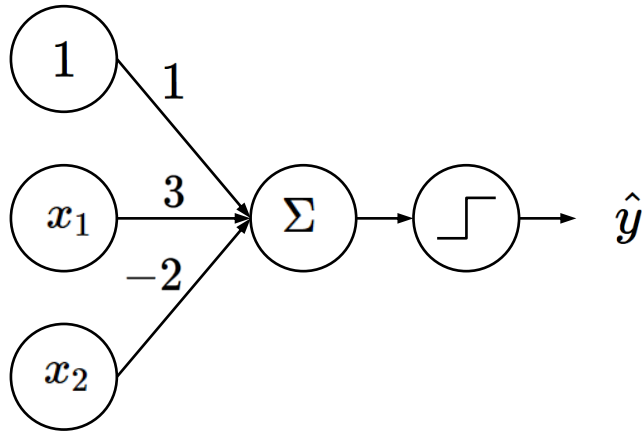




(Artificial) Neural Networks: From Perceptron to MLP

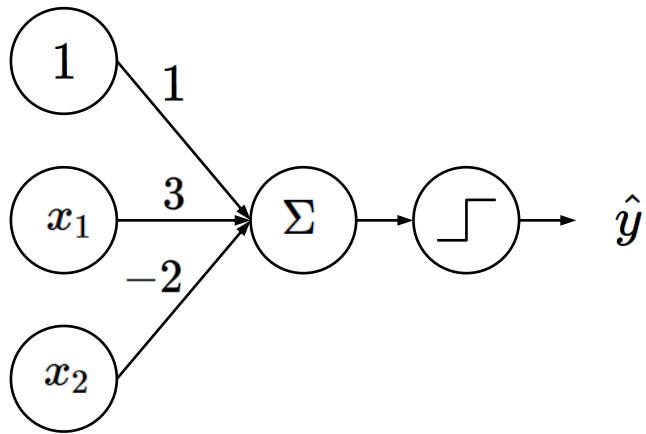
Prof. Seungchul Lee
Industrial AI Lab.

Perceptron: Example

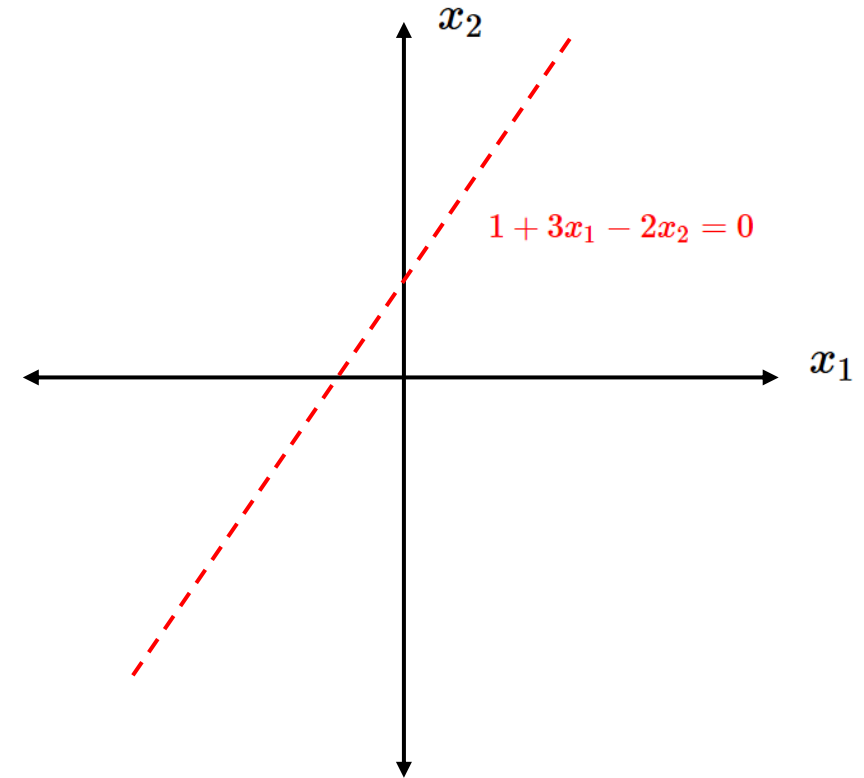


$$\begin{aligned}\hat{y} &= g(\omega_0 + X^T \omega) \\ &= g\left(1 + \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 3 \\ -2 \end{bmatrix}\right) \\ &= g(1 + 3x_1 - 2x_2)\end{aligned}$$

