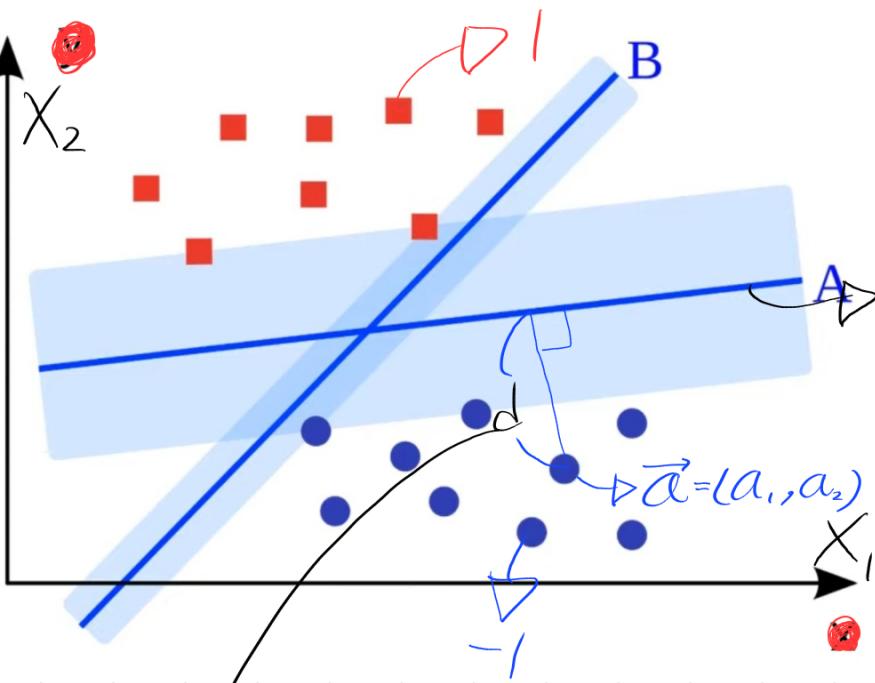


4. SVM

연세대학교 디지털애널리틱스 융합학과



$$w_1 x_1 + w_2 x_2 + b = 0$$

$$\vec{w} = (w_1, w_2)$$

$$d = \frac{|\vec{w} \cdot \vec{a} + b|}{\|\vec{w}\|_2} = \frac{|w_1 a_1 + w_2 a_2 + b|}{\sqrt{w_1^2 + w_2^2}}$$

[점과 초평면 ($\vec{w} \cdot \vec{x} + b = 0$) 사이의 거리]

만약 \vec{x} 가 d 차원이면

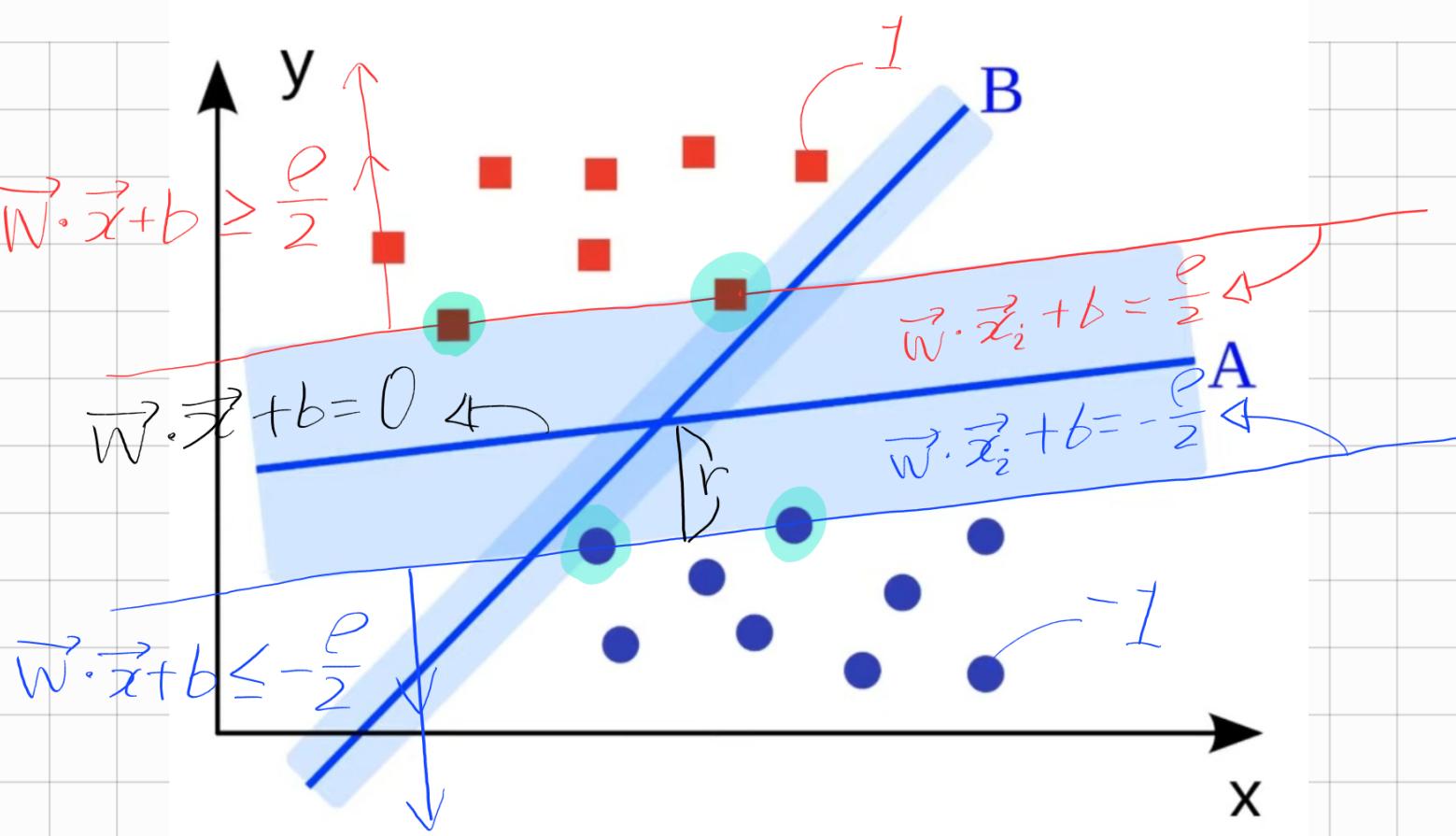
초평면은 몇 차원? $\rightarrow d-1$ 차원

Ex) $\vec{x} = (x_1, x_2) \in \mathbb{R}^2 \rightarrow 2$ 차원

$$\text{초평면} \rightarrow w_1 x_1 + w_2 x_2 + b = 0$$

$$\rightarrow x_1 = \frac{-w_2 x_2 - b}{w_1}$$

$\rightarrow 2-1$ 차원 초평면



$$y(\vec{w} \cdot \vec{x}' + b) = \frac{\rho}{2} \text{ at } S.V$$

초평면 $\vec{w} \cdot \vec{x} + b = 0$ 과 $S.V$ 사이의 거리 r ?

