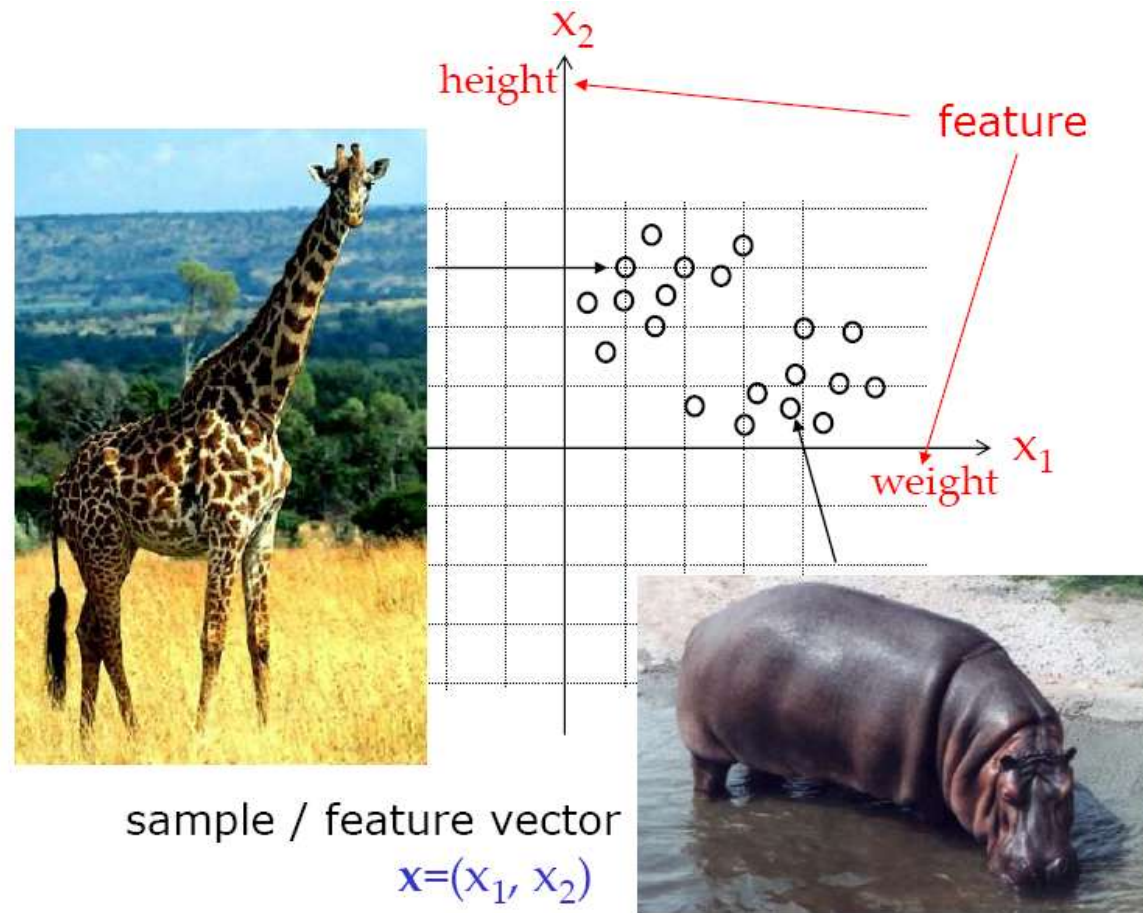


# Support Vector Machine

## Part 1

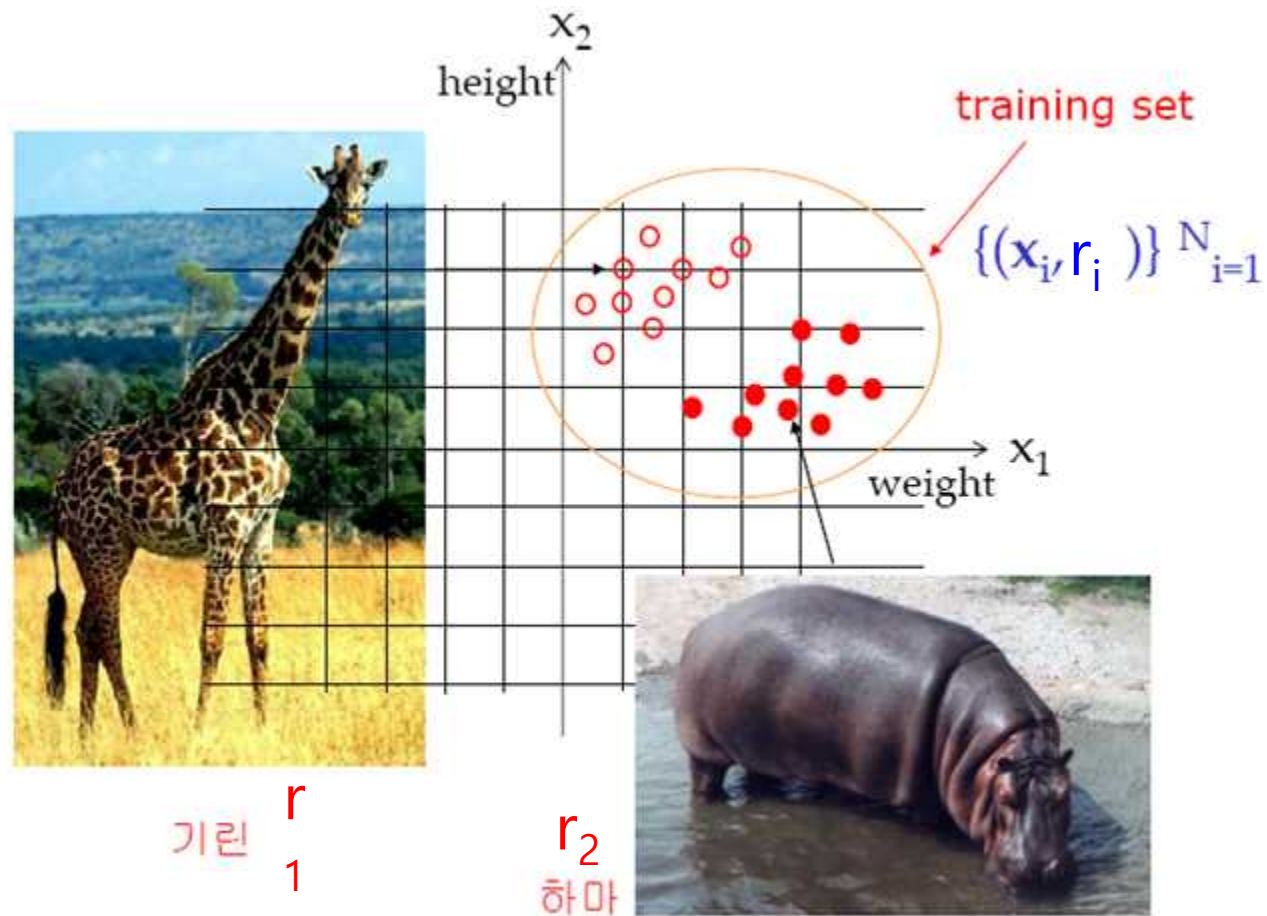
# Feature Space

- Sample