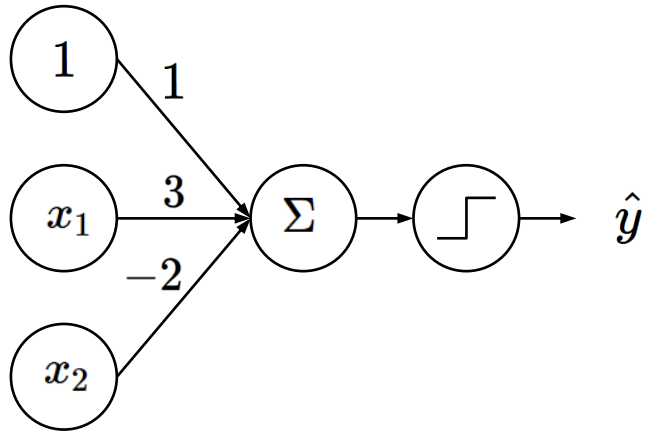
Perceptron: Example



$$\hat{y} = g(1 + 3x_1 - 2x_2)$$

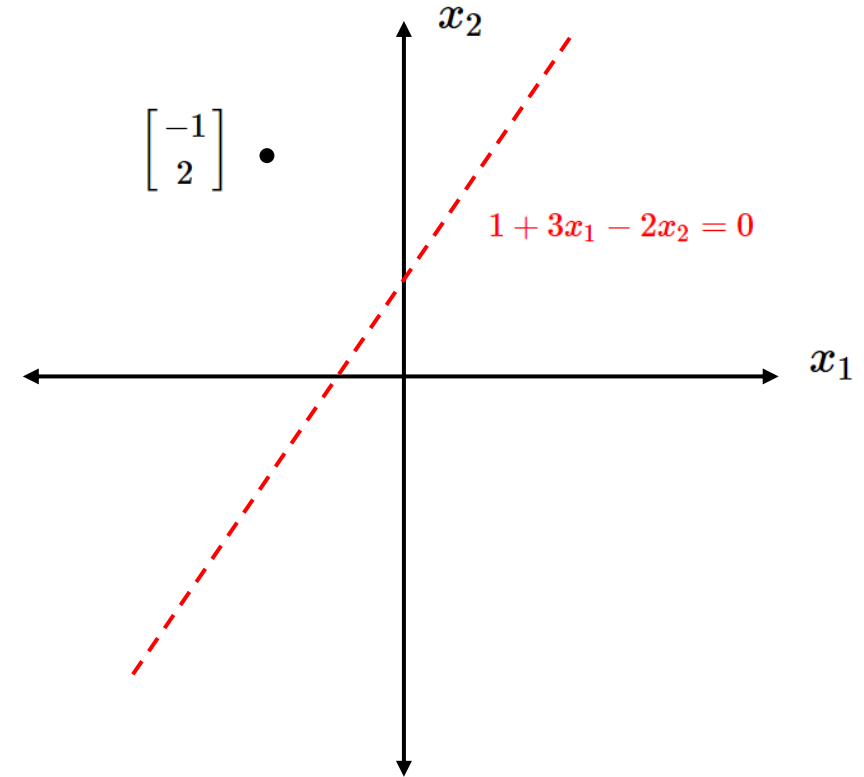


Perceptron: Example

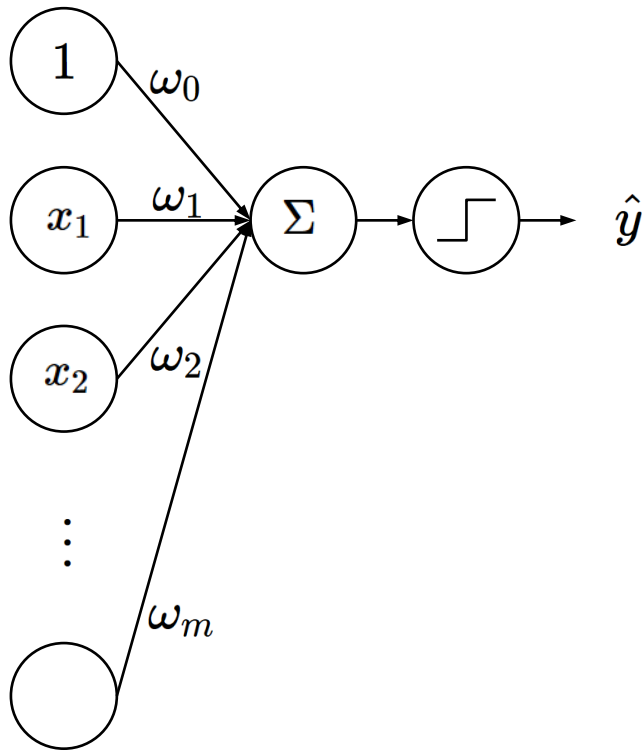


$$\hat{y} = g(1 + 3 \times (-1) - 2 \times 2) = g(-6) = -1$$

$$\hat{y} = g(1 + 3x_1 - 2x_2)$$



Perceptron: Forward Propagation



$$\hat{y} = g(\omega_0 + X^T \omega)$$

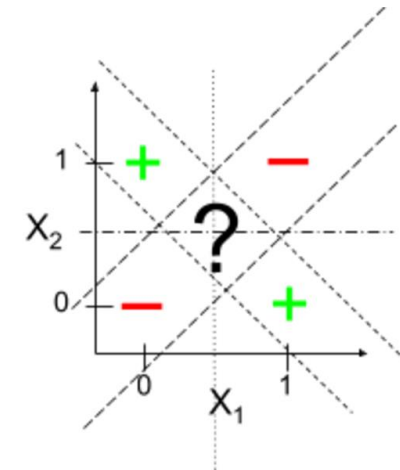
