# Support Vector Machine

## Instructor: Junghye Lee

**Department of Industrial Engineering**
**junghyelee@unist.ac.kr**

# Contents

**1** **Linear Support Vector Machine (Separable case)**

# Linear Classifier

$$x \longrightarrow \boxed{f} \longrightarrow \widehat{y}$$

- denotes +1 (class 1)
- denotes -1 (class 2)

$$f(x, w, b) = \text{sign}(w^T x + b)$$

$w^T x + b > 0$

$w^T x + b = 0$

$w^T x + b < 0$

**How would you classify this data?**

# Linear Classifier

$$x \longrightarrow \boxed{f} \longrightarrow \widehat{y}$$

- denotes +1 (class 1)
- denotes -1 (class 2)

$w^T x + b > 0$

$f(x, w, b) = \text{sign}(w^T x + b)$

**Which one is the best?**

$w^T x + b < 0$

# Linear Classifier



$$f(x, w, b) = \text{sign}(wx + b)$$

- denotes +1 (class 1)
- denotes -1 (class 2)

$wx + b > 0$

$wx + b = 0$

$wx + b < 0$

**Define the margin of a linear classifier as the width that the boundary could be increased by before hitting a data point.**

# Linear Classifier

$$x \longrightarrow \boxed{f} \longrightarrow \hat{y}$$

$$w^T x + b = 1$$

$$f(x, w, b) = \text{sign}(w^T x + b)$$

- denotes +1 (class 1)

$$w^T x + b > 0$$

- denotes -1 (class 2)

$$w^T x + b = -1$$

**Support Vectors**
are those data
points that the
margin pushes up
against.

$$w^T x + b < 0$$

$$w^T x + b = 0$$

# Linear Classifier

• denotes +1 (class 1)

○ denotes -1 (class 2)

$w^T x + b$)

**Support Vectors** are those data points that the margin pushes up against.

$w^T x + b = 0$

$w^T x + b < 0$

# Linear SVM Mathematically

**What we know:**

$$w^T x^+ + b = +1$$
$$w^T x^- + b = -1$$



"Predict Class = +1" zone

$x^+$

Margin

$x^-$

$wx + b = 1$

$wx + b = 0$

$wx + b = -1$

"Predict Class = -1" zone

$$M(\text{margin}) = \frac{2}{\|w\|_2} = \frac{2}{\sqrt{w^T w}}$$

# Linear SVM Mathematically



**Normal vector**

$w^T x + b = 1$

$w^T x + b = 0$

$w^T x + b = -1$

**What we know:** $x^+ = x^- + \lambda w$

## 1st step

$w^T x^+ + b = 1$

$w^T(x^- + \lambda w) + b = 1$

$w^T x^- + b + \lambda w^T w = 1$

$-1 + \lambda w^T w = 1$

$\therefore \lambda = \dfrac{2}{w^T w}$

## 2nd step

$$Margin = distance(x^+, x^-)$$

$$= \|x^+ - x^-\|_2$$

$$= \|x^- + \lambda w - x^-\|_2$$

$$= \|\lambda w\|_2$$

$$= \lambda \sqrt{w^T w}$$

$$= \frac{2}{w^T w} \sqrt{w^T w}$$

$$\therefore \frac{2}{\sqrt{w^T w}} = \frac{2}{\|w\|_2}$$

# Linear SVM Mathematically

1) Correctly classify all training data

$$w^T x_i^+ + b \geq +1 \qquad \text{if } y_i = +1$$

$$w^T x_i^- + b \leq -1 \qquad \text{if } y_i = -1$$

$$y_i(w^T x_i + b) \geq +1 \qquad \text{for all } i$$

2) Maximize the Margin $\dfrac{2}{\|w\|_2}$ same as minimize $\dfrac{1}{2} w^T w$

**We can formulate a Quadratic Optimization Problem and solve for $w$ and $b$**

**Minimize** $\quad \Phi(w) = \dfrac{1}{2} w^T w$

**subject to** $\quad y_i(w^T x_i + b) \geq +1$

# Solving the Optimization Problem

**Find $w$ and $b$ such that**

$\Phi(w) = \frac{1}{2} w^T w$ **is minimized;**

**and for all** $\{(x_i, y_i)\}: y_i(w^T x_i + b) \geq 1$

- Need to optimize a *quadratic* function subject to **linear** constraints.
- Quadratic optimization problems are a well-known class of mathematical programming problems, and many (rather intricate) algorithms exist for solving them.
- The solution involves constructing a **dual problem** where a **Lagrange multiplier $\alpha_i$** is associated with every constraint in the primary problem:

