# Introduction to Deep Neural Networks
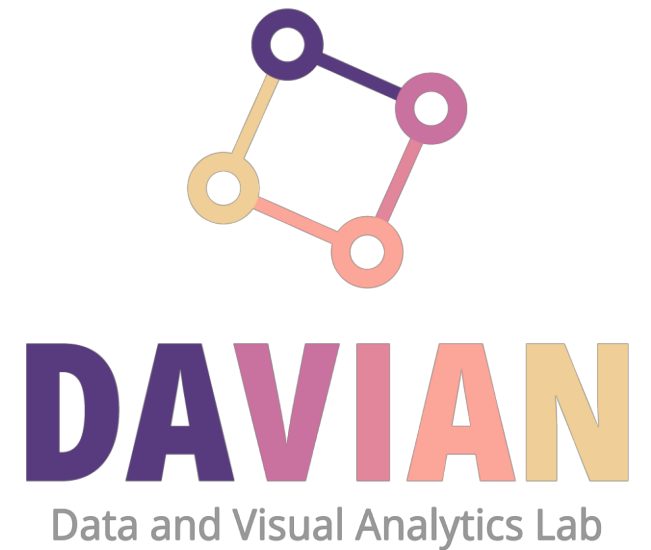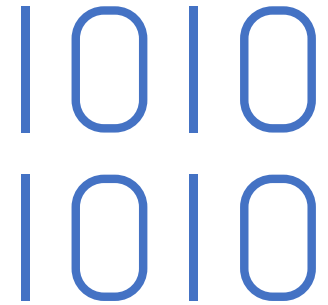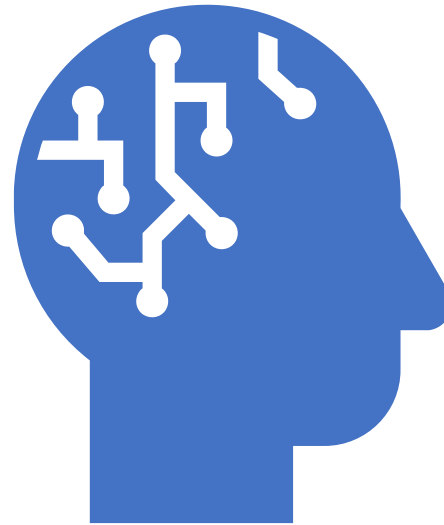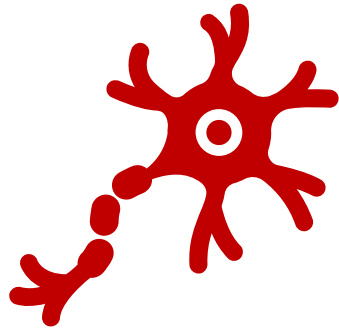
## 주재걸 교수
### KAIST 김재철AI대학원

DAVIAN
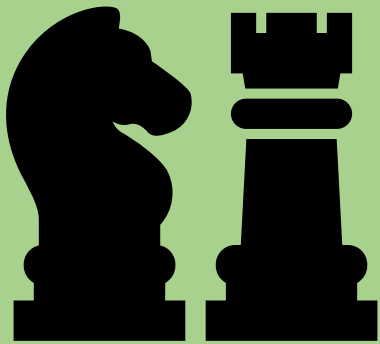Data and Visual Analytics Lab

# Deep Learning

## Deep Learning

- Deep learning refers to artificial neural networks that are composed of many layers.

