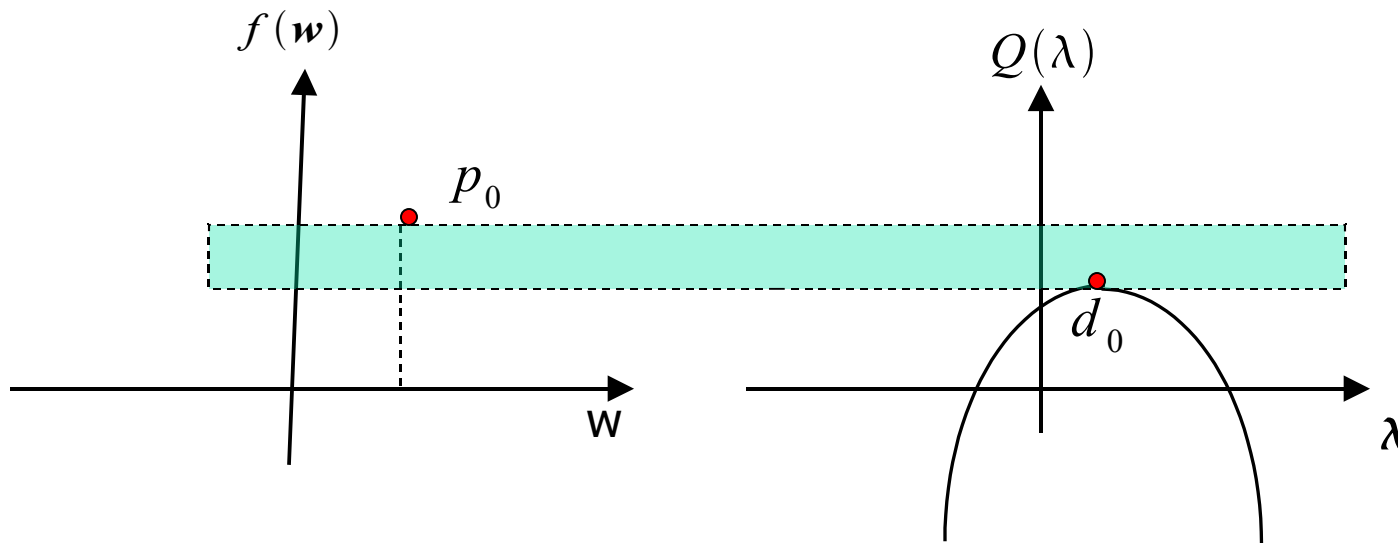
- Primal and dual space optimization:
 - Same result!



Duality Gap

$$d_0 < p_0$$

- In a general case
 - Strong duality is not true
 - “Step 1-2-3” a lower bound, not a solution





Convexity Saves the Day

- Convex function
 - Quadratic programming
- Convex set
 - Linear constraints
- No duality gap

$$\min_{w,b} \frac{1}{2} \langle \mathbf{w}^T \cdot \mathbf{w} \rangle$$
$$y_i (\langle \mathbf{w}^T \cdot \mathbf{x}_i \rangle + b) \geq 1$$

$$d_0 = p_0$$



Dual Form

- Formalize “step 1-2-3”
- H
 - Hessian matrix
 - Gram matrix
- Lambda
 - Support vector
 - Sparse

$$\max_{\lambda} Q(\lambda) = -0.5 \lambda^T H \lambda + f^T \lambda$$

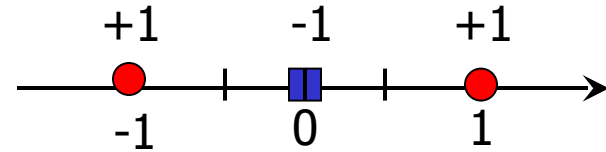
$$y^T \lambda = 0$$

$$\lambda \geq 0$$

$$\text{where, } H_{ij} = y_i y_j \langle \mathbf{x}_i^T \cdot \mathbf{x}_j \rangle$$