**Find $\alpha_1 \ldots \alpha_N$ such that**

$Q(\alpha) = \Sigma\alpha_i - \frac{1}{2}\Sigma\Sigma\alpha_i\alpha_j y_i y_j x_i^T x_j$ **is maximized and**

(1) $\Sigma\alpha_i y_i = 0$

(2) $\alpha_i \geq 0$ for all $\alpha_i$

# The Optimization Problem Solution

- The solution has the form:

$$\boldsymbol{w} = \Sigma_i \alpha_i y_i \boldsymbol{x_i} \quad b = y_i - \boldsymbol{w^T x_i} \quad \text{for any } \boldsymbol{x_i} \text{ such that } \alpha_i \neq 0$$

- Each non-zero $\alpha_i$ indicates that corresponding $\boldsymbol{x_i}$ is a support vector.

- Then the classifying function will have the form:

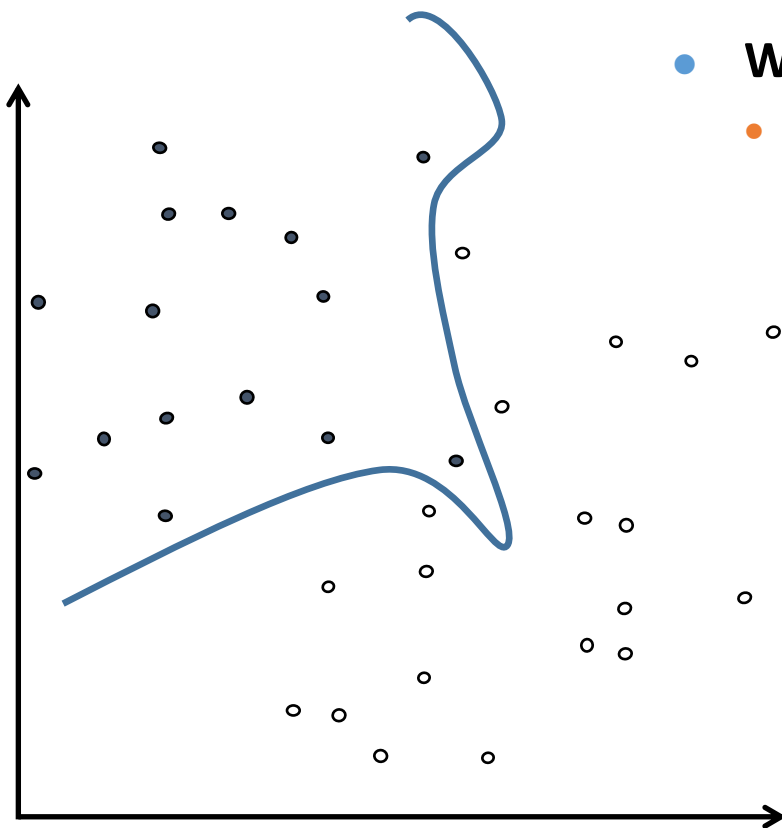$$f(\boldsymbol{x}) = \Sigma \alpha_i y_i \boldsymbol{x_i^T x} + b$$

- Notice that it relies on an *inner product* between the test point **x** and the support vectors $\boldsymbol{x_i}$ we will return to this later.

- Also keep in mind that solving the optimization problem involved computing the inner products $\boldsymbol{x_i^T x_j}$ between all pairs of training points.