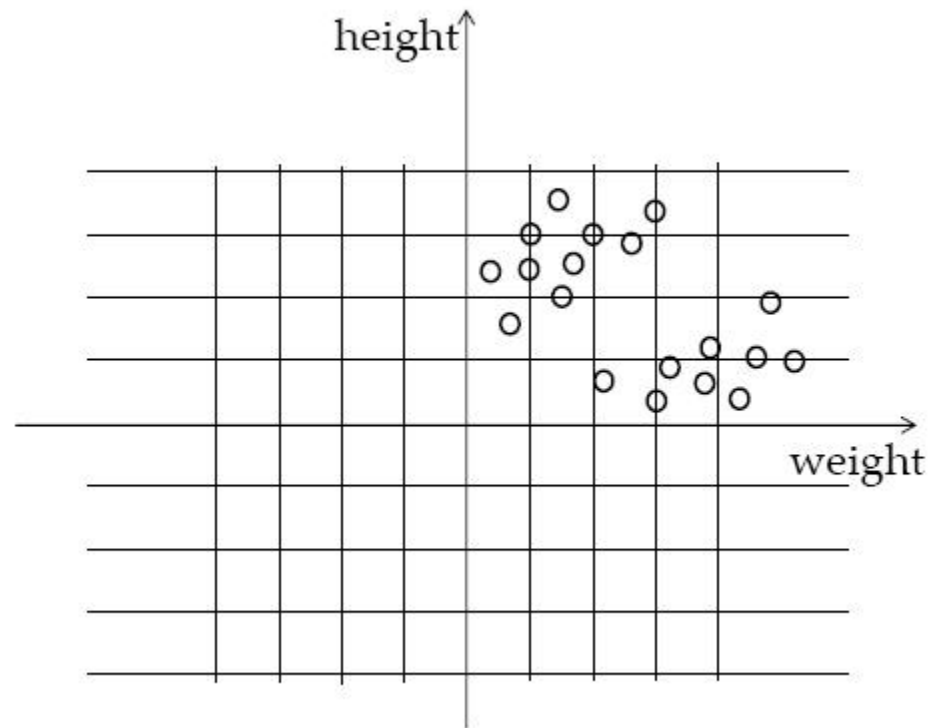
# Feature Space

- Training Set



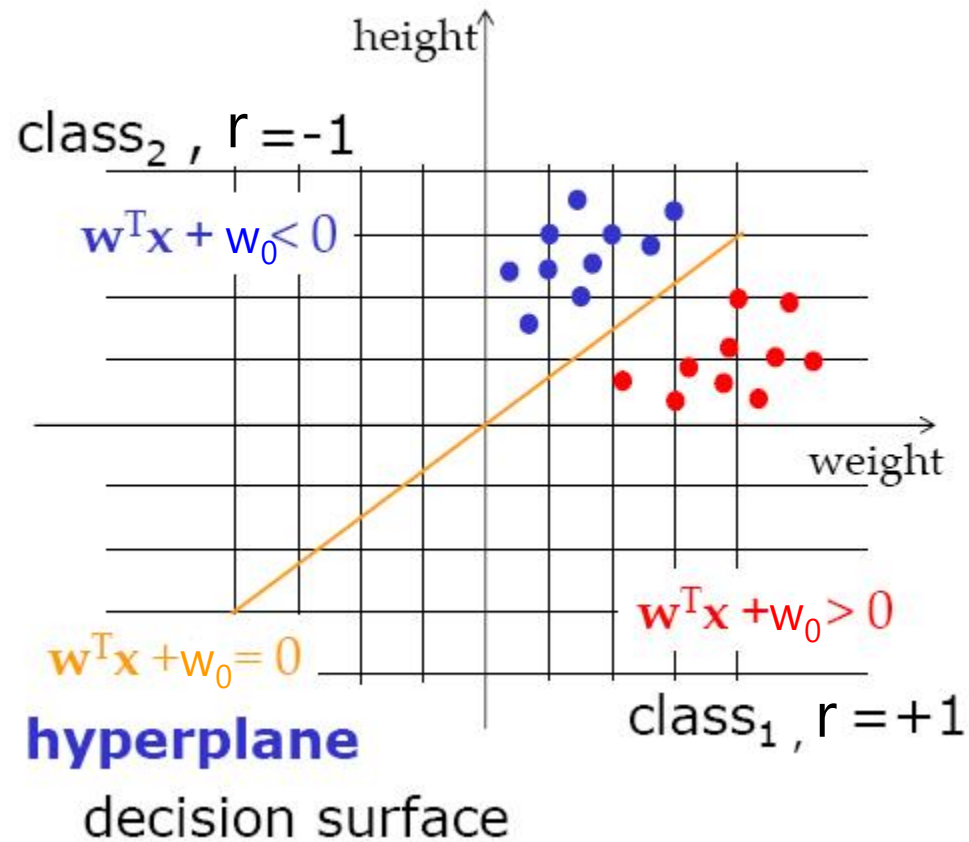
# Feature Space

- How to classify them using computer ?