# Artificial Intelligence vs. Machine Learning vs. Deep Learning

## Artificial Intelligence
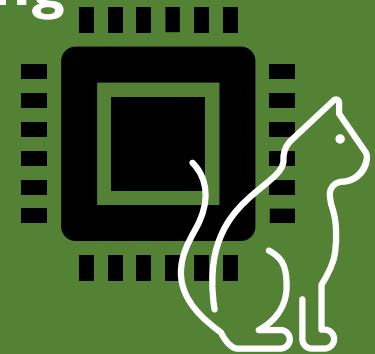Early artificial intelligence stirs excitement.

## Machine Learning
Machine learning begins to flourish

## Deep Learning
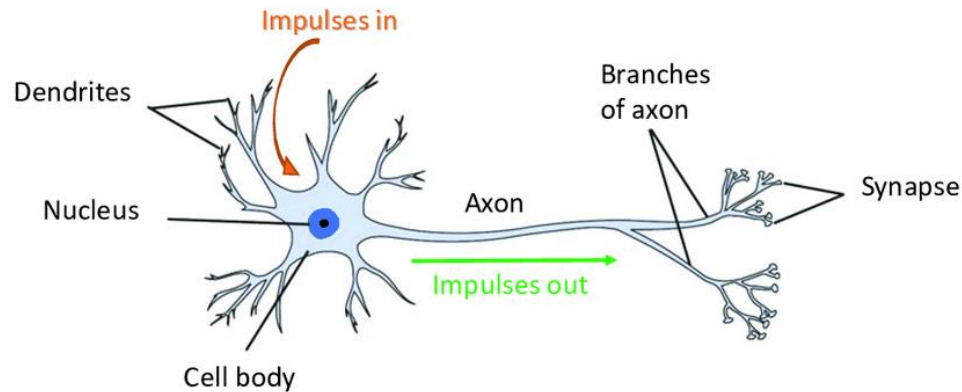Deep learning breakthroughs drive AI boom.

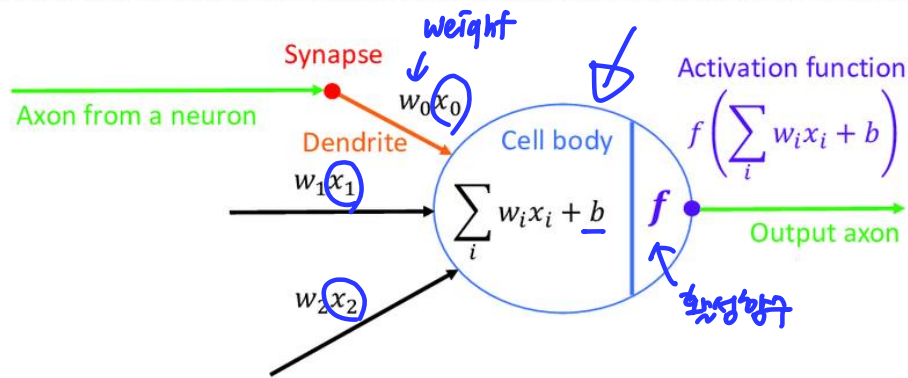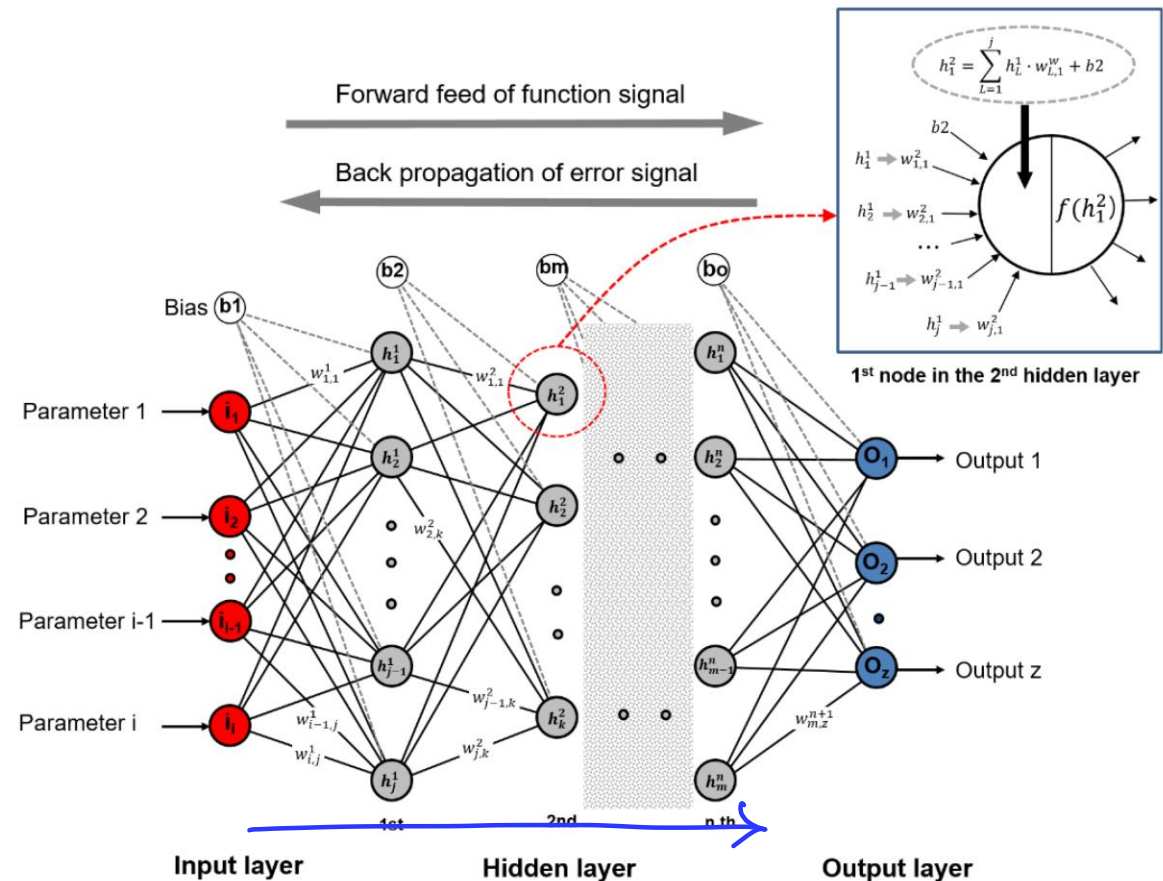1950's        1980's        2010's

# Artificial Neural Networks

- A technology that imitates neurons existing in the human brain



(a)

(b)

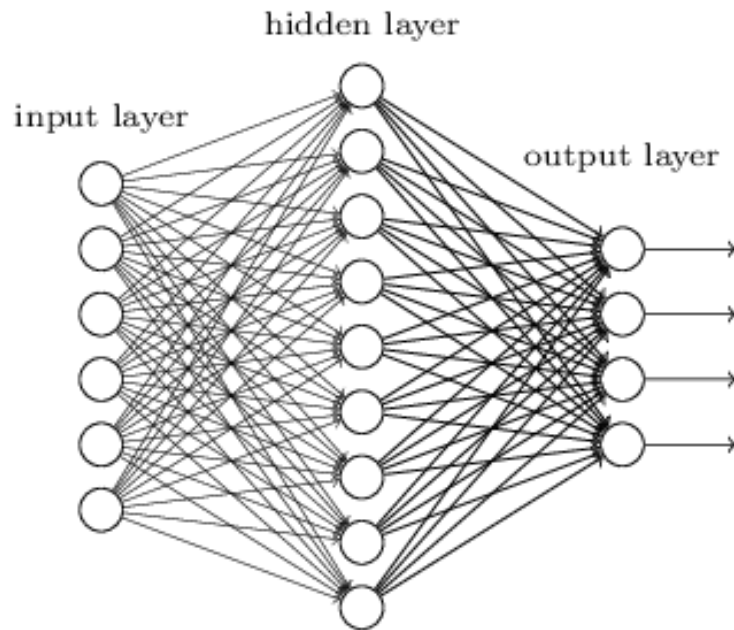https://arxiv.org/abs/1706.05933

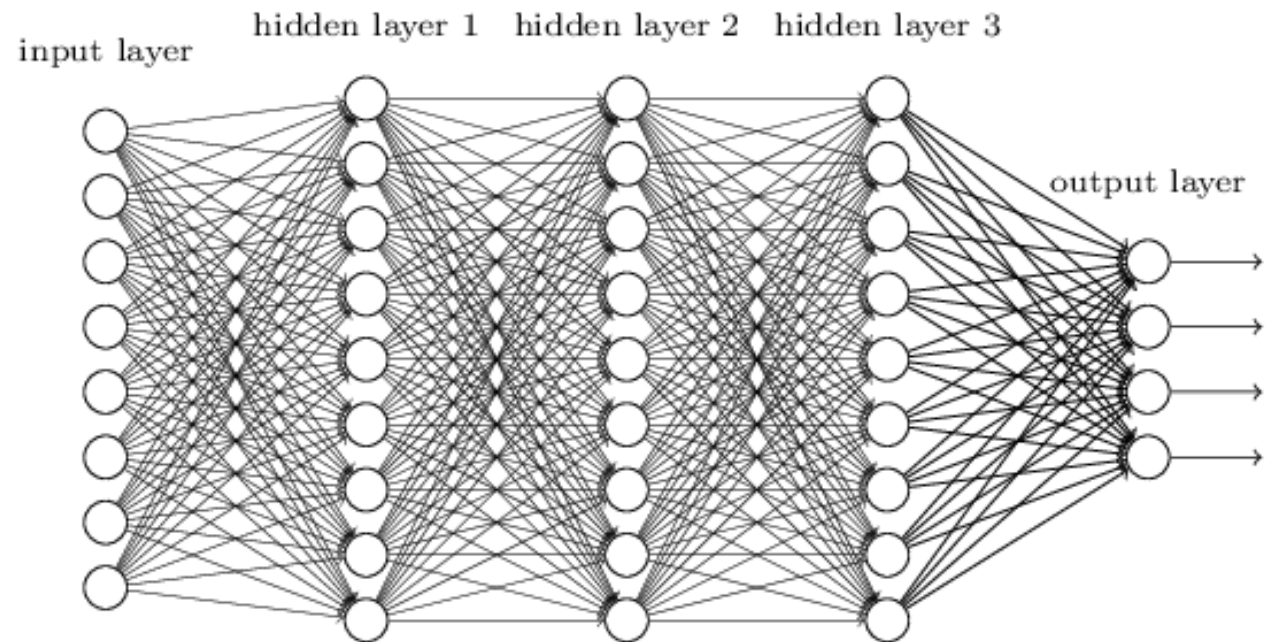https://www.mdpi.com/2076-3417/10/4/1302

# Deep Neural Network (DNN)

- DNN improves accuracy of AI technology by stacking neural network layers



"Non-deep" feedforward neural network

Deep neural network

http://neuralnetworksanddeeplearning.com/chap5.html

# Reason Why Deep Learning has been Successful

**Big Data**

**GPU Acceleration**

**Algorithm Improvements**

# Computer Vision



Object Detection

Image Synthesis

Liu, Ze et al. "Swin Transformer: Hierarchical Vision Transformer Using Shifted Windows."
Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). 2021.

Karras, Tero et al. "Analyzing and Improving the Image Quality of StyleGAN."
Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). 2020.

