

CS570 Artificial Intelligence & Machine Learning

Support Vector Machines

Kee-Eung Kim

Department of Computer Science

KAIST

Support Vector Machines (classification) → assume binary

□ Key idea: find the optimal separating hyperplane

- $\mathcal{X} = \{\mathbf{x}^t, r^t\}_t$ 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

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

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

- Equivalently,
 $\underbrace{r^t(\mathbf{w}^T \mathbf{x}^t + w_0)}_{\text{항상 만족}} \geq +1$

Margins

□ Distance from the discriminant to the closest instance on either side

□ Distance of x to the hyperplane:

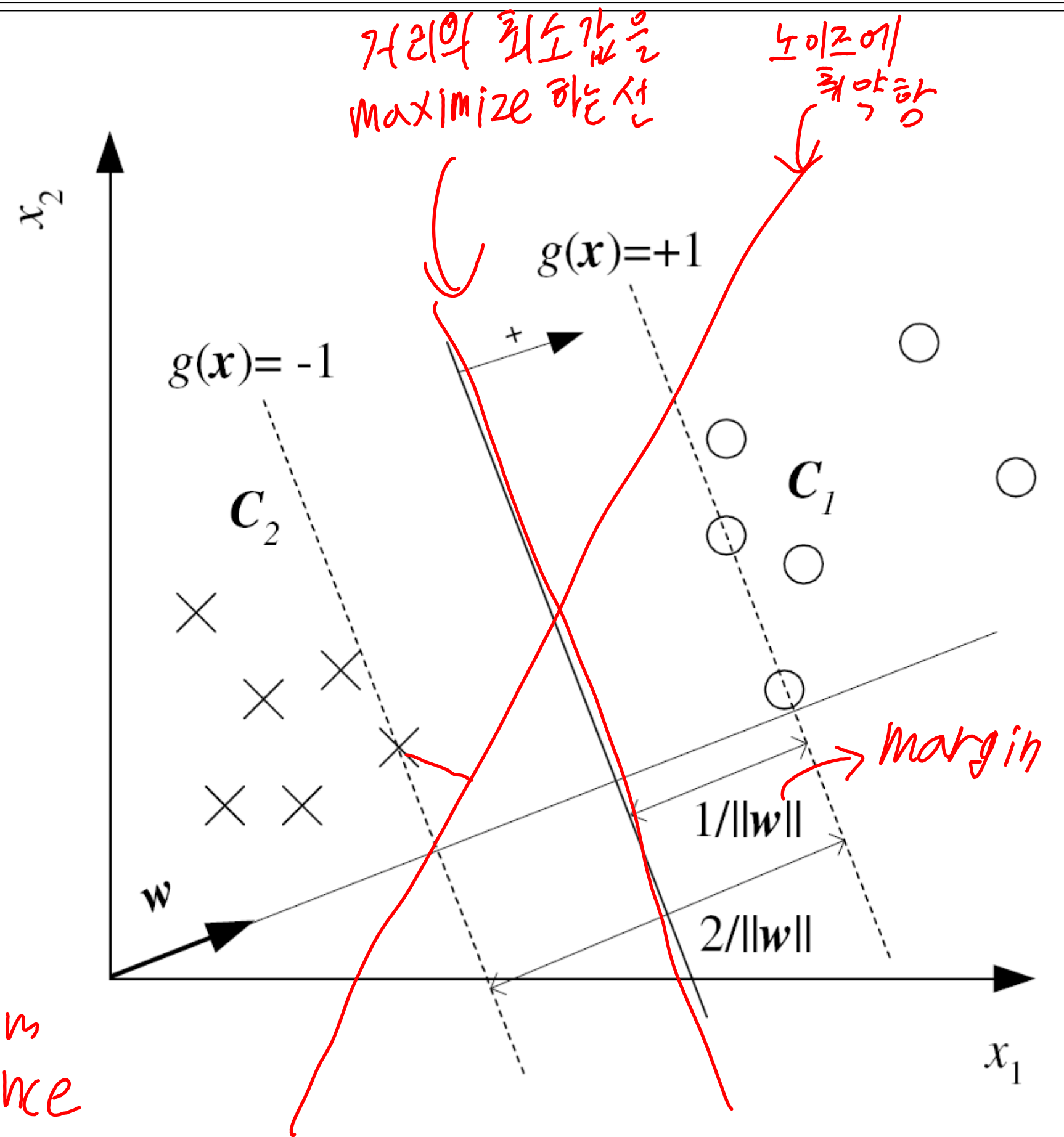
$$\frac{|w^T x^t + w_0|}{\|w\|}$$

□ Want: $\frac{r^t(w^T x^t + w_0)}{\|w\|} \geq \rho, \forall t$

Handwritten notes:
 - r^t : 부호를 유지시키기 위해
 - ρ : minimum distance
 - ρ : margin
 - ρ : 이걸 만족시키는 w 는 사실 많음