$$\begin{aligned} r &= \frac{|\vec{w} \cdot \vec{x}' + b|}{\|\vec{w}\|_2} = \frac{|y(\vec{w} \cdot \vec{x}' + b)|}{\|\vec{w}\|_2} = \frac{\frac{\rho}{2}}{\|\vec{w}\|_2} \\ &= \frac{\frac{\rho}{2} \times \frac{2}{\rho}}{\|\vec{w}\|_2 \times \frac{2}{\rho}} = \frac{1}{\|\vec{w} \times \frac{2}{\rho}\|_2} = \frac{1}{\|\vec{w}'\|_2} \end{aligned}$$

margin r 을 최대화 하는

\vec{w} 와 b 를 찾는 것이

SVM의 목표

초평면

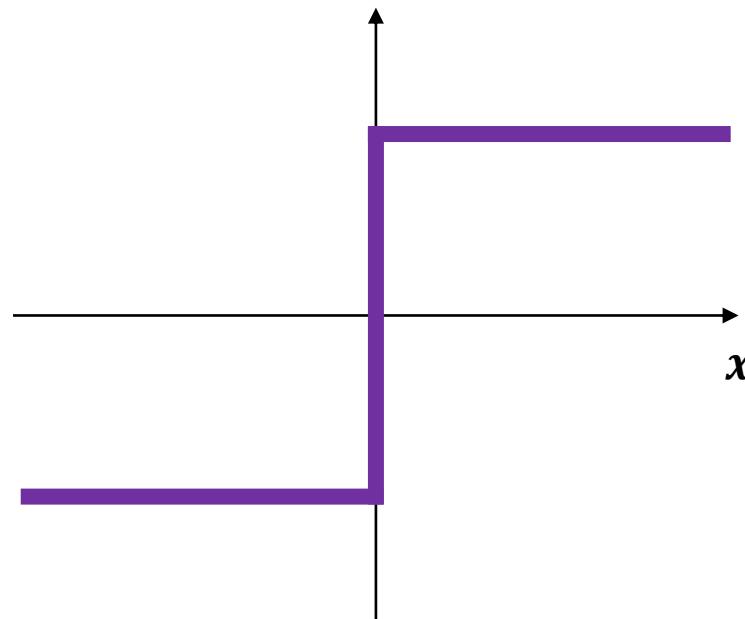
$\vec{w} \cdot \vec{x} + b = 0$ 를 찾는 것

$$\begin{aligned} \underset{\vec{w}, b}{\text{MAX}} \quad & \frac{1}{\|\vec{w}\|_2} \\ \text{s.t.} \quad & y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1 \end{aligned}$$

What is Support Vector Machine?

Perceptron

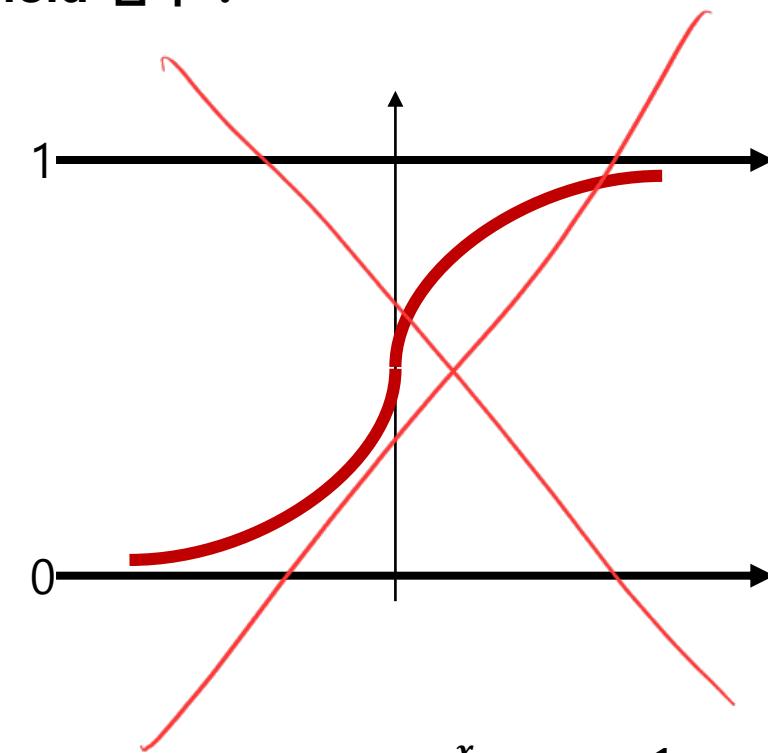
Sign 함수 :



$$sign(x) = \begin{cases} 1 & \text{for } x > 0 \\ -1 & \text{else} \end{cases}$$

→ logistic Regression

Sigmoid 함수 :



$$sigmoid(x) = \frac{e^x}{1 + e^x} = \frac{1}{1 + e^{-x}}$$

Hard Margin SVM

Primal

$$\min_{\vec{w}, b, \xi_i} \frac{1}{2} \|\vec{w}\|_2^2$$

$$s.t. \quad y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1$$

KKT 조건

만족하면

Dual

$$\Rightarrow \max_{\vec{\alpha}, \vec{\beta}} \left[\min_{\vec{w}, b, \xi_i} \frac{1}{2} \|\vec{w}\|_2^2 - \sum_{i=1}^n \alpha_i (y_i (\vec{w} \cdot \vec{x}_i + b) - 1) \right]$$

Primal

최대

K
K
T

$$① \nabla_{\vec{w}} \mathcal{G}(\vec{w}, b) = 0 \rightarrow$$

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

$$② \nabla_b \mathcal{G}(\vec{w}, b) = 0 \rightarrow$$

$$\sum_{i=1}^n \alpha_i y_i = 0$$

$$③ \alpha_i (y_i (\vec{w} \cdot \vec{x}_i + b) - 1) = 0$$

$$④ y_i (\vec{w} \cdot \vec{x}_i + b) - 1 \geq 0$$

$$⑤ \alpha_i \geq 0$$

Dual

$$\Rightarrow \max_{\vec{\alpha}, \vec{\beta}} \left[\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \vec{x}_i^T \vec{x}_j \right]$$

$$s.t. \quad \sum_{i=1}^n \alpha_i y_i = 0, \quad \alpha_i \geq 0$$

Solution

$$\Rightarrow f(\vec{x}_{\text{new}}) = \text{sign}(\vec{w} \cdot \vec{x}_{\text{new}} + b)$$