# Natural Language Processing



Machine Translation



Mail Classification

# Reinforcement Learning



Go



Atari Gane

https://www.deepmind.com/research/highlighted-research/alphago
https://www.deepmind.com/publications/playing-atari-with-deep-reinforcement-learning

# Perceptron and Neural Networks

# What is Perceptron?

$$y = f(w_0 + w_1 x_1 + w_2 x_2)$$



**Perceptron**

- One kind of neural networks

- Frank Rosenblatt devised in 1957

- Linear classifier



Impulses in

Dendrites

Branches of axon

Nucleus

Axon

Synapse

Impulses out

Cell body

(a)

- Similar with structure of a neuron

# What is Perceptron?

$$y = f(w_0 + w_1 x_1 + w_2 x_2)$$



Input: $x_1, x_2$    Output: $y$

Weights: $w_0, w_1, w_2$

# What is Perceptron?

$$y = f(w_0 + w_1 x_1 + w_2 x_2)$$



Impulses in

Dendrites

Branches of axon

Nucleus

Axon

Synapse

Impulses out

Cell body

(a)

Input: $x_1, x_2$    Output: $y$

Weights: $w_0, w_1, w_2$

# What is Perceptron?

$$y = f(w_0 + w_1 x_1 + w_2 x_2)$$



Input: $x_1, x_2$   Output: $y$

Weights: $w_0, w_1, w_2$

# What is Perceptron?
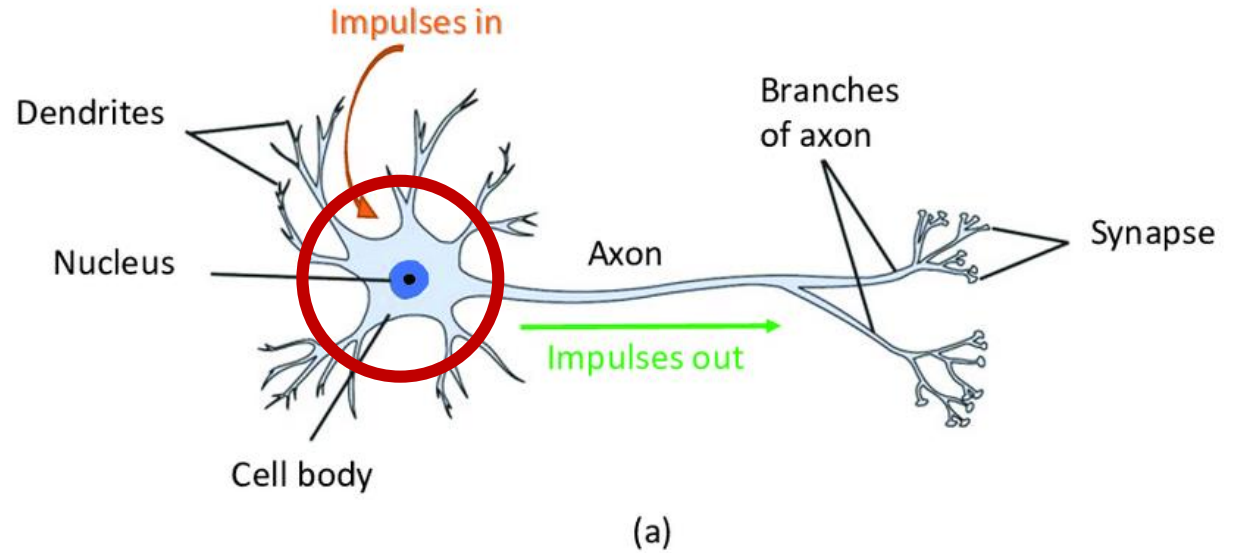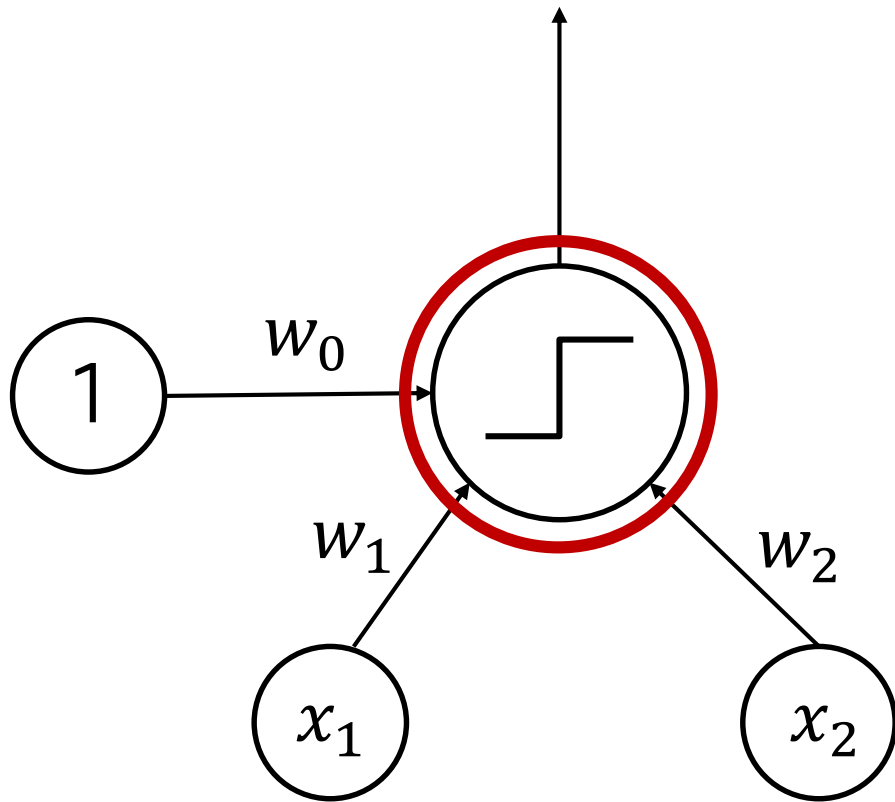
$$y = f(w_0 + w_1 x_1 + w_2 x_2)$$



