

최소화 문제 (어제) → Convex optimization.

SVM 객체함 → ε-소실,

kernel.

## Inclass 10: Support Vector Machine

[SCS4049] Machine Learning and Data Science

---

Seongsik Park (s.park@dgu.edu)

AI Department, Dongguk University

다양한 선형 모델

→ 중간 고사 준비.

# Optimization (최적화).

✓ linear regression

✓ perceptron.

✓ logistic regression

Cost.

• SSE, MSE,  $\epsilon$   
RMSE  $\epsilon$

$$-\sum_M W_n^T \phi(x_n) t_n$$

• cross entropy.

Constraint.  
(model)

$$\hat{y} = \theta^T x \\ = W^T \phi(x)$$

model.

model.  
 $\sigma(W^T x)$



grad.  
st. g.d.  
mini - ...



grad.  
descent.

Solution

• closed form

(normal equation)

• gradient descent.

↪ cost 함수 정의.

## SVM: Maximum Margin Classifier

---

constraint.



convex,  
Lagrange.

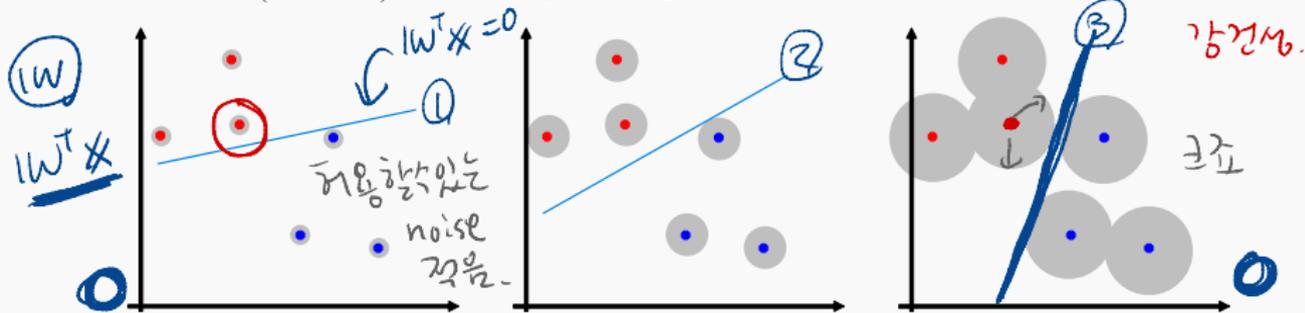


Solution

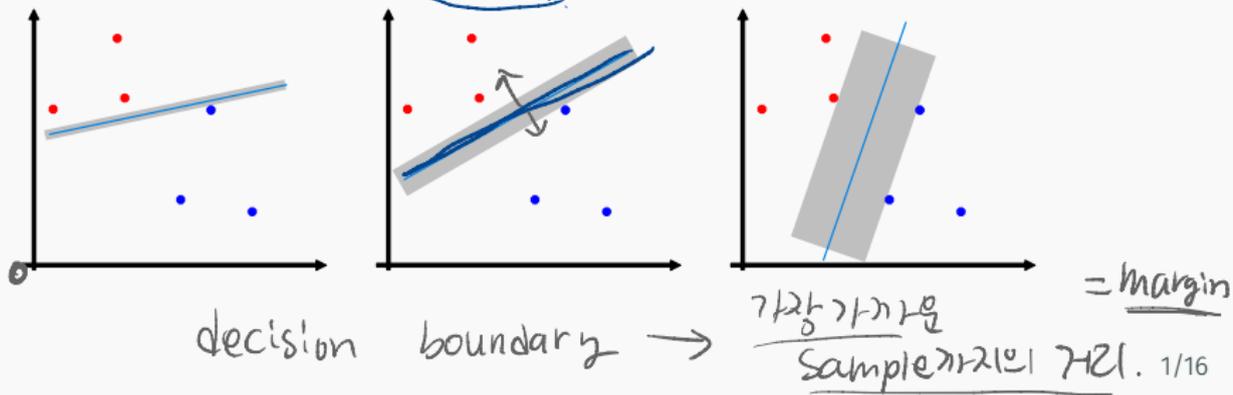
⇒ Support vector.  
(complementary  
slackness).

# What is a good decision boundary?

- 데이터 노이즈에 대한 강건성 (Robustness)
  - 노이즈(측정 오차)에 대해서 강건한 것이 좋은 모델이다.