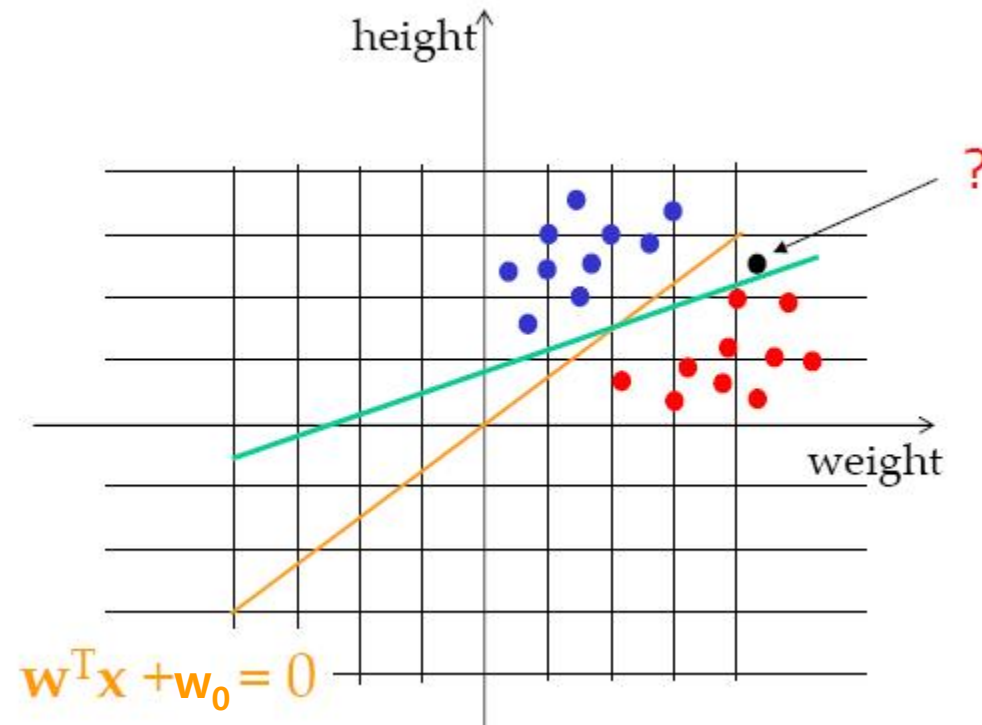
# Optimal Hyperplane

- Linear classification



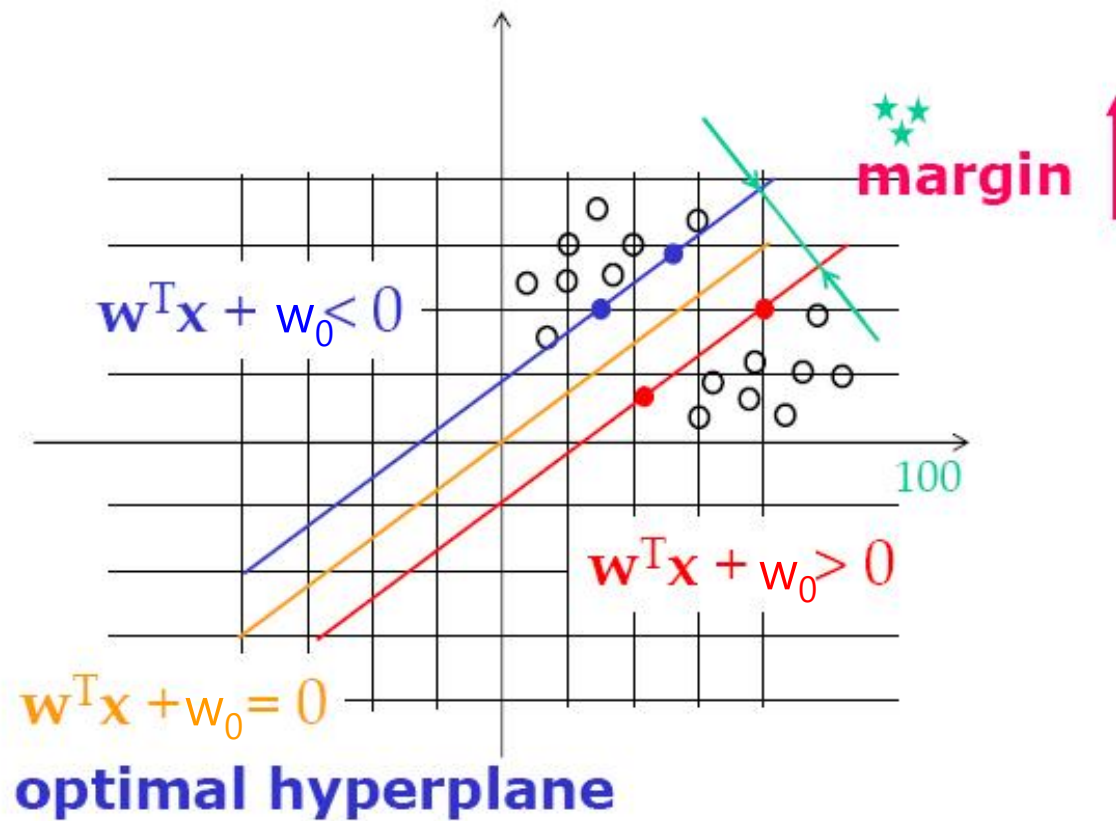
# Optimal Hyperplane

- Linear classification