Input: $x_1, x_2$    Output: $y$

Weights: $w_0, w_1, w_2$

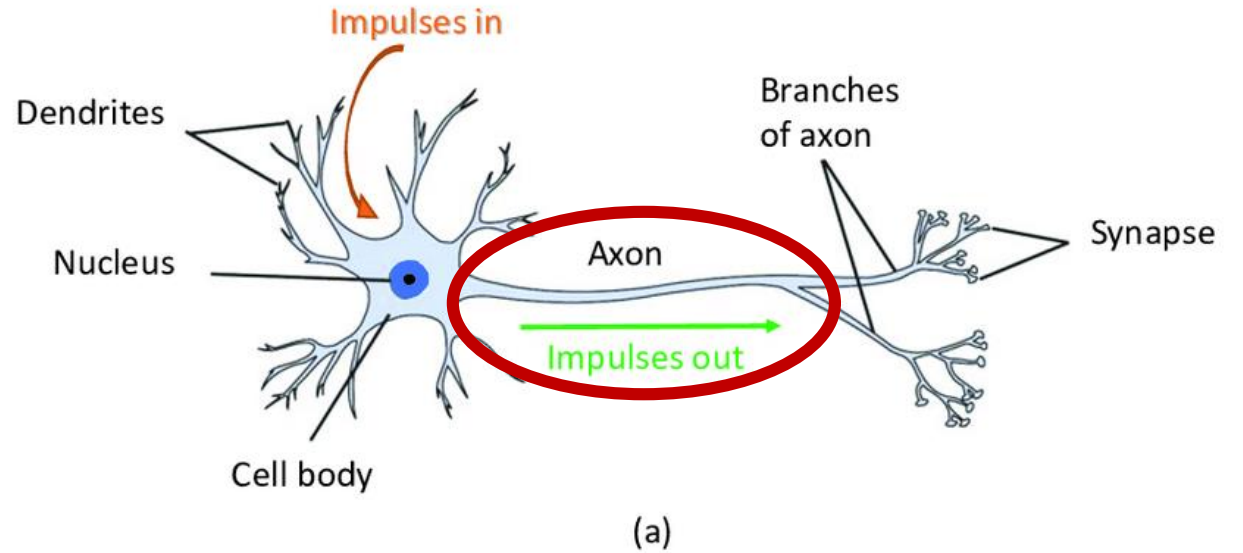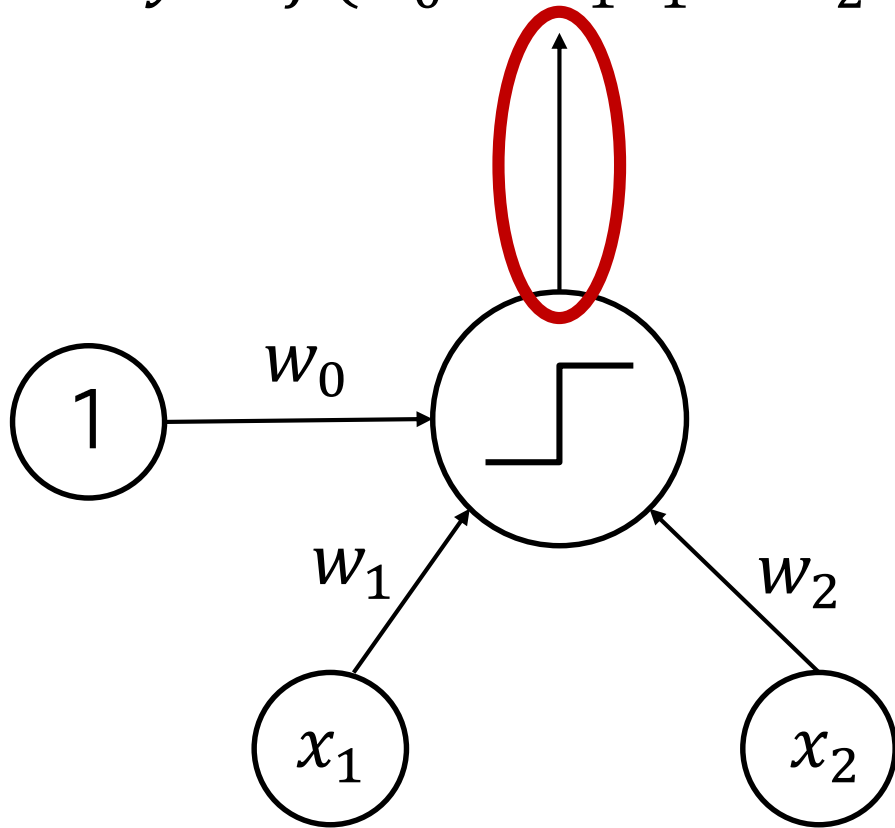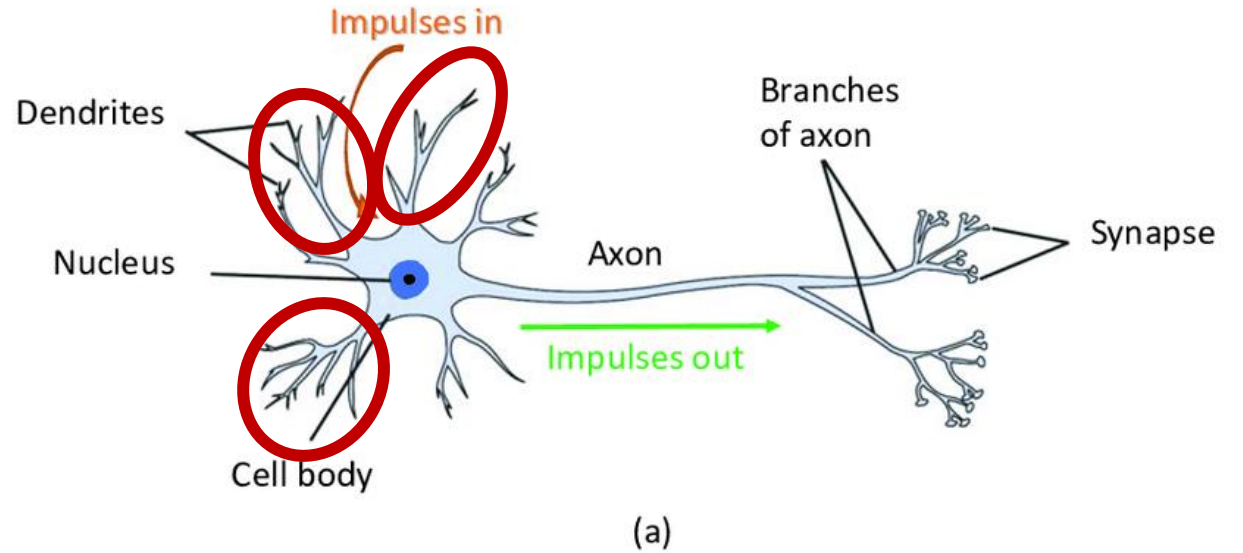# Single Layer Perceptron

$$y = f(w_0 + w_1 x_1 + w_2 x_2)$$



**Hard thresholding function**

$$y = \begin{cases} 1 & w_0 + w_1 x_1 + w_2 x_2 \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Input: $x_1, x_2$    Output: $y$

Weights: $w_0, w_1, w_2$

# Single Layer Perceptron

$w_0 x 1$

$y = f(3 + 4 \times 1 - 5 \times 0) = 1$



$w_0 = 3$

$w_1 = 4$

$w_2 = -5$

$x_1$

$x_2$

**Hard thresholding function**

$$y = \begin{cases} 1 & w_0 + w_1 x_1 + w_2 x_2 \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Input: $x_1, x_2$   Output: $y$