□ For a unique solution, fix $\rho\|w\| = 1 \rightarrow \|w\|$ 를 최소화
 and thus to maximize margin, $\rightarrow \rho(\text{margin})$ 최대화
 minimize $\|w\|$: $\min \frac{1}{2}\|w\|^2$ subject to $r^t(w^T x^t + w_0) \geq +1, \forall t$
Handwritten notes:
 - ρ 가 커짐
 - ρ 가 커짐
 - ρ 가 커짐

• Quadratic programming problem! \rightarrow ρ 가 커짐
 \rightarrow Quad prog (?)



Maximizing Margins

SVM 문제 형식화
→ margin 최대화

$$\square \min \frac{1}{2} \|\mathbf{w}\|^2 \text{ subject to } r^t(\mathbf{w}^T \mathbf{x}^t + w_0) \geq +1, \forall t$$

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

$$= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{t=1}^N \alpha^t r^t (\mathbf{w}^T \mathbf{x}^t + w_0) + \sum_{t=1}^N \alpha^t$$

$$\square \frac{\partial L_p}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{t=1}^N \alpha^t r^t \mathbf{x}^t \quad \frac{\partial L_p}{\partial w_0} = 0 \Rightarrow \sum_{t=1}^N \alpha^t r^t = 0$$

$$\square L_d = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \mathbf{w}^T \sum_t \alpha^t r^t \mathbf{x}^t - w_0 \sum_t \alpha^t r^t + \sum_t \alpha^t$$

$$= -\frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_t \alpha^t$$

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

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

- Most $\alpha^t = 0$ and only small number have $\alpha^t > 0$; \mathbf{x}^t with $\alpha^t > 0$ are the support vectors

D+1 variables
N constraints

→ 데이터셋이
큰 경우, 최적화할
구하는 데 시간이
많이 걸림

→ 그래서 변수가 많되
constraint가 적은게
나름

t, s : index to the dataset

$\mathbf{w} = \sum \alpha^t r^t \mathbf{w}^t$
→ $\alpha^t \neq 0$ 일 때만
 \mathbf{x}^t 가 contribution
→ 그래서 support SVM

optimizer가 찾아낸 대다수의 α 는 0 일 것임

hyperplane을
결정하는 가까운 점들

특정 데이터의 support
→ 확률 값이
0이 아님

KAIST

Korea Advanced Institute of Science and Technology
한국과학기술원

Soft Margins

□ If not linearly separable

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

□ Soft error $\sum_t \xi^t$

□ New objective function:

$$\min \left[\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_t \xi^t \right]$$