**2** Linear Support Vector Machine
(Non-separable case)

# Dataset with noise
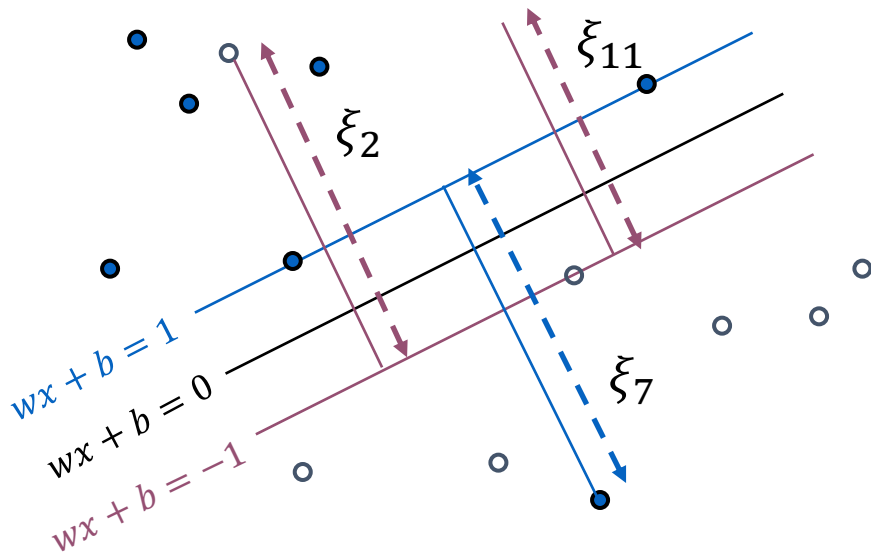
- denotes +1 (class 1)
- denotes -1 (class 2)



- **Hard Margin: So far we require all data points to be classified correctly**
  - **No training error**
- **What if the training set is noisy?**
  - **Solution 1: use very powerful kernels**

# Overfitting!

# Soft Margin Classification

*Slack variables* $\xi_i$ can be added to allow misclassification of difficult or noisy examples.



**What should our quadratic optimization criterion be?**

Minimize $\quad \dfrac{1}{2}\boldsymbol{w}^T\boldsymbol{w} + C\Sigma_{k=1}^{R}\xi_k$

# Hard Margin v.s. Soft Margin

- The old formulation:

**Find $w$ and $b$ such that**
$$\Phi(w) = \frac{1}{2} w^T w \text{ is minimized and for all } \{(x_i, y_i)\}$$
$$y_i (w^T x_i + b) \geq 1$$

- The new formulation incorporating slack variables:

**Find $w$ and $b$ such that**
$$\Phi(w) = \frac{1}{2} w^T w + C \Sigma \xi_i \text{ is minimized and for all } \{(x_i, y_i)\}$$
$$y_i (w^T x_i + b) \geq 1 - \xi_i \text{ and } \xi_i \geq 0 \text{ for all } i$$

- Parameter $C$ can be viewed as a way to control overfitting.

# Linear SVM:  Overview

- The classifier is a *separating hyperplane.*

- Most "important" training points are support vectors; they define the hyperplane.

- Quadratic optimization algorithms can identify which training points $x_i$ are support vectors with non-zero Lagrangian multipliers $\alpha_i$.

- Both in the dual formulation of the problem and in the solution training points appear only inside dot products:

**Find $\alpha_1 \dots \alpha_N$ such that**

$Q(\alpha) = \Sigma\alpha_i - \frac{1}{2}\Sigma\Sigma\alpha_i\alpha_j y_i y_j x_i^T x_j$ **is maximized and**

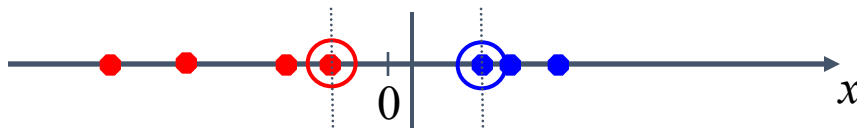(1) $\Sigma\alpha_i y_i = 0$

(2) $0 \leq \alpha_i \leq C$ for all $\alpha_i$

$$f(x) = \Sigma\alpha_i y_i x_i^T x + b$$

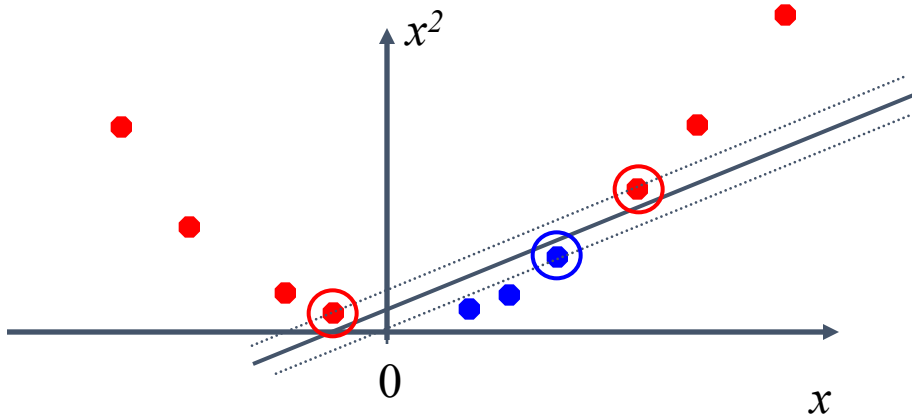**3**     **Nonlinear Support Vector Machine**

