

Support Vector Machines

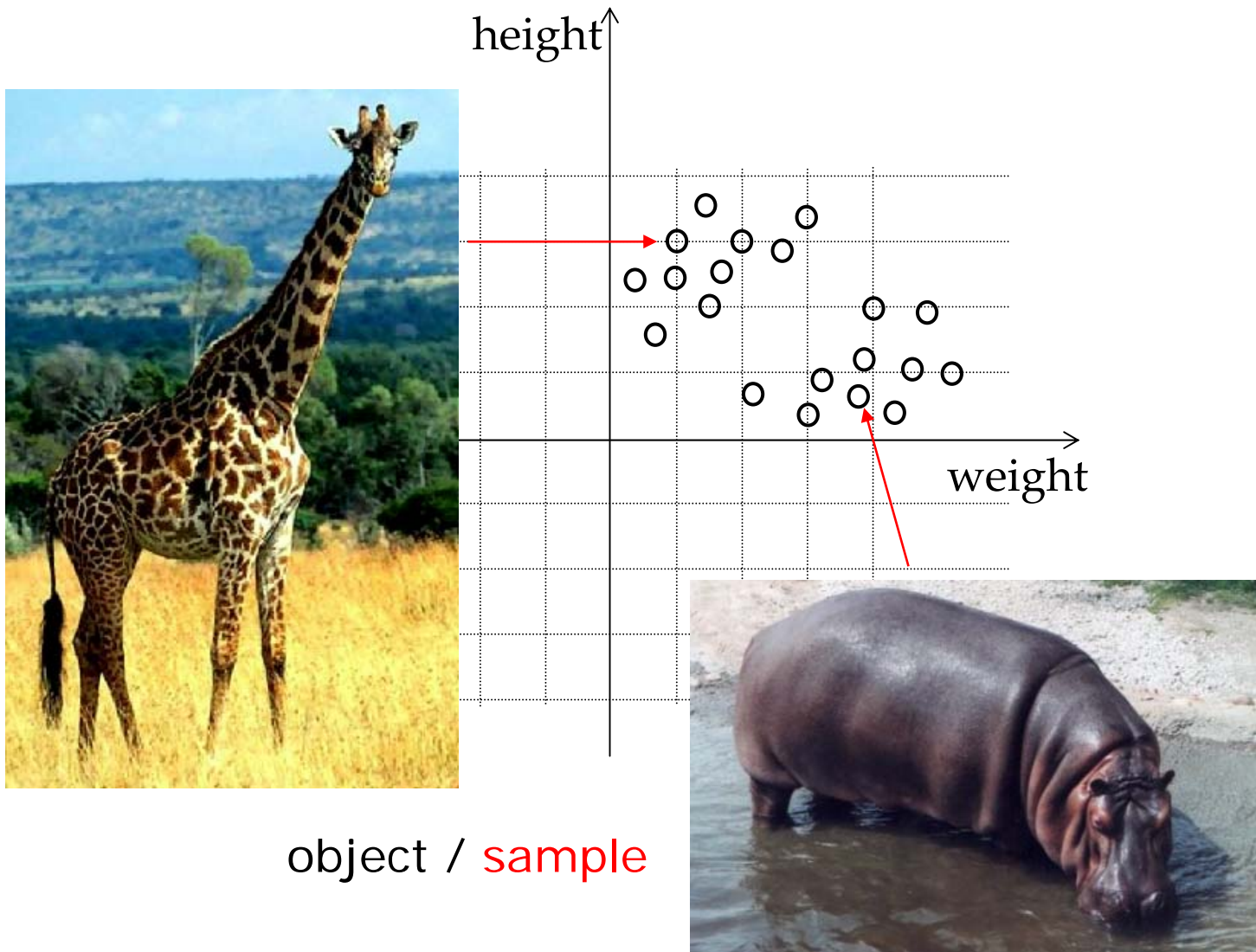
<http://cespc1.kumoh.ac.kr/~nonezero/svm-ws-cvpr.pdf>

금오공과대학교
컴퓨터공학부
고재필

Contents

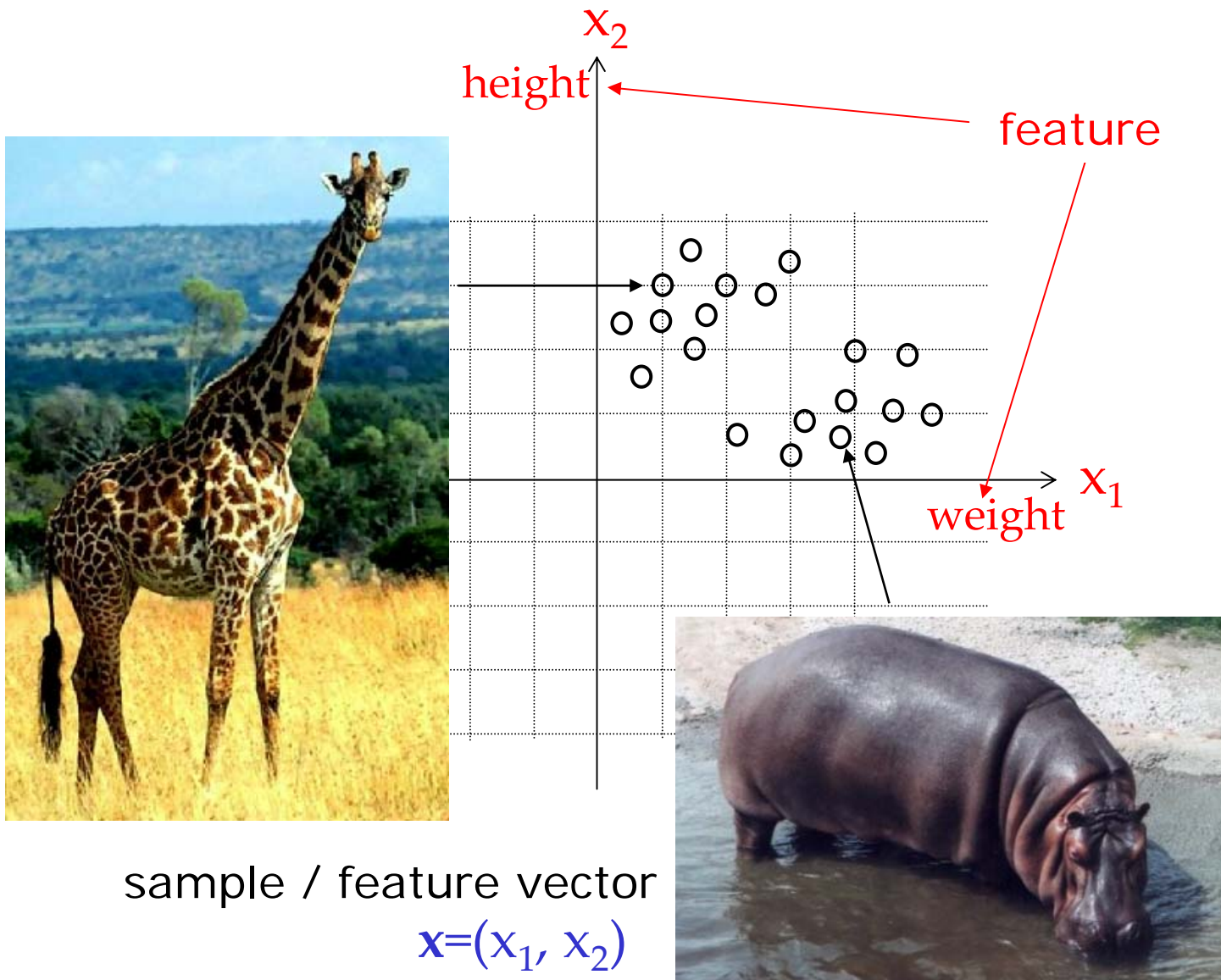
- **Introduction**
- **Optimal Hyperplane**
- **Soft-Margin SVM**
- **Nonlinear SVM with Kernel Trick**
- **Training SVM: SMO**
- **Link to Statistical Learning Theory**
- **Multiclass SVM with Classifiers Ensemble**
- **SVM Demo**
- **References**