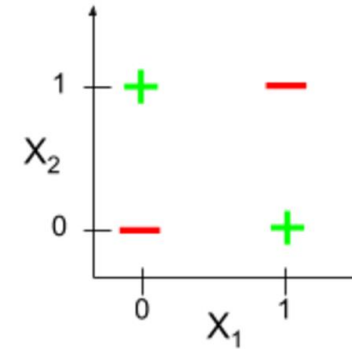
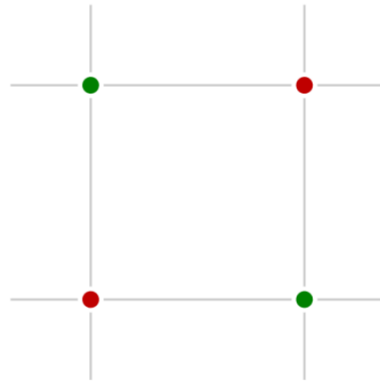
$$= g\left(\omega_0 + \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}^T \begin{bmatrix} \omega_1 \\ \vdots \\ \omega_m \end{bmatrix}\right)$$

From Perceptron to MLP

XOR Problem

- Minsky-Papert Controversy on XOR
 - Not linearly separable
 - Limitation of perceptron

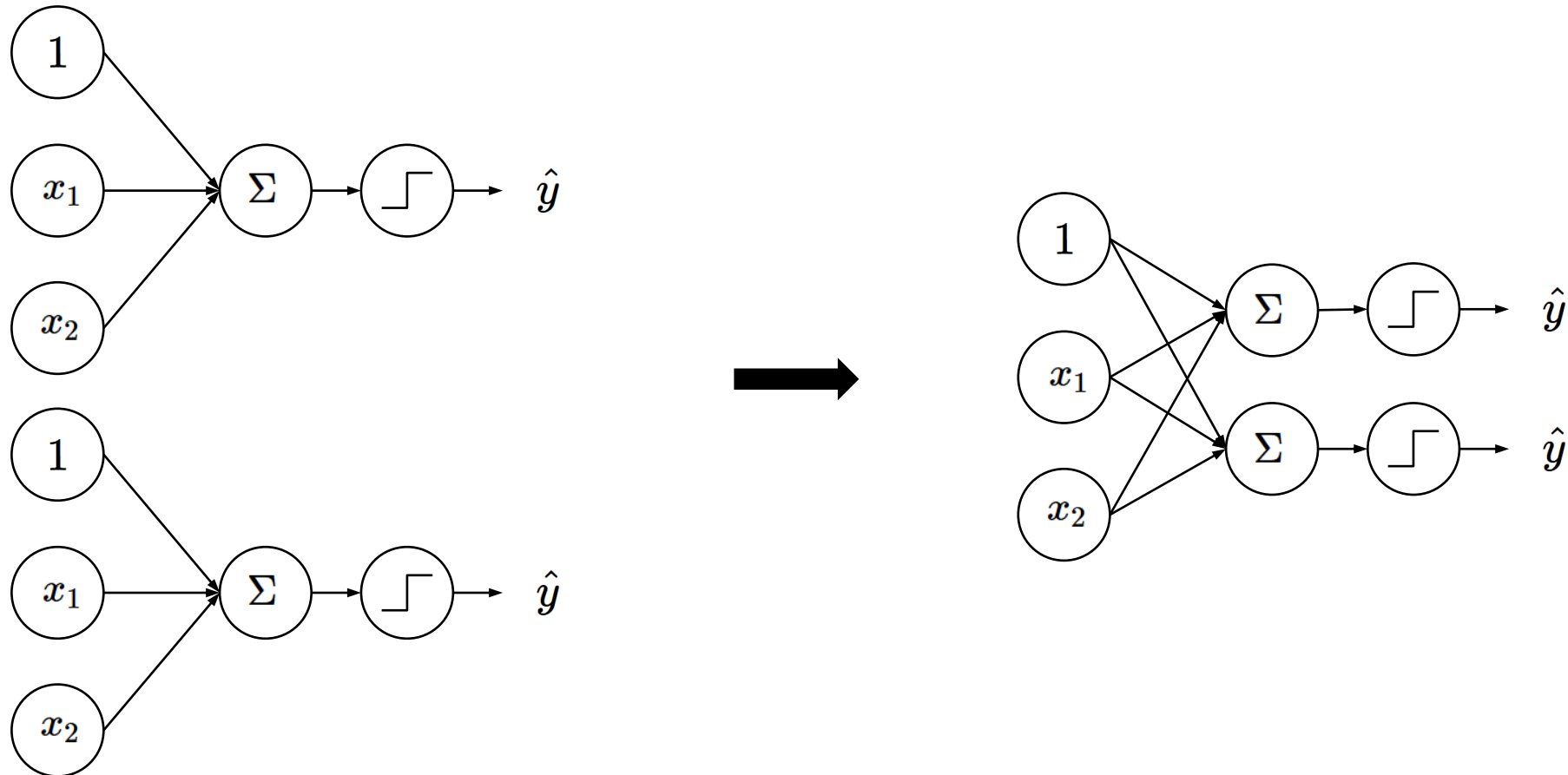
x_1	x_2	$x_1 \text{ XOR } x_2$
0	0	0
0	1	1
1	0	1
1	1	0



