

Advanced Classification Model

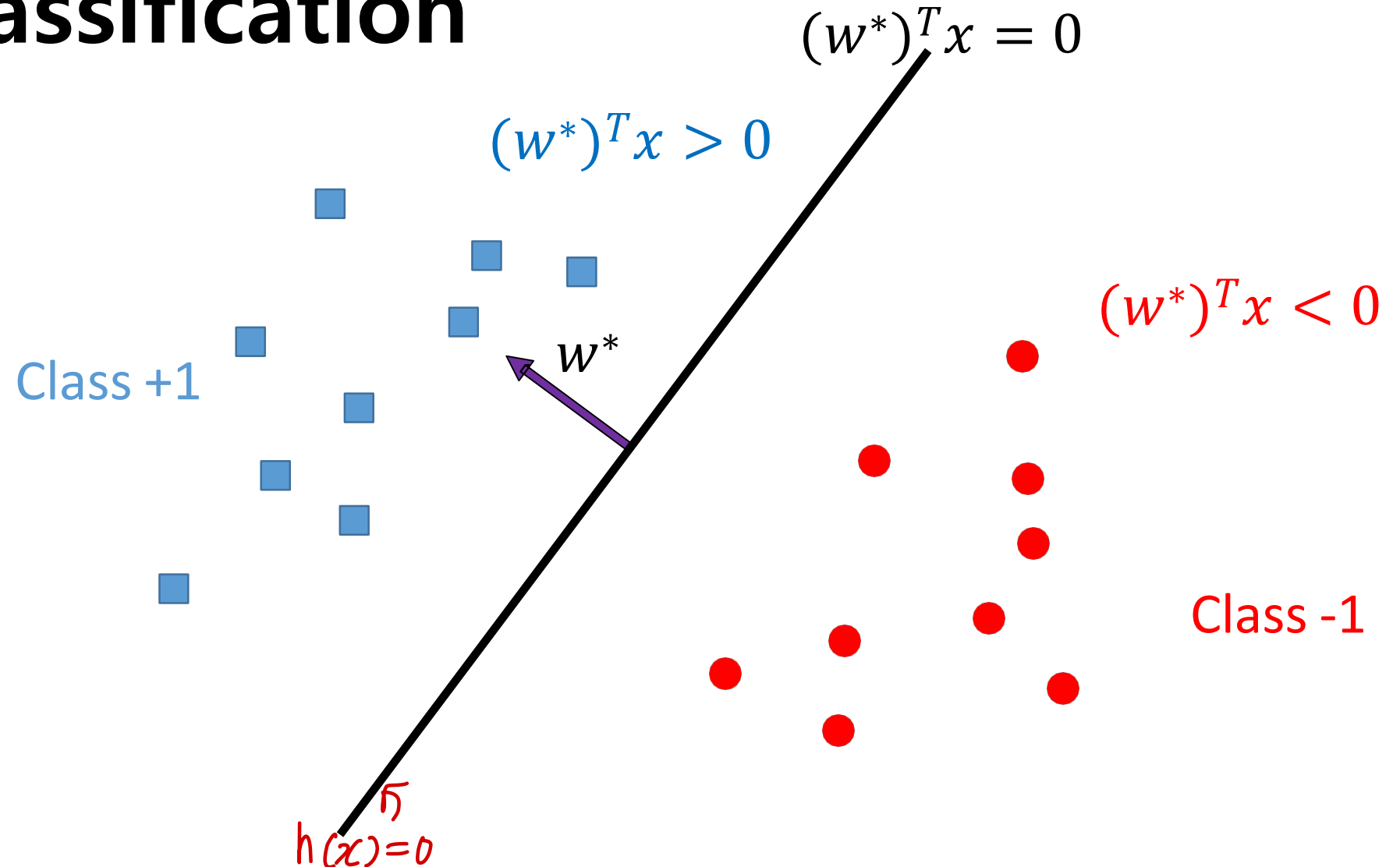
Linear and non-linear model...

Prof. Je-Won Kang
Electronic & Electrical Engineering
Ewha Womans University

Contents

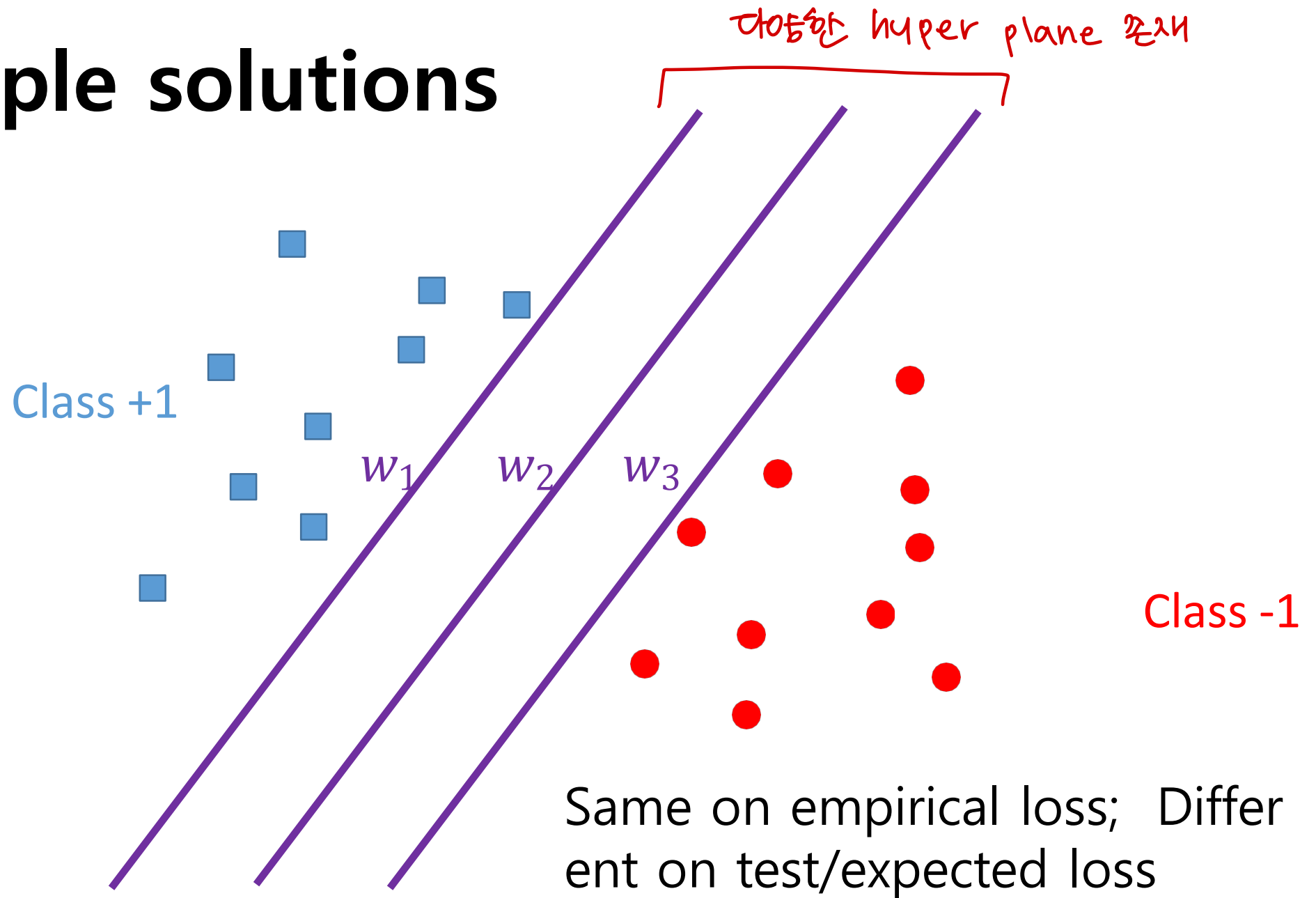
- Support vector machine (SVM)
- Neural network (NN)

Linear classification

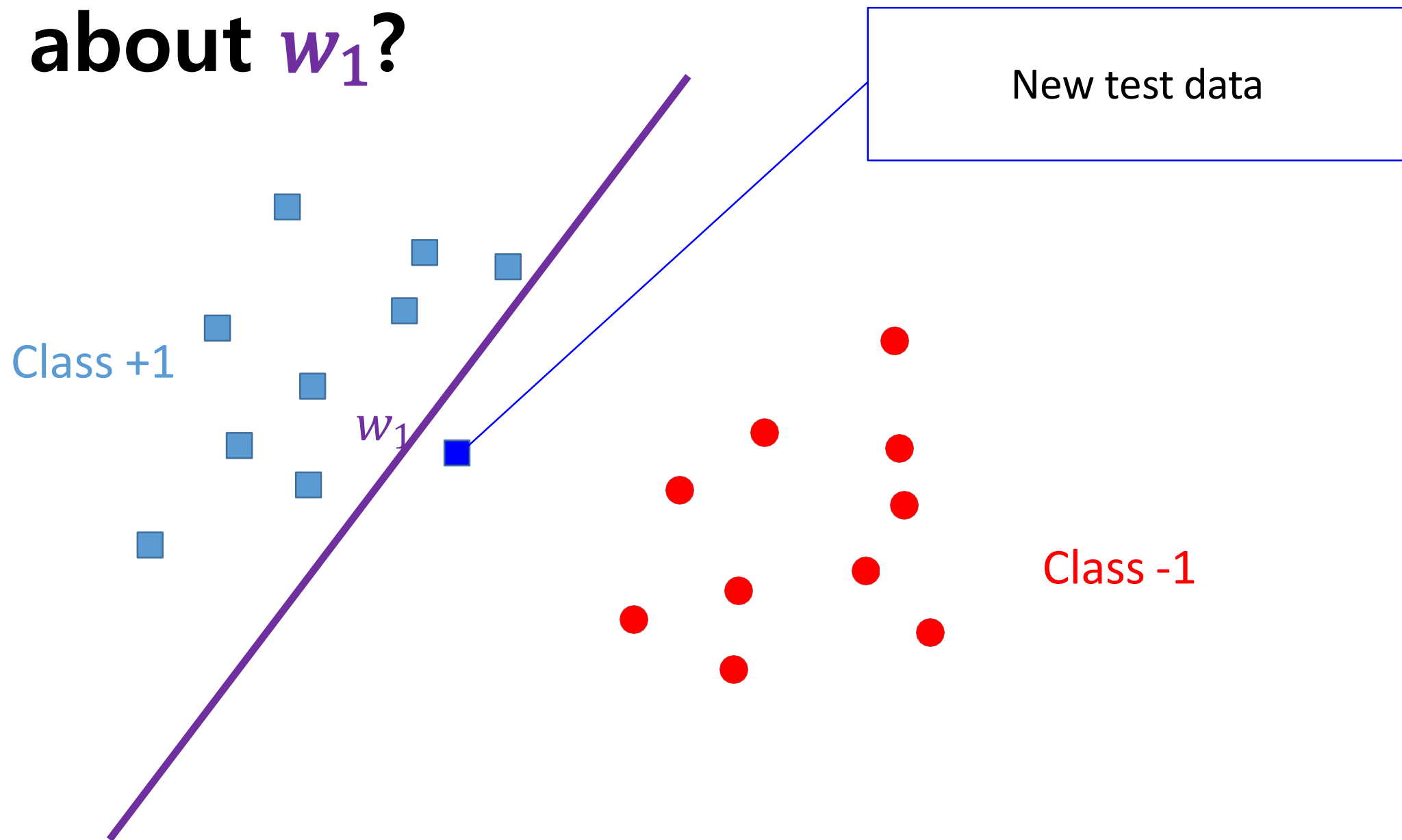


Assume perfect separation between the two classes using discriminant function w^*