f is a unit vector

Non-separable Data

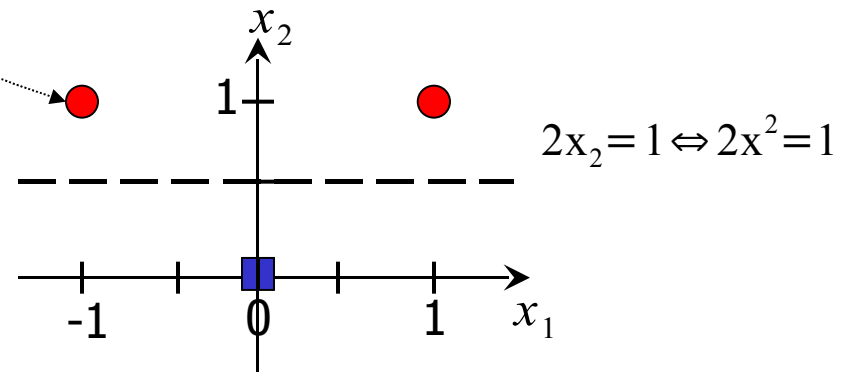
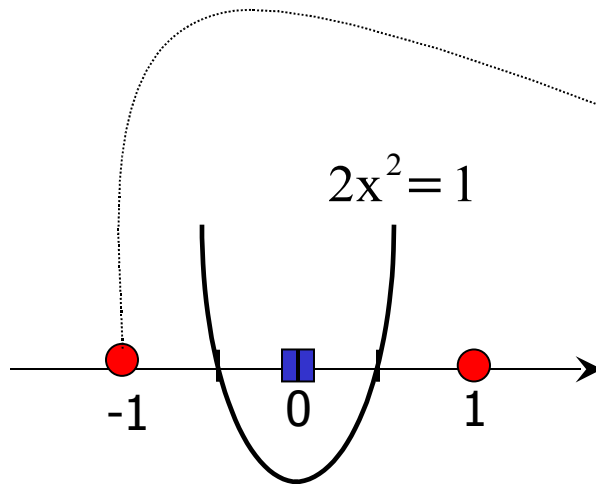


■ Solutions:

■ Nonlinear classifier

■ Increase the dimension

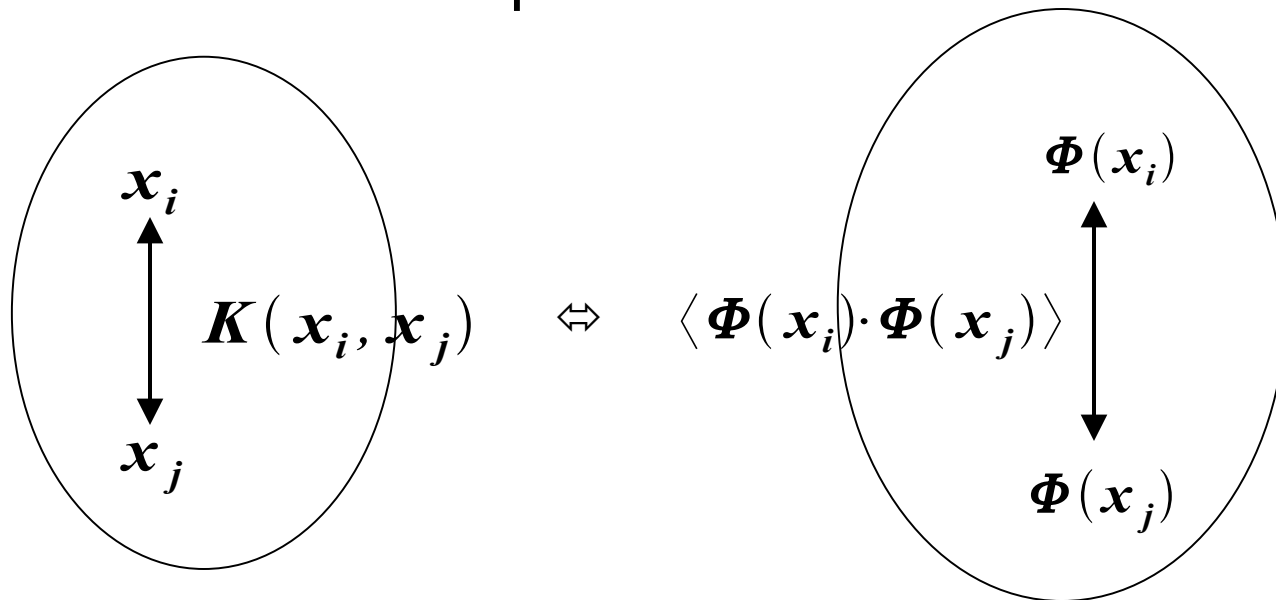
$$[x] \rightarrow \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x \\ x^2 \end{bmatrix}$$