- Single neuron = **one linear classification boundary**

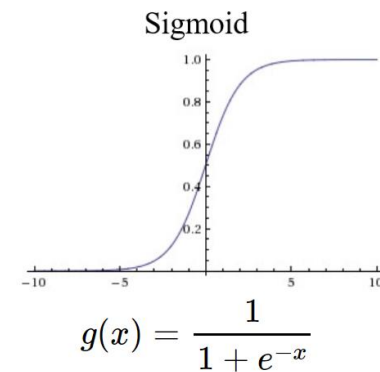
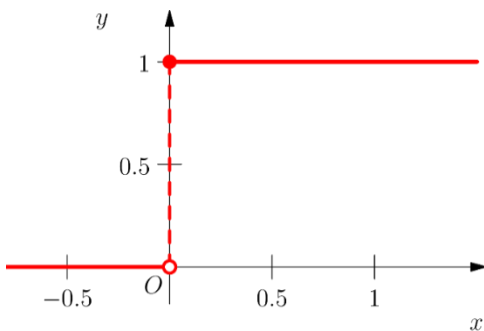
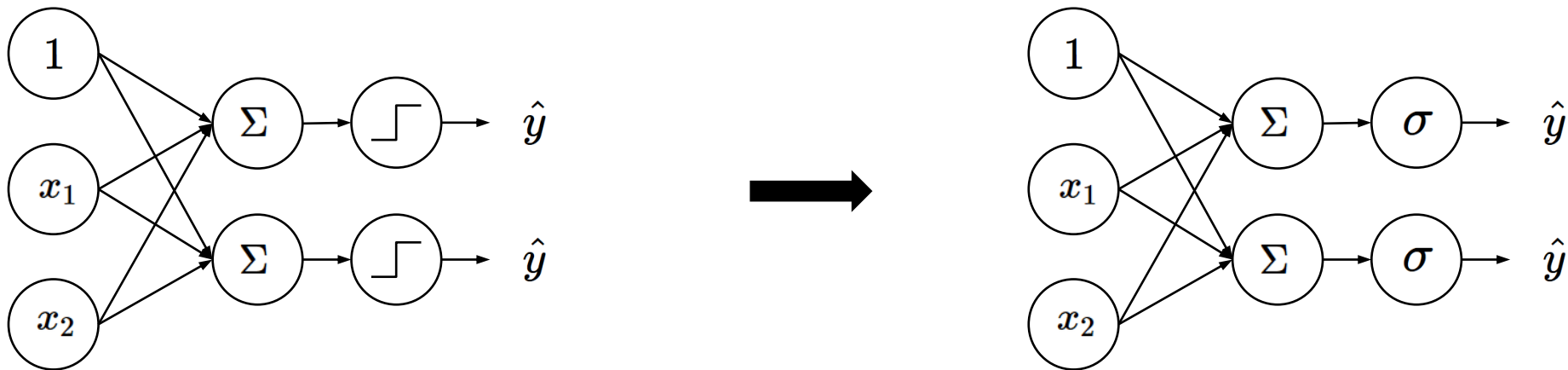
Artificial Neural Networks: MLP

- Multi-layer Perceptron (MLP) = Artificial Neural Networks (ANN)
 - Multi neurons = multiple linear classification boundaries



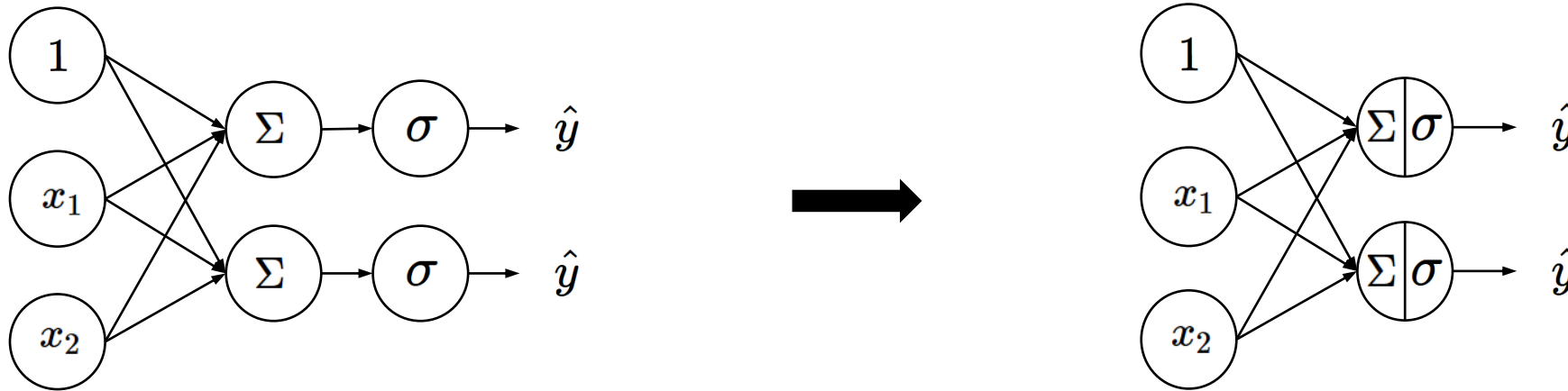
Artificial Neural Networks: Activation Function