- 여유로운 것이 더 강건하다  $\Rightarrow$  넓은 통로가 좋다  $\Rightarrow$  Large Margin Classification



# What is a good decision boundary?

→ 의사 결정은 경계의 데이터 (support vectors)에 의해서 결정됨  
 decision,

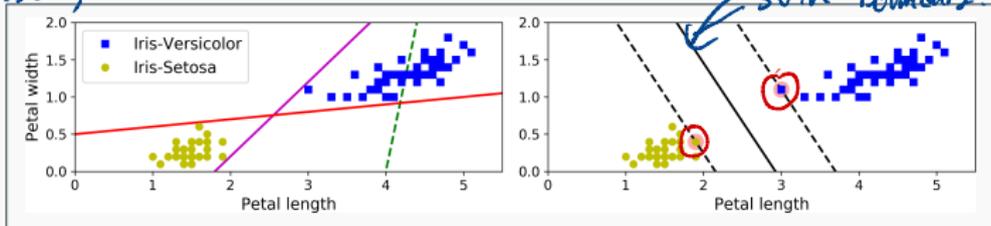


Figure 5-1. Large margin classification

Classifier  
S.V.M.  
 ↓  
 (?)

boundary ← Support vector (margin)  
 결정. ↑

→ Input feature의 scale에 민감한 support vector machine

training sample  
 중의 일부만  
 대표.

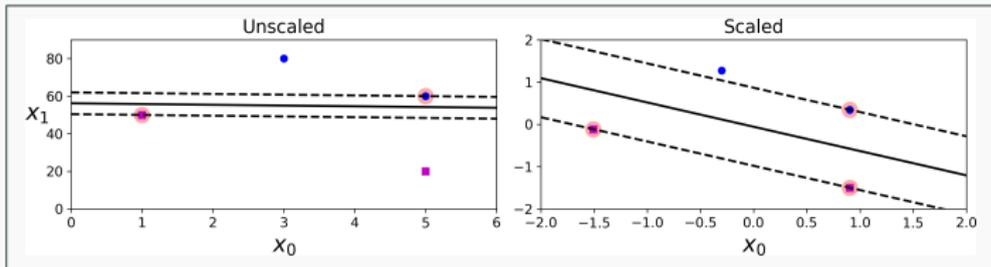
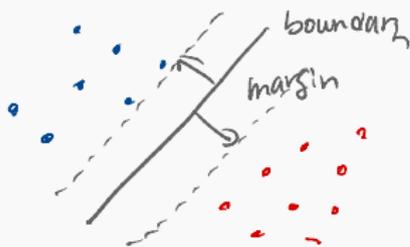


Figure 5-2. Sensitivity to feature scales

# Hard margin vs. soft margin



## Hard margin classification (hard-SVM)

- 모든 데이터들이 margin 밖에 위치하도록 boundary를 설정
- 데이터가 선형적으로 분리 가능(linearly separable)할 때만 적용 가능
- outlier에 매우 민감

## X Soft margin classification (soft-SVM)

바깥으로  
사라짐.

- margin을 가능한 넓게 하면서도 margin 안쪽으로 들어오는 것을 허용
- hyperparameter C: 클수록 좁아짐 (엄격), 작을 수록 넓어짐 (위반 허용)

Slack variable,  $\xi$  →  $\xi$ 가 클수록

저비용 데이터이지.

Advanced ML.  
1~2333 SVM.

# Hard margin vs. soft margin

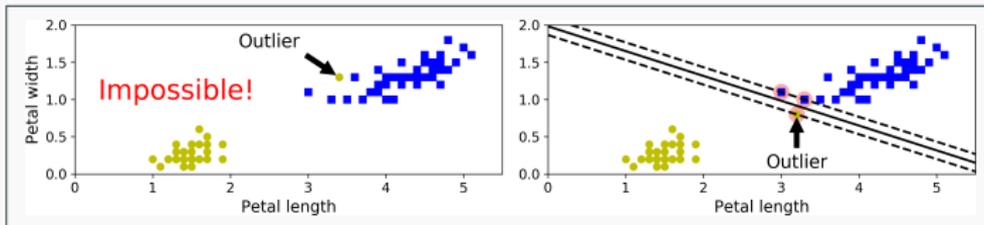


Figure 5-3. Hard margin sensitivity to outliers

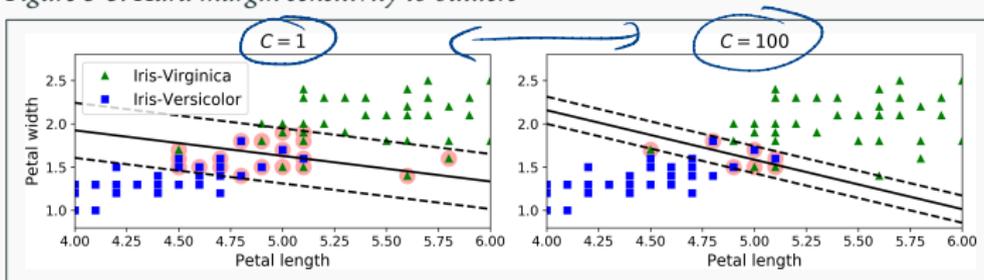


Figure 5-4. Large margin (left) versus fewer margin violations (right)