최종 Model

$$= \text{sign} \left(\sum_{i=1}^n \alpha_i y_i \vec{x}_i^T \vec{x}_{\text{new}} + b \right)$$

KKT 조건에 따른 SVM의 특징

$$③ \alpha_i(y_i(\vec{w} \cdot \vec{x}_i + b) - 1) = 0$$

★ 둘 중 하나만 0

$$⑤ \alpha_i \geq 0$$

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

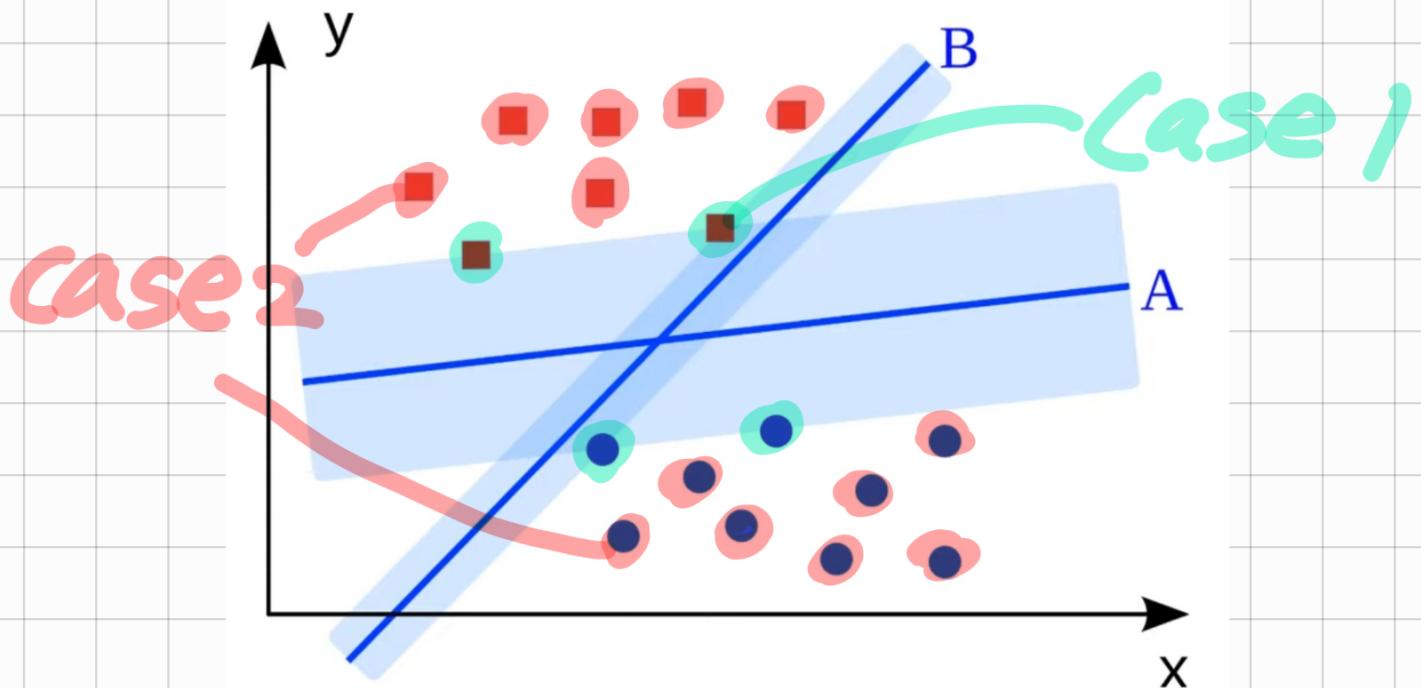
$$\text{Case 1: } \alpha_i = 0 \stackrel{①}{\Rightarrow} y_i(\vec{w} \cdot \vec{x}_i + b) - 1 = 0$$

∴ S.V이다

$$\text{Case 2: } \alpha_i > 0 \stackrel{①}{\Rightarrow} y_i(\vec{w} \cdot \vec{x}_i + b) - 1 \neq 0$$

∴ S.V가 아님

$$\begin{aligned} f(\vec{x}_{\text{new}}) &= \text{sign}(\vec{w} \cdot \vec{x}_{\text{new}} + b) \\ &= \text{sign}\left(\sum_{i=1}^n \alpha_i y_i \vec{x}_i^\top \vec{x}_{\text{new}} + b\right) \end{aligned}$$