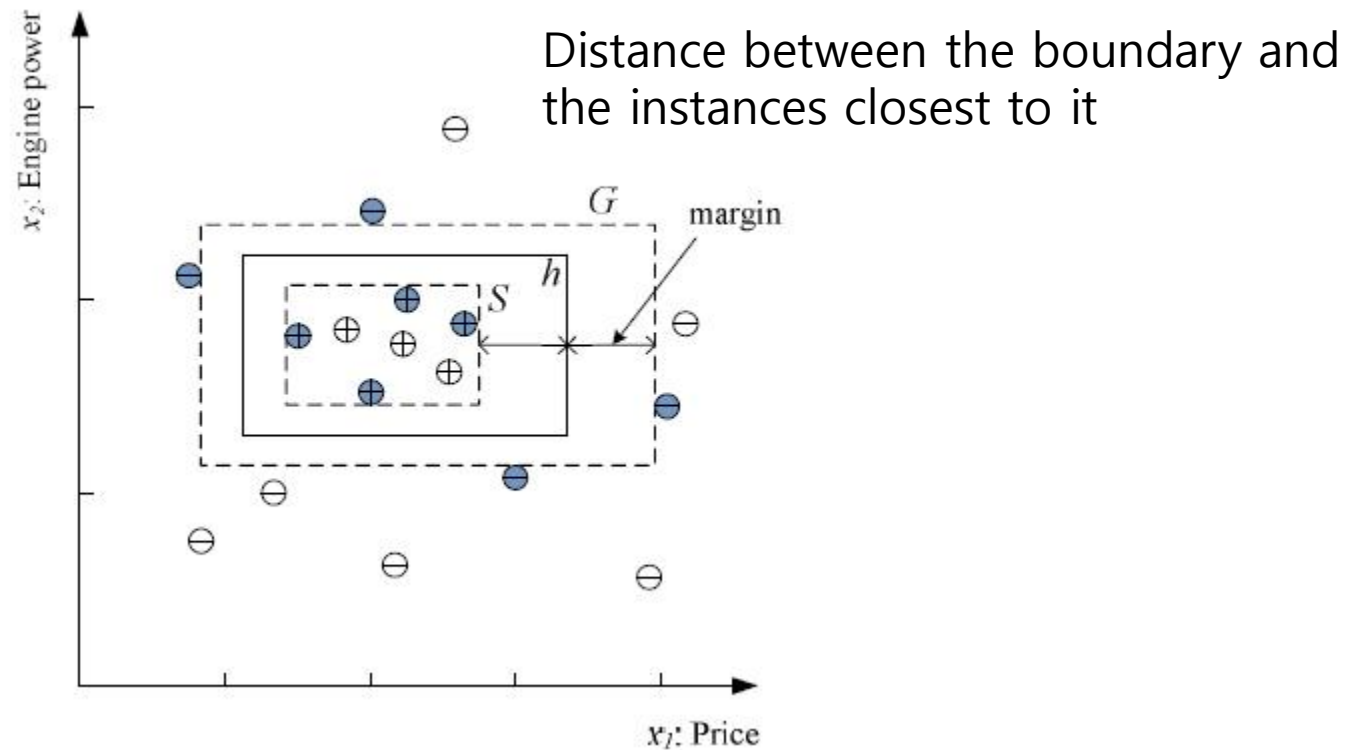
# Optimal Hyperplane

- Margin



# Optimal Hyperplane

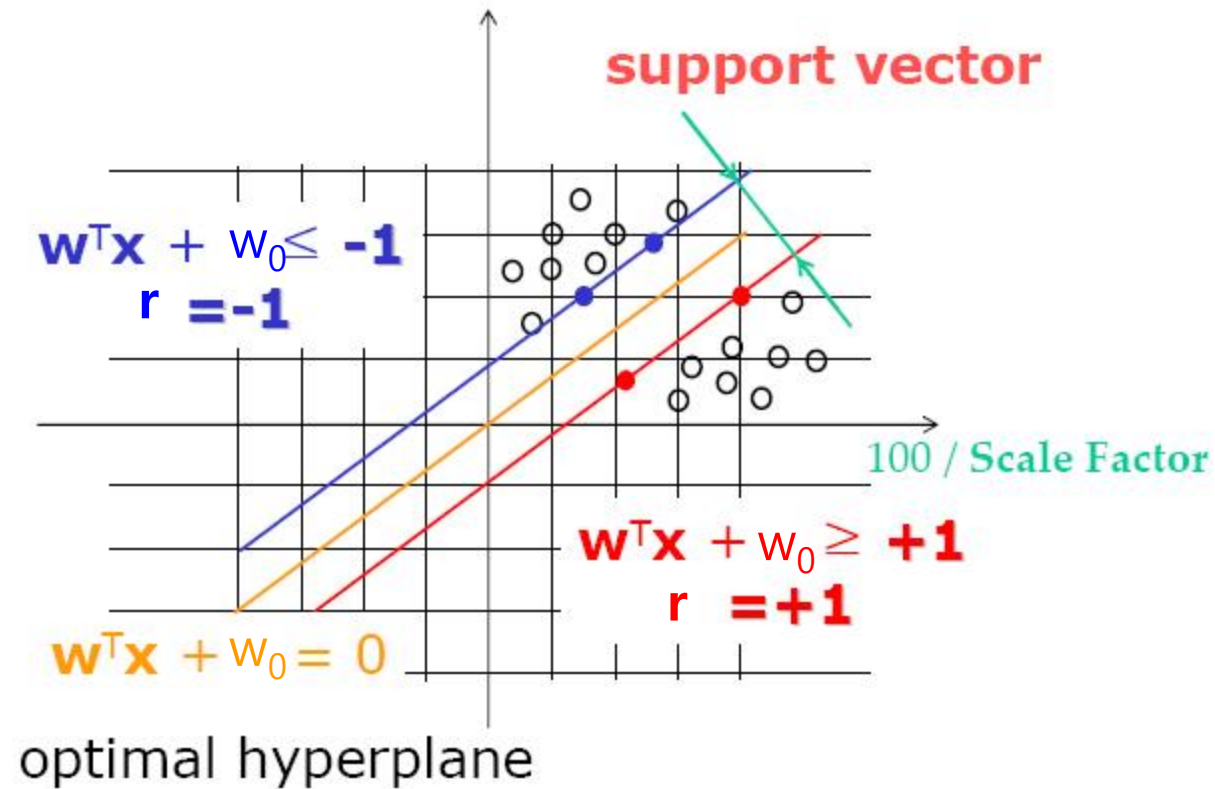
- Choose  $h$  with largest margin
  - Minimize the error function