Sample

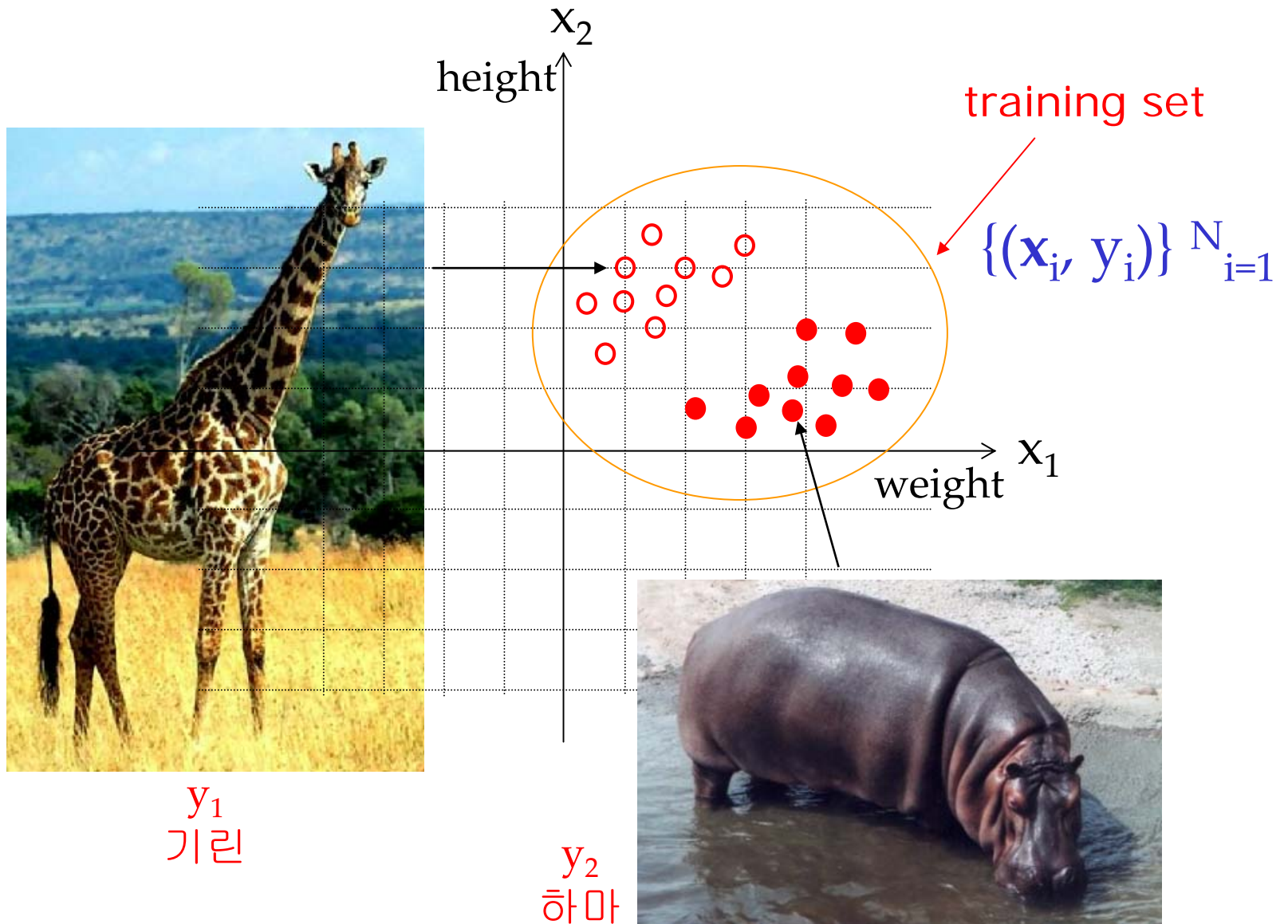


object / **sample**

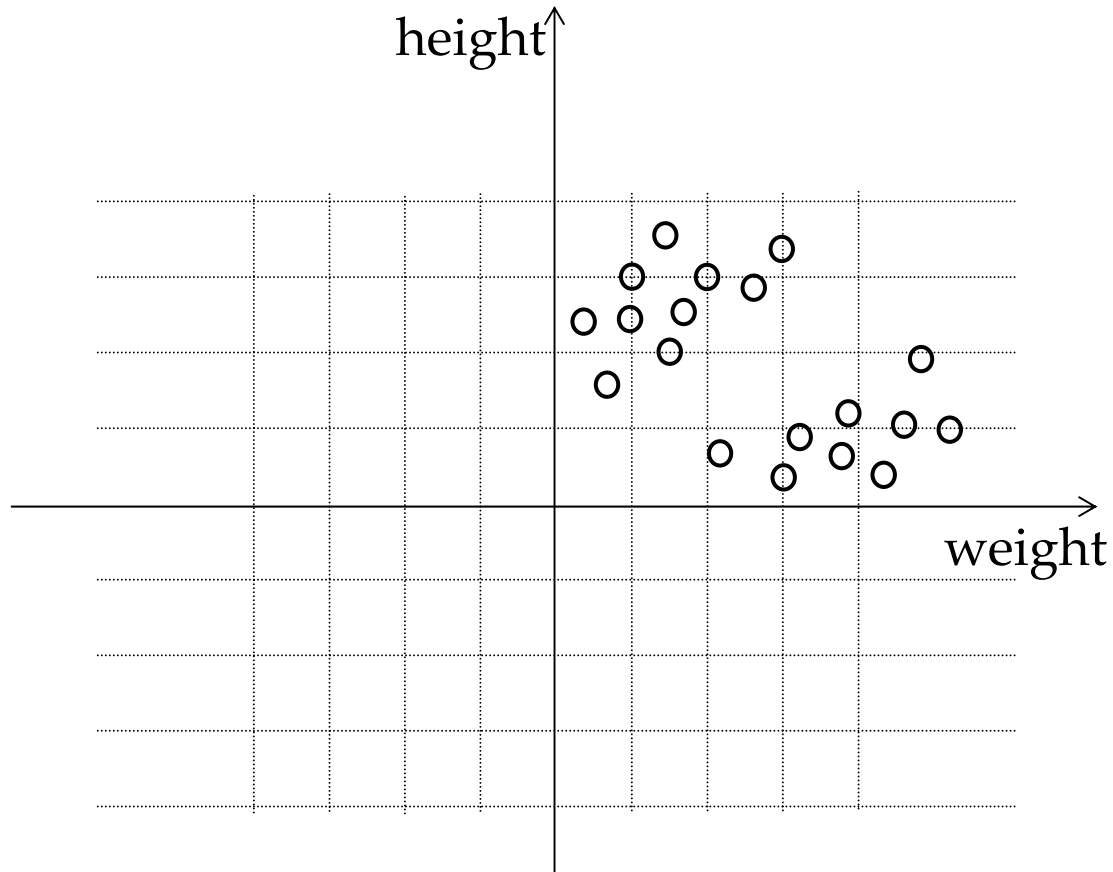
Feature



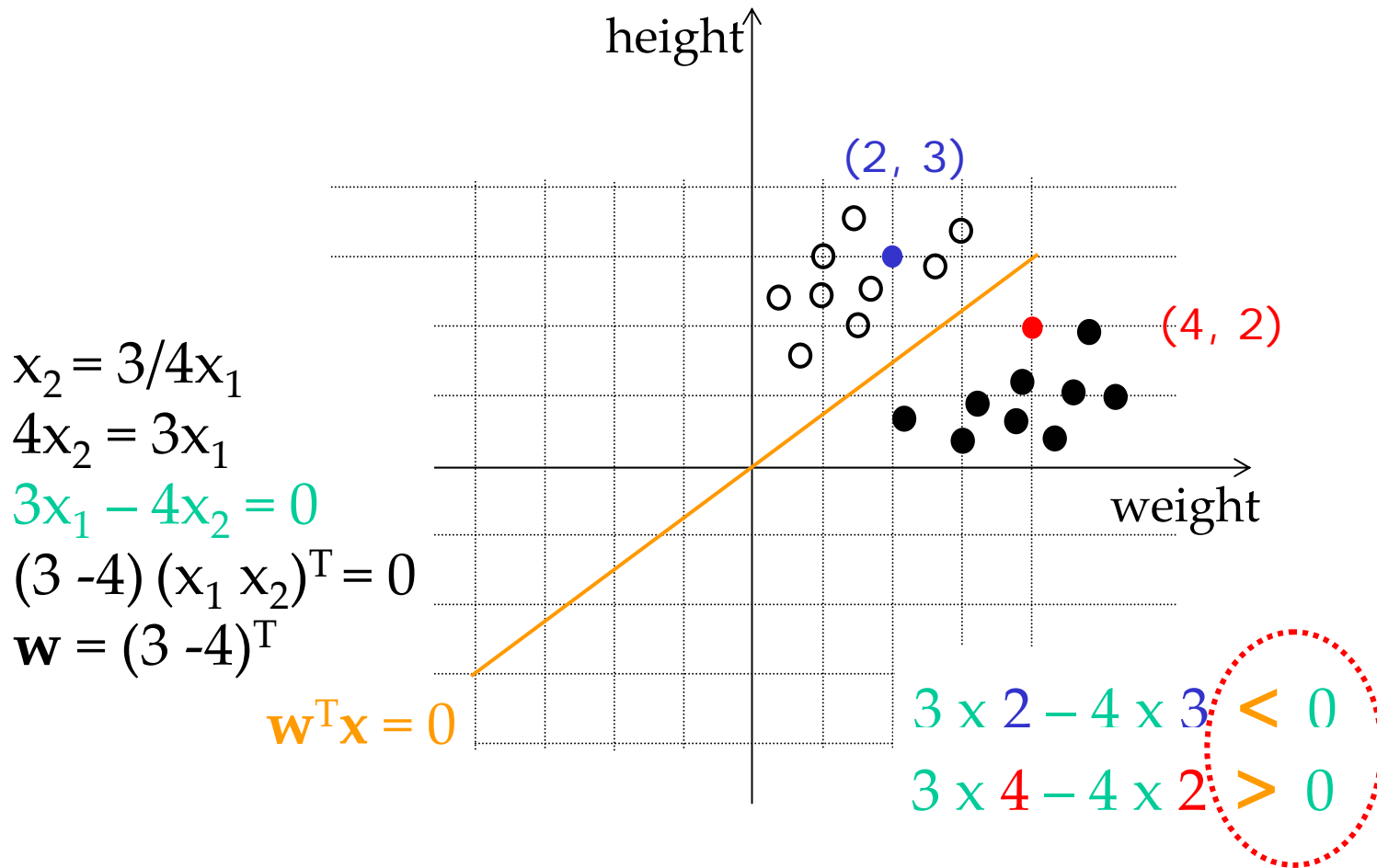
Training Set



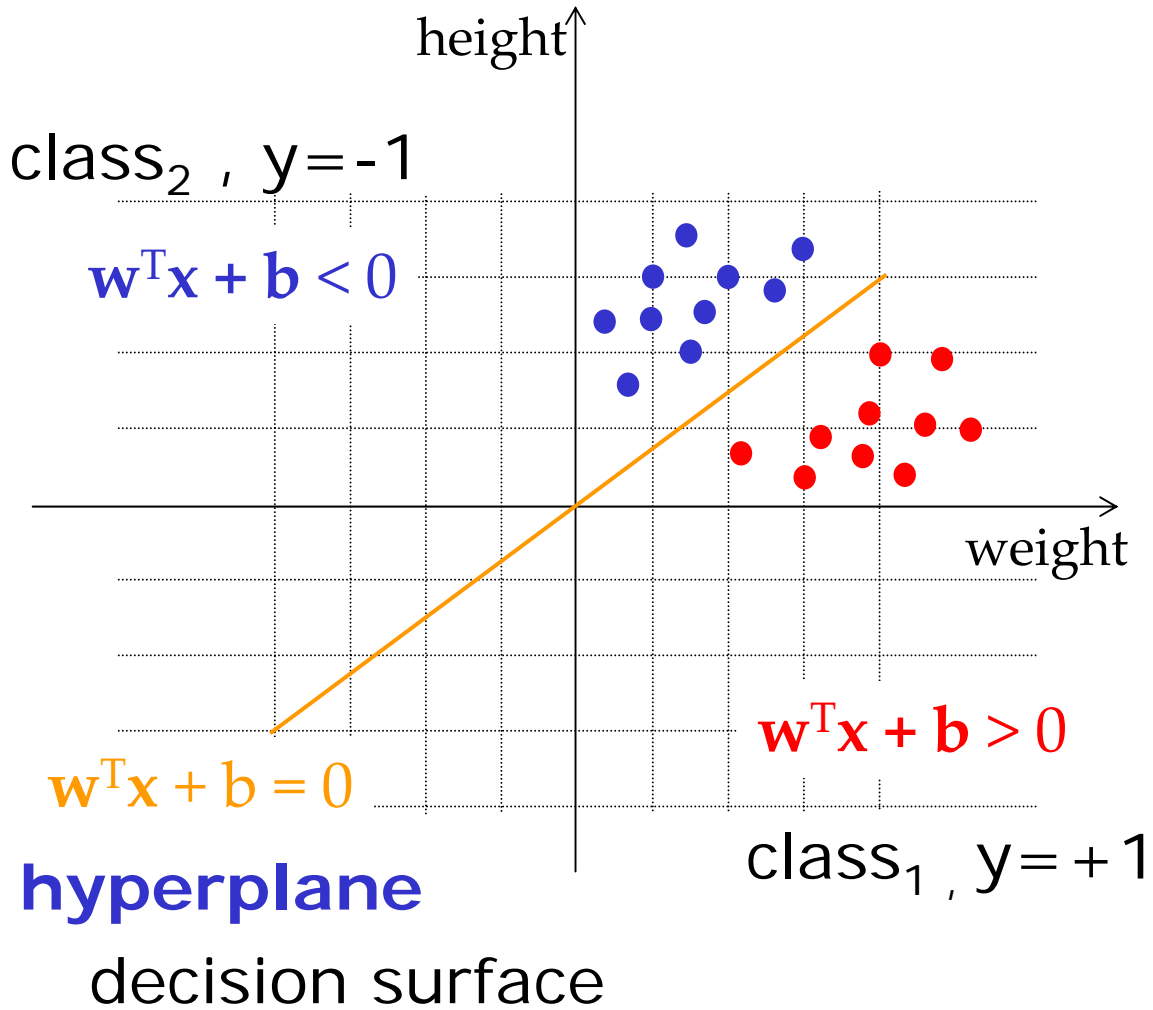
How to Classify Them Using Computer ?



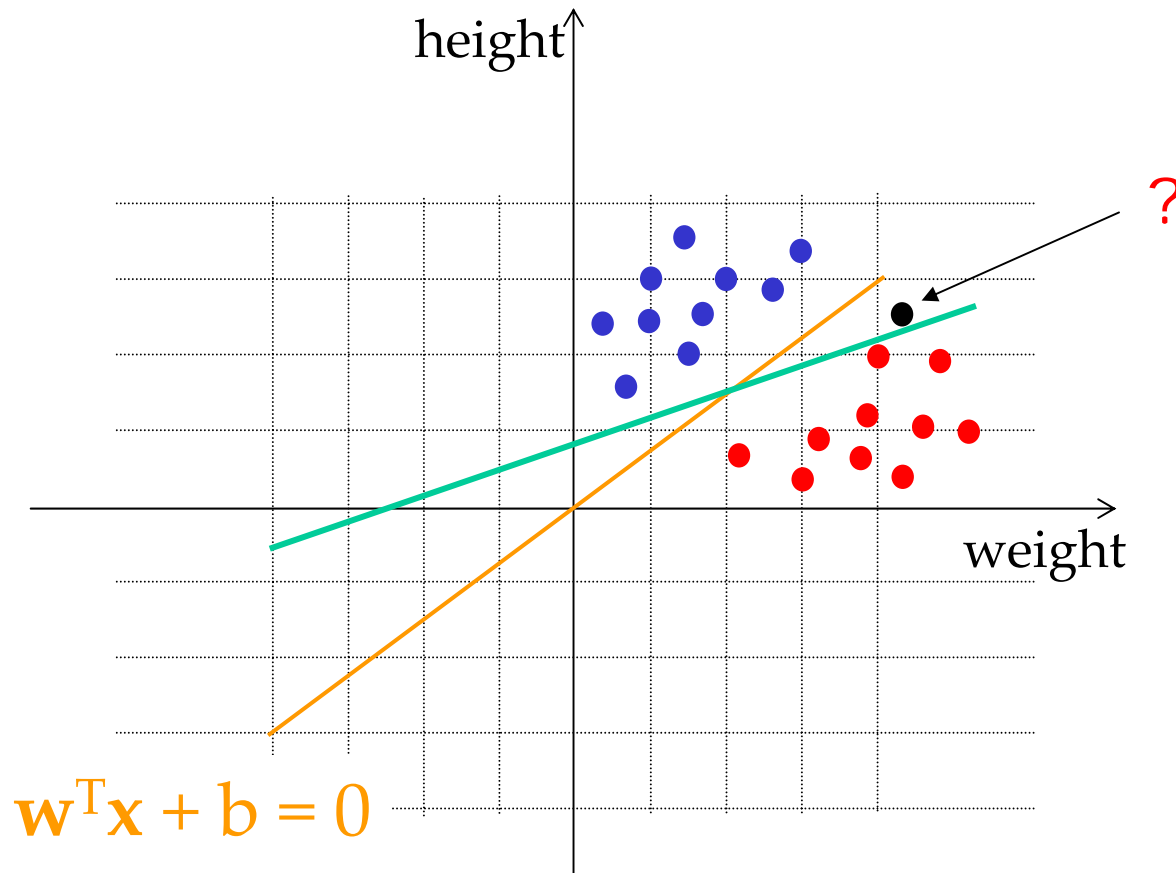
How to Classify Them Using Computer ?