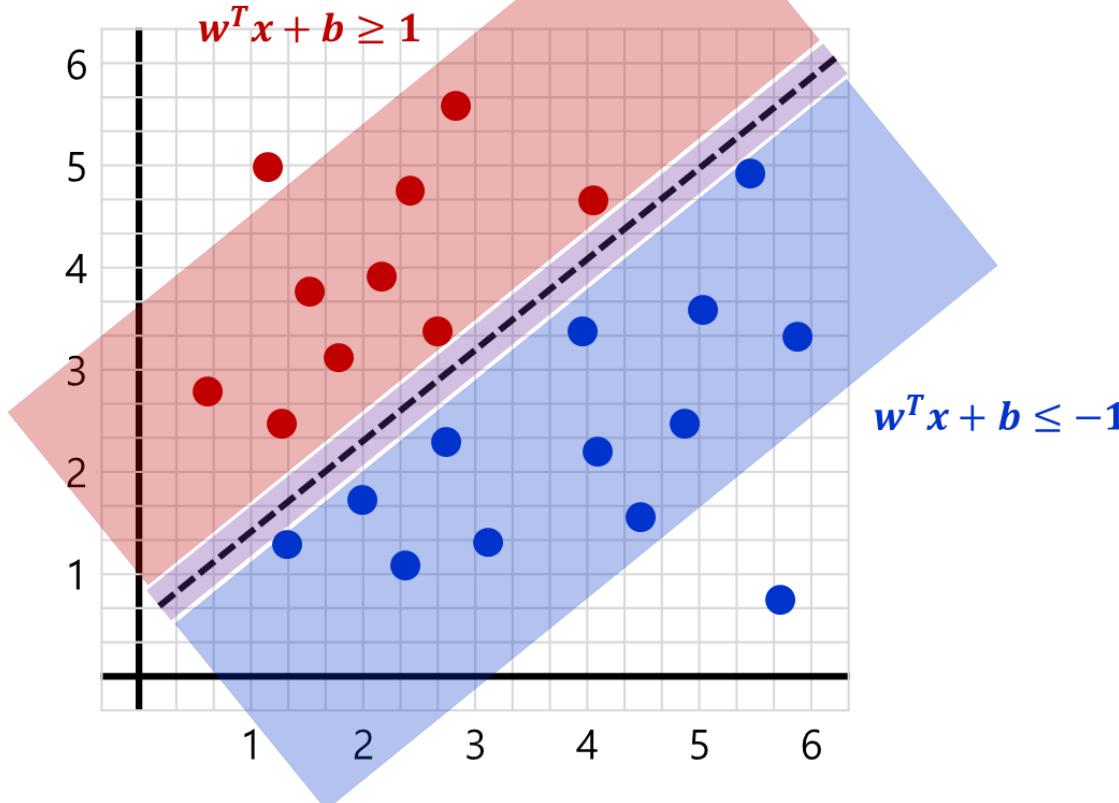


Support Vector Machine Concepts

Gradient Descent

$$\text{Hinge loss} = \max\{0, 1 - y \times \hat{y}\}$$

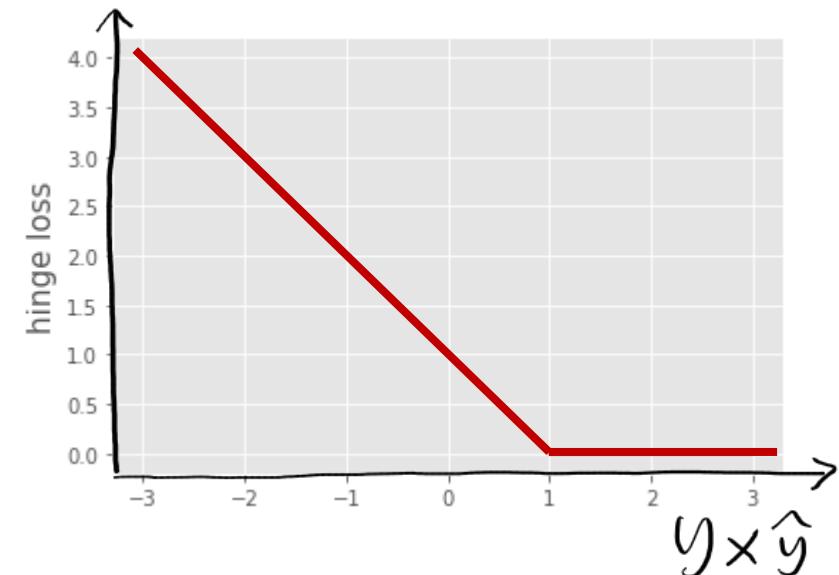
If $y = 1, w^T x + b \geq 1$ else if $y = -1, w^T x + b \leq -1$



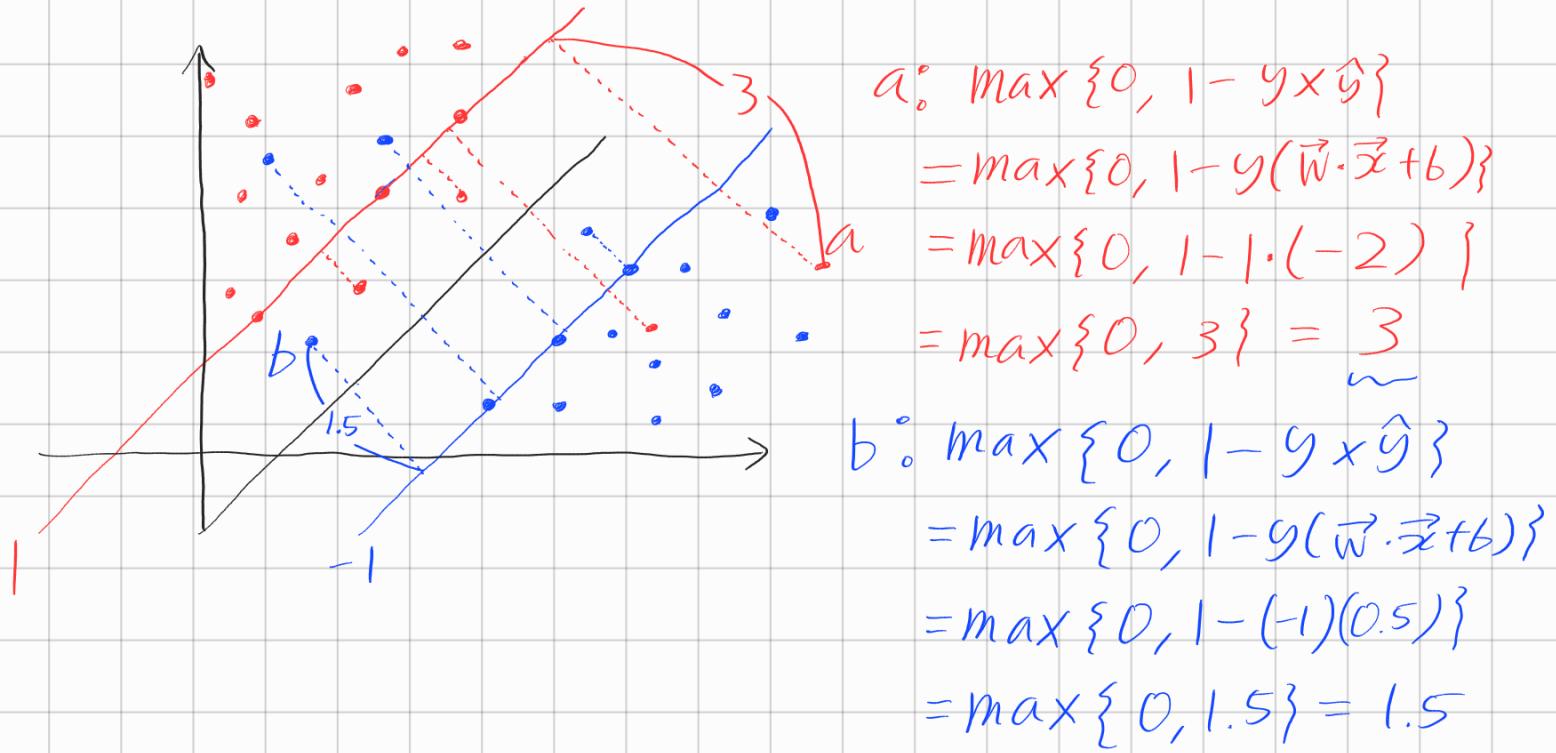
If y and \hat{y} are equal , then

$$y \times \hat{y} = y(w^T x + b) \geq 1$$

$y \times \hat{y} < 1$ means any observations exist in margin area



$$\text{Hinge Loss} = \max\{0, 1 - y \cdot \hat{y}\}$$



즉, Hinge Loss는 (\vec{x} 자신의 참값의 SV 초평면으로 부터 \vec{x} 자신의 점까지의 거리 - 1)에 비례한다.