Weights: $w_0, w_1, w_2$

# Decision Boundary in Perceptron



AND Gate

$$-0.8 + 0.5x_1 + 0.5x_2 = 0$$

Input feature space

output $x_1 = 1$

OR Gate

$$-0.3 + 0.5x_1 + 0.5x_2 = 0$$

Multi-Layer Perceptron for XOR Gate

Is it possible to solve a XOR problem using a single layer perceptron?
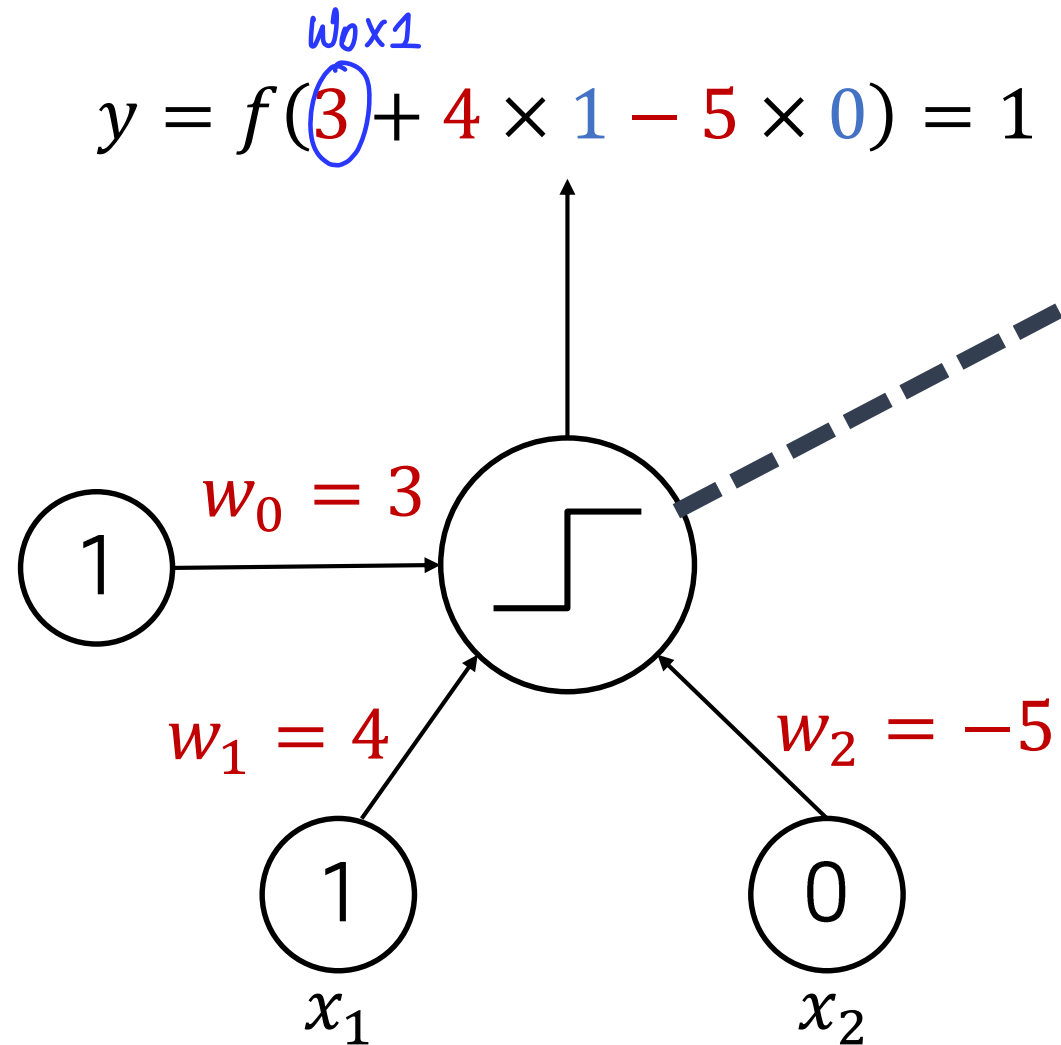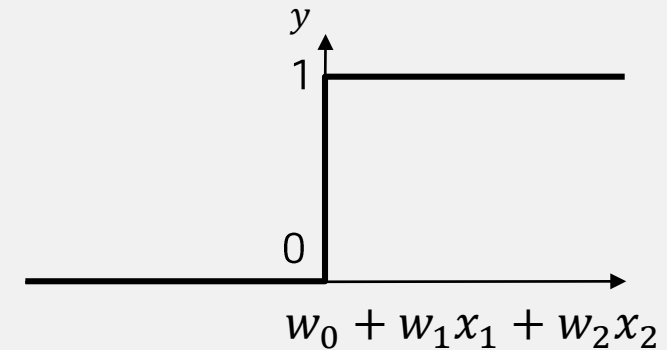→ No. Single layer perceptron can only solve linear problem. XOR problem is non-linear

$a + bx_1 + cx_2 = 0$

● : 1
▲ : 0

**XOR Gate**

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0     | 0     | 0   |
| 0     | 1     | 1   |
| 1     | 0     | 1   |
| 1     | 1     | 0   |

# Multi-Layer Perceptron

But if we use two-layer perceptron, we can solve XOR problem
→ This model is called multi-layer perceptron

# Multi-Layer Perceptron

# Multi-Layer Perceptron



$$z_1 = -0.8 + 0.5x_1 + 0.5x_2$$
$$z_2 = -0.3 + 0.5x_1 + 0.5x_2$$

$$y_i = \begin{cases} 1 & z_i \geq 0 \\ 0 & z_i < 0 \end{cases}$$

Mapping

$$z = -0.5 - 0.8y_1 + 0.8y_2$$

Mapping

$$y = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$

$$y$$

25

# Hidden Layer

2-layer Neural Network
or
1-hidden–layer Neural Network

Input layer

Hidden layer

Output layer

3-layer Neural Network
or
2–hidden–layer Neural Network

Input layer

Hidden layer 1

Hidden layer 2

Output layer

# Tensorflow Playground

https://playground.tensorflow.org/

# Forward Propagation

# Forward Propagation



$W^{(1)}$     $W^{(2)}$

$x_0$

$a_0^{(2)}$

$x_1$     $a_1^{(2)}$     $a_1^{(3)}$

$x_2$     $a_1^{(2)}$

Input layer     Hidden layer     Output layer

- $a_j^{(i)}$ : **"Activation"**
  of the $i$-th unit in the $j$-th layer

- $W^{(j)}$: **"Weight Matrix"** mapping
  from the $j$-th layer to the $(j+1)$-th layer

몇번째 노드?

layer

# Forward Propagation



$$z_1^{(2)} = W_{10}^{(3)} x_0 + W_{11}^{(1)} + W_{12}^{(1)} x_2$$

$$= \begin{bmatrix} W_{10}^{(3)} & W_{11}^{(1)} & W_{12}^{(1)} \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

$$a_1^{(2)} = g(z_1^{(2)})$$

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

**Logistic function
(Sigmoid function)**

# Forward Propagation

$$\begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \end{bmatrix} = \begin{bmatrix} W_{10}^{(3)} & W_{11}^{(1)} & W_{12}^{(1)} \\ W_{20}^{(3)} & W_{21}^{(1)} & W_{22}^{(1)} \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

$$\begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \end{bmatrix} = \begin{bmatrix} g\left(z_1^{(2)}\right) \\ g\left(z_2^{(2)}\right) \end{bmatrix}$$
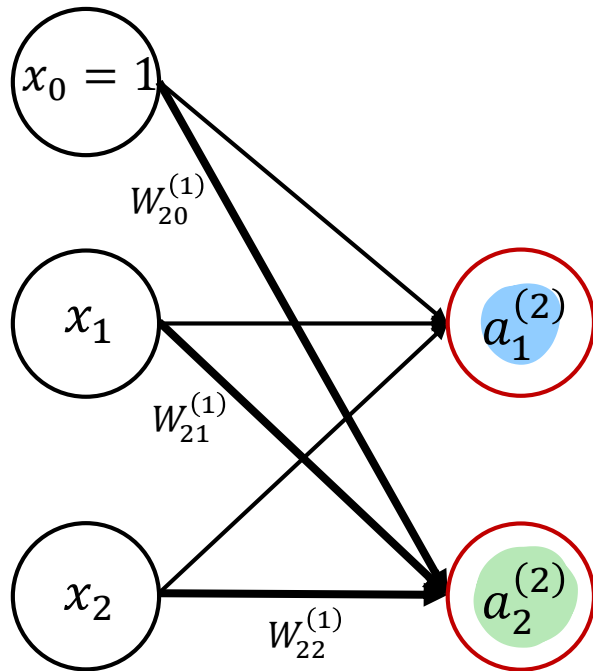
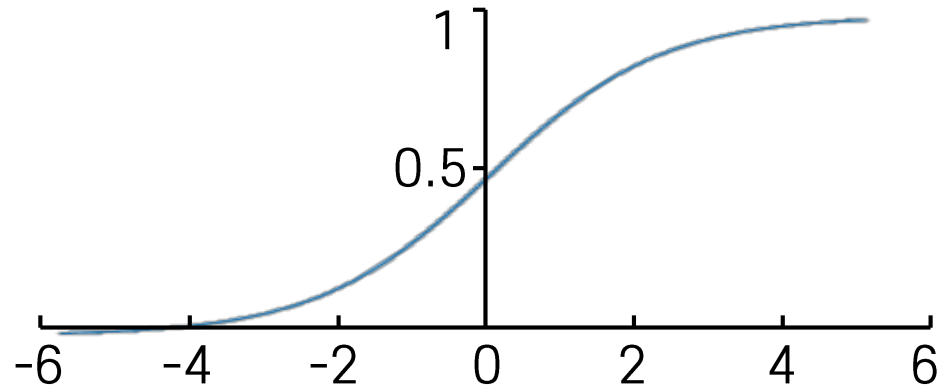$$g(x) = \frac{1}{1+e^{-x}}$$