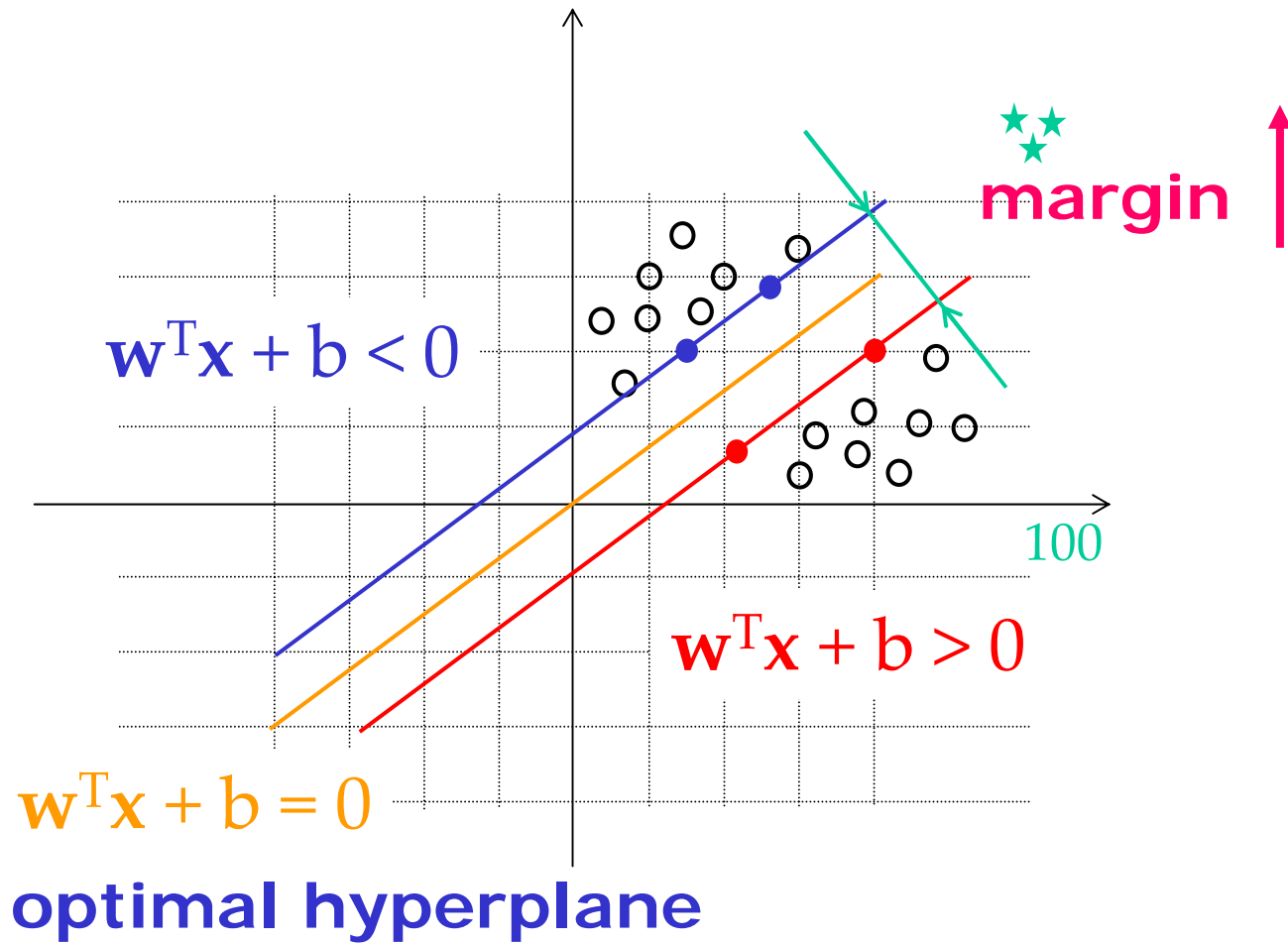
Linear Classification



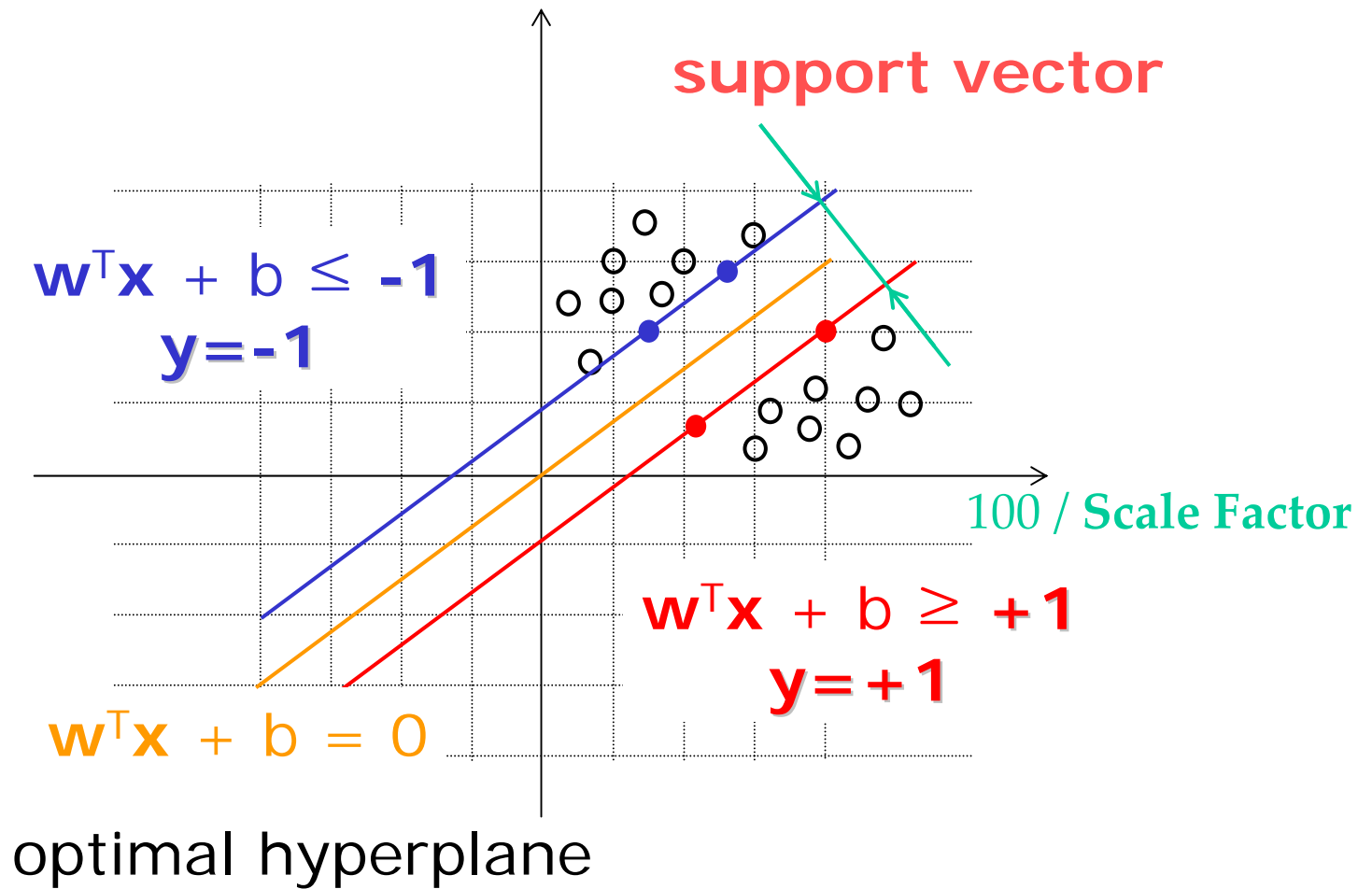
Optimal Hyperplane



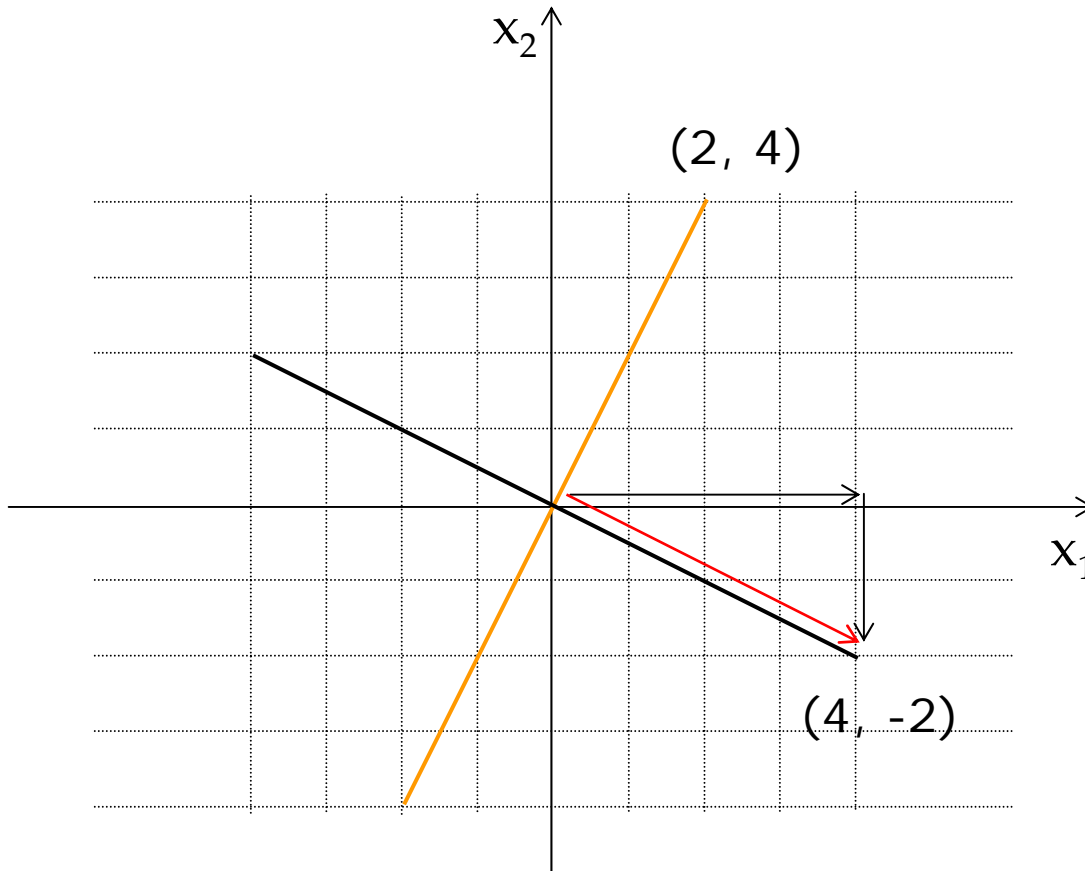
Optimal Hyperplane



Canonical Hyperplane



Normal Vector

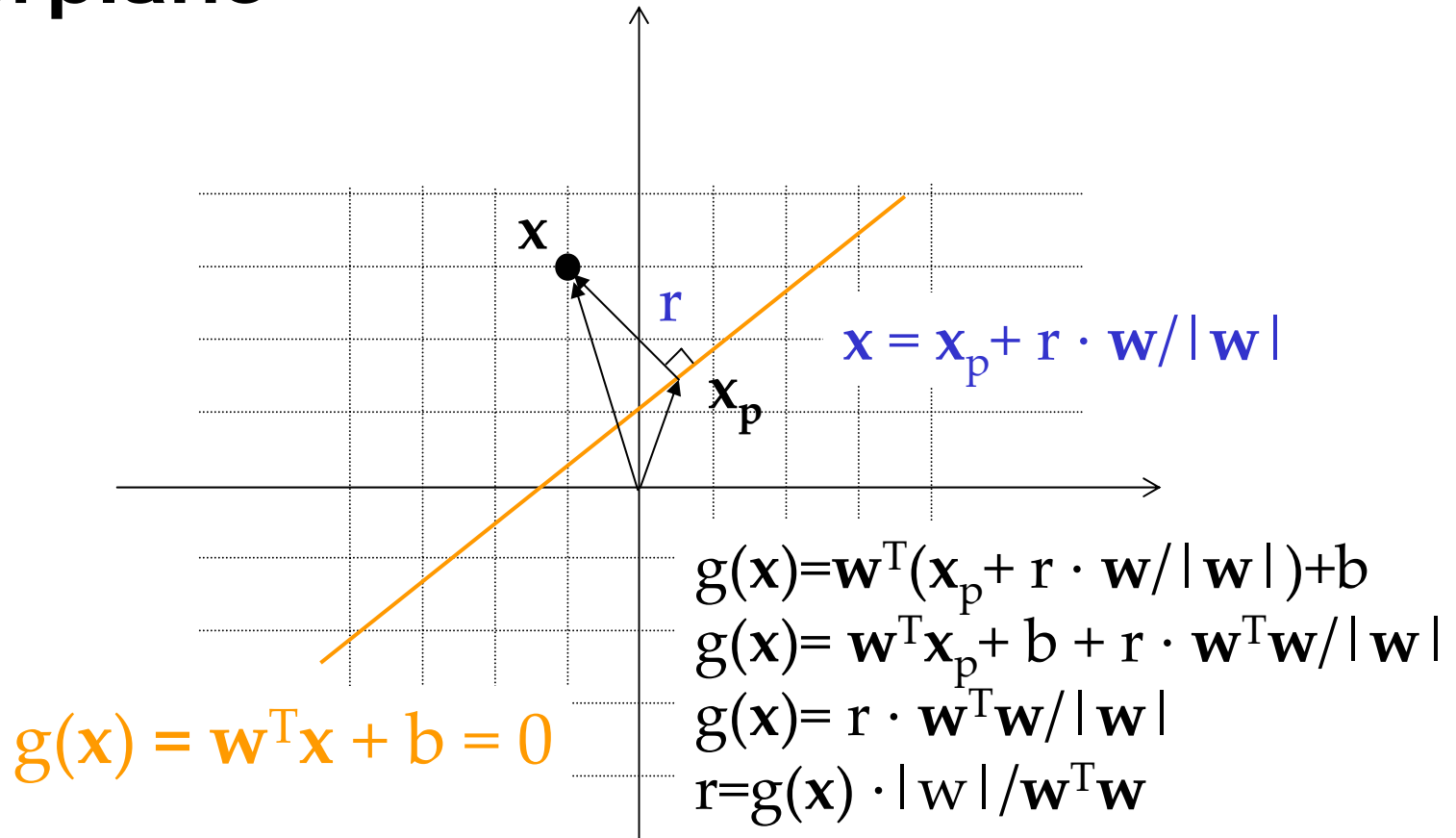


$$\begin{aligned}x_2 &= 4/2x_1 \\2x_2 &= 4x_1 \\4x_1 - 2x_2 &= 0 \\(4 \ -2) (x_1 \ x_2)^T &= 0 \\w^T x &= 0 \\w &= (4 \ -2)^T\end{aligned}$$