- Differentiable nonlinear activation function



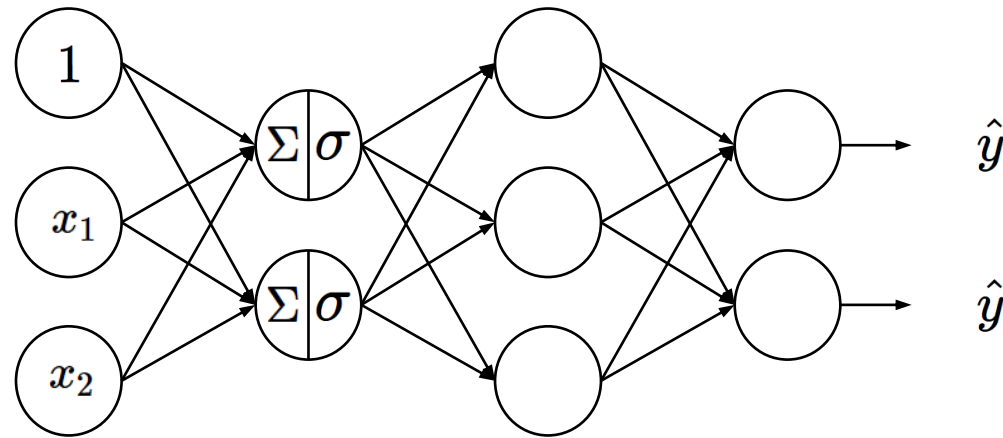
Artificial Neural Networks

- In a compact representation



Artificial Neural Networks

- A single layer is not enough to be able to represent complex relationship between input and output
⇒ perceptron with many layers and units



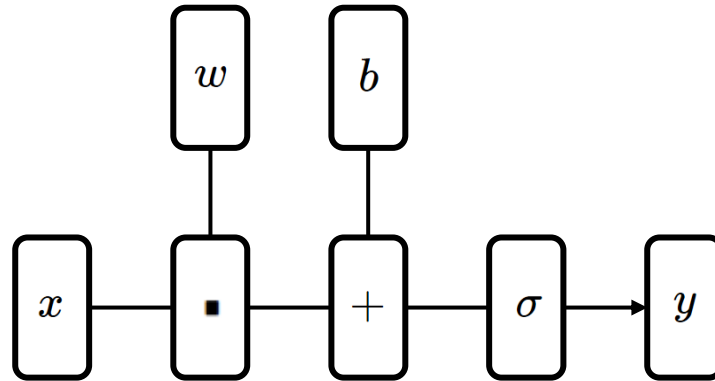
- Multi-layer perceptron
 - Features of features
 - Mapping of mappings

Another Perspective: ANN as Kernel Learning

Neuron

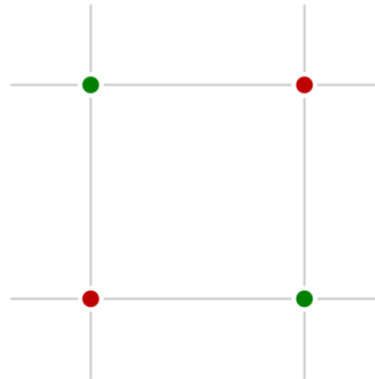
- We can represent this “neuron” as follows:

$$f(x) = \sigma(w \cdot x + b)$$



XOR Problem

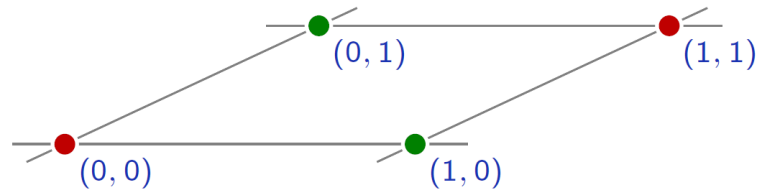
- The main weakness of linear predictors is their lack of capacity.
- For classification, the populations have to be linearly separable.