Soft  
SVM

# A brief history of SVM

Classification  
problem  
↓  
SVM

- SVM은 1992년 Boser, Guyon and Vapnik에 의해서 소개됨
- Statistical Learning Theory에 이론적 바탕을 둔 알고리즘 (Vapnik & Chervonenkis)  
*MNIST.*
- 손글씨 숫자 인식에서 뛰어난 성능을 보이면서 널리 쓰이게 됨
- SVM으로 1.1% Test error rate  $\approx$  신중히 설계된 신경망과 맞먹음  
*Neural Network.*
- 실용적으로 우수한 성능
- Bioinformatics, text, image recognition 등을 포함한 많은 성공 사례
- 선형/비선형 분류 뿐 아니라 회귀, outlier detection 도 수행
- 복잡한 소규모/중규모 데이터셋의 분류에 특히 잘 맞음
- Kernel 방법을 사용하는 대표적 알고리즘
- Liblinear & libsvm: Scikit-Learn 에서 liblinear 및 libsvm 을 사용

## Hard margin SVM: Theory

---

# Maximum margin classifier

training data  $\rightarrow$   $W, b$   $\stackrel{?}{\approx}$   $\frac{2}{2+1}$

We begin our discussion of support vector machines to the two-class classification problem using linear models of the form

$$\underline{y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b} \quad \tilde{\mathbf{w}}^T \mathbf{x} \quad (1)$$

where  $\phi(\mathbf{x})$  denotes a fixed feature-space transformation, and we have made the bias parameter  $b$  explicit.

The training data set comprises  $N$  input vectors  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ , with corresponding target values  $t_1, t_2, \dots, t_N$  where  $t_n \in \{-1, 1\}$ , and new data points  $\mathbf{x}$  are classified according to the sign of  $y(\mathbf{x})$

$$y(\mathbf{x}_n) \geq 0 \Leftrightarrow t_n = \begin{matrix} +1 \\ -1 \end{matrix}$$