**Logistic function
(Sigmoid function)**

# Forward Propagation



$$\begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \end{bmatrix} = \begin{bmatrix} W_{10}^{(3)} & W_{11}^{(1)} & W_{12}^{(1)} \\ W_{20}^{(3)} & W_{21}^{(1)} & W_{22}^{(1)} \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

$$\begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \end{bmatrix} = \begin{bmatrix} g\left(z_1^{(2)}\right) \\ g\left(z_2^{(2)}\right) \end{bmatrix}$$

$$z_1^{(3)} = \begin{bmatrix} W_{20}^{(3)} & W_{21}^{(1)} & W_{22}^{(1)} \end{bmatrix} \begin{bmatrix} a_0^{(2)} \\ a_1^{(2)} \\ a_2^{(2)} \end{bmatrix}$$

$$a_1^{(3)} = g\left(z_1^{(3)}\right)$$

활성화함수

# Linear Layer



$$\begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \end{bmatrix} = \begin{bmatrix} W_{10}^{(3)} & W_{11}^{(1)} & W_{12}^{(1)} \\ W_{20}^{(3)} & W_{21}^{(1)} & W_{22}^{(1)} \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$
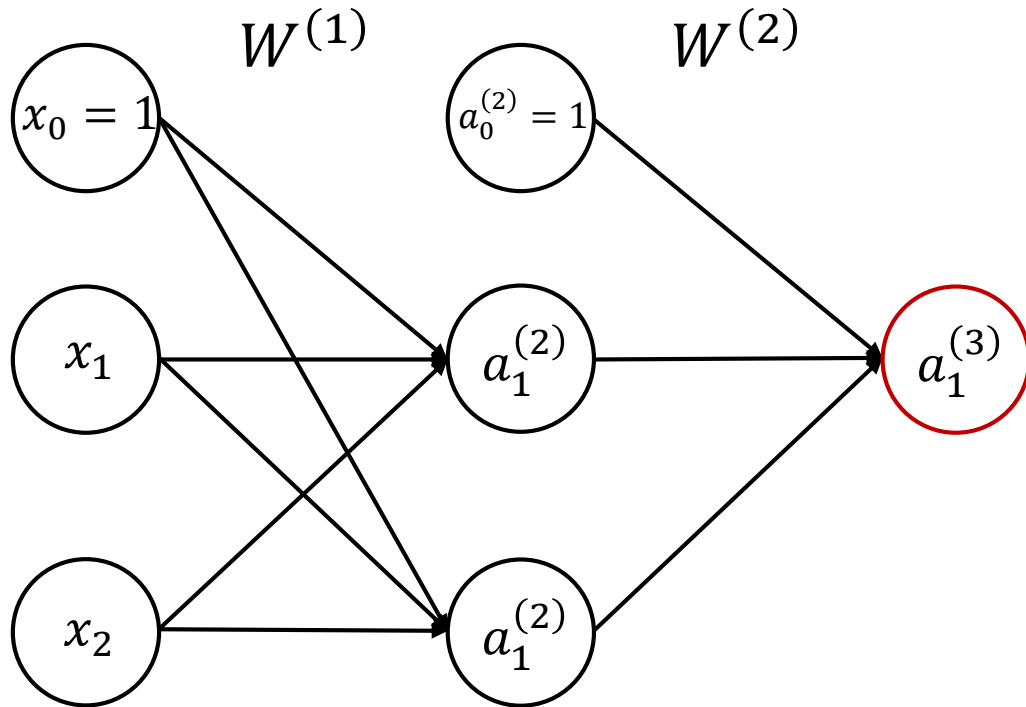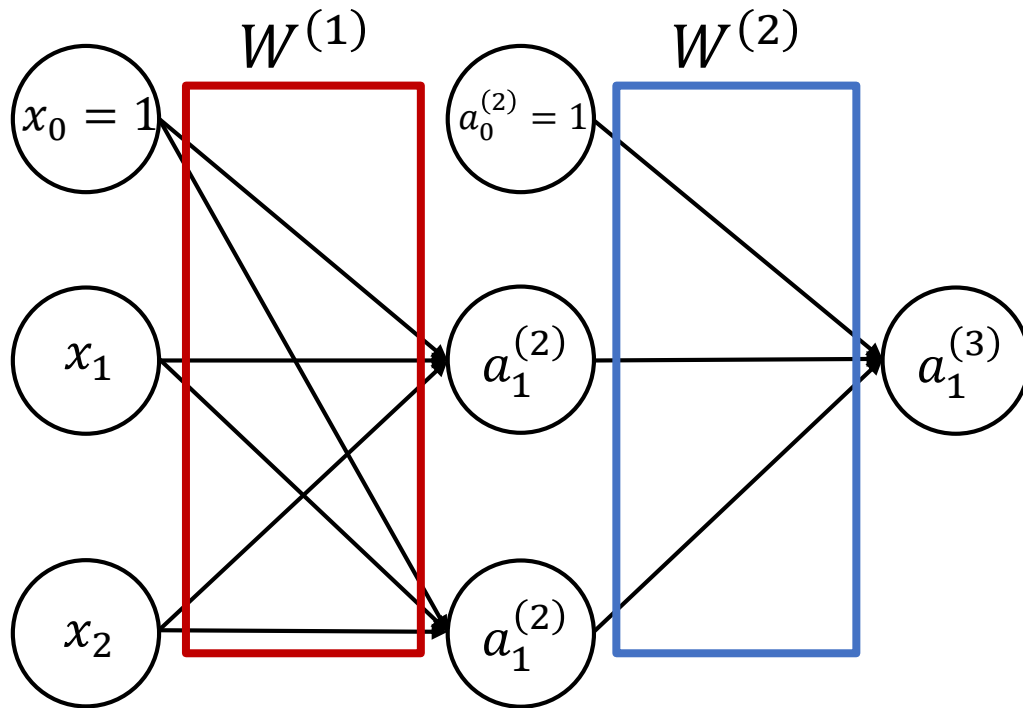
$$z_1^{(3)} = \begin{bmatrix} W_{20}^{(3)} & W_{21}^{(1)} & W_{22}^{(1)} \end{bmatrix} \begin{bmatrix} a_0^{(2)} \\ a_1^{(2)} \\ a_2^{(2)} \end{bmatrix}$$