(단, 정답인 영역으로 \hat{y} 을 예측 했을 경우는 0이 된다.)

$$\text{Hinge Loss}_i = \max\{0, 1 - y \cdot \hat{y}\}$$

$$= \max\{0, 1 - y(\vec{w} \cdot \vec{x}_i + b)\}$$

$$= \underline{\underline{\epsilon_i}}$$

Soft Margin SVM

Primal

$$\min_{\vec{w}, b, \xi_i} \frac{1}{2} \|\vec{w}\|_2^2 + C \sum_{i=1}^n \xi_i$$

$$s.t. y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0$$

KKT
만족

Hard Margin SVM

다른 부분

Dual

$$\Rightarrow \max_{\vec{\alpha}, \vec{\beta}} \left[\underbrace{\min_{\vec{w}, b, \xi_i} \frac{1}{2} \|\vec{w}\|_2^2 + C \sum_{i=1}^n \xi_i}_{\text{Primal}} - \sum_{i=1}^n \alpha_i (y_i (\vec{w} \cdot \vec{x}_i + b) - 1 + \xi_i) - \sum_{i=1}^n \beta_i \xi_i \right]$$

K
K
T

$$① \nabla_{\vec{w}} \mathcal{G}(\vec{w}, b, \xi_i) = 0 \rightarrow \vec{w} = \sum_{i=1}^n \alpha_i y_i \vec{x}_i$$

$$② \nabla_b \mathcal{G}(\vec{w}, b, \xi_i) = 0 \rightarrow \sum_{i=1}^n \alpha_i y_i = 0$$

$$③ \nabla_{\xi_i} \mathcal{G}(\vec{w}, b, \xi_i) = 0 \rightarrow C - \alpha_i - \beta_i = 0$$

$$④ \alpha_i (y_i (\vec{w} \cdot \vec{x}_i + b) - 1 + \xi_i) = 0, \quad \beta_i \xi_i = 0$$

$$⑤ y_i (\vec{w} \cdot \vec{x}_i + b) - 1 + \xi_i \geq 0, \quad \xi_i \geq 0$$

$$⑥ \alpha_i \geq 0, \quad \beta_i \geq 0$$

Hard

$$0 \leq \alpha_i \leq C$$

Dual

$$\Rightarrow \max_{\vec{\alpha}, \vec{\beta}} \left[\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \vec{x}_i \vec{x}_j^T \right]$$

$$s.t. \sum_{i=1}^n \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C$$

Solution

$$\Rightarrow f(\vec{x}_{\text{new}}) = \text{sign}(\vec{w} \cdot \vec{x}_{\text{new}} + b)$$

$$= \text{sign} \left(\sum_{i=1}^n \alpha_i y_i \vec{x}_i^T \vec{x}_{\text{new}} + b \right)$$

KKT 조건에 따른 SVM의 특징

$$\textcircled{4} \quad \alpha_i (y_i (\vec{w} \cdot \vec{x}_i + b) - 1 + \xi_i) = 0$$

$$\textcircled{3} \quad C - \alpha_i - \beta_i = 0$$

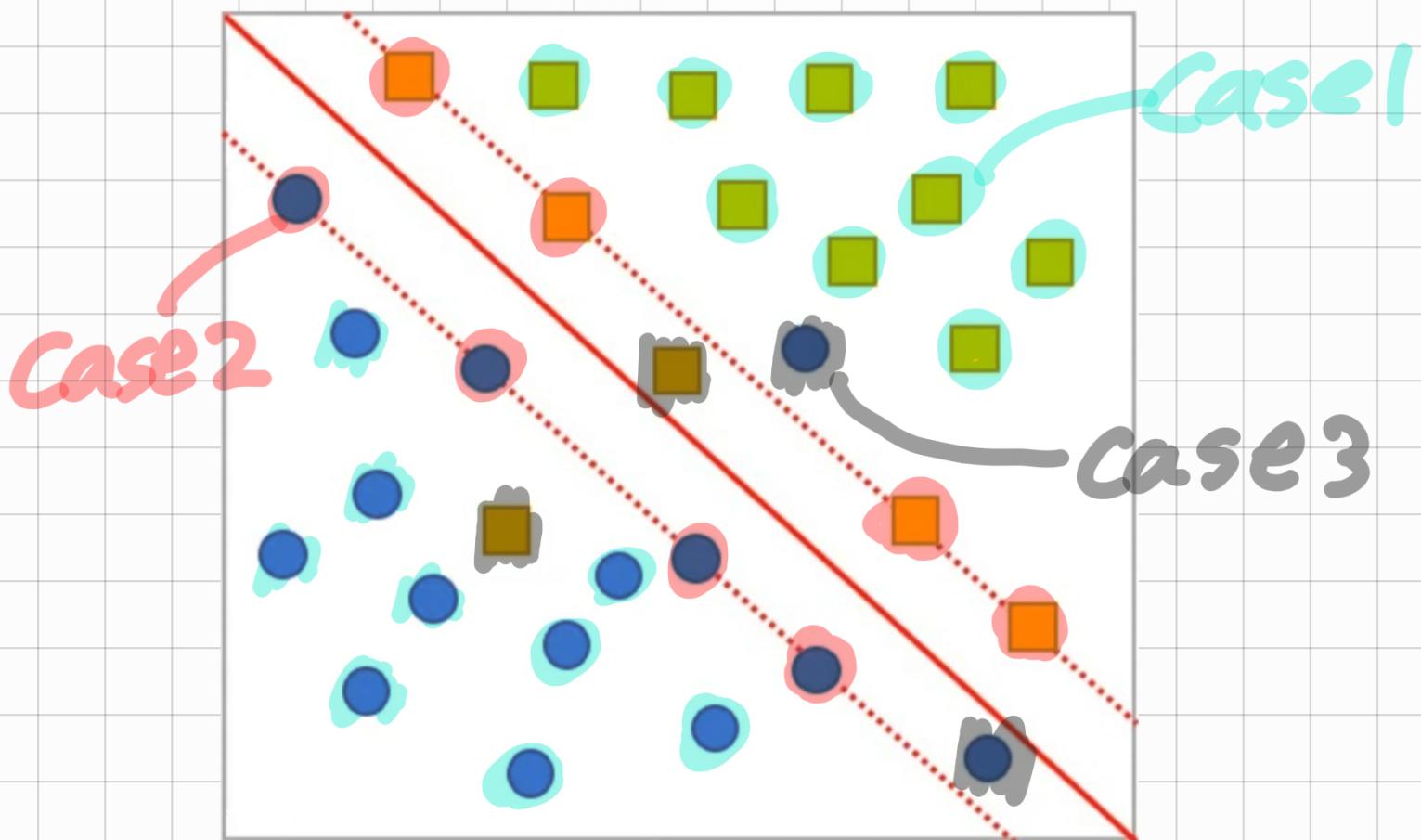
둘 중 하나만 0

$$\textcircled{4} \quad \beta_i \xi_i = 0$$

Case 1 : $\alpha_i = 0 \xrightarrow{\textcircled{3}} C = \beta_i \xrightarrow{\textcircled{2}} \xi_i = 0 \quad \boxed{y_i (\vec{w} \cdot \vec{x}_i + b) - 1 + \xi_i \neq 0}$
 $\quad \quad \quad \textcircled{1} \quad y_i (\vec{w} \cdot \vec{x}_i + b) - 1 + \xi_i \neq 0 \quad \boxed{y_i (\vec{w} \cdot \vec{x}_i + b) - 1 \neq 0}$
 $\therefore S.V \text{가 아닙니다 and 오분류가 아닙니다}$