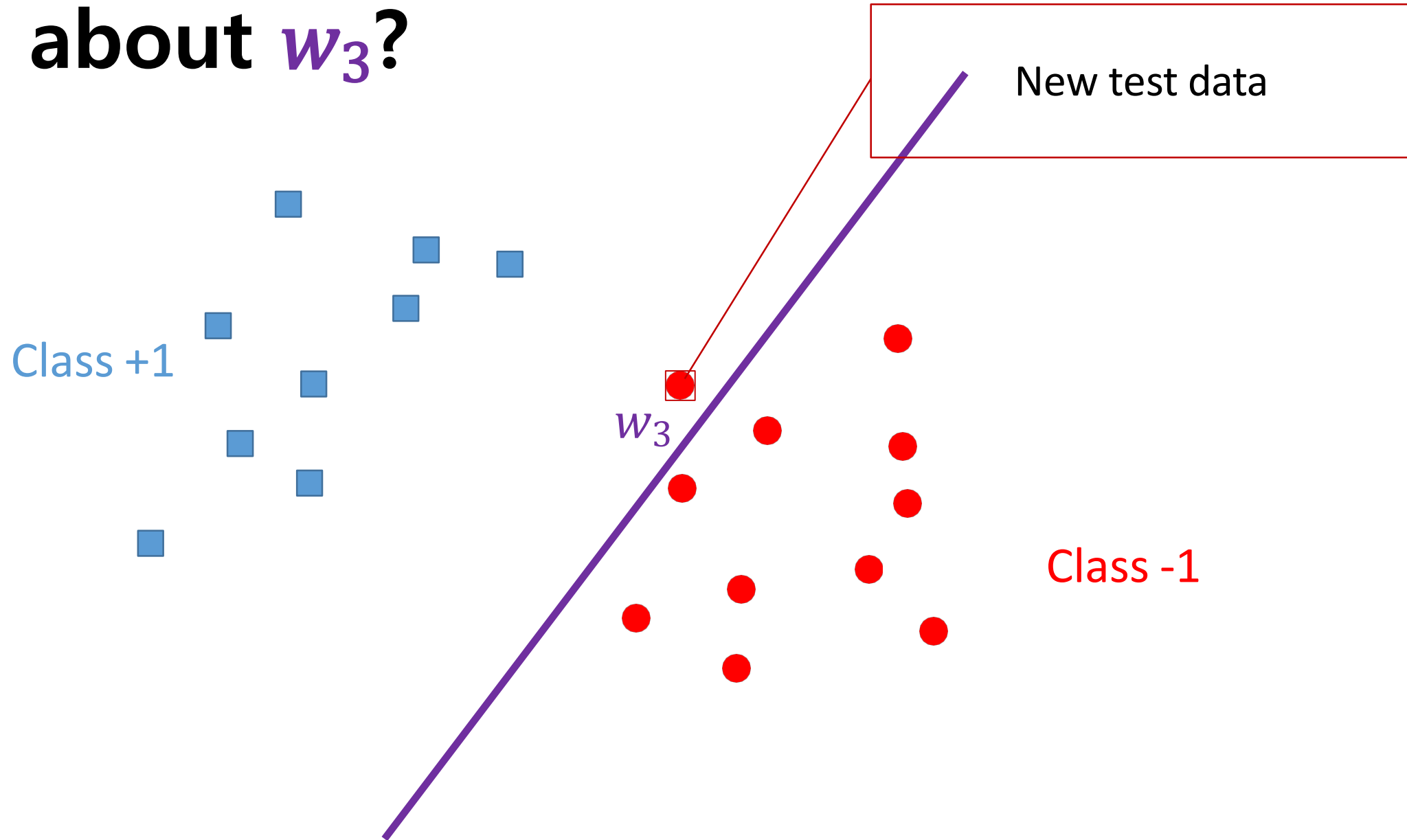
Multiple solutions



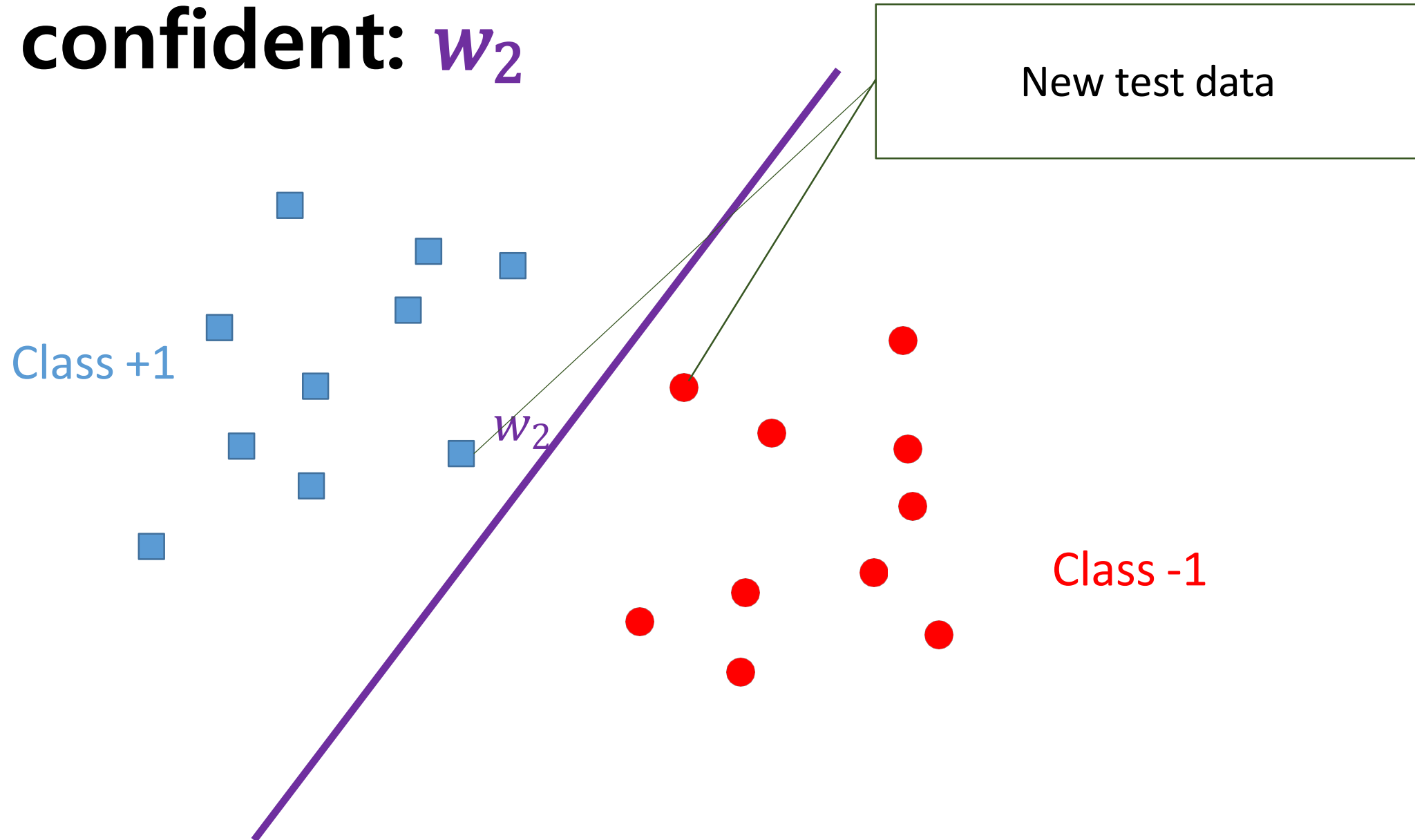
What about w_1 ?