Each layer performs linear transformation, so it is also called a **linear layer**
**Linear Layer** and **Fully-connected Layer** are the same thing.
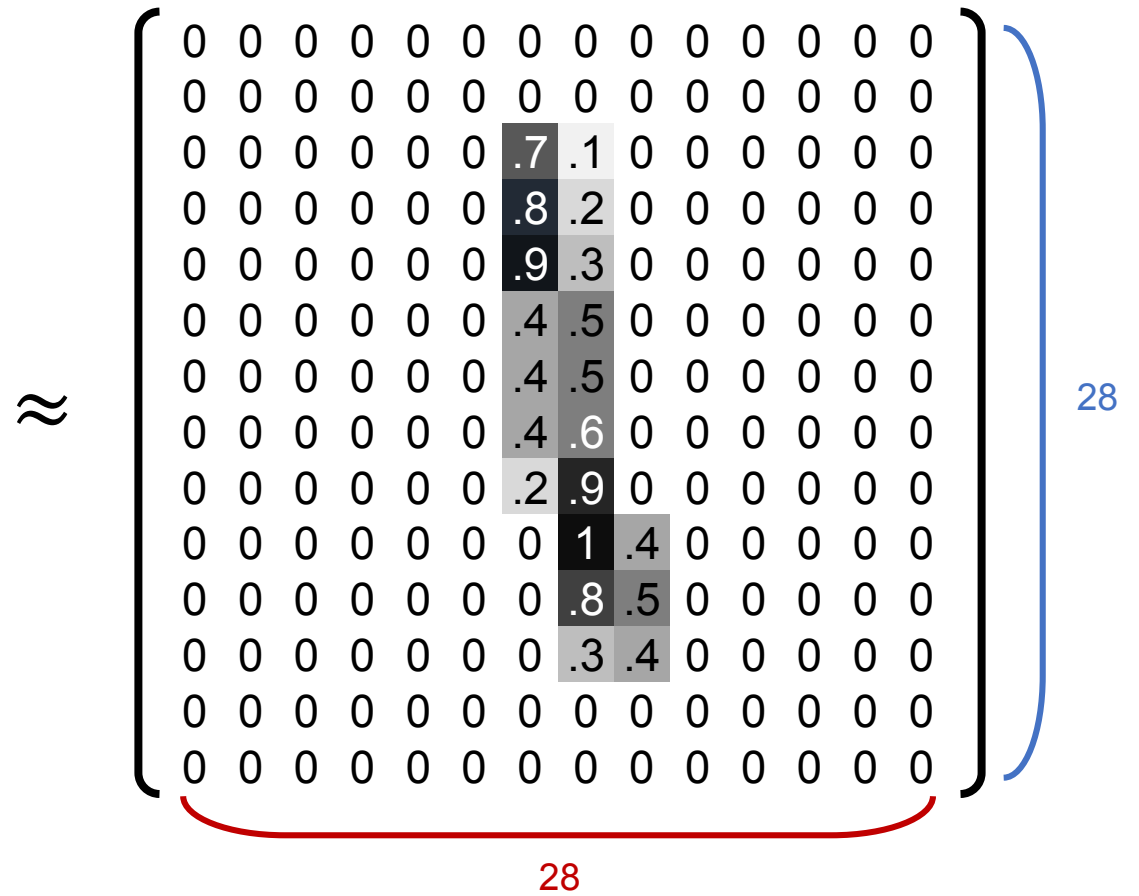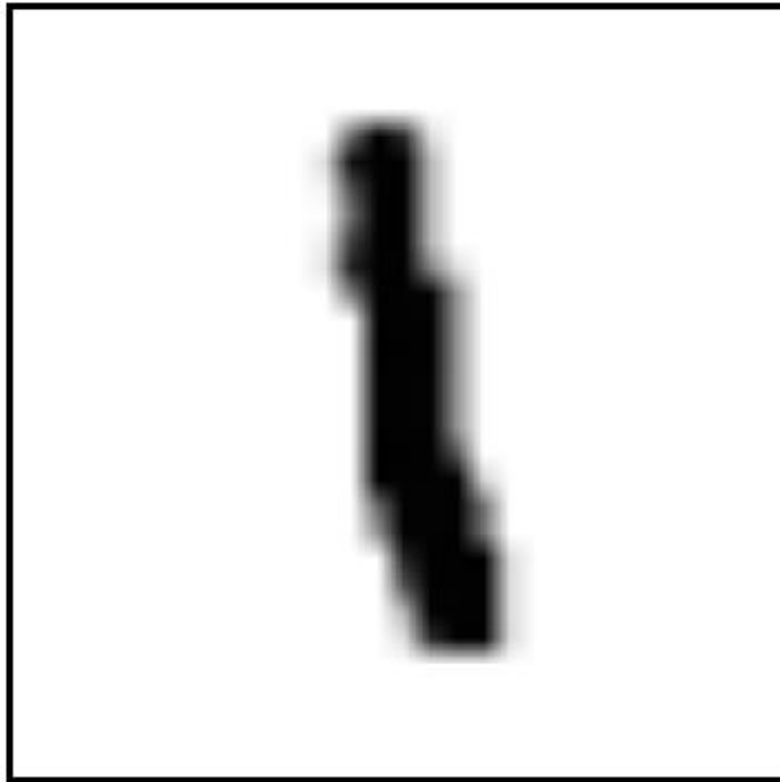
# MNIST Dataset