“xor”

Nonlinear Mapping

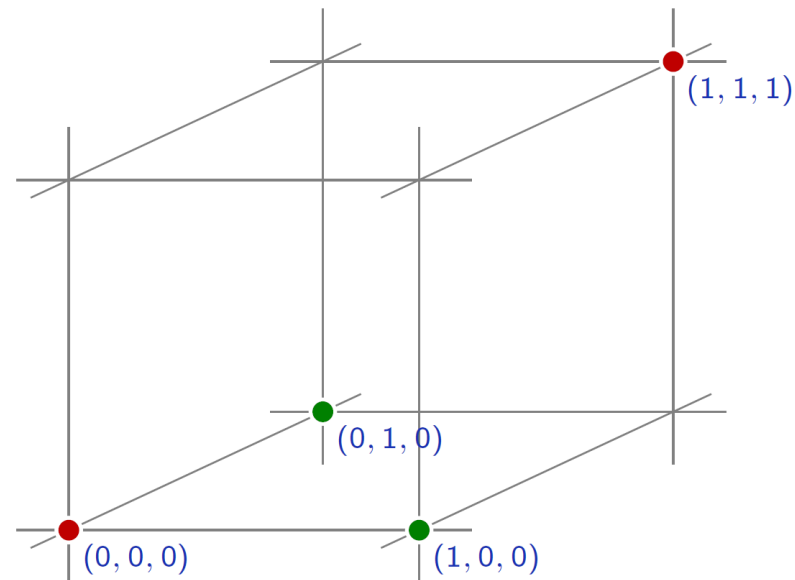
- The XOR example can be solved by pre-processing the data to make the two populations linearly separable.



Nonlinear Mapping

- The XOR example can be solved by pre-processing the data to make the two populations linearly separable.

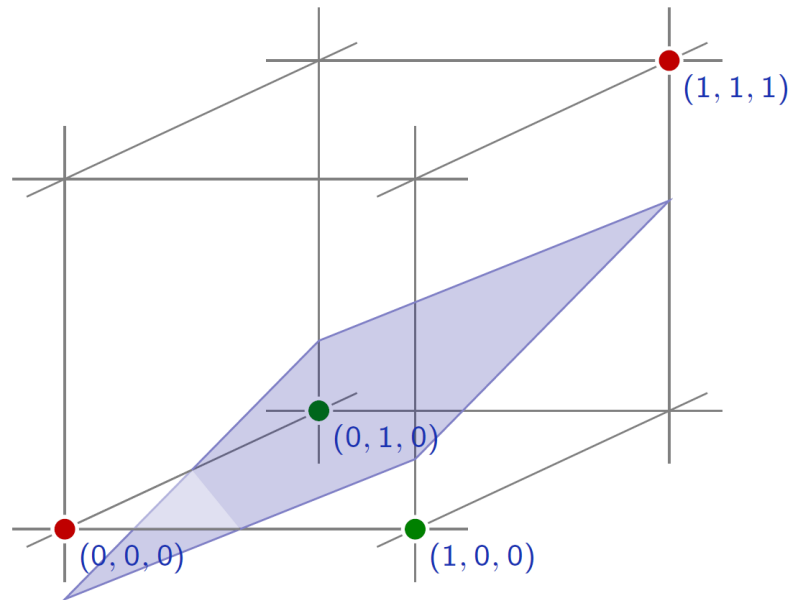
$$\phi : (x_u, x_v) \rightarrow (x_u, x_v, x_u x_v)$$



Nonlinear Mapping

- The XOR example can be solved by pre-processing the data to make the two populations linearly separable.