$$(2 \ 4)(4 \ -2)^T = 8 - 8 = 0$$

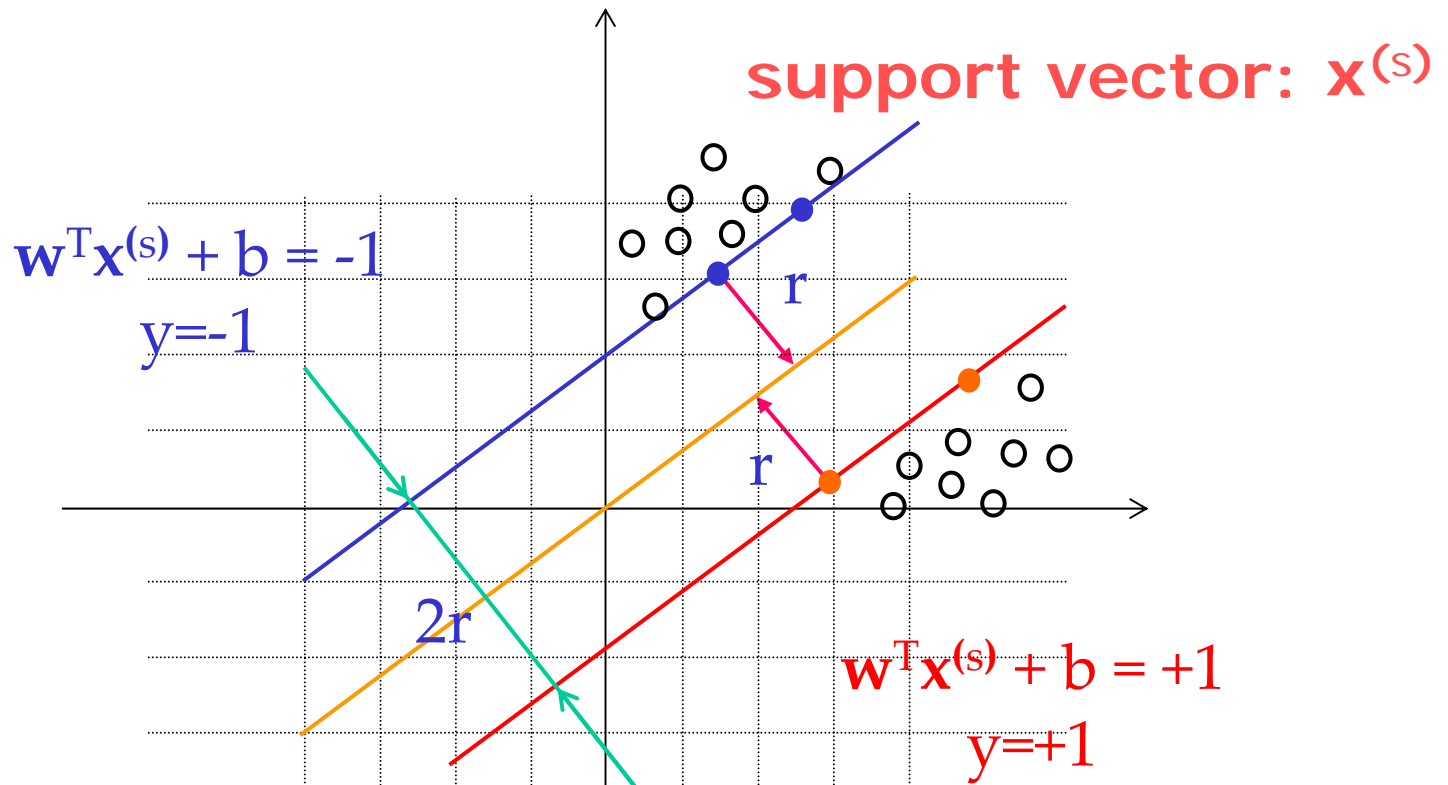
$$\mathbf{x} \cdot \mathbf{y} = \langle \mathbf{x} \ \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y} = x_1 y_1 + x_2 y_2 = |\mathbf{x}| |\mathbf{y}| \cos \theta$$

Distance from a Sample to the Optimal Hyperplane



$$r = g(\mathbf{x}) / |\mathbf{w}|$$

Margin of Separation



$$r = g(x^{(s)}) / \|w\|$$

where, $g(x^{(s)}) = w^T x^{(s)} + b = \pm 1$ for $y^{(s)} = \pm 1$

margin: $2r = 2 / \|w\| \left| \frac{1}{\|w\|} - \frac{-1}{\|w\|} \right| = \frac{2}{\|w\|}$

Finding the Optimal Hyperplane

margin: $\left| \frac{1}{\|w\|} - \frac{-1}{\|w\|} \right| = \frac{2}{\|w\|}$

$$w^T x + b \leq -1, \quad y = -1$$

$$w^T x + b \geq +1, \quad y = +1$$

Constrained Optimization Problem

(convex) objective/cost function

minimize $\frac{1}{2} w^T w$

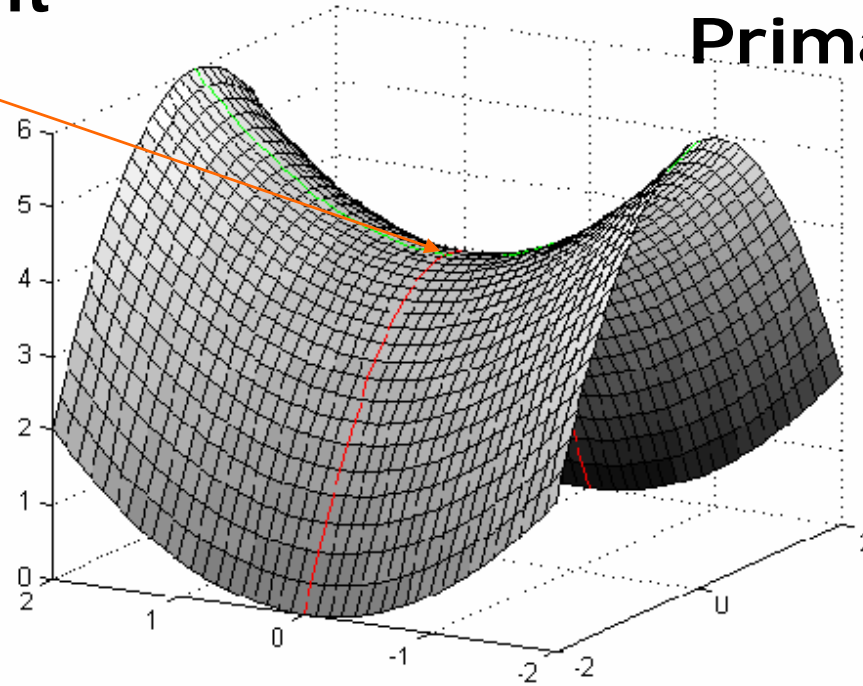
subject to $y_i (w^T x_i + b) \geq 1 \quad i = 1, 2, \dots, n$

(linear) constraints

Convex Quadratic Problem

Saddle Point

minimize
Primal problem



maximize
Dual problem

Lagrange Multipliers Method

cost function: $f(\mathbf{x})$

constraint function: $g(\mathbf{x})=0$

$$\nabla f(\mathbf{x}) = \lambda \nabla g(\mathbf{x})$$

tangent (line) ~ gradient = normal

convex / concave

↓

$$\nabla f(\mathbf{x}) = 0, g(\mathbf{x})=0$$

$$\mathbf{L}(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$$

↑
Lagrange function

↑
Lagrange multiplier

A Brief Overview of Optimization Theory

(2001년 추계 CVPR 튜토리얼)

- ▶ Theorem: $f \in C^1$ has a min. at $x^* \Rightarrow \frac{\partial f}{\partial x}(x^*) = 0$.

This condition, together with convexity of f , is also a sufficient condition.

- ▶ Example 1: $\min. f(x) = \frac{1}{2}(x_1^2 + x_2^2)$

Solution:

$$\frac{\partial f}{\partial x} = 0 \Rightarrow \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \right] = [x_1 \quad x_2] = 0 \quad \therefore x^* = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

- ▶ In a constrained min. problem,

$$f \in C^1 \text{ has a min. at } x^* \Rightarrow \frac{\partial f}{\partial x}(x^*) = 0$$

A Brief Overview of Optimization Theory

▶ Example 2:

$$\begin{aligned} \min . f(x) &= \frac{1}{2}(x_1^2 + x_2^2) \\ \text{s.t. } h(x) &= 1 - x_1 - x_2 = 0 \end{aligned}$$

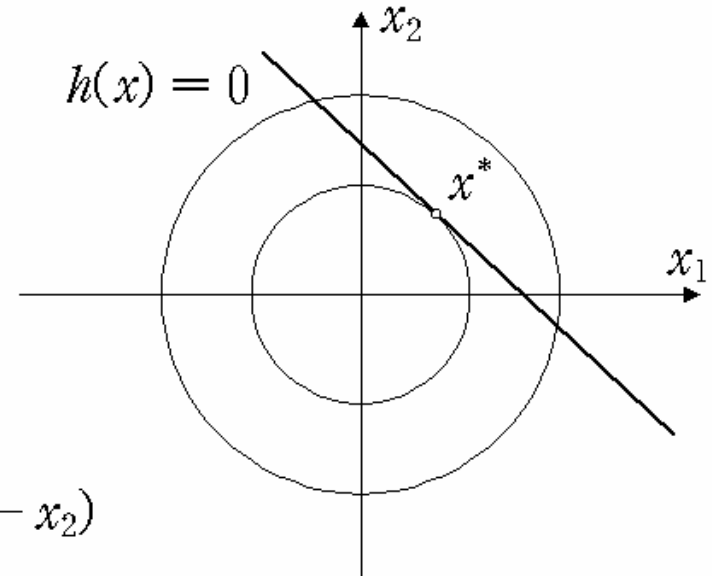
Solution: Define the Lagrange function

$$\begin{aligned} L(x, \lambda) &= f(x) + \lambda h(x) \\ &= \frac{1}{2}(x_1^2 + x_2^2) + \lambda(1 - x_1 - x_2) \end{aligned}$$

$$\frac{\partial L}{\partial x} = 0 \Rightarrow [x_1 - \lambda \quad x_2 - \lambda] = 0. \therefore x_1 = x_2 = \lambda.$$

$$\frac{\partial L}{\partial \lambda} = 0 \Rightarrow 1 - x_1 - x_2 = 0. \therefore 1 - 2\lambda = 0. \therefore \lambda = \frac{1}{2}.$$

$$\therefore x^* = \begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix}$$