# Maximum margin classifier

We shall assume that the training data set is linearly separable in feature space, so that by definition there exists at least one choice of the parameters  $\mathbf{w}$  and  $b$  such that a function satisfies  $y(\mathbf{x}_n) > 0$  for points having  $t_n = +1$  and  $y(\mathbf{x}_n) < 0$  for points having  $t_n = -1$ , so that  $t_n y(\mathbf{x}_n) > 0$  for all training data points.

~~~~~

Classification  $\frac{2}{2}$   $\frac{1}{1}$   $\frac{1}{1}$   $\frac{1}{1}$ .

$$\begin{aligned} y(\mathbf{x}_n) > 0, \quad t_n = +1. \\ y(\mathbf{x}_n) < 0, \quad t_n = -1. \end{aligned} \iff \underline{y(\mathbf{x}_n) \cdot t_n > 0.}$$

misclass

$$\underline{y(\mathbf{x}_n) \cdot t_n < 0}$$

# Maximum margin classifier: optimality criterion

Thus the distance of a point  $\mathbf{x}_n$  to the decision surface is given by

$$\min_{n=1, \dots, N} \left[ \frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|} = \frac{t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|} \right] = \text{margin} \quad (2)$$

= closest dist.

The margin is given by the perpendicular distance to the closest point  $\mathbf{x}_n$  from the data set, and we wish to optimize the parameters  $\mathbf{w}$  and  $b$  in order to maximize this distance. Thus the maximum margin solution is found by solving

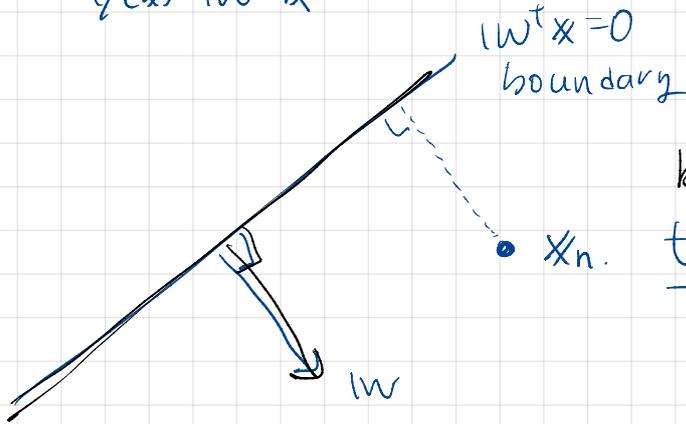
$$\text{max. margin.} \quad \arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min [t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)] \right\}$$

margin.      max. min dist.  
= margin (3)  
max. margin

where we have taken the factor  $1/\|\mathbf{w}\|$  outside the optimization over  $n$  because  $\mathbf{w}$  does not depend on  $n$ .

→ cost. function 2ger.  
= max. margin classifier.

$$y(x) = \mathbf{w}^T \mathbf{x}$$



boundary, decision surface  $\sim \mathbf{x}_n^T \mathbf{w} + b$

$$\frac{\mathbf{t}_n \cdot y(\mathbf{x}_n)}{|\mathbf{w}|} = \frac{\mathbf{t}_n \cdot (\mathbf{w}^T \mathbf{x}_n + b)}{|\mathbf{w}|}$$

# Dual problem for convex optimization

Primal optimization problem for **hard SVM**  $\rightarrow$  opt. prob.

$\rightarrow$  maximize  $\left[ \frac{1}{\|\mathbf{w}\|} \min [t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)] \right]$  (4)