$$\phi : (x_u, x_v) \rightarrow (x_u, x_v, x_u x_v)$$



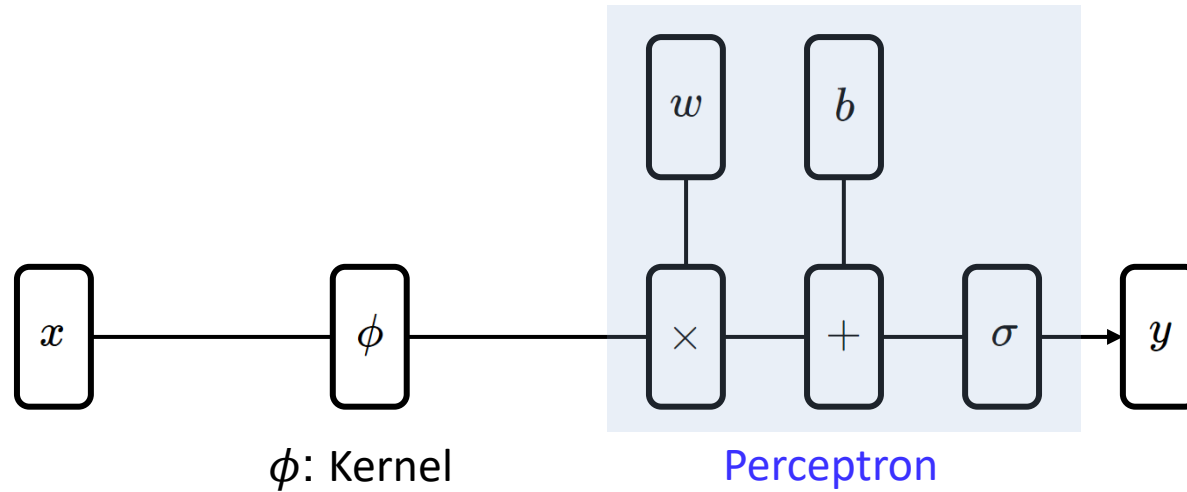
Kernel

- Often we want to capture nonlinear patterns in the data
 - nonlinear regression: input and output relationship may not be linear
 - nonlinear classification: classes may not be separable by a linear boundary
- Linear models (e.g. linear regression, linear SVM) are not just rich enough
 - by mapping data to higher dimensions where it exhibits linear patterns
 - apply the linear model in the new input feature space
 - mapping = changing the feature representation
- Kernels: make linear model work in nonlinear settings

Kernel + Neuron

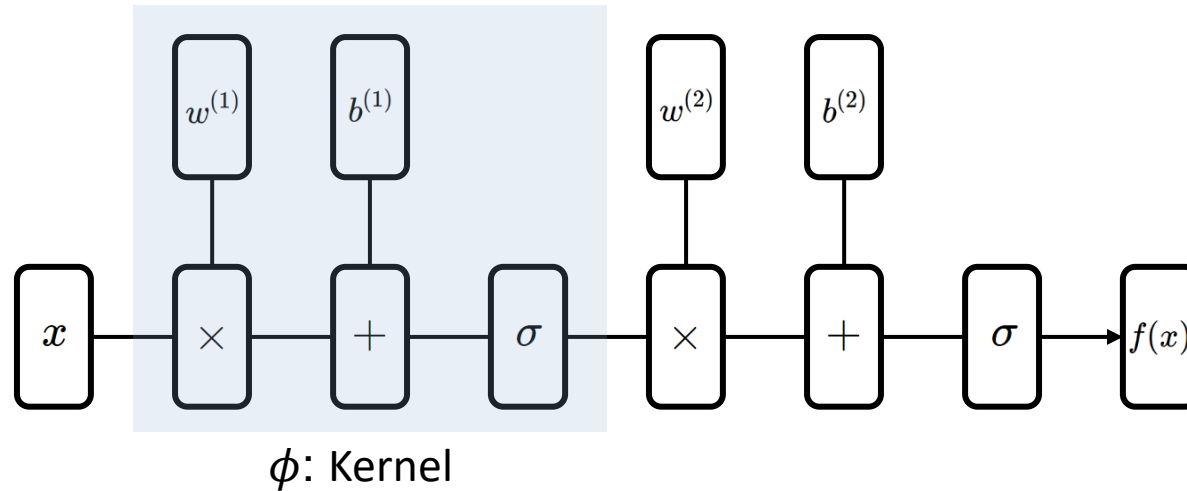
- Nonlinear mapping + neuron

$$\phi : (x_u, x_v) \rightarrow (x_u, x_v, x_u x_v)$$



Neuron + Neuron

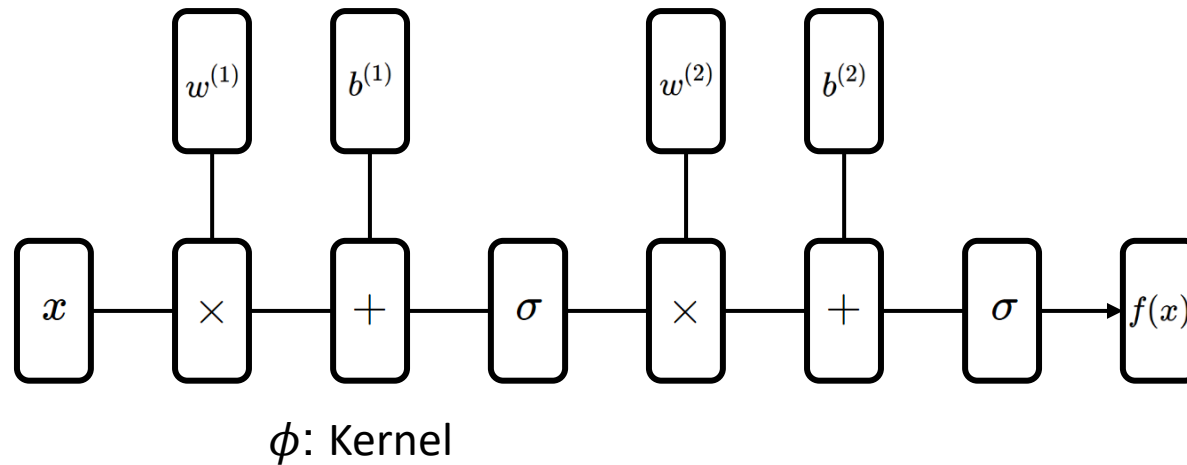
- Nonlinear mapping can be represented by another neurons