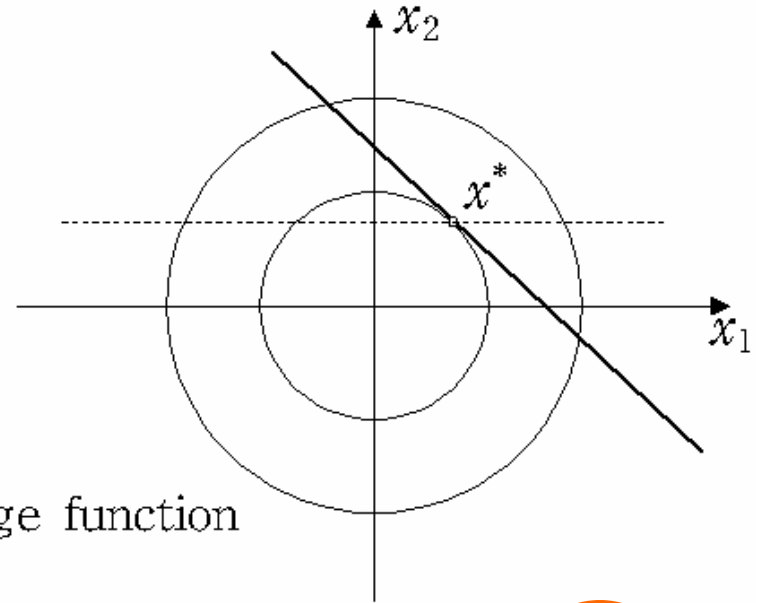
A Brief Overview of Optimization Theory

▶ Example 3:

$$\min. f(x) = \frac{1}{2}(x_1^2 + x_2^2)$$

$$\text{s.t. } h(x) = 1 - x_1 - x_2 = 0,$$

$$g(x) = \frac{3}{4} - x_2 \leq 0$$



Solution: Define the generalized Lagrange function

$$\begin{aligned} L(x, \lambda, \alpha) &\triangleq f + \lambda h + \alpha g \\ &= \frac{1}{2}(x_1^2 + x_2^2) + \lambda(1 - x_1 - x_2) + \alpha\left(\frac{3}{4} - x_2\right), \quad \alpha \geq 0 \end{aligned}$$

$$\frac{\partial L}{\partial x} = 0 \Rightarrow [x_1 - \lambda, x_2 - \lambda - \alpha] = 0. \quad \therefore x_1 = \lambda, \quad x_2 = \lambda + \alpha$$

$$\frac{\partial L}{\partial \lambda} = 0 \Rightarrow 1 - x_1 - x_2 = 0. \quad \therefore 2\lambda + \alpha = 1$$

A Brief Overview of Optimization Theory

Also, $\alpha \geq 0$ and $\frac{3}{4} - x_2 \leq 0$.

One more condition is needed to solve the problem.

Karush

→ The Kuhn-Tucker complementarity condition

$$\alpha \left(\frac{3}{4} - x_2 \right) = 0 \quad \text{i.e. } \alpha = 0 \quad \text{or} \quad x_2 = \frac{3}{4}$$

① If $\alpha = 0$, then $\lambda = \frac{1}{2}$; thus $x_1 = x_2 = \frac{1}{2} \quad \times \quad (\because x_2 \geq \frac{3}{4})$

② If $x_2 = \frac{3}{4}$, then
$$\begin{cases} \lambda + \alpha = \frac{3}{4} \\ 2\lambda + \alpha = 1 \end{cases} \therefore \begin{cases} \lambda = \frac{1}{4}, \alpha = \frac{1}{2} \\ x_1 = \frac{1}{4}, x_2 = \frac{3}{4} \end{cases}$$

$$\therefore x^* = \begin{bmatrix} 1/4 \\ 3/4 \end{bmatrix}$$

A Brief Overview of Optimization Theory

Theorem (Kuhn-Tucker Theorem)

Given an opt. prob. with convex domain $\Omega \subseteq \mathbb{R}^n$

$$\left. \begin{array}{l} \min f(x), x \in \Omega \text{ (} x \text{ is primal variable)} \\ \text{s.t. } g_i(x) \leq 0, i = 1, \dots, k \\ h_j(x) = 0, j = 1, \dots, m \end{array} \right\}$$

primal opt. prob.

with $f \in C^1$ convex, and g_i, h_j affine, the following are necessary and sufficient conditions for a point $x^* \in \Omega$ to be an opt.:

$$\text{For } L(x, \alpha, \lambda) \triangleq f(x) + \sum_{i=1}^k \alpha_i g_i(x) + \sum_{j=1}^m \lambda_j h_j(x) = f + \alpha^T g + \lambda^T h,$$

$$\exists \alpha^* \text{ and } \lambda^* \text{ s.t. } \frac{\partial L}{\partial x}(x^*, \alpha^*, \lambda^*) = 0, \frac{\partial L}{\partial \lambda}(x^*, \alpha^*, \lambda^*) = 0$$

$$g_i(x^*) \leq 0, \alpha_i^* \geq 0 \text{ for } i = 1, \dots, k,$$

$$\text{and } \underline{\alpha_i^* g_i(x^*) = 0, i = 1, \dots, k}$$

Lagrange Function for the Optimal Hyperplane (Primal Problem)

$$\begin{array}{ll} \underset{w,b}{\text{minimize}} & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{subject to} & y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, \quad i = 1, \dots, n. \end{array}$$

$$L_P(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i [y_i (\mathbf{w}^\top \mathbf{x}_i + b) - 1]$$

Solution:

$$\frac{dL_P}{d\mathbf{w}} = 0$$

$$\frac{dL_P}{db} = 0$$

$$\begin{aligned} \mathbf{w} &= \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \\ 0 &= \sum_{i=1}^n \alpha_i y_i \quad (\alpha_i \geq 0) \end{aligned}$$

KKT condition: $\alpha_i = 0$ unless $y_i (\mathbf{w}^\top \mathbf{x}_i + b) = 1$

Remark on Support Vector

$$\text{KKT Condition } \alpha_i [1 - y_i(w^T x_i + b)] = 0$$

$$\text{If } \alpha_i \neq 0 \text{ then } y_i(w^T x_i + b) = 1$$

→ x_i is **support vector** All x_i for which $\alpha_i > 0$

$$y_i(w^T x_i + b) \neq 1$$


→ $\alpha_i = 0$ x_i is **not support vector**

$$w = \sum_{i=1}^n \alpha_i y_i x_i$$

w is only related to support vectors, x_i

Lagrange Function for the Optimal Hyperplane (Dual Problem)

$$\begin{aligned} \max. \quad L_D(\alpha) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \\ \text{s.t.} \quad \sum_{i=1}^n \alpha_i y_i &= 0 \quad ; \quad \alpha_i \geq 0 \end{aligned}$$

$$-\frac{1}{2} \alpha^T Q \alpha + \alpha^T \mathbf{1}$$


Optimizing L_D only depends on the input patterns in the form of a set of **dot product** $\mathbf{x}^T \mathbf{x}$

- (+) Not depend on the **dimension** of the input pattern
- (+) Can replace dot product with **Kernel**

$$\begin{aligned}
L_P(\mathbf{w}, b, \alpha) &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i [y_i (\mathbf{w}^\top \mathbf{x}_i + b) - 1] \\
&= \frac{1}{2} \mathbf{w}^\top \mathbf{w} - \sum_{i=1}^n \alpha_i y_i \mathbf{w}^\top \mathbf{x}_i - b \sum_{i=1}^n \alpha_i y_i + \sum_{i=1}^n \alpha_i \\
&\quad \sum_{i=1}^n \alpha_i y_i = 0 \\
&\quad \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \\
&\quad \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \quad \mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \\
&= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j = -\frac{1}{2} \boldsymbol{\alpha}^\top \mathbf{Q} \boldsymbol{\alpha} + \boldsymbol{\alpha}^\top \mathbf{1} \\
&\quad \quad \quad \uparrow \\
&\quad \quad \quad n \times n (> 1000)
\end{aligned}$$

Large Scale Quadratic Problem

Support Vector Machine Classifier for Linearly Separable Case

$$f(\mathbf{x}) = \text{sgn} \left(\sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i^\top \mathbf{x} + b^* \right)$$

optimal weight

$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i$$

$$y_i (\mathbf{w}^\top \mathbf{x}_i + b) = 1$$

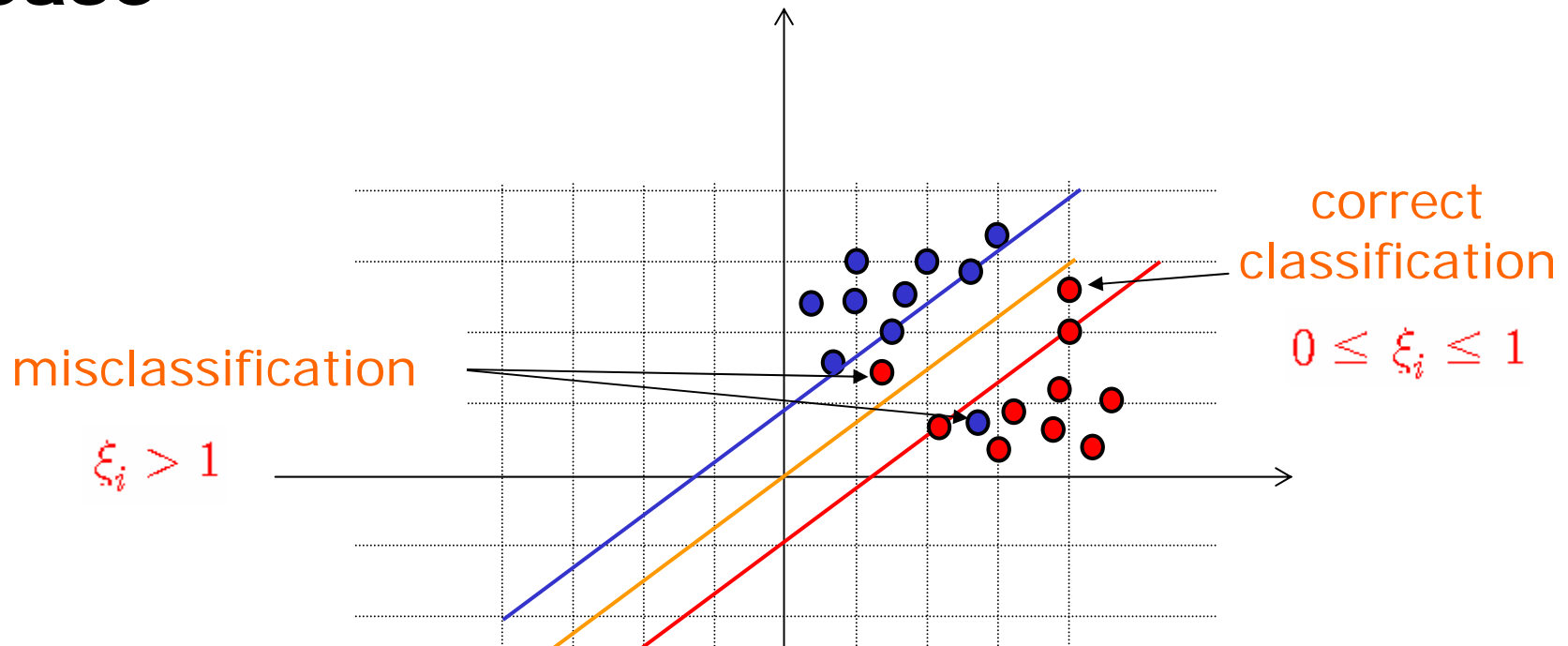
$$b_{y=+1} = 1 - \mathbf{w}^* \mathbf{x}^{(s)}$$

$$b_{y=-1} = -1 - \mathbf{w}^* \mathbf{x}^{(s)}$$

optimal bias

$$b^* = \frac{1}{2} (b_{y=+1} + b_{y=-1})$$

Optimal Hyperplane for Non-Separable Case



impossible to find a separating hyperplane

give them up as errors while minimizing the probabilities of classification error averaged over the training set

Soft Margin Technique

Problem: Can't satisfy $y_i[\mathbf{w}^\top \mathbf{x}_i + b] \geq 1$ for all i

Adopting Slack Variable

$$y_i (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$$

$$\mathbf{w}^\top \mathbf{x}_i + b \geq +1 - \xi_i \quad \text{for } y_i = +1,$$

$$\mathbf{w}^\top \mathbf{x}_i + b \leq -1 + \xi_i \quad \text{for } y_i = -1,$$

$$\xi_i \geq 0 \quad k = 1, 2, \dots, n$$

Minimizing Errors

$$\xi_i > 1 \quad \sum_{i=1}^n I(\xi_i > 1) = \# \text{ errors} \quad \min. \sum_{i=1}^n I(\xi_i > 1)$$

For QP, replace $I(\xi_i > 1)$ by ξ_i

Lagrange Function for Soft Margin Tech.

$$\begin{aligned} \underset{w, b, \xi}{\text{minimize}} \quad & L_P(\mathbf{w}, \xi) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{subject to} \quad & y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i \\ & \xi_i \geq 0 \end{aligned}$$

penalize errors

C



penalize complexity

$$\begin{aligned} L(\mathbf{w}, b, \alpha, \xi) &= \frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_{i=1}^m \xi_i \\ &\quad - \sum_{i=1}^m \alpha_i [y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i] - \sum_{i=1}^m r_i \xi_i \\ &\quad \star \alpha_i \geq 0 \text{ and } r_i \geq 0 \end{aligned}$$