Case 2 : $0 < \alpha_i < C \xrightarrow{\textcircled{3}} 0 < \beta_i < C \xrightarrow{\textcircled{2}} \xi_i = 0 \quad \boxed{y_i (\vec{w} \cdot \vec{x}_i + b) - 1 + \xi_i = 0}$
 $\quad \quad \quad \textcircled{1} \quad y_i (\vec{w} \cdot \vec{x}_i + b) - 1 + \xi_i = 0$
 $\therefore S.V \text{입니다} \rightarrow \text{당연히 오분류가 아닙니다}$

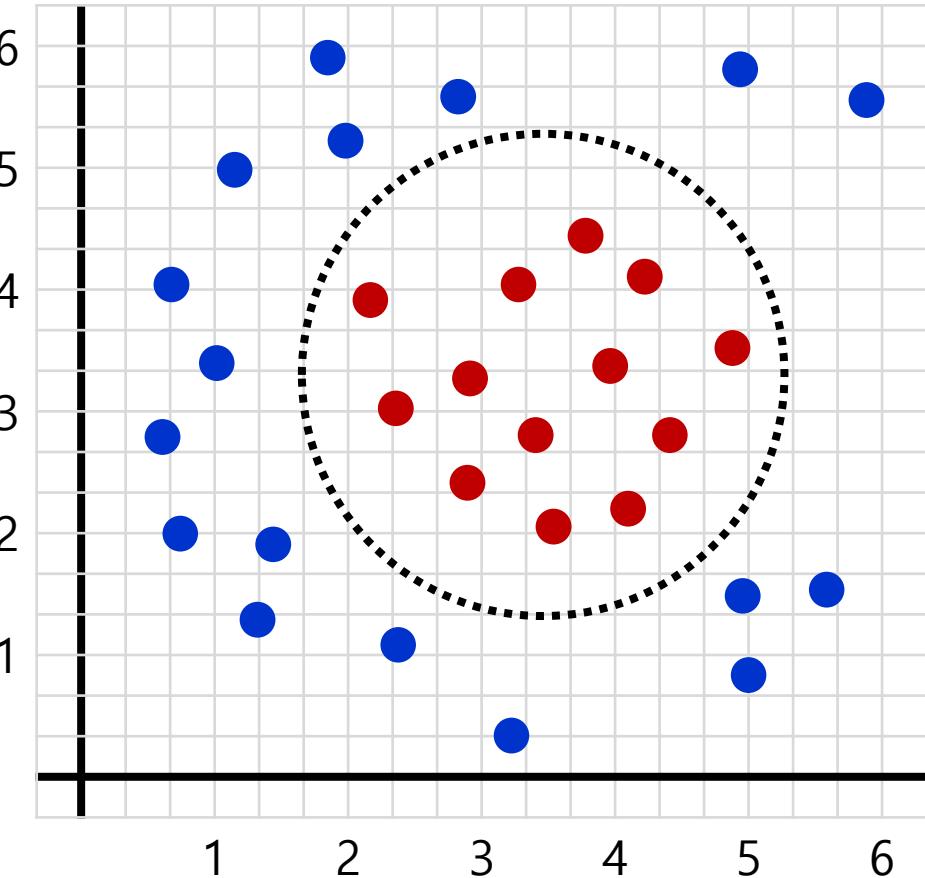
Case 3 : $\alpha_i = C \xrightarrow{\textcircled{3}} \beta_i = 0 \xrightarrow{\textcircled{2}} \xi_i > 0 \quad \boxed{y_i (\vec{w} \cdot \vec{x}_i + b) = 1 - \xi_i \neq 1}$
 $\quad \quad \quad \textcircled{1} \quad y_i (\vec{w} \cdot \vec{x}_i + b) - 1 + \xi_i = 0$
 $\therefore S.V \text{가 아닙니다 and 오분류입니다}$



Kernel Function

Non Linear SVM

Non-Linear Support Vector Machine

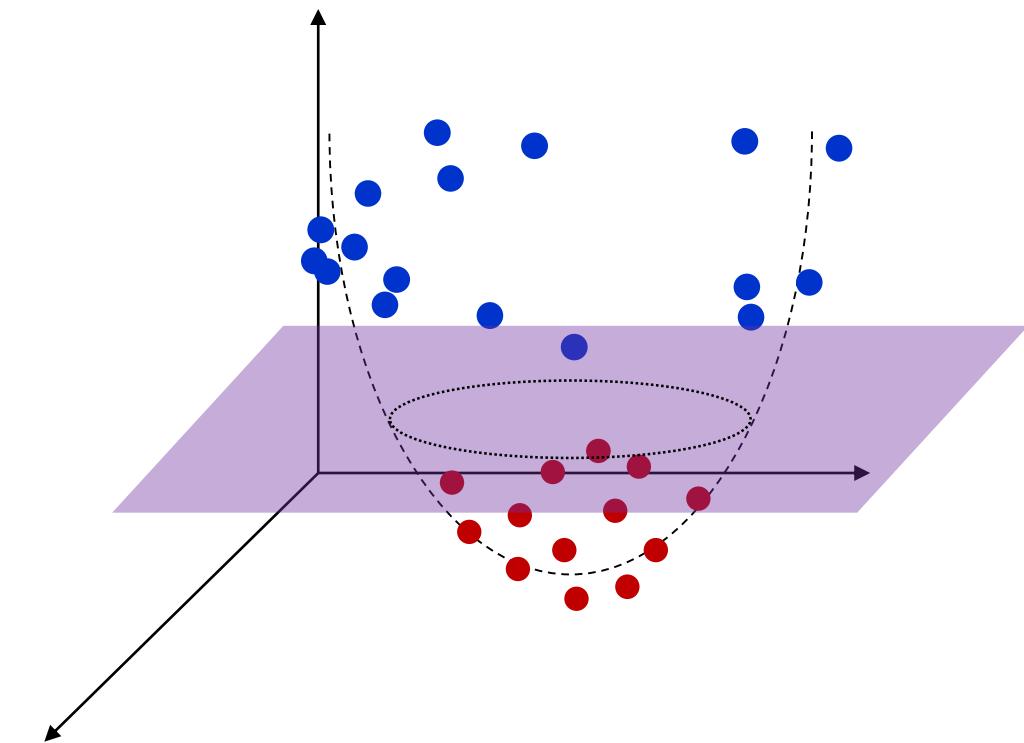


$\vec{x} \in \mathbb{R}^2$ 2차원

혹시
→ 3차원

$\phi(\vec{x}) \in \mathbb{R}^3$

$$\Phi(x)$$



Non-Linear SVM의 특징 함수

[Primal]

$$\min_{\vec{w}, b, \xi_i} \frac{1}{2} \|\vec{w}\|_2^2 + C \sum_{i=1}^n \xi_i$$

Hard Margin SVM 과
다른 부분

$$\text{s.t. } y_i (\vec{w} \cdot \phi(\vec{x}_i) + b) \geq 1 - \xi_i, \quad \xi_i \geq 0$$