# Optimal Hyperplane

- Support vector



# Optimization Problem

$$\mathcal{X} = \{\mathbf{x}^t, r^t\}_t \text{ where } r^t = \begin{cases} +1 & \text{if } \mathbf{x}^t \in C_1 \\ -1 & \text{if } \mathbf{x}^t \in C_2 \end{cases}$$

find  $\mathbf{w}$  and  $w_0$  such that

$$\mathbf{w}^T \mathbf{x}^t + w_0 \geq +1 \text{ for } r^t = +1$$

$$\mathbf{w}^T \mathbf{x}^t + w_0 \leq -1 \text{ for } r^t = -1$$

which can be rewritten as

$$r^t (\mathbf{w}^T \mathbf{x}^t + w_0) \geq +1$$

(Cortes and Vapnik, 1995; Vapnik, 1995)

# Optimization Problem

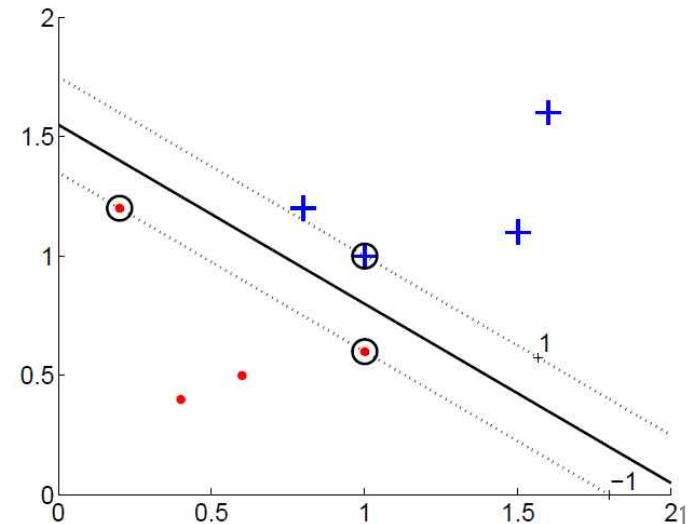
- Margin

- Distance from the discriminant to the closest instances on either side

- Distance of  $x^t$  to the hyperplane is
- Margin for support vectors

$$\frac{|\mathbf{w}^T \mathbf{x}^t + w_0|}{\|\mathbf{w}\|}$$

- Maximize  $t \left| \frac{1}{\|\mathbf{w}\|} - \frac{-1}{\|\mathbf{w}\|} \right| = \frac{2}{\|\mathbf{w}\|}$   
⇒ Minimize the  $\|\mathbf{w}\|$



# Optimization Problem

- Constrained optimization problem

$$\min \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{s.t. } r^t (\mathbf{w}^T \mathbf{x}^t + w_0) \geq +1, \forall t$$

: Convex cost function

: Linear constraints

- Standard quadratic optimization problem
- **Convex optimization problem**
  - : Local solution  $\Rightarrow$  Global optimal solution

# Optimization Problem

- Lagrange Multipliers Method

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

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

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

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:

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

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

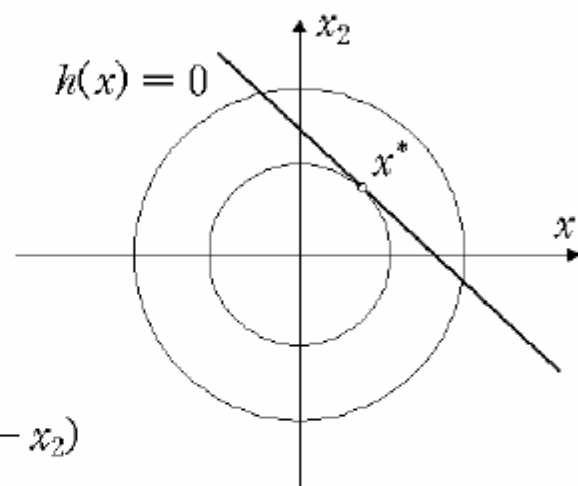
Solution: Define the Lagrange function

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

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

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

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



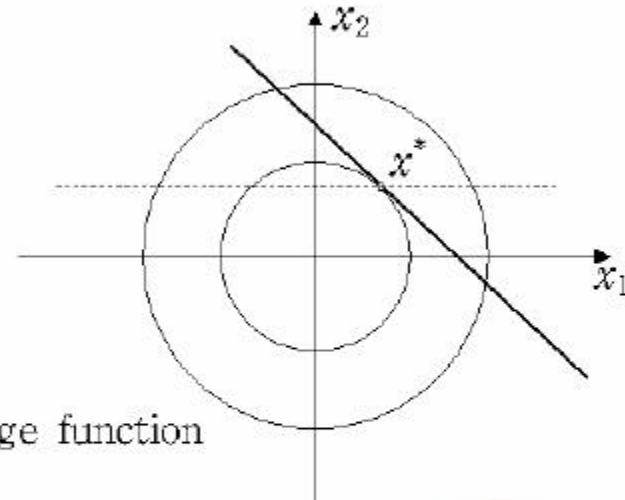
# A Brief Overview of Optimization Theory

► 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

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

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

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



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  $x_1 = x_2 = \frac{1}{2} \times (\because x_2 \geq \frac{3}{4})$

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  $\Omega \subseteq \mathbb{R}^n$

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

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.:

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

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

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

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

# Optimization Problem

- Problem

$$\begin{aligned} \min & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} & r^t (\mathbf{w}^T \mathbf{x}^t + w_0) \geq +1, \forall t \end{aligned}$$