**MNIST (Modified National Institute of Standards and Technology)**
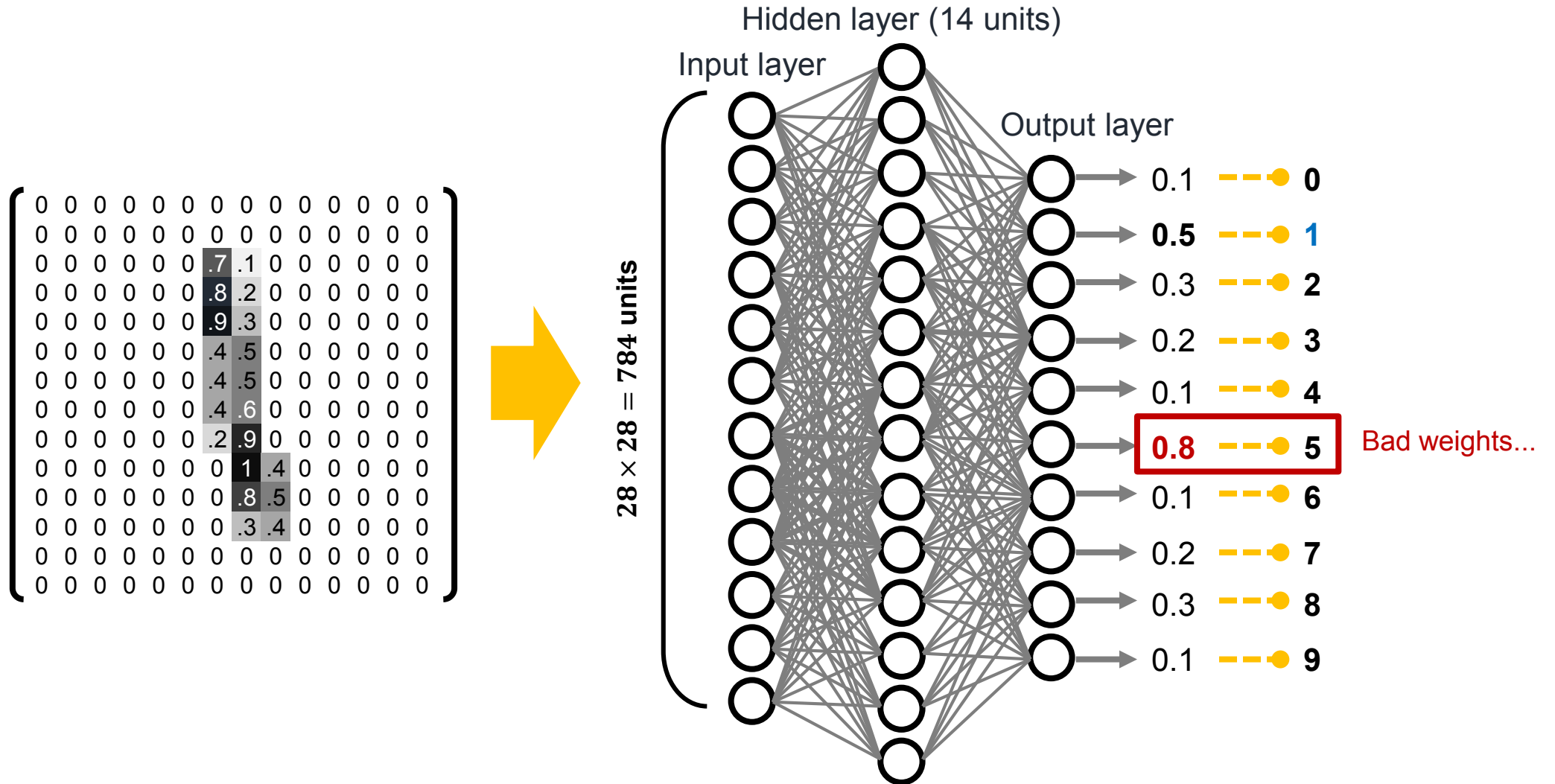
28 × 28



- Handwritten digits from 0 to 9
  - 55,000 training examples
  - 10,000 testing examples

- Each image has been preprocessed
  - Digits are center-aligned
  - Digit size is rescaled to similar size
  - Each image has fixed size of $28 \times 28$
  → Real number matrix from 0.0 to 1.0

# Example of MNIST

35

# MNIST Classification Model

# MNIST Classification Model



**Squared Error:** **1.19**    Let's find **the weight** which minimizes error!

ground truth 결과물
loss function

Softmax Layer (Softmax Classifier)

# Problem of Sigmoid Outputs and Mean-Squared Error Loss



- Because of **sigmoid outputs**, Prediction $\in (0,1)$ & Target $\in \{0,1\}$

→ **Upper limits exist on loss and gradient magnitude with MSE Loss**

$$\max L = \max_{y_i \in \{0,1\}, \widehat{y_i} \in (0,1)} \sum_{i=1}^{n} (\widehat{y_i} - y_i)^2 < 1, \qquad \max \left| \frac{\partial L}{\partial \widehat{y_i}} \right| < 2$$

- In addition, a better output would be a sum-to-one probability vector over multiple possible classes.

# Softmax Layer (or Softmax Classifier) for Multi-Class Classification



- The softmax layer applies a monotonically increasing, exponential function to a logit vector:

  - Map the value in $(-\infty, \infty)$ to $(0, \infty)$

  - Preserve the order of values

- Calculate the relative proportions with respect to the sum of these positive values, resulting in a sum-to-one probability vector

# Softmax Loss (or Cross-Entropy Loss)



The $\mathbf{\mathit{i}}$-th image

$W^\top$   $x_i$   $p_c$

$c \in \mathbf{Classes} = \{\mathbf{1}, \mathbf{2}, \mathbf{3}\}$

Groundtruth class
Cat := 2

$y_i$

$y_i \in \mathbf{Class}$

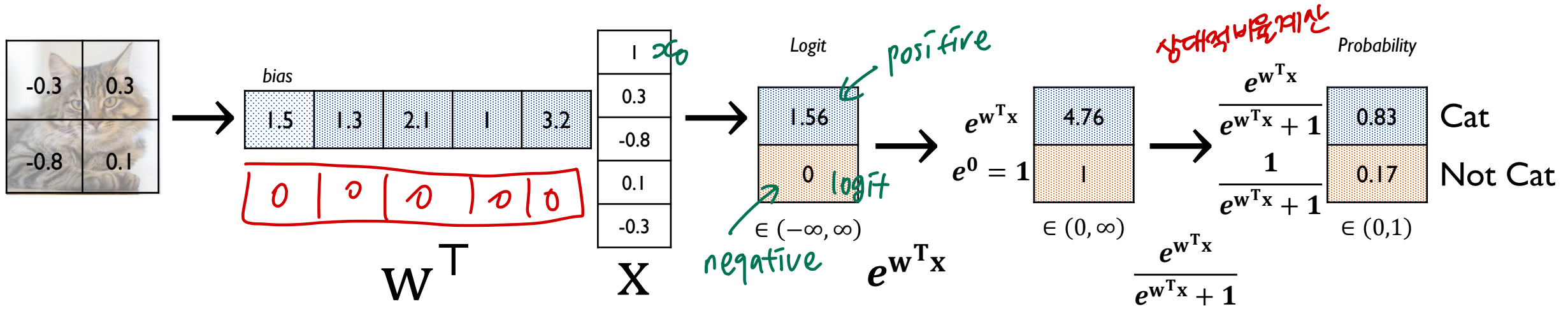- **Softmax loss**, also known as **cross-entropy loss** or **negative log-likelihood** (NLL) loss used for training a softmax classifier is defined as

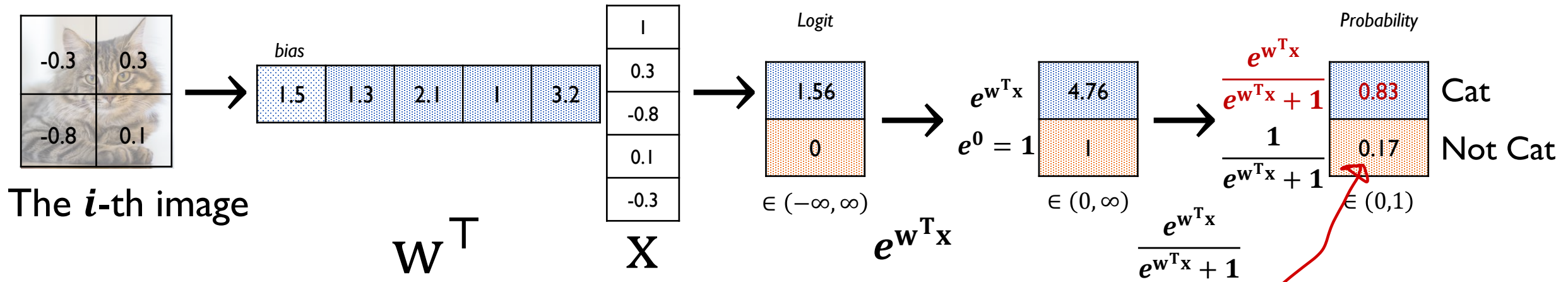$$L = -\sum_{c=1}^{C} y_c \log(\hat{p}_c) = -\log(\hat{p}_{y_i})$$

ground truth vector

$-\log x$

# Logistic Regression as a Special Case of Softmax Classifier



- **Logistic regression** ➔ Softmax classifier whose logit for a negative class is set as a constant value of 0.

- Logistic regression is used for a binary classification.

- The softmax classifier can also be used for two classes by using the matrix $W$ with two columns, i.e., using the twice the number of parameters of a logistic regression.

# Logistic Regression as a Special Case of Softmax Classifier



- Binary cross-entropy (BCE) loss for logistic regression is defined as

$$L = -\sum_{c=1}^{2} y_c \log(\hat{p}_c) = -y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i)$$

$\leftarrow$ BCE Loss

where $y_i = 1$ for a positive class, e.g., **cat**, and $y_i = 0$ for a negative class, e.g., **not cat**, and