Hard Margin SVM 과 같은 방법으로

초기화를 하기도 면

[Dual]

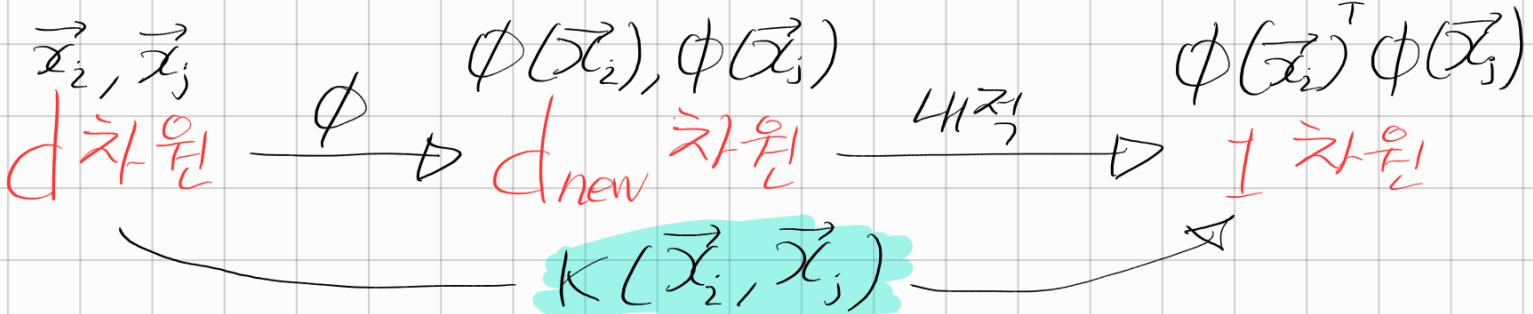
$$\Rightarrow \max_{\vec{\alpha}, \vec{\beta}} \left[\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \phi(\vec{x}_i)^T \phi(\vec{x}_j) \right]$$

$$\text{s.t. } \sum_{i=1}^n \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C$$

[Solution]

$$\begin{aligned} \Rightarrow f(\vec{x}_{\text{new}}) &= \text{sign}(\vec{w}^T \phi(\vec{x}_{\text{new}}) + b) \\ \text{초기 model} &= \text{sign} \left(\sum_{i=1}^n \alpha_i y_i \phi(\vec{x}_i)^T \phi(\vec{x}_{\text{new}}) + b \right) \\ &= \text{sign} \left(\sum_{i=1}^n \alpha_i y_i K(\vec{x}_i, \vec{x}_{\text{new}}) + b \right) \end{aligned}$$

여기 $K(\vec{x}_i, \vec{x}_j) = \phi(\vec{x}_i)^T \phi(\vec{x}_j)$
가 되는 함수 K 를 찾아 대입



Kernel Function

Kernel Trick 을 통해 연산량을 줄여줌

Non-Linear Support Vector Machine

1) Linear : $k(x_i, x_j) = x_i^T x_j$ $d \rightarrow d$
Mapping $\Phi: x \rightarrow \Phi(x)$, where $\Phi(x)$ is x itself

2) Polynomial of power p : $k(x_i, x_j) = (1 + x_i^T x_j)^p$ $d \rightarrow d+p$ C_P
Mapping $\Phi: x \rightarrow \Phi(x)$, where $\Phi(x)$ has $\binom{d+p}{p}$ dimensions (d is the dimension of data)

3) Gaussian (radial-basis function) : $k(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$ $d \rightarrow \infty$
Mapping $\Phi: x \rightarrow \Phi(x)$, where $\Phi(x)$ is infinite-dimensional

x 에 대한 mapping 차원이 증가
→ 비선형성이 증가함으로 복잡한 문제에 대한 분류 성능 ↑,
but, overfitting의 위험도 또한 ↑

4) Sigmoid : $k(x_i, x_j) = \tanh(\kappa \cdot x_i^T x_j + c)$, $\kappa > 0$ and $c < 0$ $d \rightarrow \infty$

Kernel Function

Non-Linear Support Vector Machine

2) Polynomial of power p : $k(x_i, x_j) = (1 + x_i^T x_j)^p$

Mapping $\Phi: x \rightarrow \Phi(x)$, where $\Phi(x)$ has $\binom{d+p}{p}$ dimensions (d is the dimension of data)

Ex)

$$\text{if } degree = 2, x_i = (x_{i_1}, x_{i_2}), x_j = (x_{j_1}, x_{j_2}) \quad 2 \xrightarrow{\phi} 2+2 C_2 = 6$$

$$k(x_i, x_j) = (1 + x_i^T x_j)^2 = (1 + x_{i_1}x_{j_1} + x_{i_2}x_{j_2})^2 = 1 + 2x_{i_1}x_{j_1} + 2x_{i_2}x_{j_2} + 2x_{i_1}x_{j_1}x_{i_2}x_{j_2} + {x_{i_1}}^2{x_{j_1}}^2 + {x_{i_2}}^2{x_{j_2}}^2$$

$$= (1, \sqrt{2}x_{i_1}, \sqrt{2}x_{i_2}, \sqrt{2}x_{i_1}x_{i_2}, {x_{i_1}}^2, {x_{i_2}}^2) \cdot (1, \sqrt{2}x_{j_1}, \sqrt{2}x_{j_2}, \sqrt{2}x_{j_1}x_{j_2}, {x_{j_1}}^2, {x_{j_2}}^2)$$

Gaussian (RBF) kernel

$$\vec{x}, \vec{x}' \in \mathbb{R}^d$$

$$\underbrace{\begin{matrix} d \\ \text{Dim} \end{matrix}}_{\vec{x}} \xrightarrow{\phi(\vec{x})} \infty \text{ Dim} \xrightarrow{\frac{\phi(\vec{x})^\top \phi(\vec{x}')}{4\pi\sigma^2}} \underbrace{\begin{matrix} 1 \\ \text{Dim} \end{matrix}}$$

$$K(\vec{x}, \vec{x}') = \exp\left(-\frac{(\vec{x} - \vec{x}')^2}{\sigma^2}\right)$$

$$= \left[\frac{1}{\exp((\vec{x} - \vec{x}')^2)} \right]^{\frac{1}{\sigma^2}} = \left[\frac{1}{\exp(\vec{x}^\top \vec{x} + \vec{x}'^\top \vec{x}' - 2\vec{x}^\top \vec{x}')} \right]^{\frac{1}{\sigma^2}}$$

$$= \left[\frac{1}{\exp(\vec{x}^\top \vec{x}) \times \exp(\vec{x}'^\top \vec{x}') \times \underbrace{1 / \exp(2\vec{x}^\top \vec{x}')}_{\star}} \right]^{\frac{1}{\sigma^2}}$$