- Lagrange function (Primal)

$$\begin{aligned} L_p &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{t=1}^N \alpha^t [r^t (\mathbf{w}^T \mathbf{x}^t + w_0) - 1] \\ &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{t=1}^N \alpha^t r^t (\mathbf{w}^T \mathbf{x}^t + w_0) + \sum_{t=1}^N \alpha^t \end{aligned}$$

- Solution

$$\begin{aligned} \frac{\partial L_p}{\partial \mathbf{w}} = 0 &\Rightarrow \mathbf{w} = \sum_{t=1}^N \alpha^t r^t \mathbf{x}^t \\ \frac{\partial L_p}{\partial w_0} = 0 &\Rightarrow \sum_{t=1}^N \alpha^t r^t = 0 \end{aligned}$$

$$\text{KKT condition} \Rightarrow \alpha^t \left( 1 - r^t (\mathbf{w}^T \mathbf{x}^t + w_0) \right) = 0$$

# Support Vector Machine

- Support Vector Machine (SVM) (Vapnik 1995)

$$\text{KKT condition : } \alpha^t \left( 1 - r^t (w^T x^t + w_0) \right) = 0$$

$$\text{If } \alpha^t \neq 0 \text{ then } r^t (w^T x^t + w_0) = 1$$

$$\Rightarrow x^t \text{ is support vector}$$

$$\text{If } r^t (w^T x^t + w_0) \neq 1 \text{ then } \alpha^t = 0$$

$$\Rightarrow x^t \text{ is not support vector}$$

$$w = \sum_{t=1}^N \alpha^t r^t x^t$$

:  $w$  is only related to **support vectors**  $x^t$

Most  $\alpha^t$  are 0 and only a small number have  $\alpha^t > 0$ ; they are the **support vectors**

# Support Vector Machine

- Dual Problem

$$\begin{aligned}L_d(\alpha^t) &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{t=1}^N \alpha^t [r^t (\mathbf{w}^T \mathbf{x}^t + w_0) - 1] \\&= \frac{1}{2} (\mathbf{w}^T \mathbf{w}) - \mathbf{w}^T \sum_t \alpha^t r^t \mathbf{x}^t - w_0 \sum_t \alpha^t r^t + \sum_t \alpha^t \\&= -\frac{1}{2} (\mathbf{w}^T \mathbf{w}) + \sum_t \alpha^t\end{aligned}$$

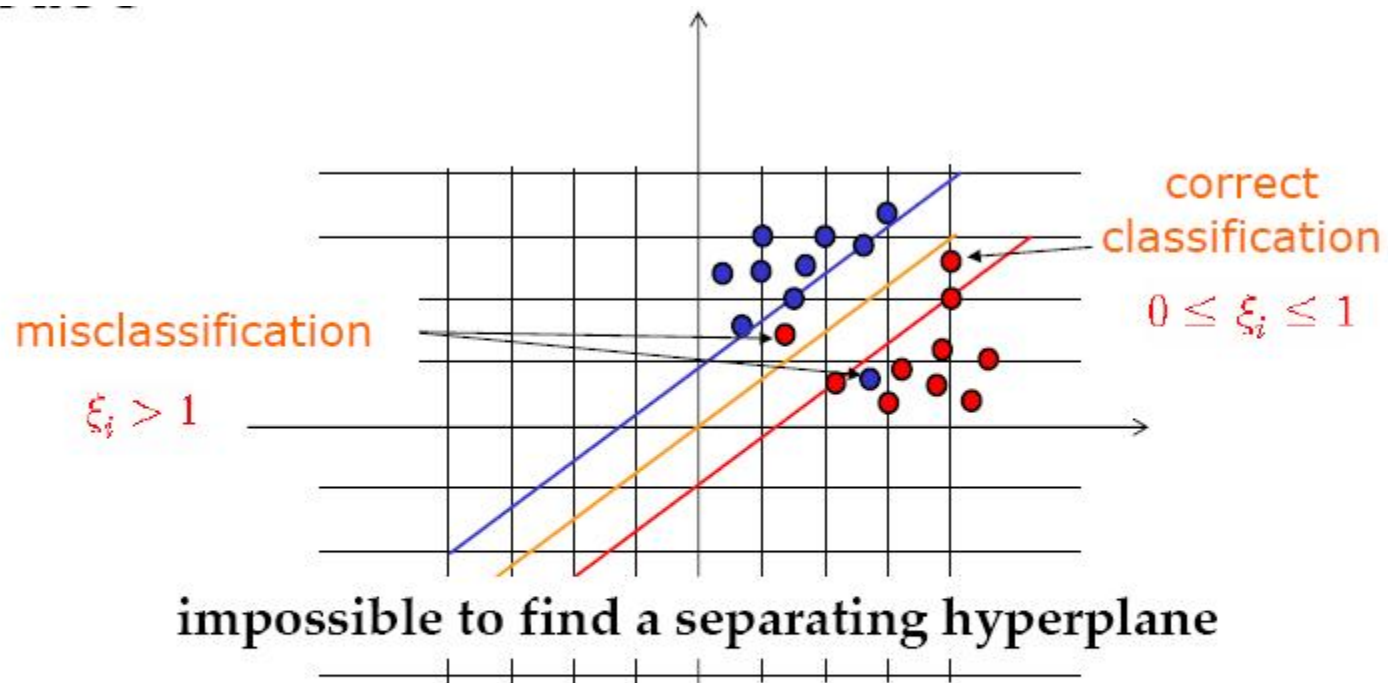
$$= -\frac{1}{2} \sum_t \sum_s \alpha^t \alpha^s r^t r^s (\mathbf{x}^t)^T \mathbf{x}^s + \sum_t \alpha^t$$