# Non-linear SVM

- Datasets that are linearly separable with some noise work out great:



- But what are we going to do if the dataset is just too hard?
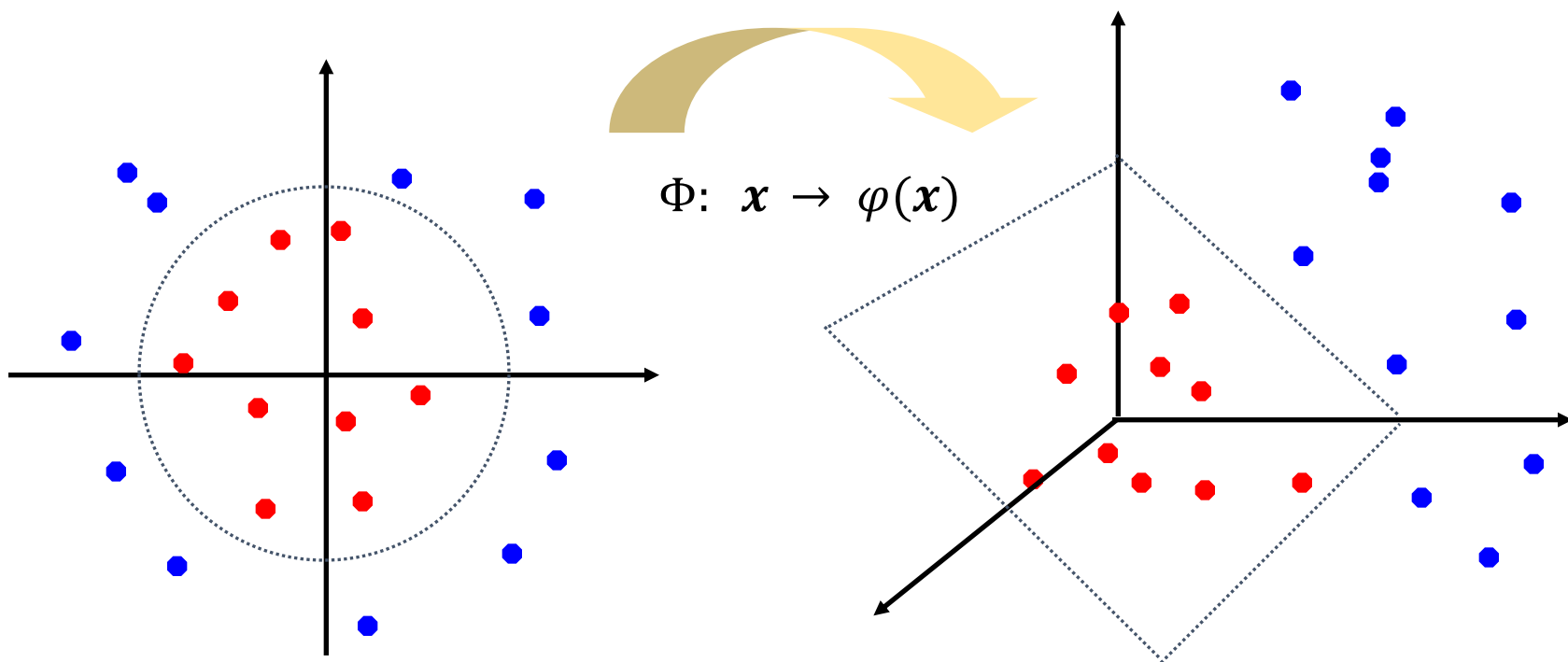


- How about mapping data to a higher-dimensional space:

# Non-linear SVM:  Feature spaces

- General idea:   the original input space can always be mapped to some higher-dimensional feature space where the training set is separable:

$$\Phi: \; \boldsymbol{x} \; \rightarrow \; \varphi(\boldsymbol{x})$$

# Kernel Trick

- The linear classifier relies on dot product between vectors $K(x_i, x_j) = x_i^T x_j$

- If every data point is mapped into high-dimensional space via some transformation $\Phi:\ x \to \varphi(x)$, the dot product becomes:
$$K(x_i, x_j) = \varphi(x_i)^T \varphi(x_j)$$

- A *kernel function* is some function that corresponds to an inner product in some expanded feature space.