Taylor expansion

$$\exp(2\vec{x}^\top \vec{x}') = \sum_{k=0}^{\infty} \frac{(2\vec{x}^\top \vec{x}')^k}{k!} = \sum_{k=0}^{\infty} \frac{2^k (\sum_{i=1}^d x_i x'_i)^k}{k!}$$

$$= 1 + 2(\sum_{i=1}^d x_i x'_i) + \frac{2^2 (\sum_{i=1}^d x_i x'_i)^2}{2!} + \frac{2^3 (\sum_{i=1}^d x_i x'_i)^3}{3!} + \frac{2^4 (\sum_{i=1}^d x_i x'_i)^4}{4!} + \frac{2^5 (\sum_{i=1}^d x_i x'_i)^5}{5!} + \dots$$

$$\Rightarrow (\sum_{i=1}^d x_i x'_i)^\infty \text{ 와 같이 } \infty \text{ 차원 까지 확장된다.}$$

RBF의 hyper parameter은 σ^2 (분산)으로

σ 가 클수록 모델의 복잡도가 증가

Sigmoid kernel은 Taylor expansion을 사용하여 ∞ 차원으로 확장

$$\rightarrow \tan(k\vec{x}_i^\top \vec{x}_j + c) \text{에서 } (c \text{가 } \infty \text{ 수록})$$

k 가 클수록 모델의 복잡도 증가, 모델의 복잡도 감소

Kernel Function의 기 위한 충분 조건은
Mercer's condition을 만족하는 것

Mercer's condition

1. It is symmetric

$$\Rightarrow K(x_i, x_j) = K(x_j, x_i)$$

2. The matrix

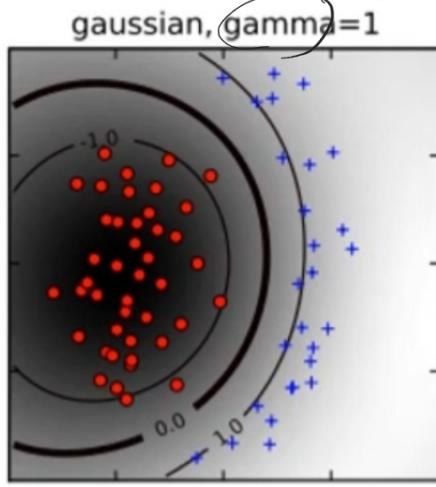
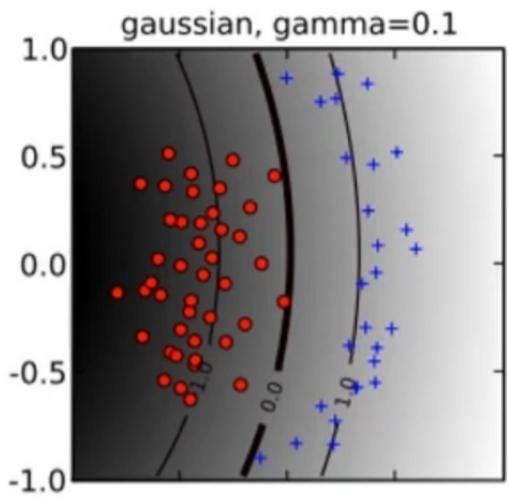
$$M = \begin{bmatrix} K(\vec{x}_1, \vec{x}_1) & K(\vec{x}_1, \vec{x}_2) & \cdots & K(\vec{x}_1, \vec{x}_n) \\ K(\vec{x}_2, \vec{x}_1) & K(\vec{x}_2, \vec{x}_2) & \cdots & K(\vec{x}_2, \vec{x}_n) \\ \vdots & \vdots & \ddots & \vdots \\ K(\vec{x}_n, \vec{x}_1) & K(\vec{x}_n, \vec{x}_2) & \cdots & K(\vec{x}_n, \vec{x}_n) \end{bmatrix}$$

is positive semi-definite

$$\Rightarrow \exists v, v^T M v \geq 0 \text{ for all } v \in (\mathbb{R}^n - \vec{0})$$

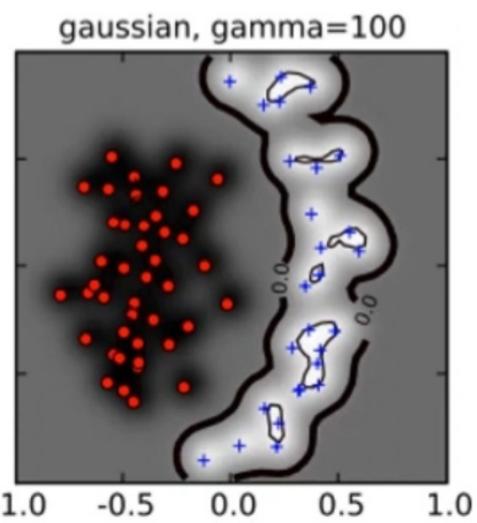
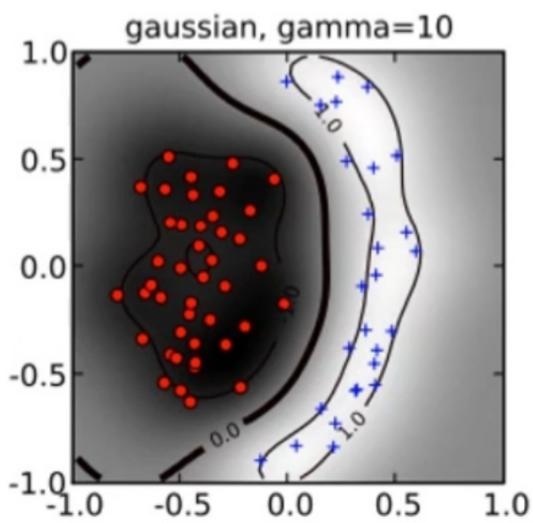
Support Vector Machine의 hyper parameter

$$\frac{1}{\sigma^2}$$



← 분산이 클 때
경계면이 단순

분산 ↓



← 분산이 작을 때
경계면이 복잡

분산 ↑

위의 목적함수 $\frac{1}{2} \vec{W}^T \vec{W} + C \sum_{i=1}^n \xi_i$

분산 ↓

margin 늘어남

오분류에 관대해짐

작을수록

분산 ↑

margin 줄여짐

오분류에 엄격해짐

C
큰수록