$\rightarrow$  subject to  $t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 0$  (5)

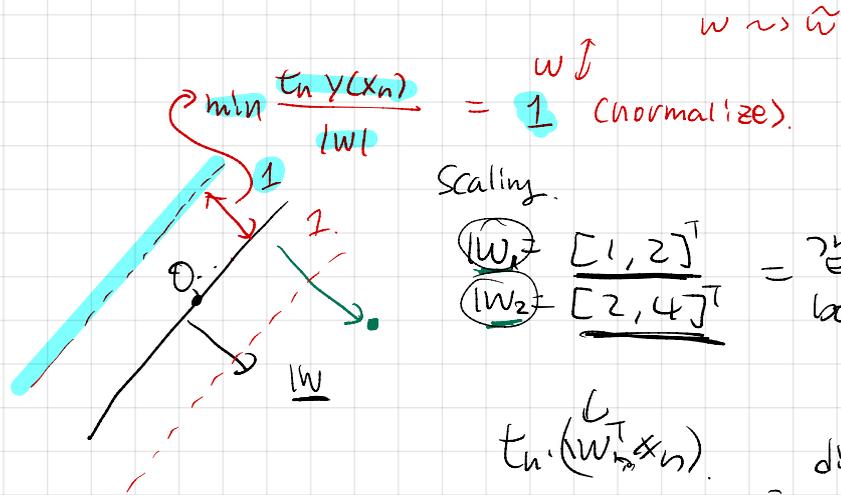
Normalize  $t_n y(x_n) \geq 0. \Rightarrow$  모든 sample이 자기 영역에 위치하도록. (6)

Equivalently,  $\frac{1}{\|\mathbf{w}\|} \rightarrow \frac{1}{2} \|\mathbf{w}\|^2$

$\max \frac{1}{\|\mathbf{w}\|}, \min \|\mathbf{w}\|, \text{ minimize } \frac{1}{2} \|\mathbf{w}\|^2 \rightarrow$  convex opt. (7)

$\min \|\mathbf{w}\|^2, \min \frac{1}{2} \|\mathbf{w}\|^2$  subject to  $t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1$  (8)

Direct solution of this optimization problem would be very complex, so we shall convert it into an equivalent problem that is much easier to solve.



Scaling.

$$\begin{matrix} |w_1| & [1, 2]^T \\ |w_2| & [2, 4]^T \end{matrix} = \begin{matrix} \text{같은} \\ \text{boundary} \end{matrix}$$

$$\frac{t_n (w^T x_n)}{|w|} = \begin{matrix} \text{dist.} \\ \text{같은} \\ \text{class} \end{matrix}$$

# Lagrangian function with constraint

In order to solve this constrained optimization problem, we introduce **Lagrange multipliers**  $a_n \geq 0$ , with one multiplier  $a_n$  for each of the constraints, giving the Lagrangian function

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n \{t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - 1\} \quad (9)$$

*Handwritten notes: "inequality" above the brace, and a wavy underline under the denominator 2.*

where  $\mathbf{a} = (a_1, \dots, a_N)^T$ . Note the minus sign in front of the Lagrange multiplier term, because we are minimizing with respect to  $\mathbf{w}$  and  $b$ , and maximizing with respect to  $\mathbf{a}$ .

Setting the derivatives of  $L(\mathbf{w}, b, \mathbf{a})$  with respect to  $\mathbf{w}$  and  $b$  equal to zero, we obtain the following two conditions

Convex opt.

$$\begin{aligned} \rightarrow \mathbf{w} &= \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n) & \leftarrow a_n = 0 \Rightarrow & (10) \\ & & & \text{여기는 아무 영향도 안} \end{aligned}$$
$$\begin{aligned} \rightarrow 0 &= \sum_{n=1}^N a_n t_n & a_n \neq 0 \Rightarrow & \text{영향} = 0. & (11) \\ & & & \text{Support vector.} \end{aligned}$$

*Handwritten notes: "Support vector" below equation (11).*