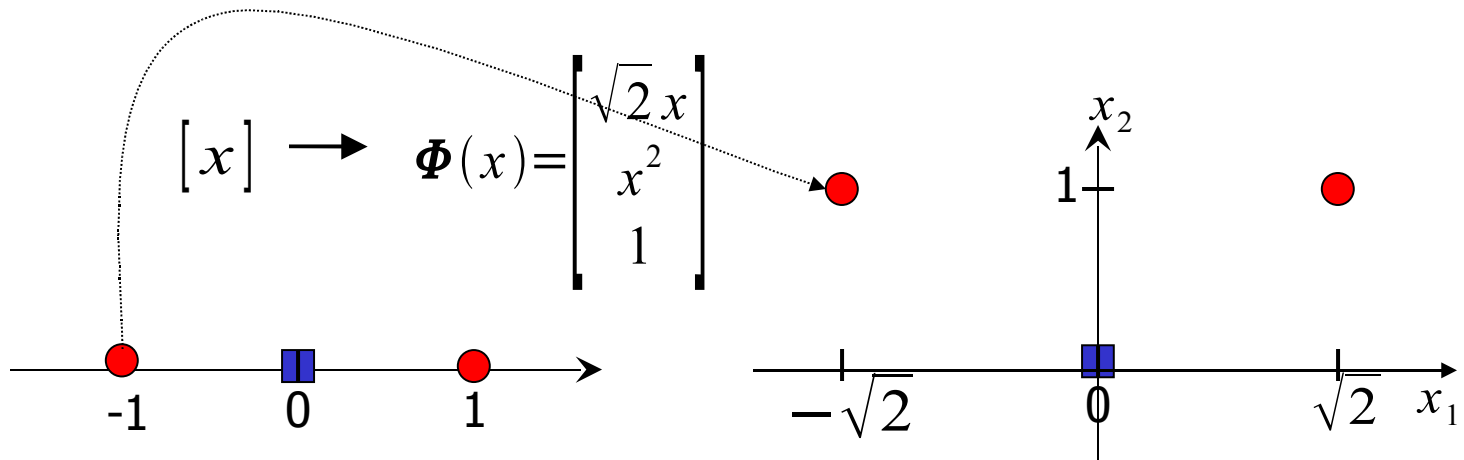


Kernel Trick

- Kernel function
 - in the original space
- Inner product
 - In the feature space with increased dimension



Example



$$\mathbf{K}(x_i, x_j) = (x_i x_j + 1)^2$$

$$\langle \Phi(x_i) \cdot \Phi(x_j) \rangle = 2x_i x_j + x_i^2 x_j^2 + 1 = (x_i x_j + 1)^2 = \mathbf{K}(x_i, x_j)$$



Curse of Dimensionality

- Primal space
 - Makes optimization much harder
- Dual space
 - Can be avoided

$$\min_{w, b} \frac{1}{2} \langle \Phi^T(\mathbf{w}) \cdot \Phi(\mathbf{w}) \rangle$$
$$y_i (\langle \Phi^T(\mathbf{w}) \cdot \Phi(\mathbf{x}_i) \rangle + b) \geq 1$$

$$\max_{\lambda} Q(\lambda) = -0.5 \lambda^T \mathbf{H} \lambda + \mathbf{f}^T \lambda$$

$$\mathbf{y}^T \lambda = 0$$

$$\lambda \geq 0$$

$$\text{where } H_{ij} = y_i y_j \mathbf{K}(\mathbf{x}_i, \mathbf{x}_j)$$

\mathbf{f} is a unit vector



Mercer Condition

- Dual form is convex
 - H is P.S.D.
 - Kernel must be P.S.D.

$$Q(\lambda) = -0.5 \lambda^T \mathbf{H} \lambda + \mathbf{f}^T \lambda$$

where, $H_{ij} = y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$

- Mercer kernels
 - Polynomial
 - Gaussian

$$K(\mathbf{x}, \mathbf{y}) = [\langle \mathbf{x}^T \mathbf{y} \rangle + 1]^p$$

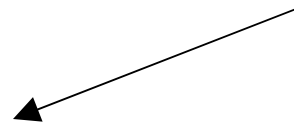
$$K(\mathbf{x}, \mathbf{y}) = e^{-\frac{(\mathbf{x}-\mathbf{y})^T \Sigma^{-1} (\mathbf{x}-\mathbf{y})}{2}}$$



Why Is SVM So Popular?

- Works very well
- Fool-proof
 - Only 2 kernels to choose from
 - Very few hand-crafted parameters
- Error bound easy to get

$$\frac{|SV|}{N}$$



$$E_{test} = E_{train} + E_{generalization}$$



Disadvantages

- Slow
 - Training: quadratic programming
 - SMO, etc.
 - Hardware
 - Testing: depend on number of SV
- Big
 - Curse of sample size: H matrix
 - Online SVM