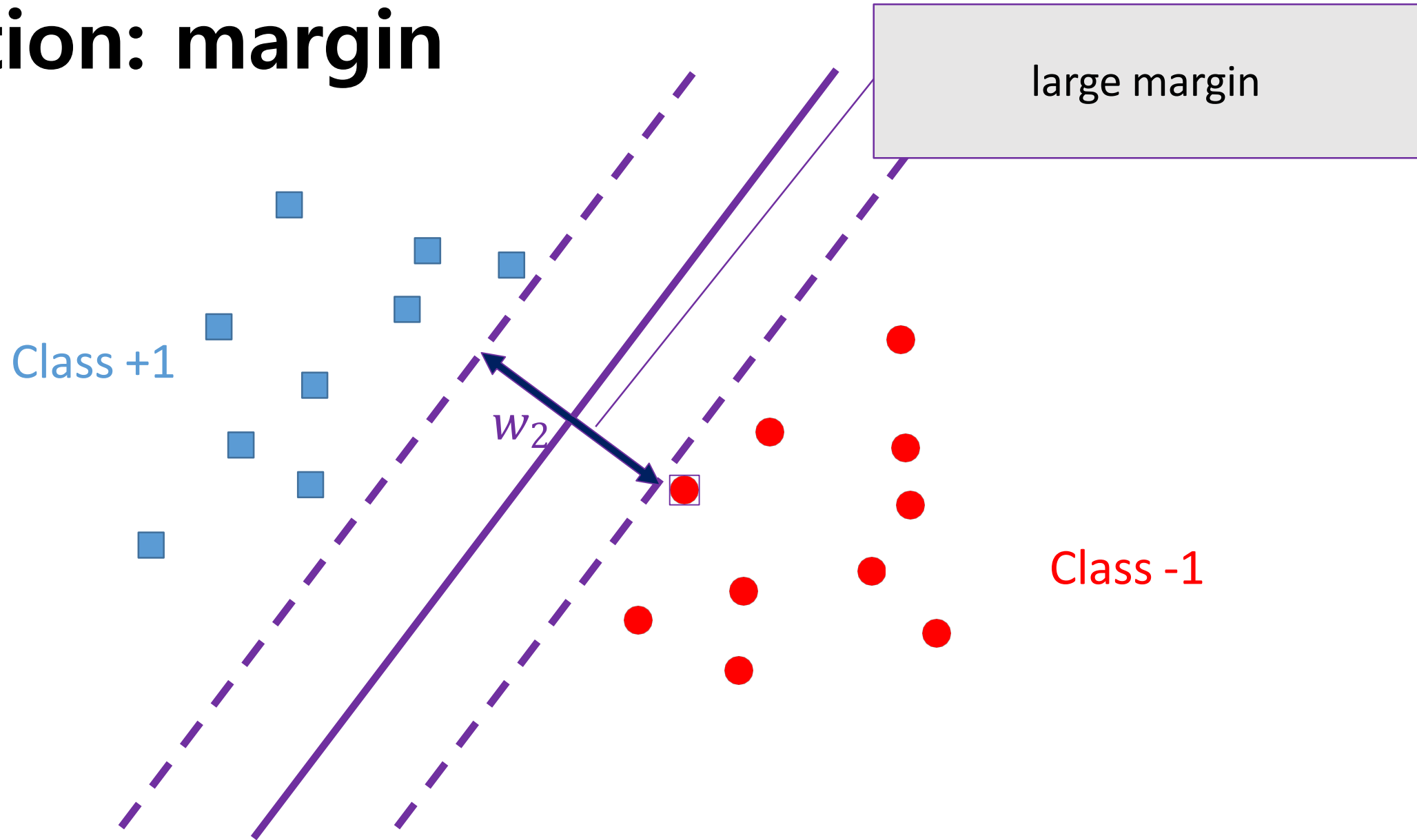
What about w_3 ?



Most confident: w_2



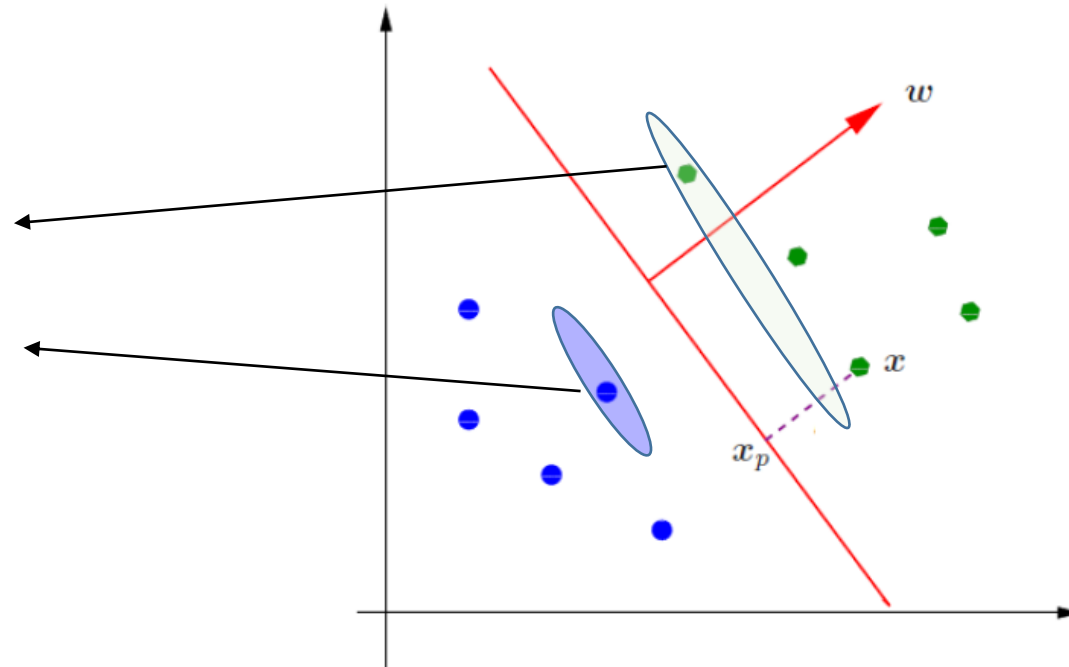
Intuition: margin



Support Vector Machine

- Choose the linear separator (hyperplane) with the largest **margin** on either side
 - Maximum margin hyperplane with **support vectors**
 - Robust to outliers

Support vector : an instance with the minimum margin, which will be the most sensible data points to affect the performance

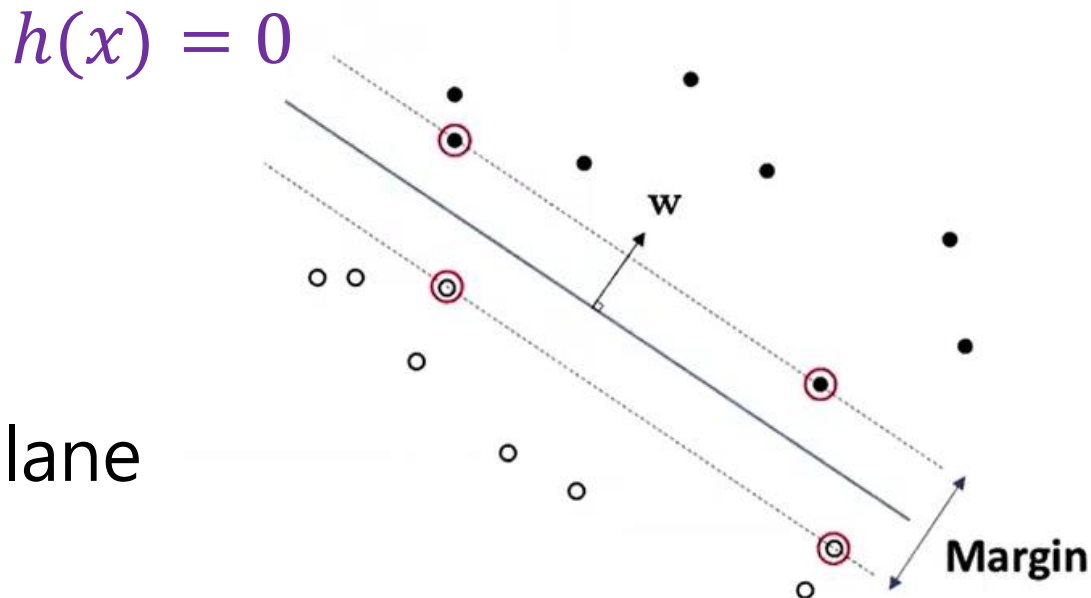


$$h(x) = w^T x + b$$

Margin

- Twice the distance from the hyperplane to the nearest instance on either side
- w is orthogonal to the hyperplane

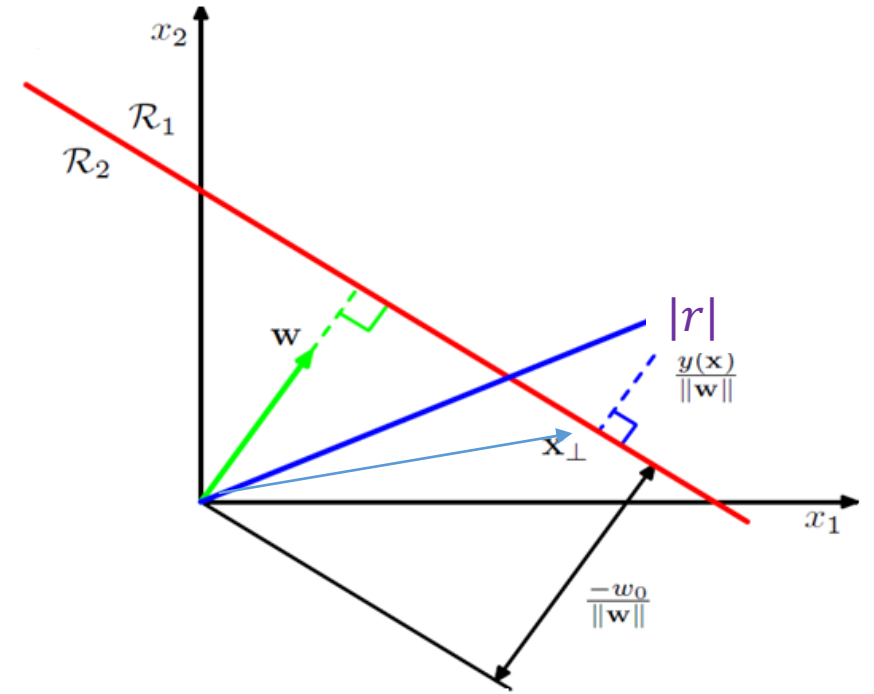
- Lemma : x has distance $\frac{|h_w b(x)|}{\|w\|}$
to the hyperplane $h_{wb}(x) = w^T x + b = 0$



Margin distance

Proof:

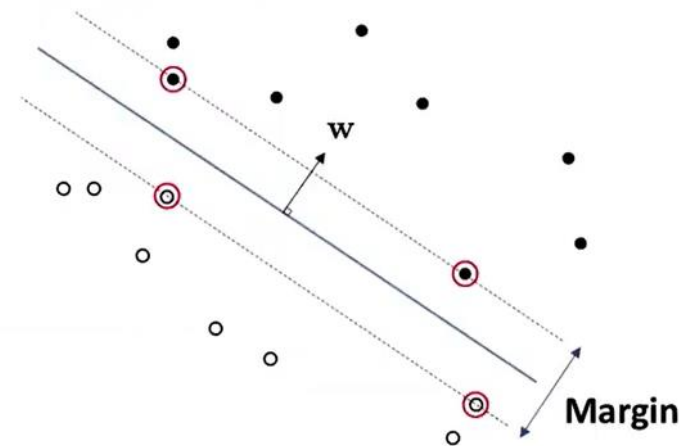
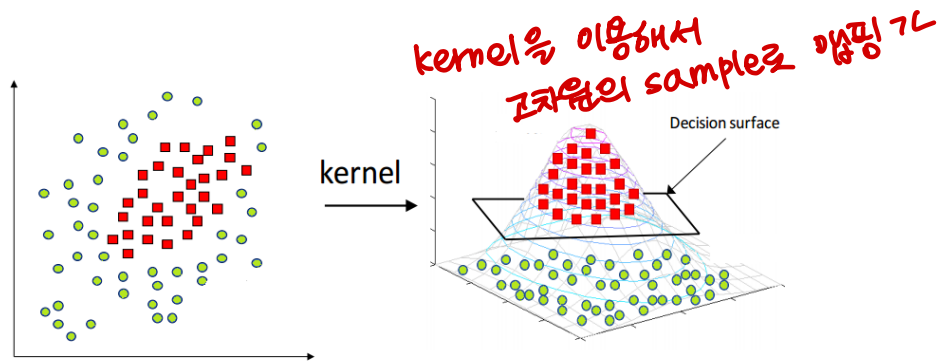
- Let $x = x_{\perp} + r \frac{w}{\|w\|}$, then $|r|$ is the distance
- Multiply both sides by w^T and add b
- Left hand side: $w^T x + b = h_{w,b}(x)$
- Right hand side: $w^T x_{\perp} + r \frac{w^T w}{\|w\|} + b = 0 + r \|w\|$