subject to  $\sum_t \alpha^t r^t = 0$  and  $\alpha^t \geq 0, \forall t$

$$\alpha^{t*} = \arg \min_{\alpha^t} L_d$$

# Soft Margin Hyperplane

- Non-separable case



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

# Soft Margin Hyperplane

- Soft Error

- Problem: can't satisfy  $r^t(\mathbf{w}^T \mathbf{x}^t + w_0) \geq 1, \forall t$
- Relaxing the equation using slack variables  $\xi^t \geq 0$

$$r^t(\mathbf{w}^T \mathbf{x}^t + w_0) \geq 1 - \xi^t$$

$\xi^t = 0$ : no problem

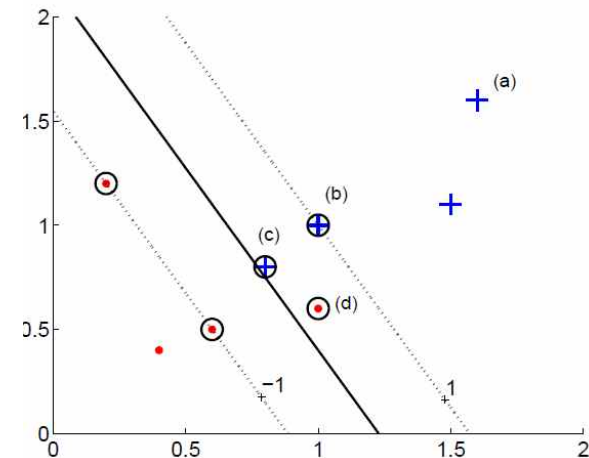
$0 < \xi^t < 1$ : correct classification

$\xi^t \geq 1$ : misclassification

- Soft error

$$\sum_t \xi^t$$

$\#\{\xi^t > 1\}$ : number of misclassification points



# Soft Margin Hyperplane

- Cost Function

$$\text{minimize } \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_t \xi^t$$

- Lagrange function

$$L_p = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_t \xi^t - \sum_t \alpha^t [r^t (\mathbf{w}^T \mathbf{x}^t + w_0) - 1 + \xi^t] - \sum_t \mu^t \xi^t$$

- Solution

$$\frac{\partial L_p}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{t=1}^N \alpha^t r^t \mathbf{x}^t$$

$$\frac{\partial L_p}{\partial w_0} = 0 \Rightarrow \sum_{t=1}^N \alpha^t r^t = 0$$

$$\frac{\partial L_p}{\partial \xi^t} = 0 \Rightarrow C - \alpha^t - \mu^t = 0$$

$$0 \leq \alpha^t \leq C \quad (\because \mu^t \geq 0)$$

The optimal value of  $C$  is determined **experimentally**



# $\nu$ (nu) -SVM

- Replacing  $C$  with the parameter  $\nu$  (Scholkopf, 2000)

$$\begin{aligned} \min & \frac{1}{2} \|\mathbf{w}\|^2 - \nu \rho + \frac{1}{N} \sum_t \xi^t \\ \text{subject to} & \\ & r^t (\mathbf{w}^T \mathbf{x}^t + w_0) \geq \rho - \xi^t, \xi^t \geq 0, \rho \geq 0 \end{aligned}$$

- $\rho$  : optimization variable (scales the margin),  
margin size =  $\rho / \|\mathbf{w}\|$
- $\nu \in [0,1]$  : parameter  
lower bound on the fraction of support vectors  
upper bound on the fraction of instances having margin errors

$$\begin{aligned} \sum_t \alpha^t r^t &= 0 \\ 0 &\leq \alpha^t \leq \frac{1}{N} \\ \sum_t \alpha^t &\leq \nu \end{aligned}$$

# Nonlinear SVM

- Nonlinear SVM

- Basis function

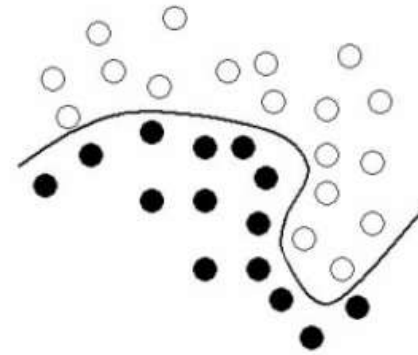
$$\mathbf{z} = \boldsymbol{\varphi}(\mathbf{x}) \quad , z_j = \phi_j(\mathbf{x}), j = 1 \dots, k$$

$$g(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}) = \sum_{j=1}^k w_j \phi_j(\mathbf{x})$$

- Solution

$$\mathbf{w} = \sum_t \alpha^t r^t \boldsymbol{\varphi}(\mathbf{x}^t)$$

$$\begin{aligned} g(\mathbf{x}) &= \mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}) = \sum_t \alpha^t r^t \boldsymbol{\varphi}(\mathbf{x}^t)^T \boldsymbol{\varphi}(\mathbf{x}) \\ &= \sum_t \alpha^t r^t K(\mathbf{x}^t, \mathbf{x}) \end{aligned}$$



# Nonlinear SVM

## 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  $k(\mathbf{x}, \mathbf{x}_i) = \exp(-\|\mathbf{x} - \mathbf{x}_i\|^2 / 2\sigma^2)$