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{w}} &= \mathbf{w} - \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i = 0 \\ \frac{\partial L}{\partial b} &= \sum_{i=1}^m \alpha_i y_i = 0 \end{aligned}$$

$$\star \frac{\partial L}{\partial \xi_i} = C - \alpha_i - r_i = 0$$

Dual Problem for Soft Margin Tech.

$$\begin{aligned} \max. \quad & L_D(\boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \quad ; \quad \boxed{0 \leq \alpha_i \leq C} \end{aligned}$$

$$\alpha_i \geq 0 \text{ and } r_i \geq 0$$

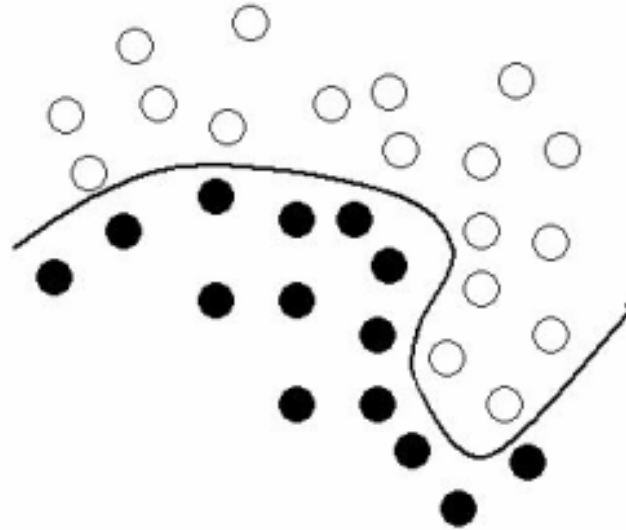
$$\frac{\partial L}{\partial \xi_i} = C - \alpha_i - r_i = 0$$

$$0 < \alpha_i < C \quad \textit{non-bound} \text{ pattern}$$

$$\alpha_i = 0 \text{ or } \alpha_i = C \quad \textit{bound} \text{ pattern}$$

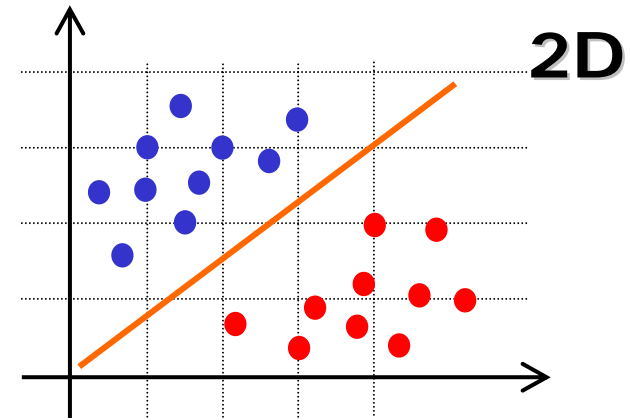
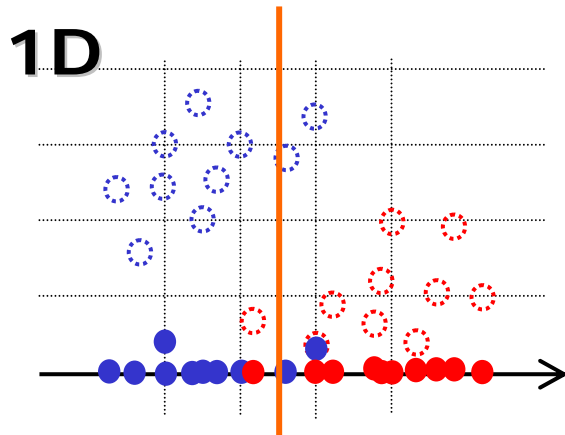
The optimal value of C is determined experimentally, it cannot be readily related to the characteristics of the dataset or model

Nonlinear Support Vector Machines



What if the problem is not linear ?

Feature Space



linearly non-separable on 1D
(**low-dimensional
input space**)

linearly separable on 2D
(**high-dimensional
feature space**)

$$\mathbf{x}_i \rightarrow \phi(\mathbf{x}_i)$$

$$\mathbf{x}_i \cdot \mathbf{x}_j \rightarrow \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$$

Dual Problem and Classifier in the Feature Space (not feasible)

$$\begin{aligned} \max. \quad L_D(\boldsymbol{\alpha}) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \leftarrow \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j) \\ \text{s.t.} \quad \sum_{i=1}^n \alpha_i y_i &= 0 \quad ; \quad 0 \leq \alpha_i \leq C \end{aligned}$$

$$f(\mathbf{x}) = \text{sgn} \left(\sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i^\top \mathbf{x} + b^* \right)$$

- The feature space $\phi(\mathbf{x}_i)$ can be huge or infinite !
- The phi feature has the form of **inner product**

Kernel

Kernel: a function k that takes 2 variables and computes a scalar value (a kind of **similarity**)

$$k(\mathbf{x}, \mathbf{y}) = (\Phi(\mathbf{x}) \cdot \Phi(\mathbf{y}))$$

Kernel Matrix: $m \times m$ matrix \mathbf{K} with elements $\mathbf{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$.

Standard Kernels

Polynomial Kernel $k(\mathbf{x}, \mathbf{x}_i) = (\mathbf{x}^T \mathbf{x}_i + 1)^d$

Radial Basis Function Kernel $k(\mathbf{x}, \mathbf{x}_i) = \exp(-\|\mathbf{x} - \mathbf{x}_i\|^2 / 2\sigma^2)$

Sigmoid Kernel $k(\mathbf{x}, \mathbf{x}_i) = \tanh(\beta_0 \mathbf{x}^T \mathbf{x}_i + \beta_1)$

Mercer's Condition

Valid Kernel functions should satisfy **Mercer's Condition**

For any $g(x)$ for which: $\int g(x)^2 dx < \infty$

It must be the case that: $\int K(x, x') g(x) g(x') dx dx' \geq 0$

A criteria is that the kernel should be **positive semi-definite**

Theorem: If a kernel is **positive semi-definite** i.e.:

$$\sum_{i,j} K(x_i, x_j) c_i c_j \geq 0$$

$\{c_1, \dots, c_n\}$ are real numbers

Then, there **exists a function** $\phi(x)$ defining an inner product of possibly higher dimension i.e.:

$$K(x, y) = \phi(x) \cdot \phi(y)$$

Dual Problem and Classifier with Kernel (Generalized Inner Product SVM)

$$\begin{aligned} \max. \quad L_D(\boldsymbol{\alpha}) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \leftarrow k(\mathbf{x}_i, \mathbf{x}_j) \\ \text{s.t.} \quad \sum_{i=1}^n \alpha_i y_i &= 0 \quad ; \quad 0 \leq \alpha_i \leq C \end{aligned}$$

$$f(\mathbf{x}) = \text{sgn} \left(\sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i^\top \mathbf{x} + b^* \right) \leftarrow k(\mathbf{x}, \mathbf{x}_i)$$

- Kernel Trick works without the mapping

$$\mathbf{x}_i \rightarrow \phi(\mathbf{x}_i)$$

Example: XOR Problem (small scale QP problem)

- Training set: $\{ \underset{\mathbf{x}_1}{(-1 \ -1; -1)}, \underset{\mathbf{x}_2}{(-1 \ +1; +1)}, \underset{\mathbf{x}_3}{(+1 \ -1; +1)}, \underset{\mathbf{x}_4}{(+1 \ +1; -1)} \}$
- Let kernel $k(\mathbf{x}, \mathbf{x}_i) = (1 + \mathbf{x}^T \mathbf{x}_i)^2$, $\mathbf{x} = (x_1, x_2)^T$, $\mathbf{x}_i = (x_{i1}, x_{i2})^T$
- Then $k(\mathbf{x}, \mathbf{x}_i) = (1 + \mathbf{x}^T \mathbf{x}_i)^2 = (1 + x_1 x_{i1} + x_2 x_{i2})^2$
 $= 1 + x_1^2 x_{i1}^2 + 2x_1 x_2 x_{i1} x_{i2} + x_2^2 x_{i2}^2 + 2x_1 x_{i1} + 2x_2 x_{i2}$
- A Mapping: $\varphi(\mathbf{x}) = (1, x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2)^T$

- Kernel Matrix:

$$\Phi = \begin{bmatrix} 1 & 1 & 1 & -\sqrt{2} & -\sqrt{2} & \sqrt{2} \\ 1 & 1 & 1 & -\sqrt{2} & \sqrt{2} & -\sqrt{2} \\ 1 & 1 & 0 & \sqrt{2} & -\sqrt{2} & -\sqrt{2} \\ 1 & 1 & 1 & \sqrt{2} & \sqrt{2} & \sqrt{2} \end{bmatrix}$$

$$K = \begin{bmatrix} 9 & 1 & 1 & 1 \\ 1 & 9 & 1 & 1 \\ 1 & 1 & 9 & 1 \\ 1 & 1 & 1 & 9 \end{bmatrix} 4 \times 4$$

Example: XOR Problem

- Dual Problem:

$$L_D(\alpha) = \alpha_1 + \alpha_2 + \alpha_3 + \alpha_4 - (9\alpha_1^2 - 2\alpha_1\alpha_2 - 2\alpha_1\alpha_3 + 2\alpha_1\alpha_4 + 9\alpha_2^2 + 2\alpha_2\alpha_3 - 2\alpha_2\alpha_4 + 9\alpha_3^2 - 2\alpha_3\alpha_4 + 9\alpha_4^2)/2$$

- Optimizing L_D : $9\alpha_1 - 2\alpha_2 - \alpha_3 + \alpha_4 = 1$, $-\alpha_1 + 9\alpha_2 + \alpha_3 - \alpha_4 = 1$
 $-\alpha_1 + \alpha_2 + 9\alpha_3 - \alpha_4 = 1$, $\alpha_1 - \alpha_2 - \alpha_3 + 9\alpha_4 = 1$

- Optimal Lagrange multipliers: $\alpha_1 = \alpha_2 = \alpha_3 = \alpha_4 = 1/8 > 0$

All the samples are support vectors

- Optimal \mathbf{w} :

$$\mathbf{w} = \sum \alpha_i y_i \boldsymbol{\varphi}(\mathbf{x}_i) = (1/8)(-1) \boldsymbol{\varphi}(\mathbf{x}_1) + (1/8)(+1) \boldsymbol{\varphi}(\mathbf{x}_2) + (1/8)(+1) \boldsymbol{\varphi}(\mathbf{x}_3) + (1/8)(-1) \boldsymbol{\varphi}(\mathbf{x}_4) = (0 \ 0 \ 0 \ 0 \ 0 \ -1\sqrt{2})^T$$

- Optimal Hyperplane:

$$\begin{aligned} \underline{f}(\mathbf{x}) &= \text{sgn}(\mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x})) \\ &= \text{sgn}(-x_1 x_2) \end{aligned}$$

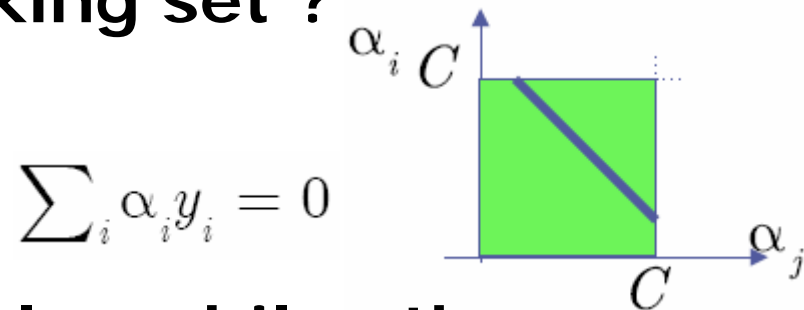
Sequential Minimal Optimization (Platt '98)

$$\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j k(x_i, x_j) = -\frac{1}{2} \alpha^T Q \alpha + \alpha^T \mathbf{1}$$

subject to $\sum_i \alpha_i y_i = 0 \quad \& \quad \alpha_i \in [0, C]$

In case of large scale QP problem

- divide a large QP into a series of smaller QP sub-problems and optimize them sequentially
- What is the smallest working set ?
- Update just 2 Lagrange multipliers at a time
- Updating subset of variables while others fixed will also converge globally



Sequential Minimal Optimization

- Write dual prob. as a function of just 2 alphas

$$L_D \propto \alpha_i + \alpha_j - \frac{1}{2} \left(K_{ii} \alpha_i^2 + 2K_{ij} \alpha_i \alpha_j + K_{jj} \alpha_j^2 \right) - h_i \alpha_i - h_j \alpha_j$$

$$\text{subject to : } y_i \alpha_i + y_j \alpha_j + \sum_{t \neq i, j} y_t \alpha_t = 0 \quad \text{and } \alpha_i, \alpha_j \in [0, C]$$