Optimization

- Optimal weight w and bias b
 - Classifies points correctly as well as achieves the largest possible margin
 - **Hard margin SVM** – assumes linear separability
 - Soft margin SVM – extends to non-separable cases
 - Nonlinear transform & kernel trick

어느정도
에러 허용



Optimization

constraints: linearly separable; **hard-margin linear SVM**

$$h(x) = w^T x + b \geq 1 \text{ for } y = 1$$

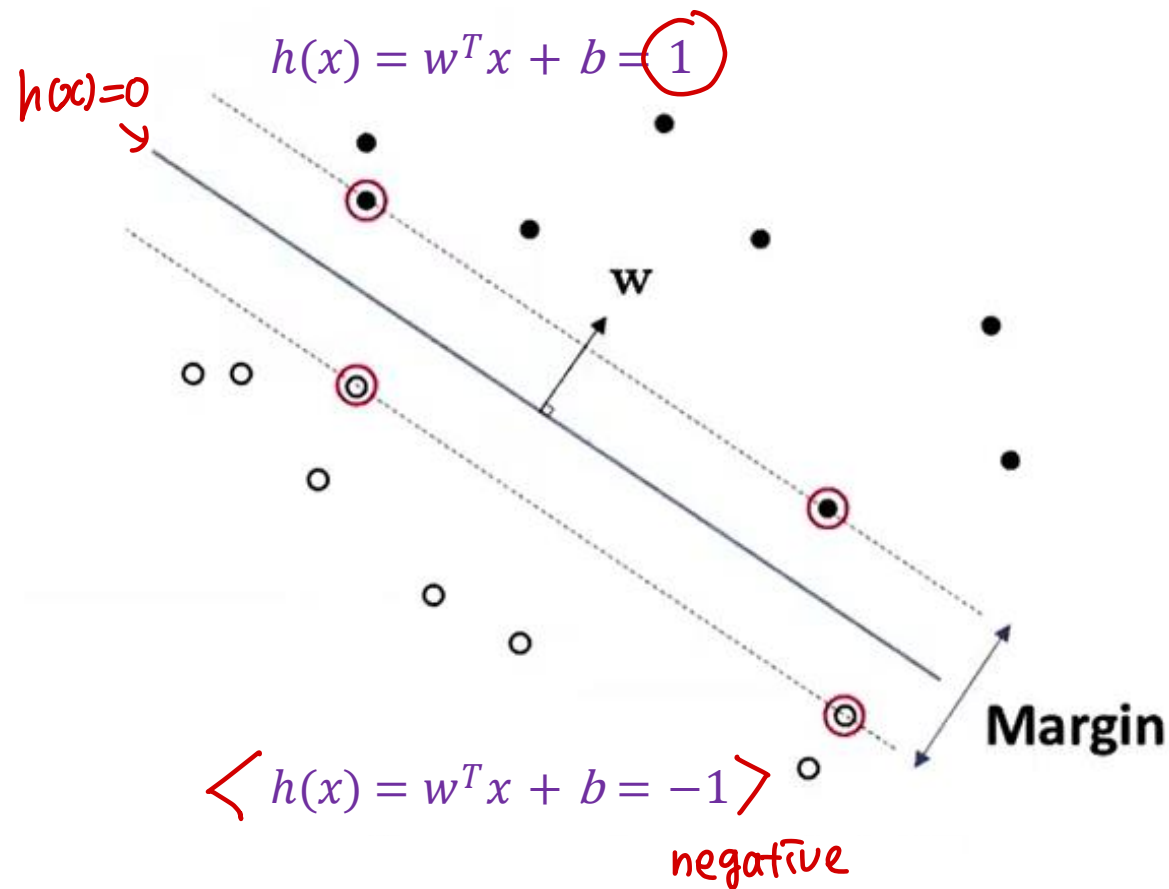
$$h(x) = w^T x + b \leq -1 \text{ for } y = -1$$



$$\textcircled{y}(w^T x + b) \geq 1 \text{ for all samples}$$



summary sample constraints



Optimization

objective function: linearly separable; hard-margin linear SVM

Distance from a support vector to the hyper plane

$$\frac{w^T x + b}{\|w\|} = \frac{\pm 1}{\|w\|} \longrightarrow \frac{2}{\|w\|} \quad \text{Maximize the margin}$$

SVM Primal problem

Equivalently, we want to minimize $\|w\|^2$ ^{계산의 편의상}

Subject to $y_i(w^T x_i + b) \geq 1$ for all samples

Constrained optimization problem, solved by convex quadratic programming

Can be converted into an unconstrained optimization problem and dual form

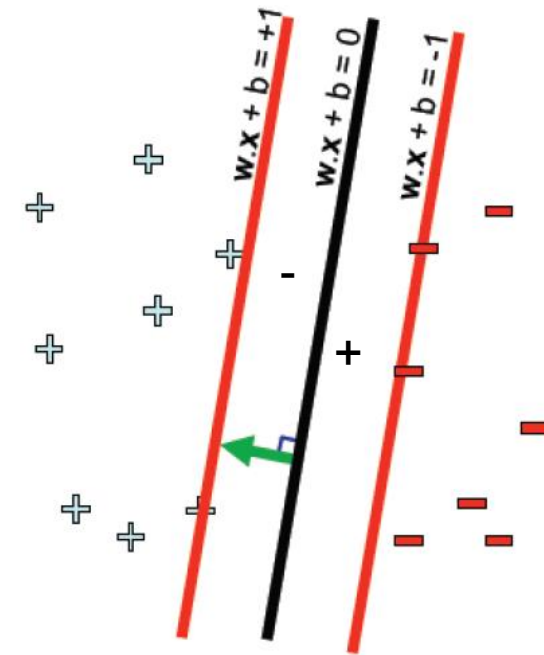
Support Vector Machine

linearly separable; soft-margin linear SVM

* kernel 함수

↳ linearly separable 하지 않은 데이터 샘플들이 있다고 할 때, 그 차이를 높여 linearly separable 하게 만드는 과정

- What if the data is not linearly separable? e.g., there are outliers or noisy measurements, or the data is slightly non-linear.
- Assign a slack variable to each instance $\xi_i \geq 0$, which can be thought of distance from the separating hyperplane if an instance is misclassified and 0 otherwise.
- Minimize $\frac{1}{2} \|w\|^2 + C \sum \xi_i$
subject to $y_i(w^T x_i + b) \geq 1 - \xi_i$



$\alpha=1$ for scale invariance

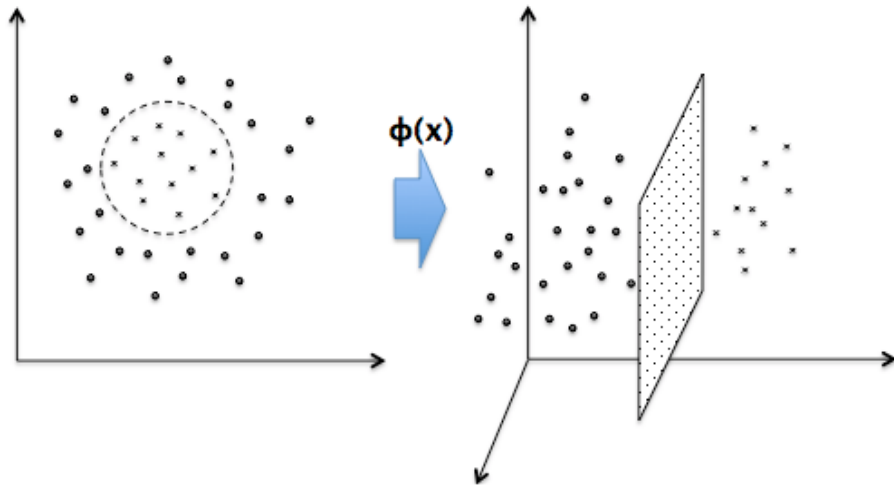
Problem of SVM

- What if the data samples are not linearly separable?
- Linearity in w

$$\sum w^T x \longrightarrow \sum w^T \phi(x)$$

Both linear to x and w

linear to w which allows non linear version of x while maintaining the analytic capability



Nonlinear transforms can be applied to the input to extend the hard-margin SVM to non-linear problem