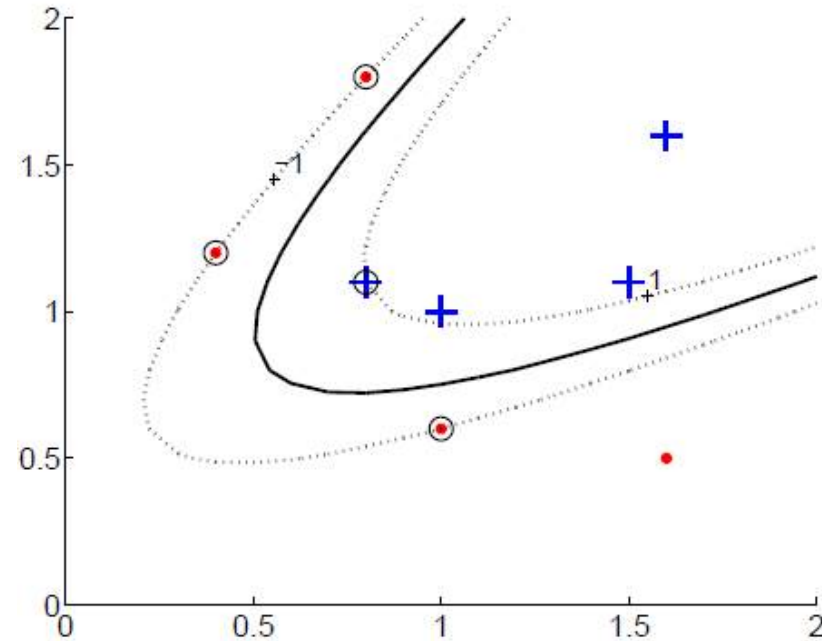
Sigmoid Kernel  $k(\mathbf{x}, \mathbf{x}_i) = \tanh(\beta_0 \mathbf{x}^T \mathbf{x}_i + \beta_1)$

# Nonlinear SVM

- Polynomials of degree  $q$ .

$$K(\mathbf{x}^t, \mathbf{x}) = (\mathbf{x}^T \mathbf{x}^t + 1)^q$$

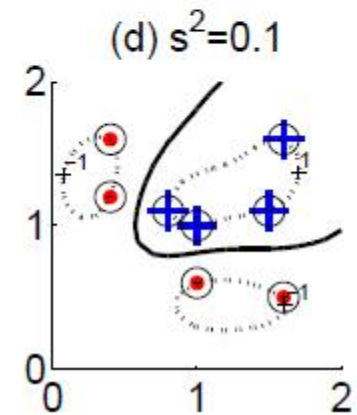
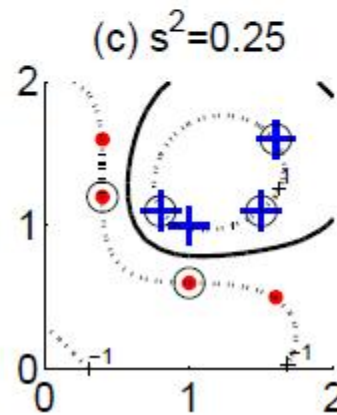
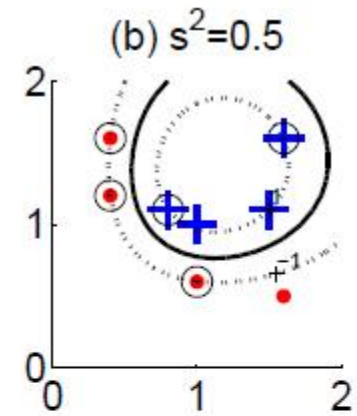
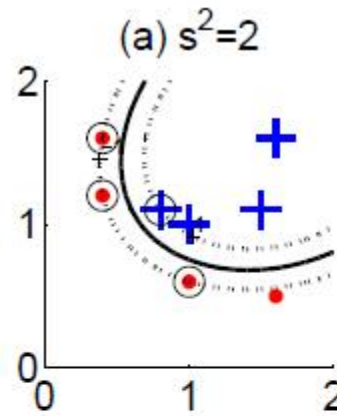
$$\begin{aligned} K(\mathbf{x}, \mathbf{y}) &= (\mathbf{x}^T \mathbf{y} + 1)^2 \\ &= (x_1 y_1 + x_2 y_2 + 1)^2 \\ &= 1 + 2x_1 y_1 + 2x_2 y_2 + 2x_1 x_2 y_1 y_2 + x_1^2 y_1^2 + x_2^2 y_2^2 \\ \phi(\mathbf{x}) &= [1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2, x_1^2, x_2^2]^T \end{aligned}$$



# Nonlinear SVM

- Radial-basis functions:

$$K(\mathbf{x}^t, \mathbf{x}) = \exp\left[-\frac{\|\mathbf{x}^t - \mathbf{x}\|^2}{2s^2}\right]$$



# Nonlinear SVM

## 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:

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

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

# Nonlinear SVM

## 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  $\mathbf{w}$ :

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

- Optimal Hyperplane:  $f(\mathbf{x}) = \text{sgn}(\mathbf{w}^T \varphi(\mathbf{x}))$   
 $= \text{sgn}(-x_1 x_2)$

# Multiclass SVM

- Multiclass Classification
  - SVM : binary classification
  - Multiclass SVM: **multiple binary classification**
- One-to-rest vs. one-to-one
  - 1) One-to-rest
    - 1 번째 class 와 1 클래스를 제외한 나머지 M-1 클래스로 이진 분류
    - M 개의 이진 분류 SVM 을 학습하여 사용
  - 2) One-to-one
    - M 개의 클래스 중 2개를 선택하여 이진 분류
    - $M(M-1)/2$  개의 이진 분류 SVM 을 학습하여 사용



## 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)	<a href="http://www.ai.mit.edu">http://www.ai. mit.edu</a>
LIBSVM	C.C. Chang, C.H. Lin (National Taiwan Univ.)	C++, Java	Python, R, Matlab, Perl	SMO(Platt), SVMLight(Joachims)	<a href="http://www.csie.ntu.edu.tw/~libsvm">http://www.csie. ntu.edu.tw/~libsvm</a>
SVMLight	T. Joachims, (Univ. of Dortmund)	C	Solaris, Linux, IRIX, Windows NT	T. Joachims	<a href="http://www.svmlight.joachims.org">http://www.svmlight. joachims.org</a>
SVM Torch	R. Collobert, (IDIAP, Switzerland)	C, C++	Windows	R. Collobert	<a href="http://www.idap.ch/learning/SVMTorch.html">http://www.idap.ch /learning/SVMTorch.html</a>

# SVM Demo

<https://jgreitemann.github.io/svm-demo>

# SVM Summary

- Optimal separating hyperplane : max. margin
- Global solution (convexity)
- Analytic solution
- Model selection problem: kernel selection
- Classification & Regression 문제에 모두 적용