- Update Rules

- no numerical part
- memory for buffering E

$$S = y_i y_j$$

$$L = \max \left(0, \alpha_j + S \alpha_i - \frac{1}{2} (S + 1) C \right)$$

$$E1 = \sum_t \alpha_t y_t k(x_i, x_t) + b - y_i$$

$$H = \min \left(C, \alpha_j + S \alpha_i - \frac{1}{2} (S - 1) C \right)$$

$$E2 = \sum_t \alpha_t y_t k(x_j, x_t) + b - y_j$$

$$\alpha_j^{NEW} = \alpha_j + \frac{y_j (E1 - E2)}{k(x_i, x_i) + k(x_j, x_j) - 2k(x_i, x_j)} \quad \text{clipped inside } [L, H]$$

$$\alpha_i^{NEW} = \alpha_i + S \left(\alpha_j - \alpha_j^{NEW} \right)$$

Link to Statistical Learning Theory (Vapnik '95)

Learn f from training set is to minimize the following:

$$R[f] = \int \frac{1}{2} |f(\mathbf{x}) - y| dP(\mathbf{x}, y) \quad \text{Expected Risk}$$

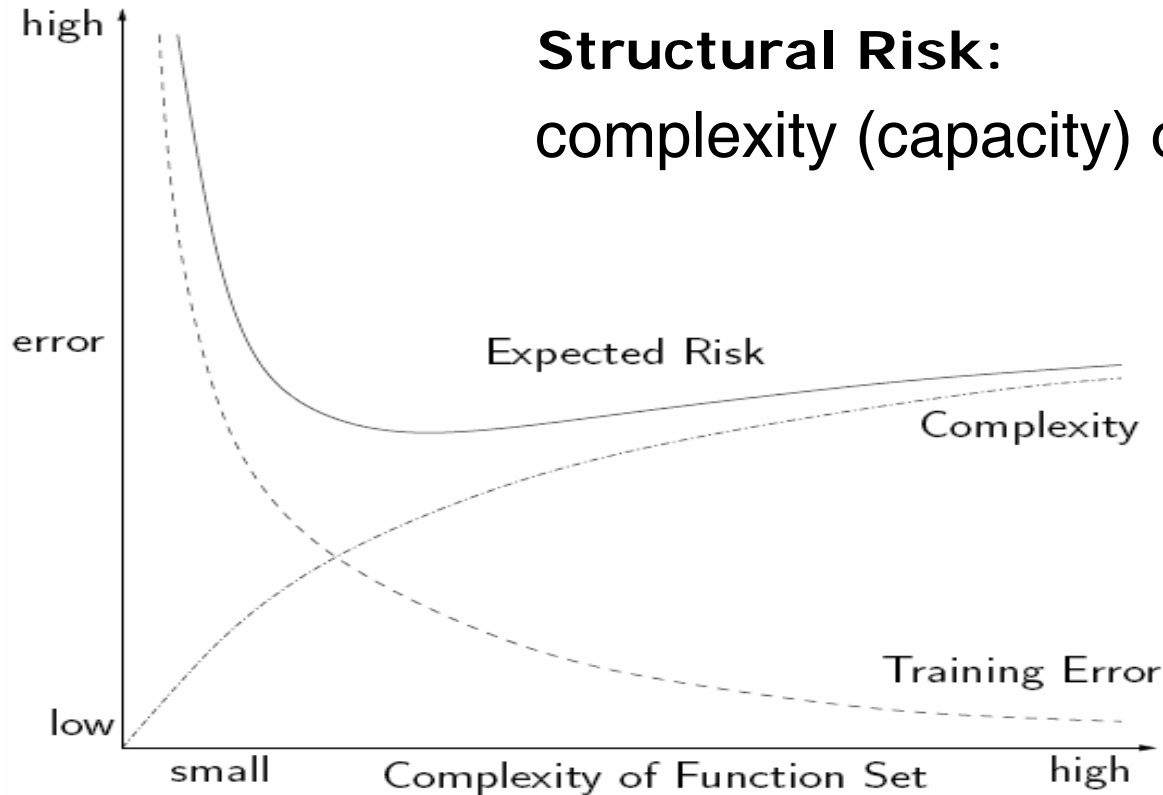
\uparrow
 $\{\pm 1\}$ Unknown

Minimize instead the average risk over the training set:

$$R_{\text{emp}}[f] = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} |f(\mathbf{x}_i) - y_i| \quad \text{Empirical Risk}$$

A Risk Bound

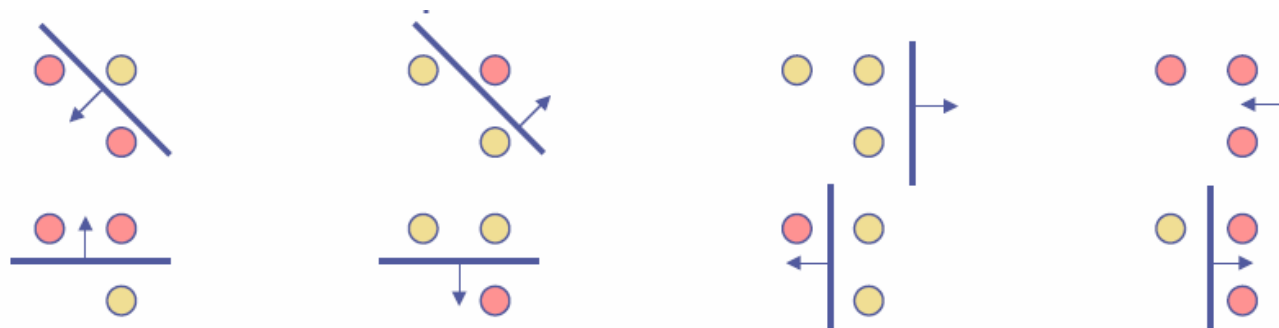
$$R[f] \leq \underbrace{R_{\text{emp}}[f]}_{\text{minimize}} + \underbrace{\sqrt{\frac{h(\ln \frac{2m}{h} + 1) - \ln(\eta/4)}{m}}}_{\text{minimize}}$$



Structural Risk:
complexity (capacity) of func. set

Vapnik-Chervonenkis (VC) Dimension

- Maximum number of points that can be labeled in all possible way
- VC dimension of linear classifiers in N -dimensions is $h=N+1$



Lines(dichotomies) can **shatter** 3 points in 2d

- Measure of Complexity of Function Set

Minimizing VC dim. → Minimizing Complexity

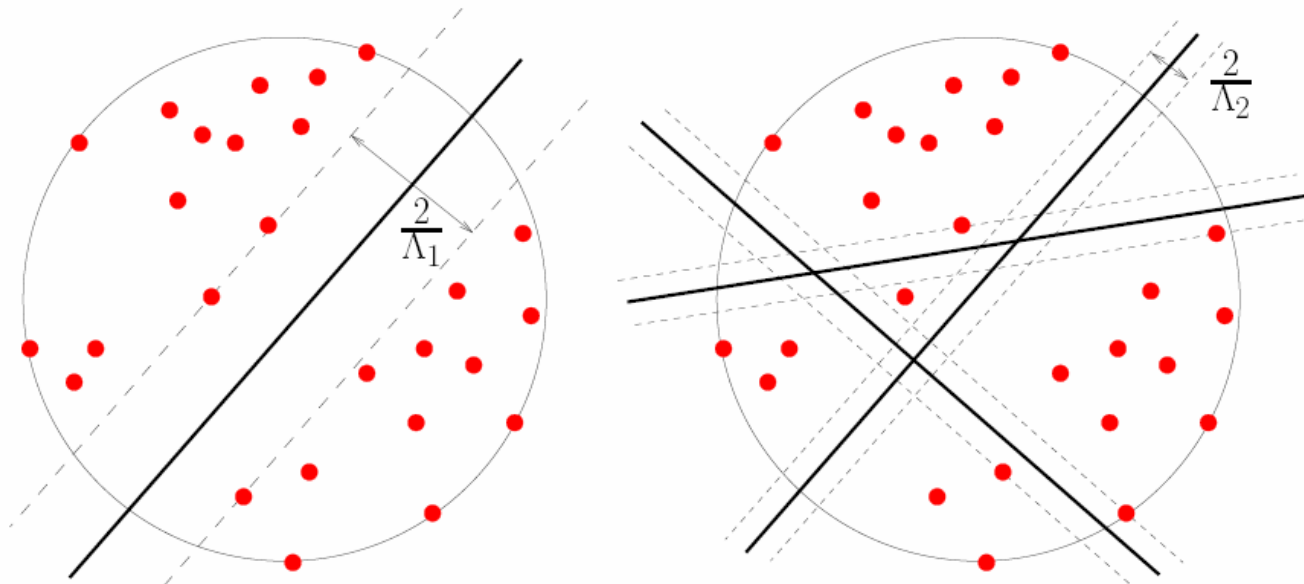
VC Dimension of Margin Hyperplanes

- VC dimension satisfies the following:

$$h \leq \min(R^2 \Lambda^2 + 1, N + 1)$$

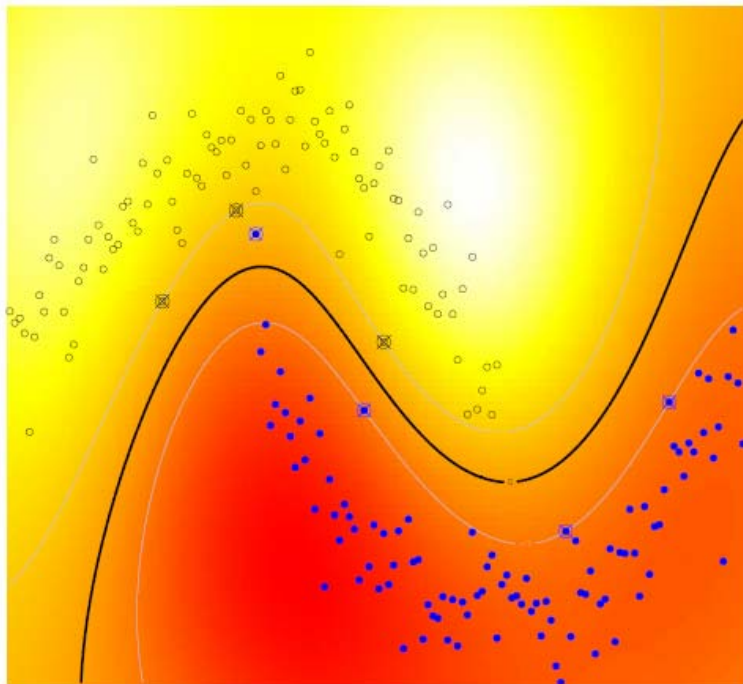
R is the smallest sphere containing a set of points

$$\|\mathbf{w}\| \leq \Lambda$$

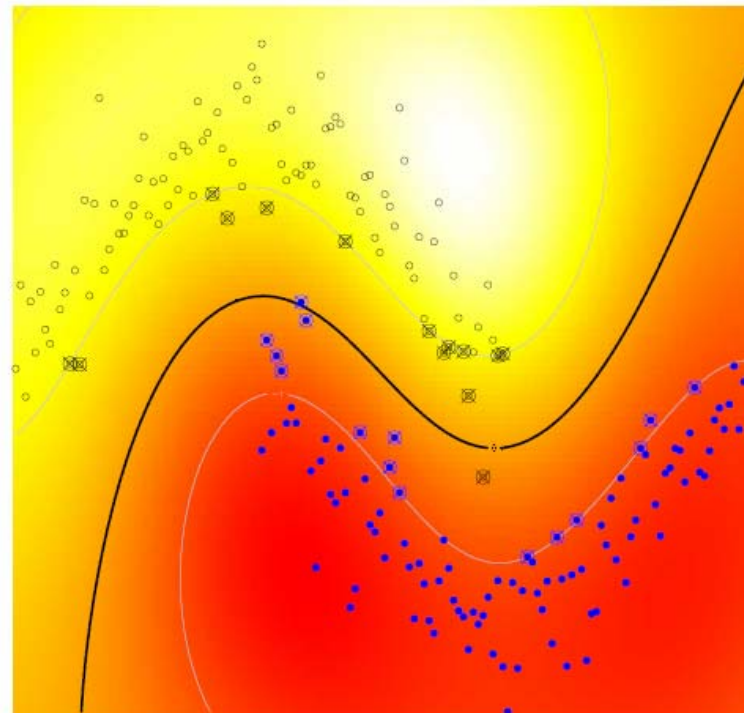


Maximizing Margin \rightarrow Minimizing VC dim.