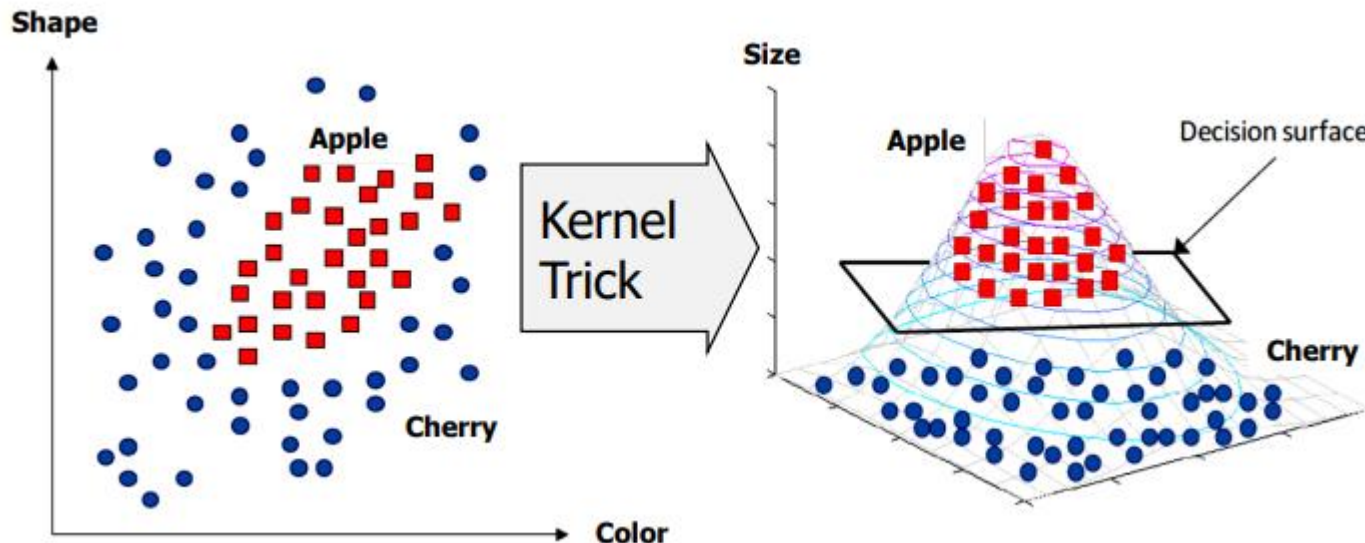
Transform input into new space and use linear model
(kernel transformation)

Support Vector Machine

not linearly separable; Kernel Trick

- Data is not linearly separable in the input space

- Data is linearly separable in the feature space obtained by a kernel



Some commonly used kernels

① Polynomial:

$$K(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y} + 1)^p$$

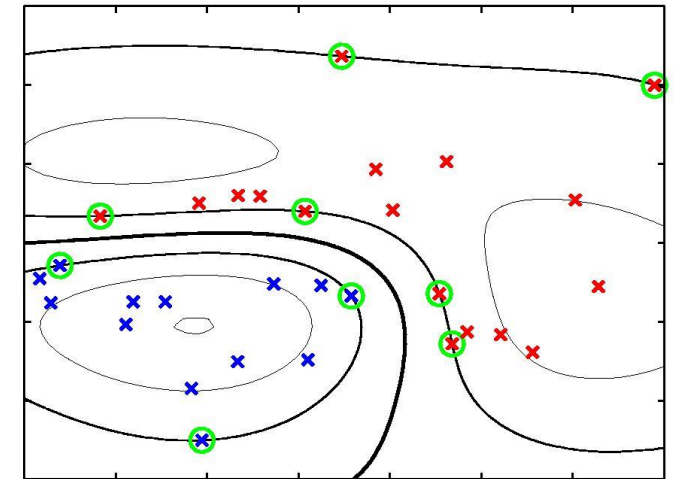
② Gaussian radial basis function (RBF)

$$K(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x} - \mathbf{y}\|^2 / 2\sigma^2}$$

③ Hyperbolic tangent (multilayer perceptron kernel)

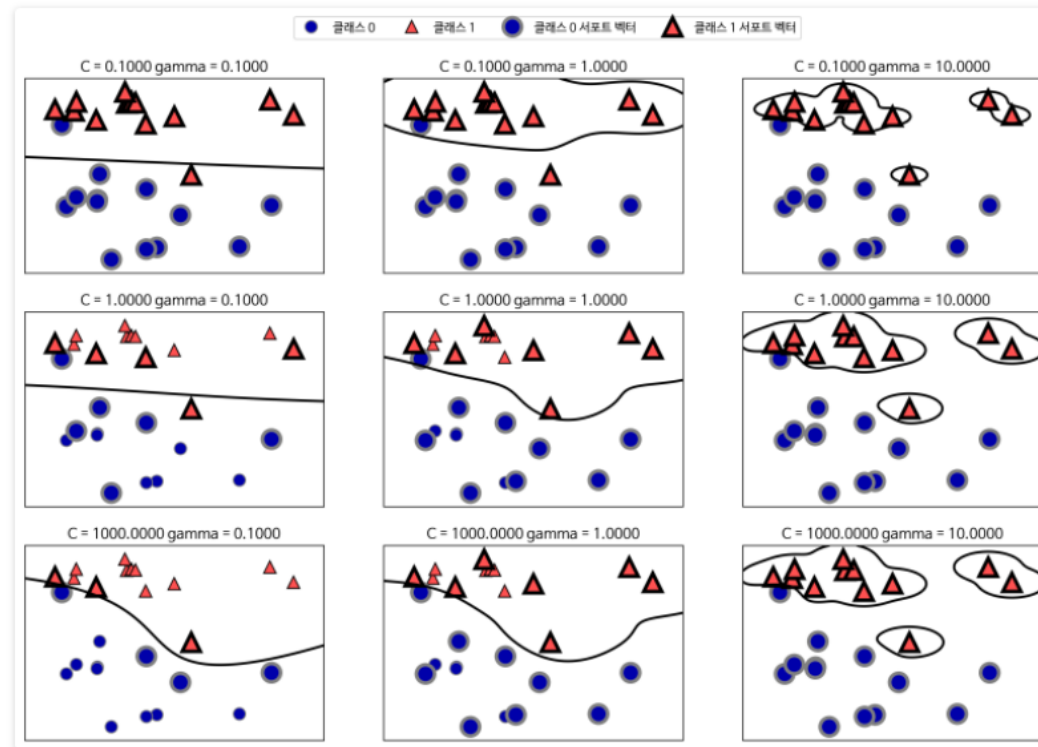
$$K(\mathbf{x}, \mathbf{y}) = \tanh(k \mathbf{x} \cdot \mathbf{y} - \delta)$$

Parameters that the user must choose



Radial-basis function (RBF) kernel

- Radial-basis function kernel



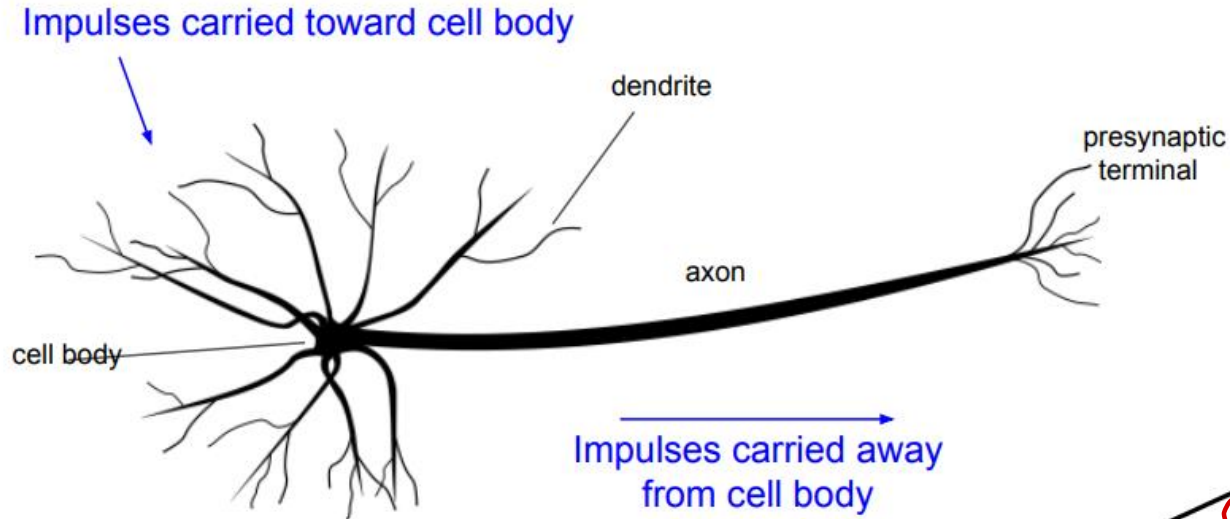
$$K(\mathbf{x}, \mathbf{x}') = \exp \left(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2 \right)$$

Large γ : sharp Gaussian \rightarrow overfit
Small γ : not enough linearity

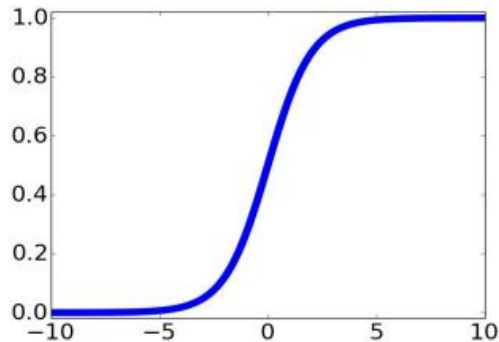
Artificial neural network (ANN)