# Lagrangian function with constraint

Eliminating  $\mathbf{w}$  and  $b$  from  $L(\mathbf{w}, b, \mathbf{a})$  using these conditions then gives the *dual representation* of the maximum margin problem in which we maximize

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m) \quad (12)$$

with respect to  $\mathbf{a}$  subject to the constraints

$$a_n \geq 0, \quad n = 1, \dots, N \quad (13)$$

$$\sum_{n=1}^N a_n t_n = 0. \quad (14)$$

Here the kernel function is defined by  $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$ .

# Prediction for a new sample: support vector machine

$$x_n^T x$$

x. 라는 벡터의  
x\_n 방향으로의 성분



In order to classify new data points using the trained model, we evaluate the sign of  $y(x)$ . This can be expressed in terms of the parameter  $\{a_n\}$  and the kernel function by substituting for  $w$  to give

x\_n 들이 모든 입력.

$$y(x) = \sum_{n=1}^N a_n t_n k(x, x_n) + b \quad (15)$$

↓

$$y(x) = \underline{w}^T x + b$$

kernel.

$x_n$ , support vector  $\vec{x}$   
( $a_n > 0$ ).

$$= \sum a_n t_n \boxed{x_n^T x} + b. \geq 0.$$

↳ support vector  $\vec{x}$ 이고  $x$ 가 있을 때

$$y(x) = \sum a_n t_n \phi(x_n)^T \phi(x) + b.$$

$$K(x, x_n).$$

$\Rightarrow$  두개의 vector의

유사성 비교하는  $\vec{x}_1, \vec{x}_2$ .

kernel function.

# KKT condition: complementary slackness

We show that a constrained optimization of this form satisfies the **Karush-Kuhn-Tucker (KKT)** conditions, which in this case require that the following three properties hold

dual.  $\lambda$ . multi  $\geq 0$ . ②  $\rightarrow a_n \geq 0$  (16)

primal. const. ①  $\frac{t_n y(\mathbf{x}_n) - 1 \geq 0}{\text{Chord SUM}}$  (17)

Complementary slackness ③.  $a_n \{t_n y(\mathbf{x}_n) - 1\} = 0$  (18)

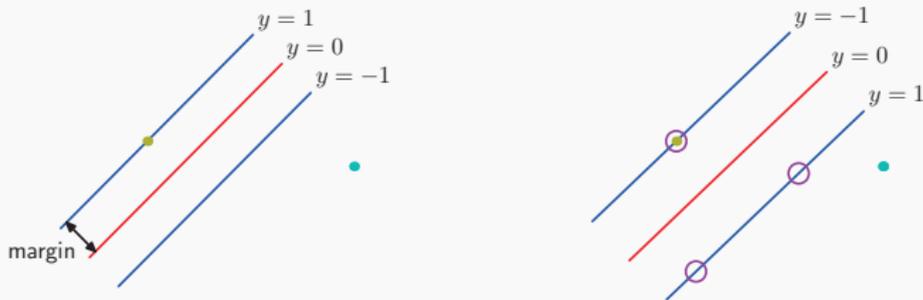
Thus for every data point, either  $a_n = 0$  or  $t_n y(\mathbf{x}_n) = 1$ . Any data point for which  $a_n = 0$  will not appear in the sum and hence plays no role in making predictions for new data points. The remaining data points are called *support vectors*, and because they satisfy  $t_n y(\mathbf{x}_n) = 1$ , they correspond to points that lie on the maximum margin hyperplanes in feature space.

$a_n = 0$ , ~~support~~  $\Rightarrow t_n y(\mathbf{x}_n) - 1 \neq 0$ , margin  $\neq 0$ , simplex.