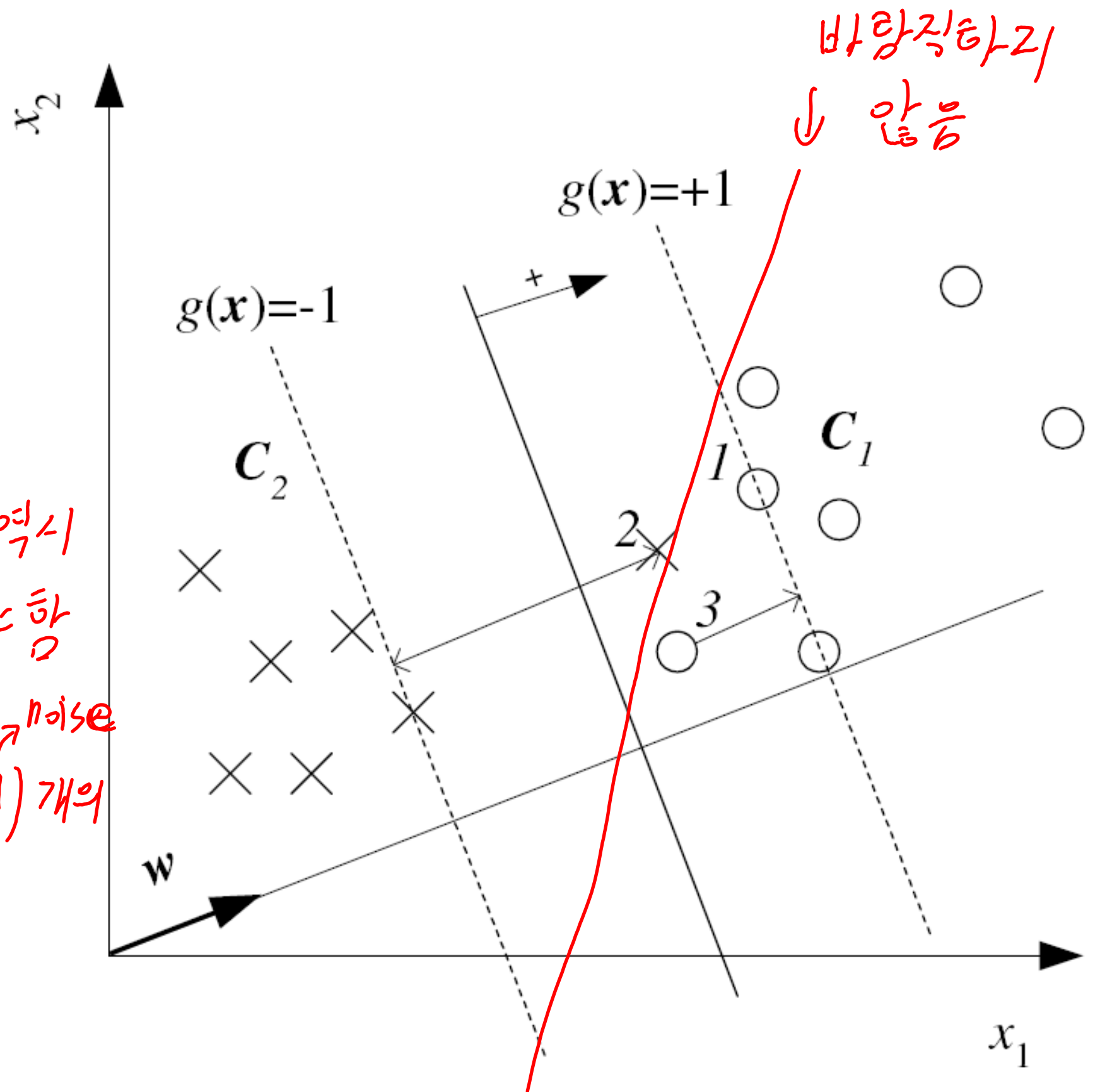
subject to

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

$$\xi^t \geq 0, \forall t$$

이런 예시 역시
최소화해야 함
noise
0+1+(N)개의
변수

tunable



□ New primal is

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

Kernel Machines

□ Preprocess input \mathbf{x} by basis functions

- Suppose $\mathbf{z} = \varphi(\mathbf{x})$
- Prepare transformed training set $\mathcal{Z} = \{\varphi(\mathbf{x}^t), r^t\}$
- Linear model in space \mathcal{Z} is nonlinear model in space \mathcal{X}

$$g(\mathbf{z}) = \mathbf{w}^T \mathbf{z} \quad g(\mathbf{x}) = \mathbf{w}^T \varphi(\mathbf{x}) \quad \text{key idea}$$

□ SVM on the transformed space \mathcal{Z}

- $\mathbf{w} = \sum_t \alpha^t r^t \mathbf{z}^t = \sum_t \alpha^t r^t \varphi(\mathbf{x}^t)$
- $g(\mathbf{x}) = \mathbf{w}^T \varphi(\mathbf{x}) = \sum_t \alpha^t r^t \varphi(\mathbf{x}^t)^T \varphi(\mathbf{x})$ $\rightarrow K(\mathbf{x}^t, \mathbf{x})$ (핵심)
- $g(\mathbf{x}) = \sum_t \alpha^t r^t K(\mathbf{x}^t, \mathbf{x})$

□ Kernel functions K

- **Polynomials of degree q :** $K(\mathbf{x}^t, \mathbf{x}) = (\mathbf{x}^T \mathbf{x}^t + 1)^q$
 - $K(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y} + 1)^2 = (x_1 y_1 + x_2 y_2 + 1)^2$
$$= 1 + 2x_1 y_1 + 2x_2 y_2 + 2x_1 x_2 y_1 y_2 + x_1^2 y_1^2 + x_2^2 y_2^2$$
$$\varphi(\mathbf{x}) = [1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2, x_1^2, x_2^2]^T$$
- **Radial-basis functions:** $K(\mathbf{x}^t, \mathbf{x}) = \exp[-\|\mathbf{x}^t - \mathbf{x}\|^2 / \sigma^2]$
- **Sigmoid functions:** $K(\mathbf{x}^t, \mathbf{x}) = \tanh(2\mathbf{x}^T \mathbf{x}^t + 1)$

SVM for Regression ☆ r^t 가 연속 값

□ Assume a linear model (possibly kernelized)