non-linear classification model

↑
nonlinear 계층, DNN의 기본

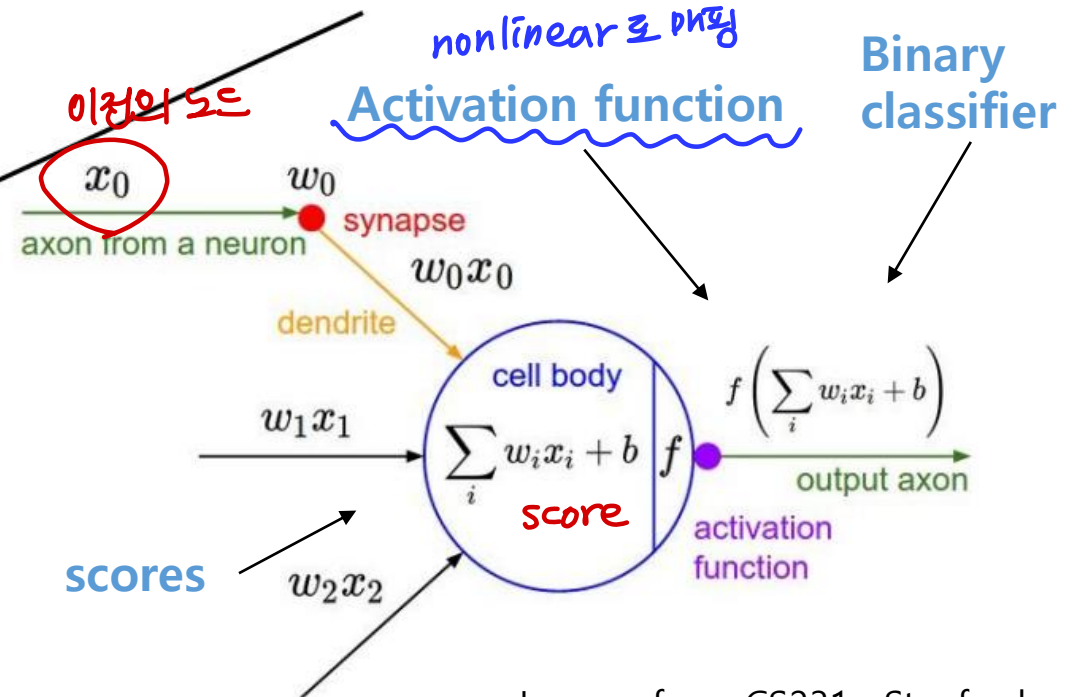


This image by Felipe Peruchio is licensed under CC-BY 3.0



sigmoid activation function

$$\frac{1}{1 + e^{-x}}$$

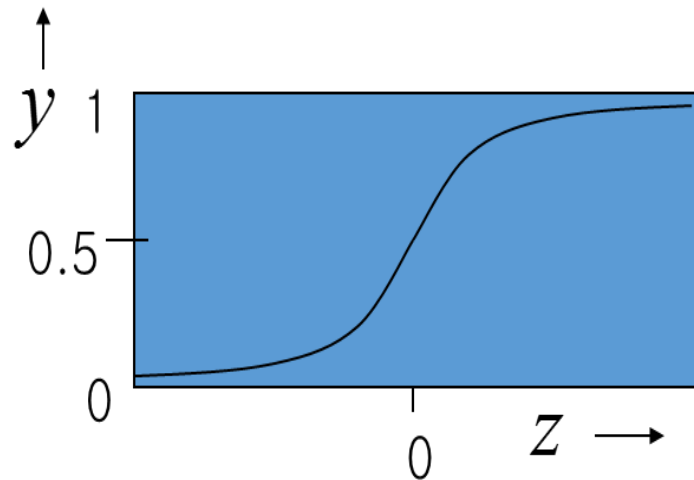


Artificial neural network (ANN) activation functions

- **Sigmoid neurons** give a real-valued output that is a smooth and bounded function of their total input.
- **Non-linearity due to the activation functions**

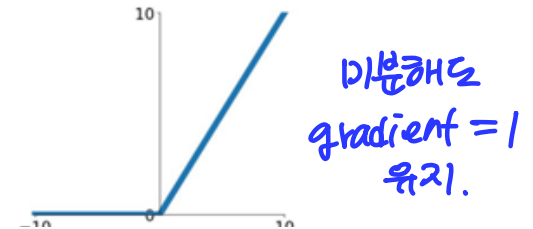
Sigmoid 함수이지 x

$$z = b + \sum_i x_i w_i$$
$$y = \frac{1}{1 + e^{-z}}$$

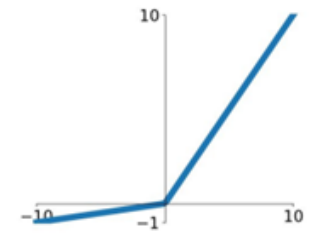


Activation functions

① **ReLU**
 $\max(0, x)$

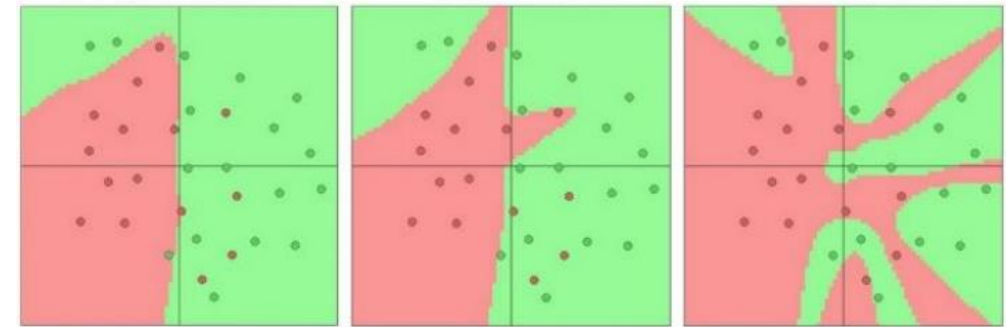
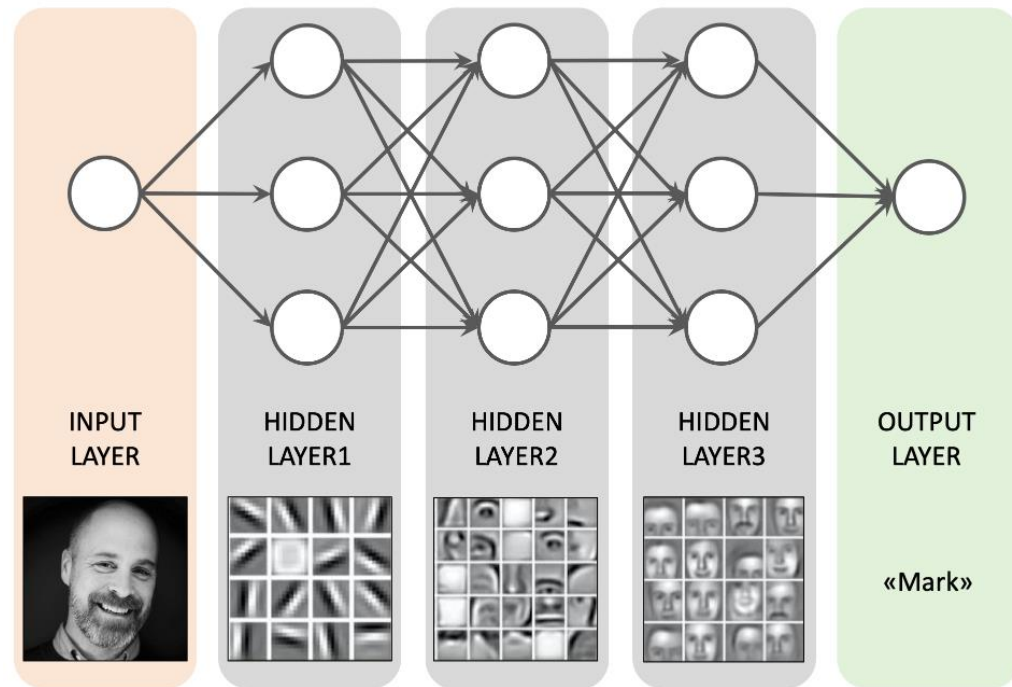


② **Leaky ReLU**
 $\max(0.1x, x)$



Artificial neural network (ANN)

deep neural network

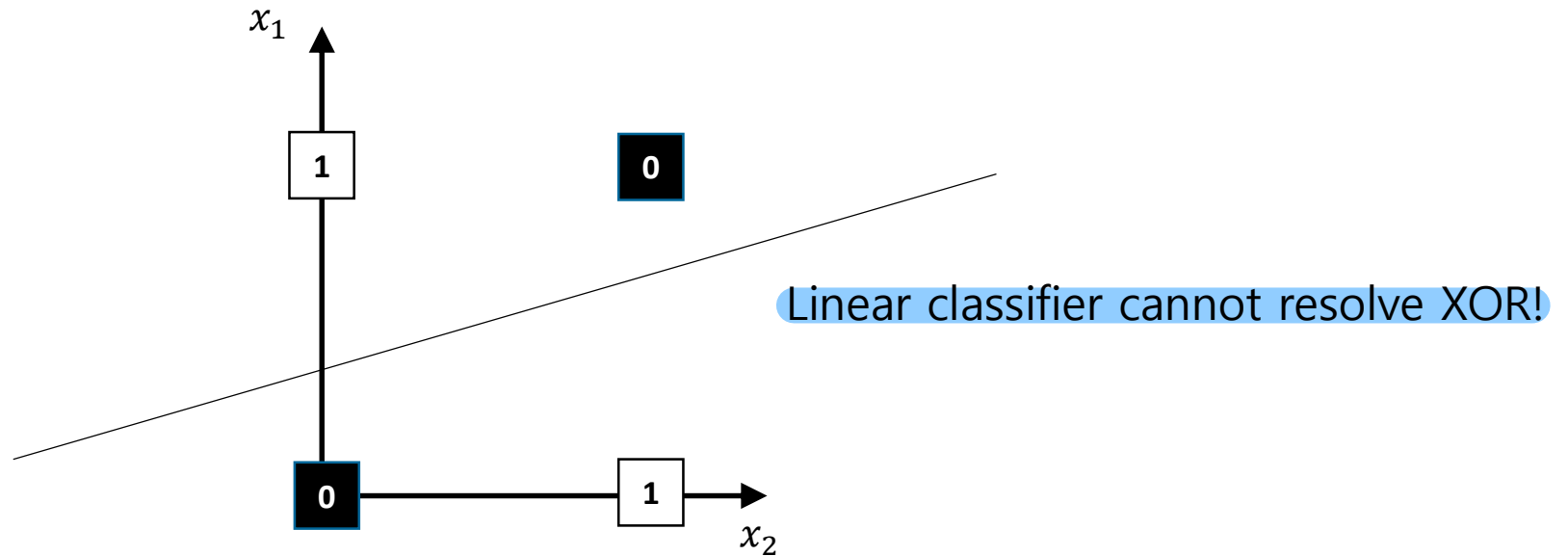


- Can represent more complex (non-linear) boundaries with increasing neurons

Artificial neural network (ANN)