$a_n > 0$ , **Support**  $\Rightarrow t_n y(\mathbf{x}_n) - 1 = 0$ , margin  $= 0$ ,  $\frac{0}{101}$ ,  $\frac{13}{16}$

# KKT condition: complementary slackness

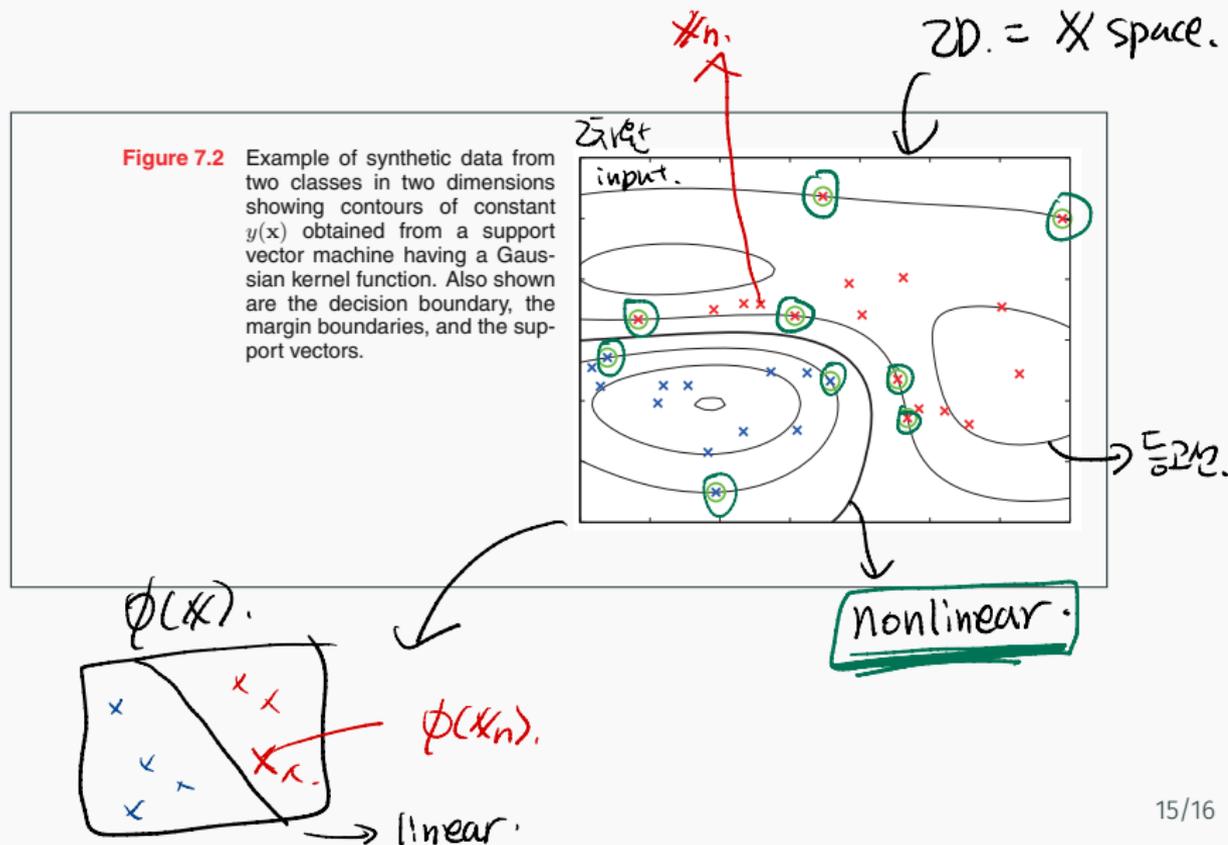
$\rightarrow \bar{\lambda}_1 \bar{x}_1$  margin = 1.



**Figure 7.1** The margin is defined as the perpendicular distance between the decision boundary and the closest of the data points, as shown on the left figure. Maximizing the margin leads to a particular choice of decision boundary, as shown on the right. The location of this boundary is determined by a subset of the data points, known as support vectors, which are indicated by the circles.

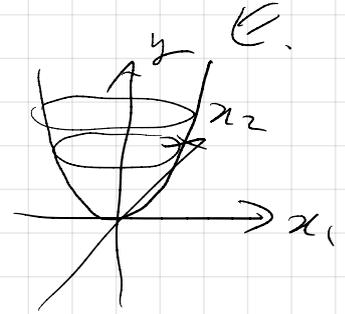
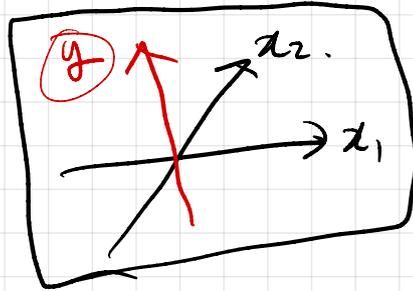
# KKT condition: complementary slackness

**Figure 7.2** Example of synthetic data from two classes in two dimensions showing contours of constant  $y(x)$  obtained from a support vector machine having a Gaussian kernel function. Also shown are the decision boundary, the margin boundaries, and the support vectors.



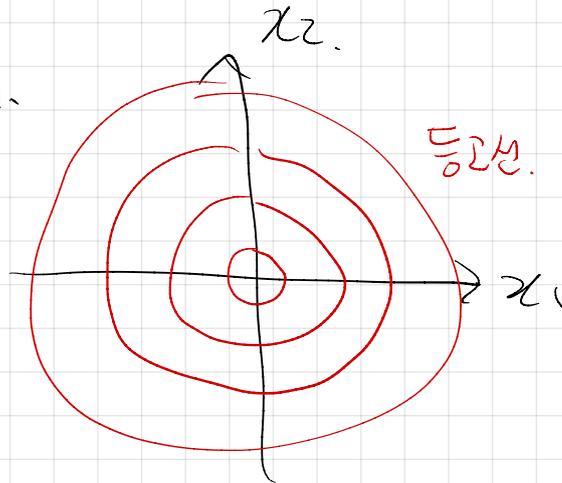
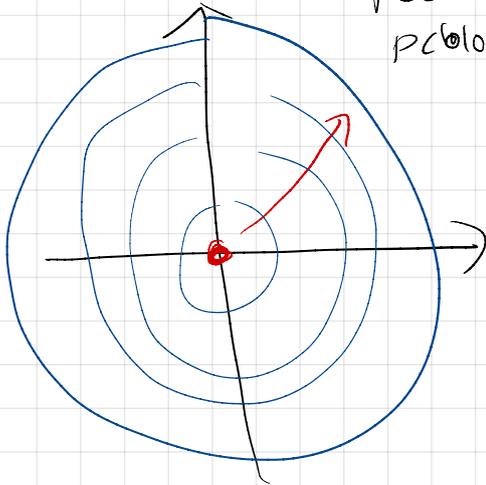
pcolormesh.

$y(x) \geq 0$ .  
2D  $\rightarrow$  1D.



$$y = x_1^2 + x_2^2$$

pcolor  
pcolormesh.



# Kernel Method

---

# Kernel method

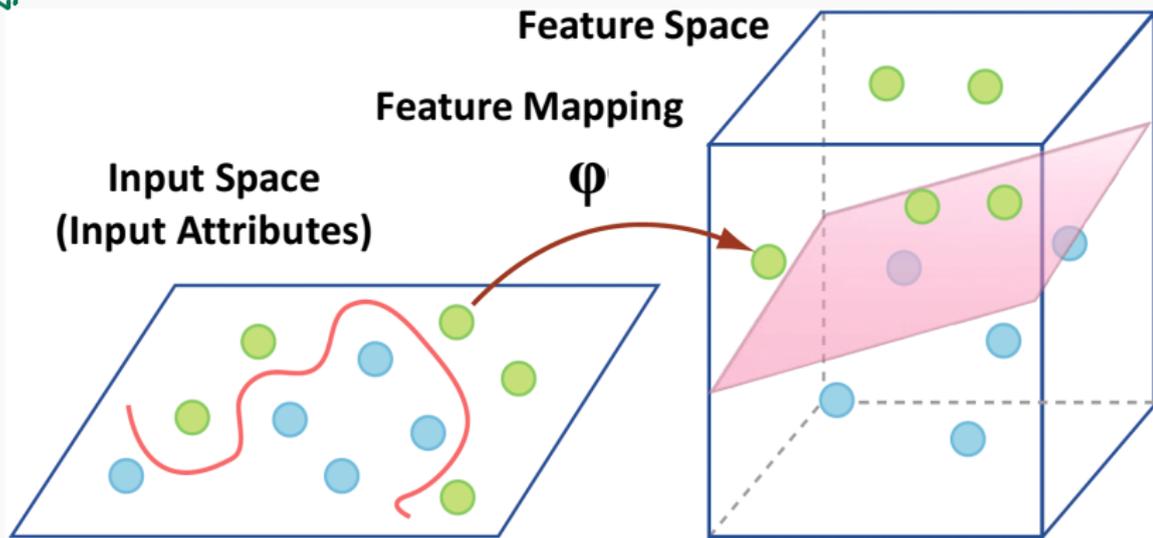


Image by MIT OpenCourseWare

**Kernel function:** Kernel function이란 어떤 feature space에서 feature의 inner product와 동등한 함수를 말함. Feature mapping  $\phi(\mathbf{x})$ 에 대한 kernel function은 다음과 같이 정의할 수 있음.

$$\text{예 1)} \quad k(\mathbf{x}, \mathbf{x}') = \underline{\phi(\mathbf{x})^T \phi(\mathbf{x}')}$$
 (19)

**Similarity:** Kernel function  $k(\mathbf{x}, \mathbf{x}')$ 는  $\mathbf{x}$ 와  $\mathbf{x}'$  간의 similarity로 볼 수 있음.

또한 일반적으로 feature mapping을 통한  $\phi(\mathbf{x})^T \phi(\mathbf{x}')$ 에 비해  $k(\mathbf{x}, \mathbf{x}')$ 의 계산이 훨씬 간단함. 예를 들어

$$k(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^T \mathbf{x}')^2$$
 (20)

$$\phi(\mathbf{x}) = [ 1 \quad x_1^2 \quad \sqrt{2}x_1x_2 \quad x_2^2 \quad \sqrt{2}x_1 \quad \sqrt{2}x_2 ]$$
 (21)

# Construction of kernel

X 안함, 시험X.

임의의 kernel function 대해서  $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$ 가 성립하는 feature mapping이 존재하는지 검사하는 것은 매우 까다로움.