Results for Gaussian Kernel

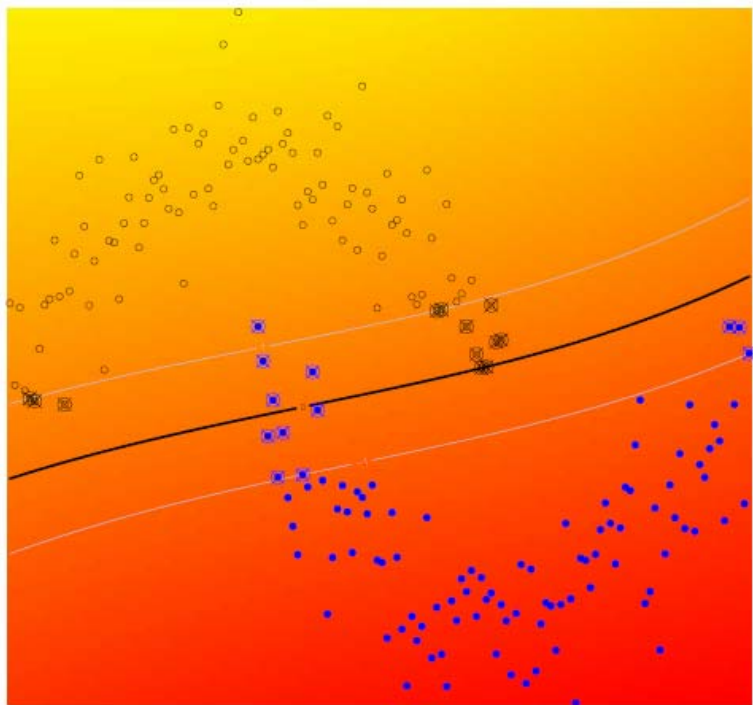


$$\sigma = 0.5, C = 50$$

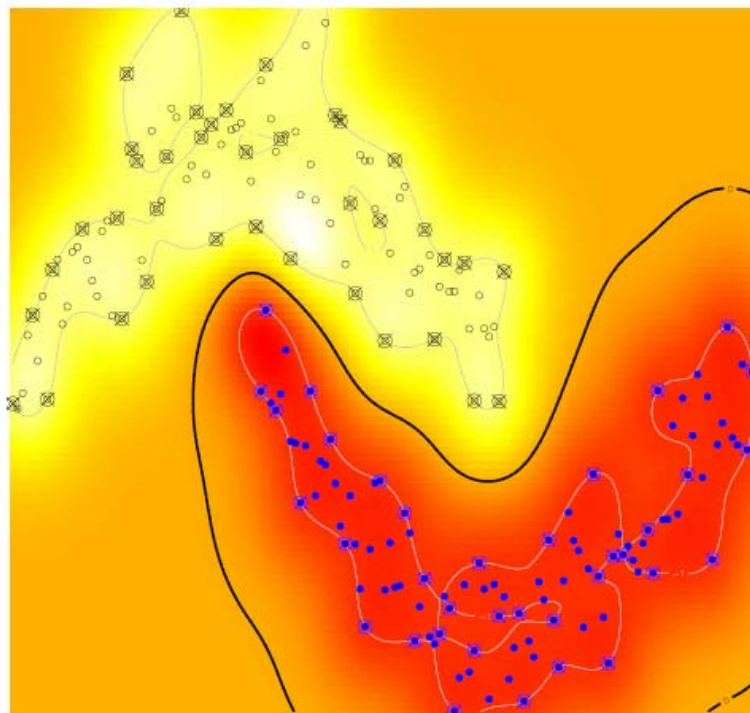


$$\sigma = 0.5, C = 1$$

Results for Gaussian Kernel



$$\sigma = 0.02, C = 50$$



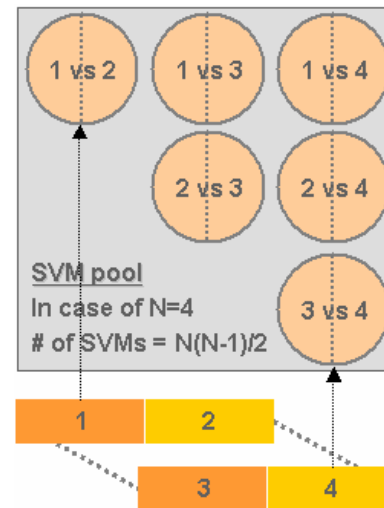
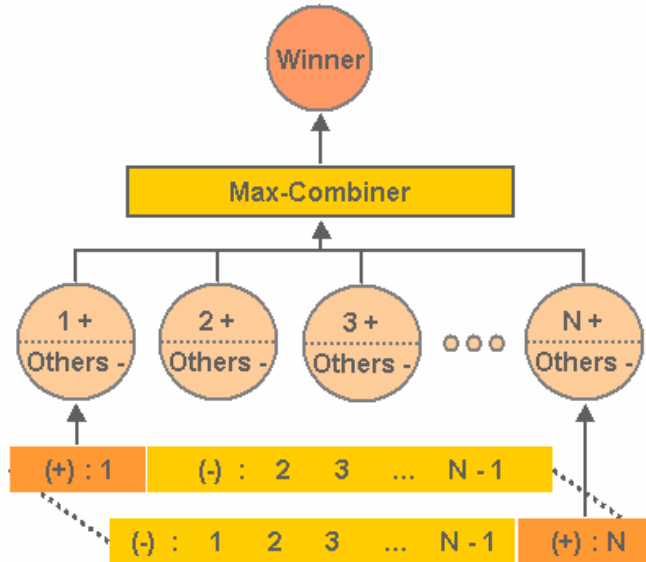
$$\sigma = 10, C = 50$$

Summary

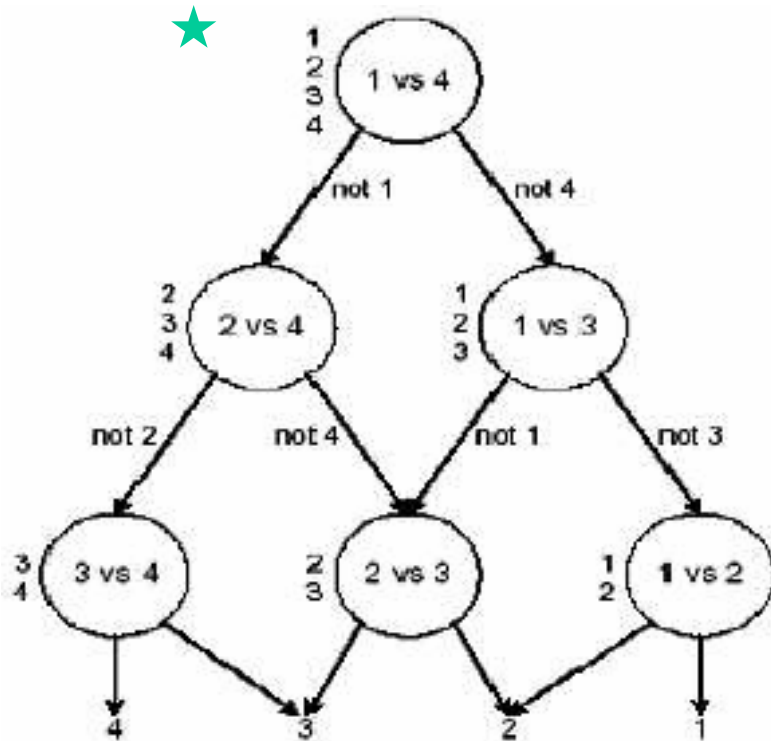
- **Optimal Separating Hyperplane (Margin)**
- **Global Minimum Solution (Convexity)**
- **Only SVs are Relevant (Sparseness)**
- **Automatically selects SVs; # of SVs can be considered as # of hidden units of MLP**
- **Model selection problem; kernel selection**
- **Training speed and method for a large training set**
- **Binary classifier**

Multiclass Support Vector Machines (Multiple Classifiers System)

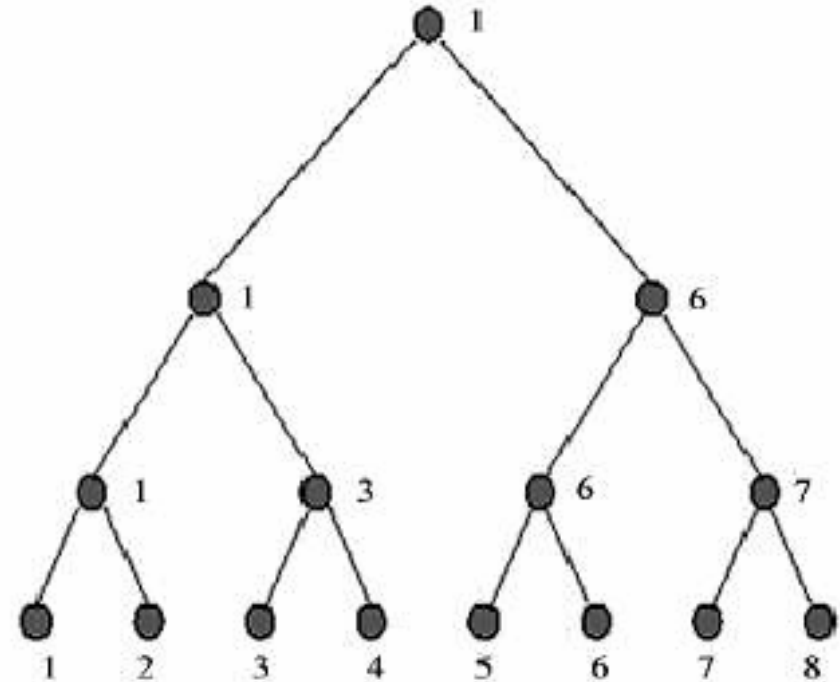
- Ensemble of binary support vector machines
- Categorized by Coding / Decoding
- OPC (one-per-class), PWC(pair-wise coupling), ECOC (error correcting output coding)



Tree-based Methods



(a) Example of top-down tree structure (DDAG)



(b) Example of bottom-up tree structure (Pairwise SVM)

Open Software

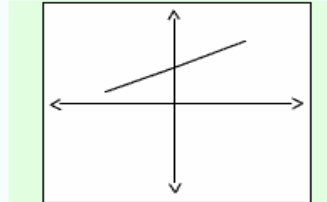
Software	Developer	Language	Environment	Algorithms	URL
SVMFu	R. Rifkin M. Nadermann (MIT)	C++	Unix-like system	Osuna <i>et al.</i> , SMO(Platt)	http://www.ai. mit.edu
LIBSVM	C.C. Chang, C.H. Lin (National Taiwan Univ.)	C++, Java	Python, R, Matlab, Perl	SMO(Platt), SVMLight(Joachims)	http://www.csie. ntu.edu.tw/~libsvm
SVMLight	T. Joachims, (Univ. of Dortmund)	C	Solaris, Linux, IRIX, Windows NT	T. Joachims	http://www.svmlight. joachims.org
SVM Torch	R. Collobert, (IDIAP, Switzerland)	C, C++	Windows	R. Collobert	http://www.idap.ch /learning/SVMTorch.html

Demo

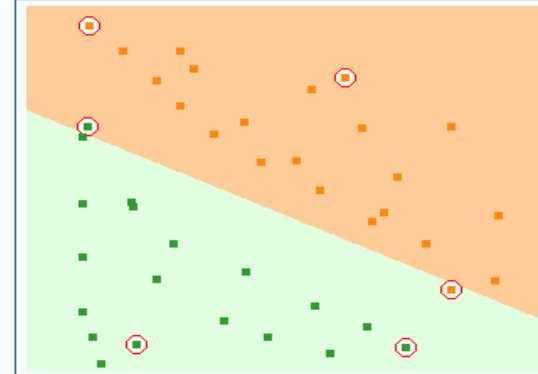
Kernel Options

Kernel :

Formula : $(s \cdot a \cdot b + c)$



SVM Data Display



Detailed SVM Options

Stopping Criteria (Epsilon) :

Coefficient of the Error Term (C) :

Use Shrinking :

Cache In Mega Bytes :

SVM Results Display

```
#Iterations : 5  
rho : [-1,0748244104614357]  
probA : -  
probB : -  
number of support vectors : 6
```

SVM Applet...

Developed for :
EE-583 Pattern
Recognition

Developed by :
Hakan Serçe, 2005

This applet demonstrates
SVM (Support Vector

- <http://www.eee.metu.edu.tr/~alatan/Courses/Demo/AppletSVM.html>
- <http://www.csie.ntu.edu.tw/~cjlin/libsvm/#GUI>

Applications

- **Biometrics**
- **Object Detection and Recognition**
- **Character Recognition**
- **Information and Image Retrieval**
- **Other Applications**

References

S.Haykin, Neural Networks A Comprehensive Foundation, Prentice Hall, 1999. Chap. 6

Burges, C. J. C., "A Tutorial on Support Vector Machines for Pattern Recognition," Data Mining and Knowledge Discovery, Vol. 2, pp. 121~167, 1998.

John C. Platt, "Sequential Minimal Optimization: A Fast Algorithm for Training Support Vector Machines," Microsoft Research, Technical Report MSR-TR-98-14, 1998

C.-W. Hsu and C.-J. Lin. "A Comparison of Methods for Multi-Class Support Vector Machines," IEEE Transactions on Neural Networks, 13, 415-425, 2002.

Books

Vladimir N. Vapnik, *The Nature of Statistical Learning Theory*, Springer-Verlag, 1995.

Nello Cristianini and John Shawe-Taylor, *An Introduction to Support Vector Machines and Other Kernel-based Learning Methods*, Cambridge University Press, 2000.

Bernhard Scholkopf and Alexander J. Smola, *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*, MIT Press, 2001.

Links

<http://videlectures.net/>

<http://www.kernel-machines.org/>

감사합니다