- $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$

그리고 예외적인 예외 함수를 쓰는 법

→ $\frac{0}{2}$ 만쓰면 \checkmark formulation 이 더 좋음

□ Use ε -sensitive error function (instead of squared error function)

$$Err(r^t, f(\mathbf{x}^t)) = \begin{cases} 0 & \text{if } |r^t - f(\mathbf{x}^t)| < \epsilon \\ |r^t - f(\mathbf{x}^t)| - \epsilon & \text{otherwise} \end{cases}$$

□ Problem formulation:

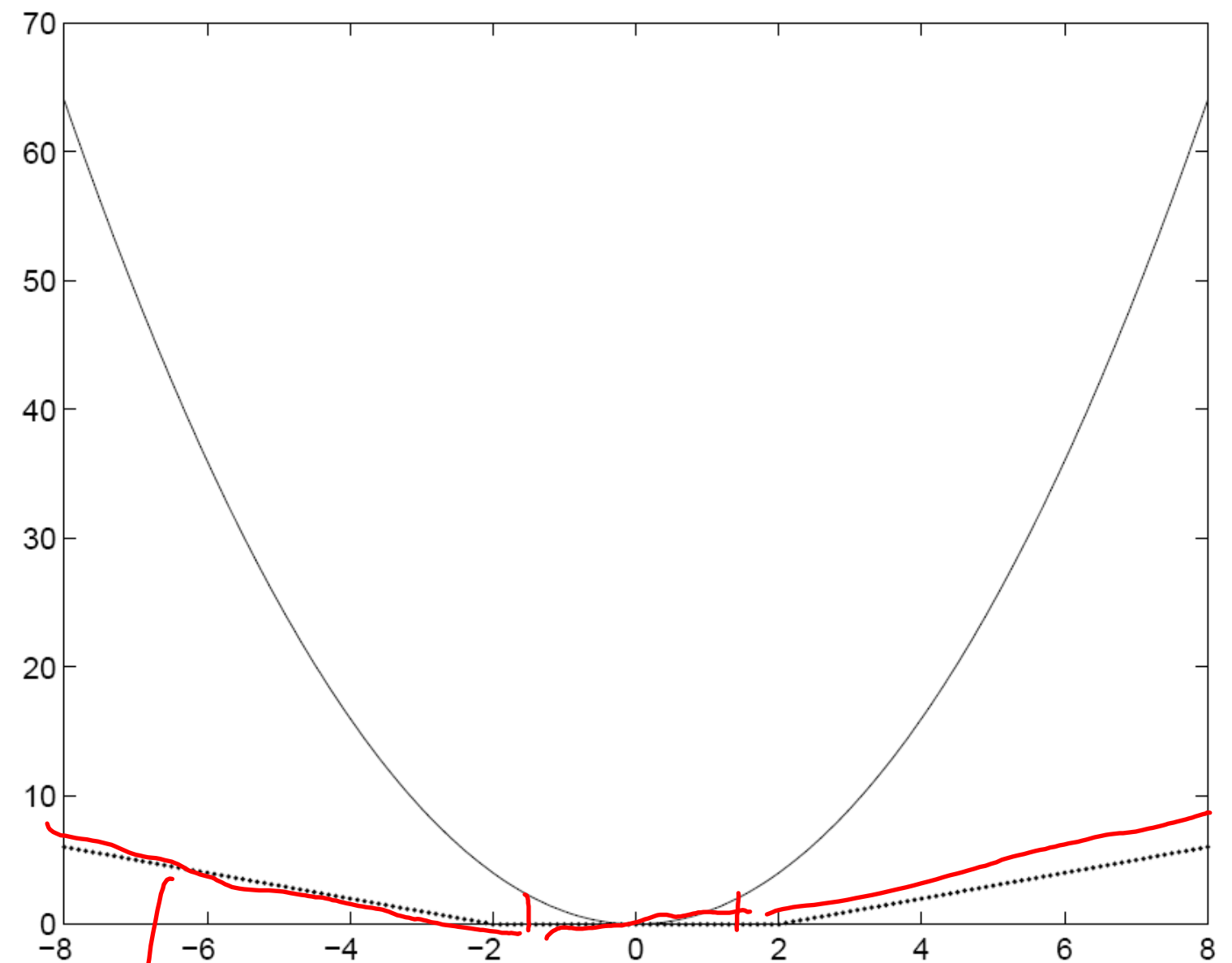
$$\min \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_t (\xi_+^t + \xi_-^t) \right\}$$

subject to

$$r^t - (\mathbf{w}^T \mathbf{x} + w_0) \leq \epsilon + \xi_+^t$$

$$(\mathbf{w}^T \mathbf{x} + w_0) - r^t \leq \epsilon + \xi_-^t$$

$$\xi_+^t, \xi_-^t \geq 0$$



Err가 선형 증가 Err=0