## Mercer's theorem

Every positive semi-definite symmetric function is a kernel.

다음과 같이 정의되는 Gram matrix  $K$ 가 positive semi-definite symmetric matrix면 function  $k$ 는 kernel function임.

$$\mathbf{K} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix} \quad (22)$$

# Gaussian radial basis function (RBF) kernel

많이 쓰는 kernel.

Gaussian RBF kernel

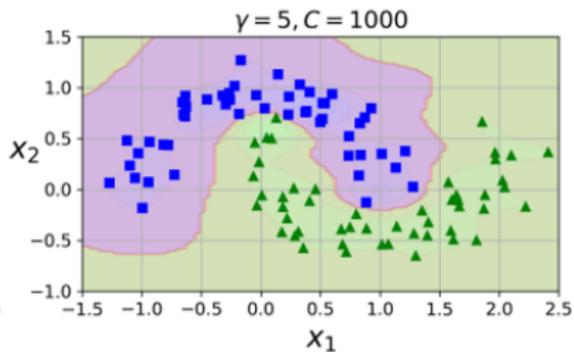
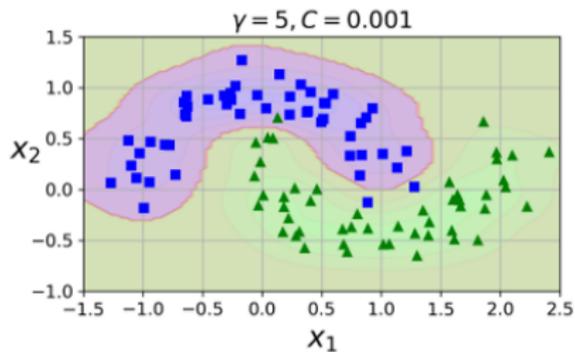
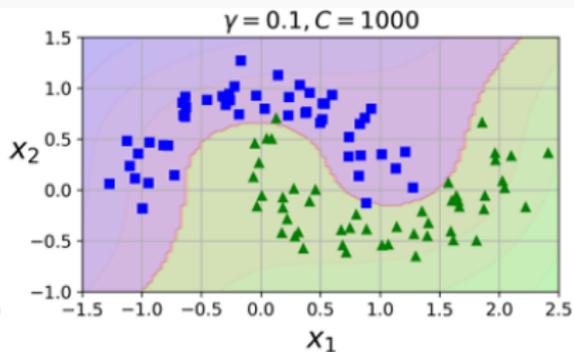
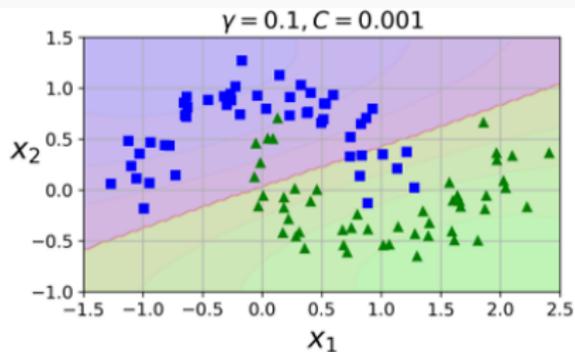
$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right) = \exp(-\underbrace{\gamma}_{\sim} \underbrace{\|\mathbf{x} - \mathbf{x}'\|}_{\sim}^2) \quad (23)$$

유사도:  $\mathbf{x}, \mathbf{x}'$  가까움  $\rightarrow \uparrow$   
멀면  $\rightarrow \downarrow$ .

- kernel value는  $\mathbf{x}'$ 가 landmark  $\mathbf{x}$ 와 가까울수록 큰 값
- $\gamma$  값이 커지는 경우 landmark에서 멀어질수록 kernel 값이 급격하게 작아짐 = 영향력의 범위가 좁아짐

# Hyperparameter of soft SVM-RBF

X



# SVM: strength and limitations

장점  $\left[ \begin{array}{l} \text{Convex.} \\ \text{opt.} \end{array} \right] \rightarrow \text{global. 최적화.}$

G.D.  $\rightarrow$  local 문제.

- 강력하고 우수한 이론을 바탕으로 함
- 많은 블랙박스 알고리즘과는 대조적으로 비교적 직관적인 해석과 이해가 가능  $\hookrightarrow$  설명, NN
- 학습이 상대적으로 쉬움  $\rightarrow$  Conv. opt.
- 신경망처럼 지역 최적값에 빠지는 일이 없음
- 학습 시간이 차원에 의존하지 않으며 kernel trick 덕분에 고정된 입력에만 의존
- 과적합이 잘 조절되는 경향
- 많은 분야에서 신경망 및 기타 알고리즘과 필적하는 성능
- 데이터가 작은 조건이나 고차원 공간에서도 잘 일반화

# SVM: strength and limitations

Cost  $\leftarrow$  max. margin.

hard SVM  $\rightarrow$  linearly separable.

Soft SVM  $\rightarrow$  완한.



다른 알고리즘.

단점

- ① 노이즈에 민감: 비교적 적은 수의 잘못된 label로 성능이 심각하게 악화
  - 커널 함수의 선택/구축하는 방법에 대한 정리된 원칙이 없음
  - hyperparameter C의 적정값을 정하기 위한 원칙이 없음
- ② 컴퓨터의 메모리와 계산 시간 측면에서 비용이 높은 편이며 multiclass에서 더욱더 심화됨

ECoC.

Support vector 기법.

- 항상 들고 다녀야 함

-  $\forall$ .  $\rightarrow$  S.V. 들고 kernel은 전부 비교.

# Appendix

---

## Reference and further reading

- “Chap 7 | Sparse Kernel Machines” of C. Bishop, Pattern Recognition and Machine Learning
- “Chap 5 | Support Vector Machines” of A. Geron, Hands-On Machine Learning with Scikit-Learn, Keras & TensorFlow
- “Chap 4 | Convex Optimization Problems”, “Chap 5 | Duality” of S. Boyd, Convex Optimization
- “Lecture 6 | Support Vector Machines” of Kwang Il Kim, Machine Learning (2019)