multilayer perceptron

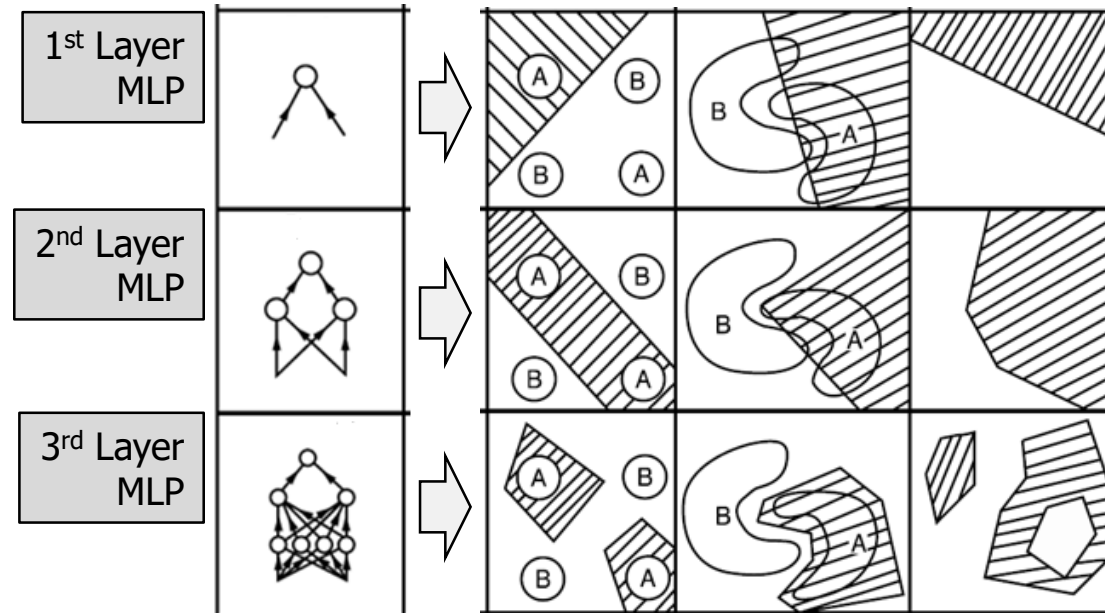
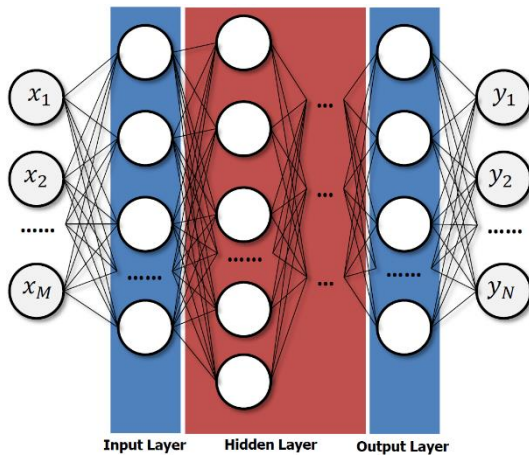
x_1	x_2	XOR
0	0	0
0	1	1
1	0	1
1	1	0



Mathematically proven by
Prof. Marvin Minsky at MIT (1969)

Artificial neural network (ANN) multilayer perceptron

- Multilayer Perceptron (MLP)
 - Proposed by Prof. Marvin Minsky at MIT (1969)
 - Can solve XOR Problem



점점복잡



$$W = \begin{pmatrix} 5 \\ 5 \end{pmatrix}, b = -8$$

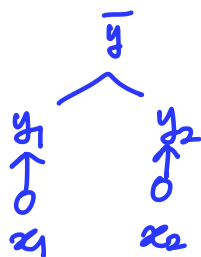
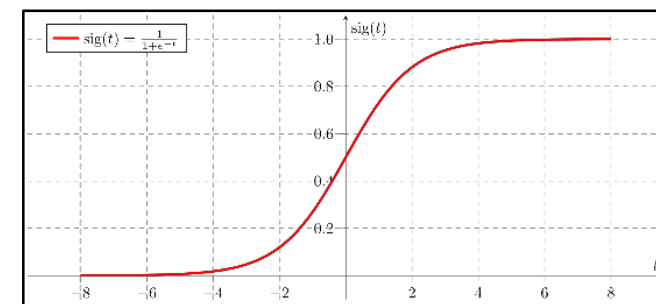


$$W = \begin{pmatrix} -7 \\ -7 \end{pmatrix}, b = 3$$



$$W = \begin{pmatrix} -11 \\ -11 \end{pmatrix}, b = 6$$

- $(x_1 \ x_2) = (0 \ 0)$
 - $(0 \ 0) \begin{pmatrix} 5 \\ 5 \end{pmatrix} + (-8) = -8$, i.e., $y_1 = \text{Sigmoid}(-8) \cong 0$
 - $(0 \ 0) \begin{pmatrix} -7 \\ -7 \end{pmatrix} + (3) = 3$, i.e., $y_2 = \text{Sigmoid}(3) \cong 1$
 - $(y_1 \ y_2) \begin{pmatrix} -11 \\ -11 \end{pmatrix} + (6) = -11 + 6 = -5$, i.e., $\bar{y} = \text{Sigmoid}(-5) \cong 0$



x_1	x_2	y_1	y_2	\bar{y}	XOR
0	0	0	1	0	0
0	1				1
1	0				1
1	1				0



$$W = \begin{pmatrix} 5 \\ 5 \end{pmatrix}, b = -8$$

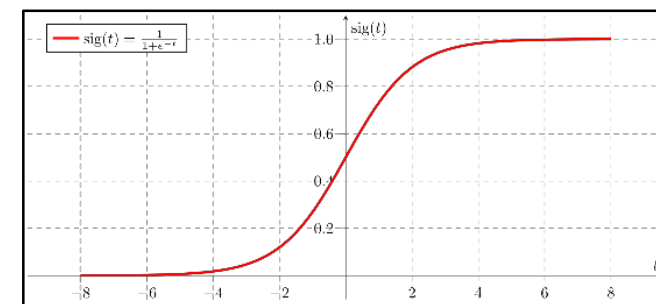


$$W = \begin{pmatrix} -7 \\ -7 \end{pmatrix}, b = 3$$



$$W = \begin{pmatrix} -11 \\ -11 \end{pmatrix}, b = 6$$

- $(x_1 \ x_2) = (0 \ 1)$
 - $(0 \ 1) \begin{pmatrix} 5 \\ 5 \end{pmatrix} + (-8) = -3$, i.e., $y_1 = \text{Sigmoid}(-3) \cong 0$
 - $(0 \ 1) \begin{pmatrix} -7 \\ -7 \end{pmatrix} + (3) = -4$, i.e., $y_2 = \text{Sigmoid}(-4) \cong 0$
 - $(y_1 \ y_2) \begin{pmatrix} -11 \\ -11 \end{pmatrix} + (6) = 6$, i.e., $\bar{y} = \text{Sigmoid}(6) \cong 1$



x_1	x_2	y_1	y_2	\bar{y}	XOR
0	0	0	1	0	0
0	1	0	0	1	1
1	0				1
1	1				0



$$W = \begin{pmatrix} 5 \\ 5 \end{pmatrix}, b = -8$$

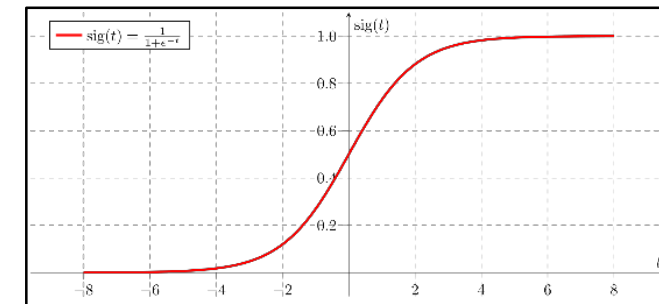


$$W = \begin{pmatrix} -7 \\ -7 \end{pmatrix}, b = 3$$



$$W = \begin{pmatrix} -11 \\ -11 \end{pmatrix}, b = 6$$

- $(x_1 \ x_2) = (1 \ 0)$
 - $(1 \ 0) \begin{pmatrix} 5 \\ 5 \end{pmatrix} + (-8) = -3$, i.e., $y_1 = \text{Sigmoid}(-3) \cong 0$
 - $(1 \ 0) \begin{pmatrix} -7 \\ -7 \end{pmatrix} + (3) = -4$, i.e., $y_2 = \text{Sigmoid}(-4) \cong 0$
 - $(y_1 \ y_2) \begin{pmatrix} -11 \\ -11 \end{pmatrix} + (6) = 6$, i.e., $\bar{y} = \text{Sigmoid}(6) \cong 1$



x_1	x_2	y_1	y_2	\bar{y}	XOR
0	0	0	1	0	0
0	1	0	0	1	1
1	0	0	0	1	1
1	1				0



$$W = \begin{pmatrix} 5 \\ 5 \end{pmatrix}, b = -8$$

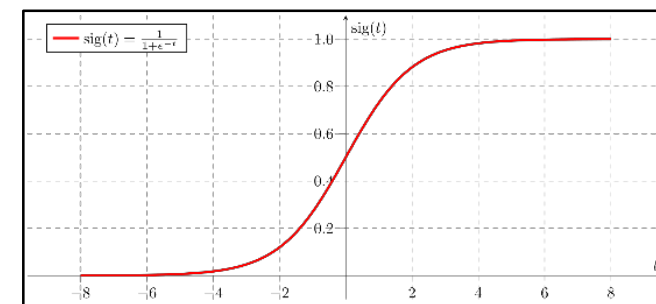


$$W = \begin{pmatrix} -7 \\ -7 \end{pmatrix}, b = 3$$



$$W = \begin{pmatrix} -11 \\ 11 \end{pmatrix}, b = 6$$

- $(x_1 \ x_2) = (1 \ 1)$
 - $(1 \ 1) \begin{pmatrix} 5 \\ 5 \end{pmatrix} + (-8) = 2$, i.e., $y_1 = \text{Sigmoid}(2) \cong 1$
 - $(1 \ 1) \begin{pmatrix} -7 \\ -7 \end{pmatrix} + (3) = -11$, i.e., $y_2 = \text{Sigmoid}(-11) \cong 0$
 - $(y_1 \ y_2) \begin{pmatrix} -11 \\ 11 \end{pmatrix} + (6) = -5$, i.e., $\bar{y} = \text{Sigmoid}(-5) \cong 0$