- Example: 2-dimensional vectors $x = [x_1\ x_2]$; let $K(x_i, x_j) = \left(1 + x_i^T x_j\right)^2$,
Need to show that $K(x_i, x_j) = \varphi(x_i)^T \varphi(x_j)$:

$K(x_i, x_j)$

$= \left(1 + x_i^T x_j\right)^2$

$= 1 + x_{i1}^2 x_{j1}^2 + 2x_{i1}x_{j1}x_{i2}x_{j2} + x_{i2}^2 x_{j2}^2 + 2x_{i1}x_{j1} + 2x_{i2}x_{j2}$

$= \left[1\, x_{i1}^2\, \sqrt{2}x_{i1}x_{i2}\, x_{i2}^2\, \sqrt{2}x_{i1}\, \sqrt{2}x_{i2}\right]^T \left[1\, x_{j1}^2\, \sqrt{2}x_{j1}x_{j2}\, x_{j2}^2\, \sqrt{2}x_{j1}\, \sqrt{2}x_{j2}\right]$

$= \varphi(x_i)^T \varphi(x_j)$, where $\varphi(x) = \left[1\, x_1^2\, \sqrt{2}x_1x_2\, x_2^2\, \sqrt{2}x_1\, \sqrt{2}x_2\right]$

# What Functions are Kernels?

- For some functions $K(x_i, x_j)$ checking that $K(x_i, x_j) = \varphi(x_i)^T \varphi(x_j)$ can be cumbersome.

- Mercer's theorem: *Every semi-positive definite symmetric function is a kernel*

- Semi-positive definite symmetric functions correspond to a semi-positive definite symmetric Gram matrix:

$$K = \begin{array}{|c|c|c|c|c|}
\hline
K(x_1, x_1) & K(x_1, x_2) & K(x_1, x_3) & \ldots & K(x_1, x_N) \\
\hline
K(x_2, x_1) & K(x_2, x_2) & K(x_2, x_3) & & K(x_2, x_N) \\
\hline
\ldots & \ldots & \ldots & \ldots & \ldots \\
\hline
K(x_N, x_1) & K(x_N, x_2) & K(x_N, x_3) & \ldots & K(x_N, x_N) \\
\hline
\end{array}$$

# Semi-positive Definiteness of a Matrix

- In linear algebra, a **symmetric** $n \times n$ real matrix $M$ is said to be **positive definite** if the scalar $z^T M z$ **is strictly positive** for every non-zero column vector $z$ of $n$ real numbers.

$$z^T M z > 0 \text{ for all } z \in \mathbb{R}^n \setminus \mathbf{0}$$

- Positive semi-definite matrices are defined similarly, except that the above scalars $z^T M z$ must be positive or zero (i.e. non-negative).

$$z^T M z \geq 0 \text{ for all } z \in \mathbb{R}^n$$

# Examples of Kernel Functions

- Linear: $K(x_i, x_j) = x_i^T x_j$

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

- Gaussian (radial-basis function): $K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$

- Sigmoid: $K(x_i, x_j) = \tanh(\beta_0 x_i^T x_j + \beta_1)$

# Non-linear SVMs Mathematically

- Dual problem formulation:

**Find** $\alpha_1 \dots \alpha_N$ **such that**

$Q(\alpha) = \Sigma\alpha_i - \frac{1}{2}\Sigma\Sigma\alpha_i\alpha_j y_i y_j K(x_i, x_j)$ **is maximized and**

(1) $\Sigma\alpha_i y_i = 0$

(2) $\alpha_i \geq 0$ for all $\alpha_i$

- The solution is:

$$f(x) = \Sigma\alpha_i y_i K(x_i, x_j) + b$$

- Optimization techniques for finding $\alpha_i$'s remain the same!

# Nonlinear SVM - Overview

- SVM locates a separating hyperplane in the feature space and classify points in that space

- It does not need to represent the space explicitly, simply by defining a kernel function

- The kernel function plays the role of the dot product in the feature space.