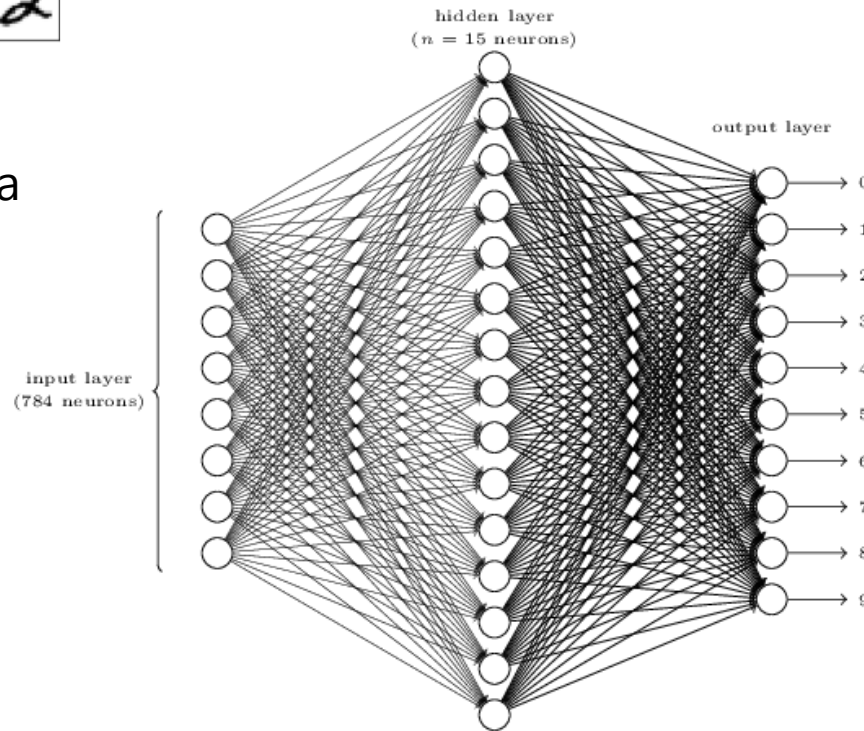


x_1	x_2	y_1	y_2	\bar{y}	XOR
0	0	0	1	0	0
0	1	0	0	1	1
1	0	0	0	1	1
1	1	1	0	0	0

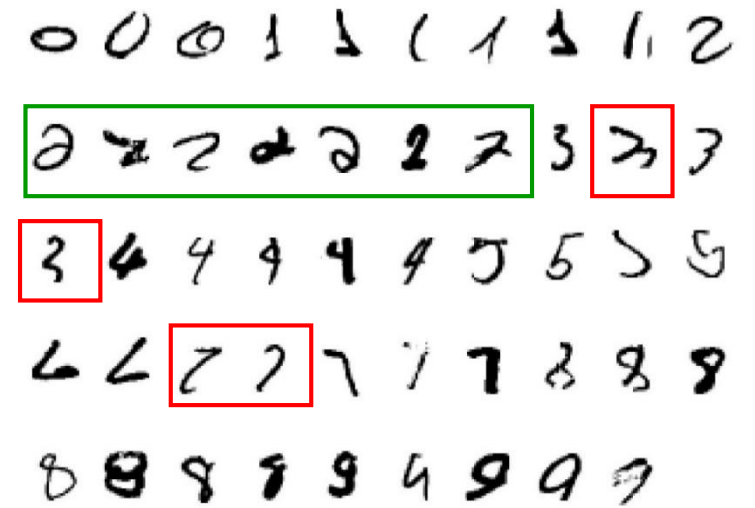
Example: MNIST data recognition

5 0 4 1 9 2

28x28 input image
-> 1x28x28 input data
vector



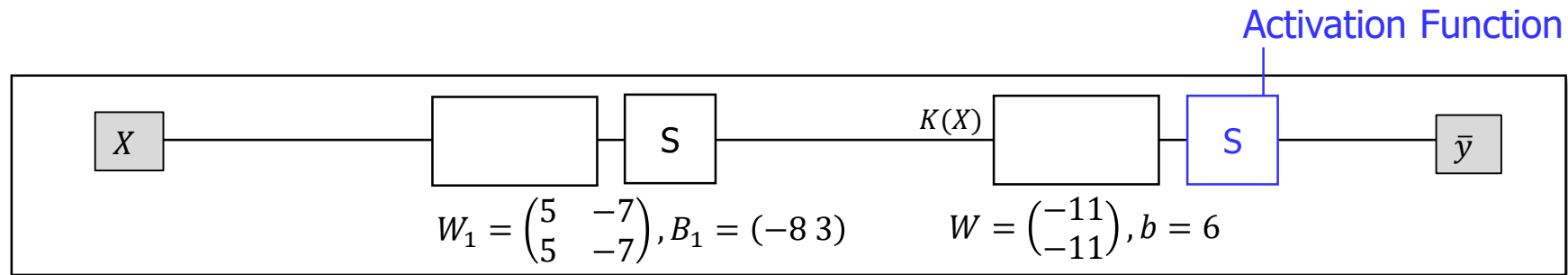
Class of 1~10 Digits



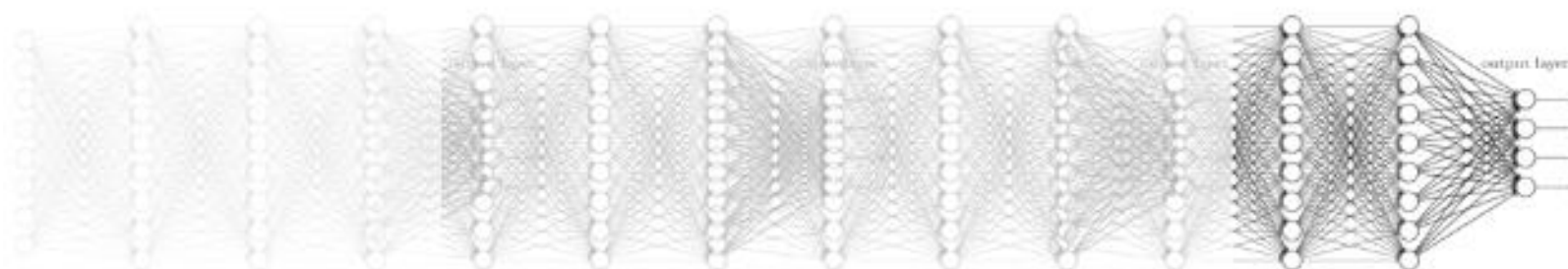
MNIST data used for Optical
Character Recognition (OCR)

Artificial neural network (ANN)

ANN for non-linear problem



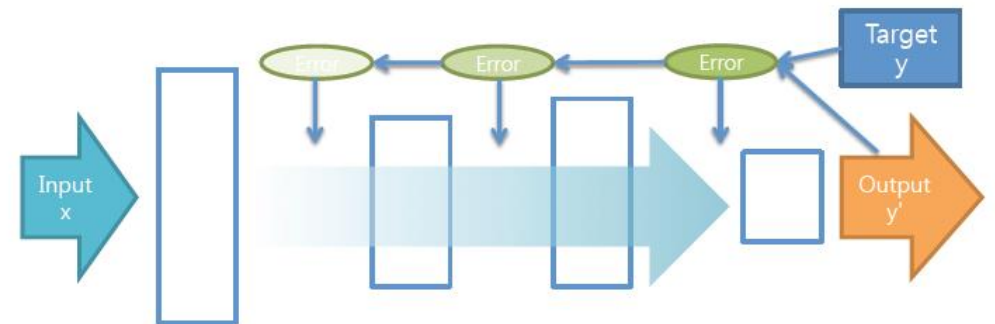
- Observation
 - There exists cases when the accuracy is low even if the # layers is high. Why?
 - Answer
 - The result of one ANN is the result of sigmoid function (**between 0 and 1**).
 - The numerous multiplication of this result converges to near zero.
 → **Gradient Vanishing Problem!!**



$$\frac{\partial \text{error}}{\partial w_1} = \frac{\partial \text{error}}{\partial \text{output}} * \frac{\partial \text{output}}{\partial \text{hidden2}} * \frac{\partial \text{hidden2}}{\partial \text{hidden1}} * \frac{\partial \text{hidden1}}{\partial w_1}$$

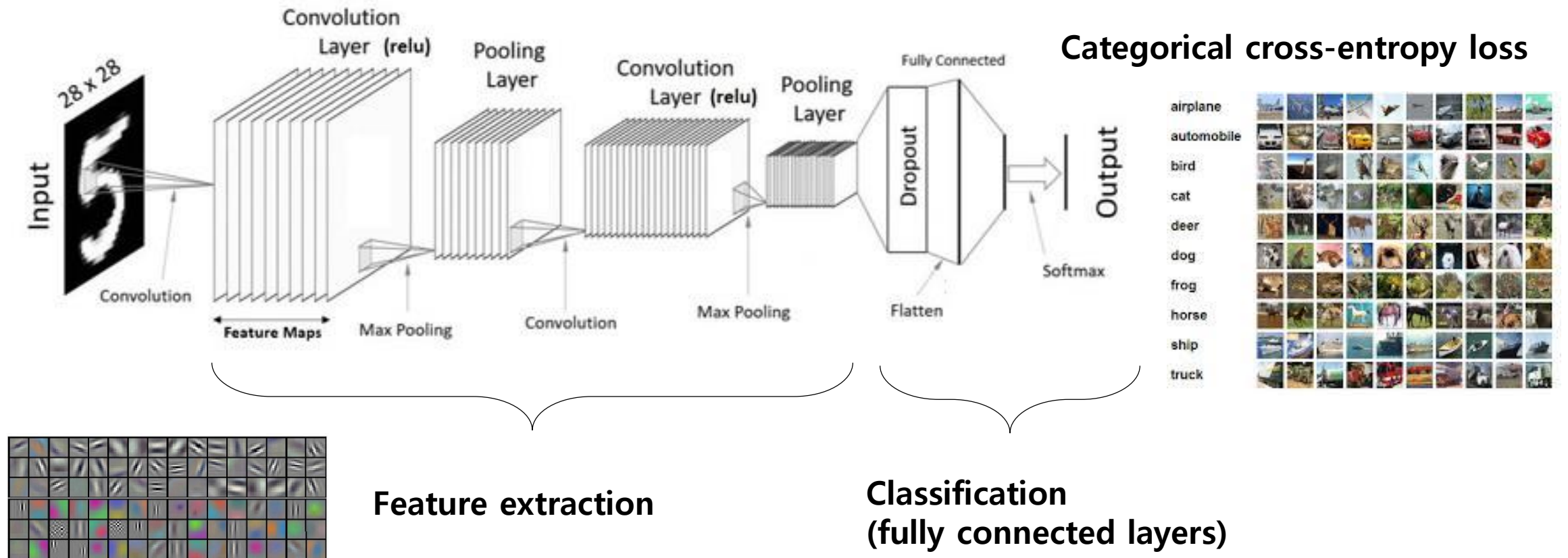
Breakthrough in Back Propagation

- Backpropagation (BP) barely changes lower-layer parameters (vanishing gradient)
- Breakthrough
 - Pre-training+ fine tuning
 - Convolutional neural networks (CNN) for reducing redundant parameters.
 - Rectified linear unit (constant gradient propagation)
 - Dropout



Convolutional neural network

State-of-the-art classification model for high-dimensional data (image, video, etc.)



Reference

- Book: Pattern Recognition and Machine Learning (by Christopher M. Bishop)
- Book: Machine Learning: a Probabilistic Perspective (by Kevin P. Murphy)
- <https://www.andrewng.org/courses/>