- Nonlinear Kernel
 - Nonlinear activation functions

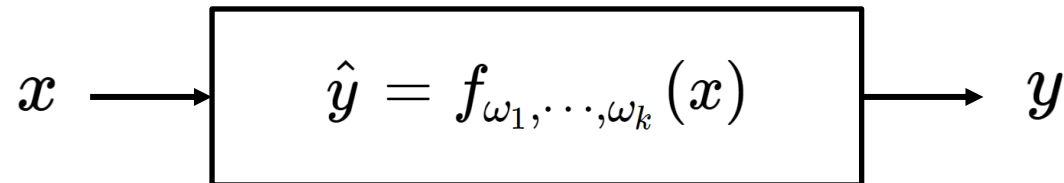
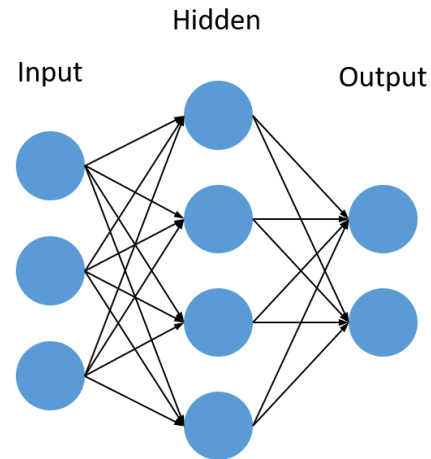
Multi Layer Perceptron

- Nonlinear mapping can be represented by another neurons
- We can generalize an MLP



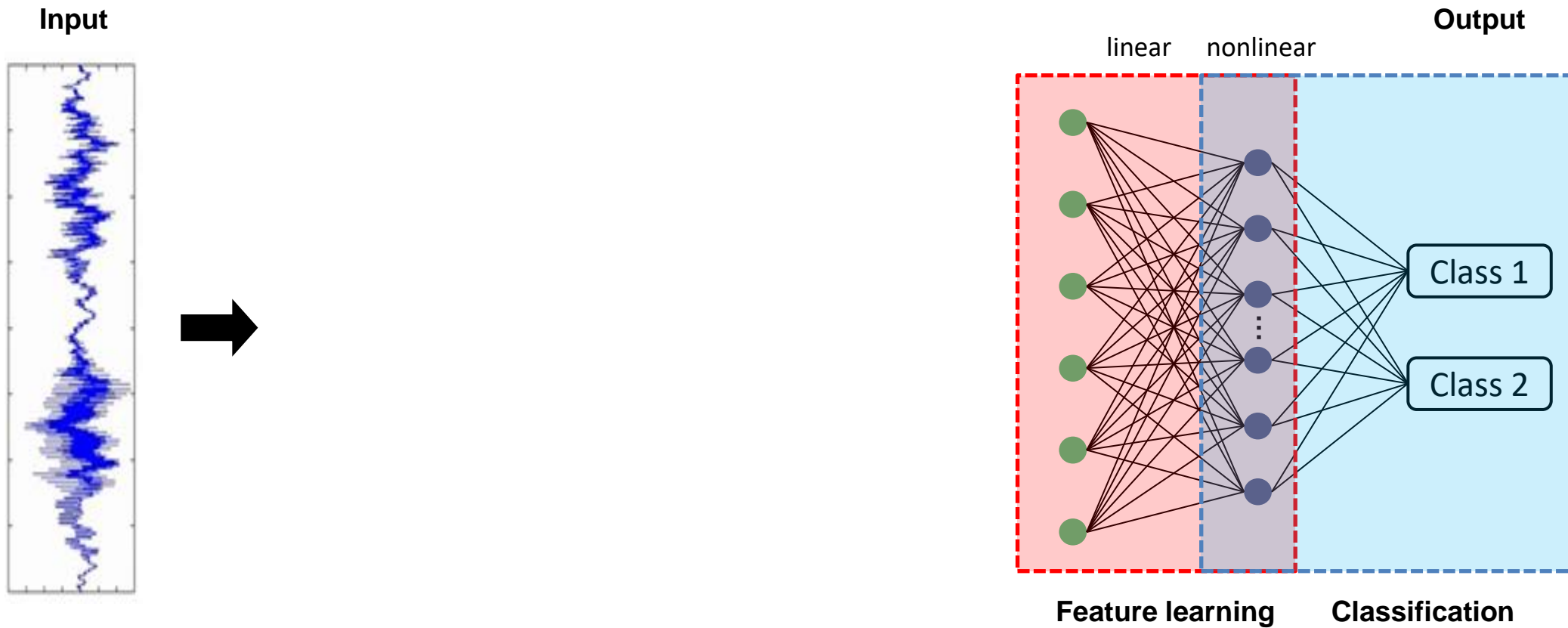
Summary

- Universal function approximator
- Universal function classifier
- Parameterized



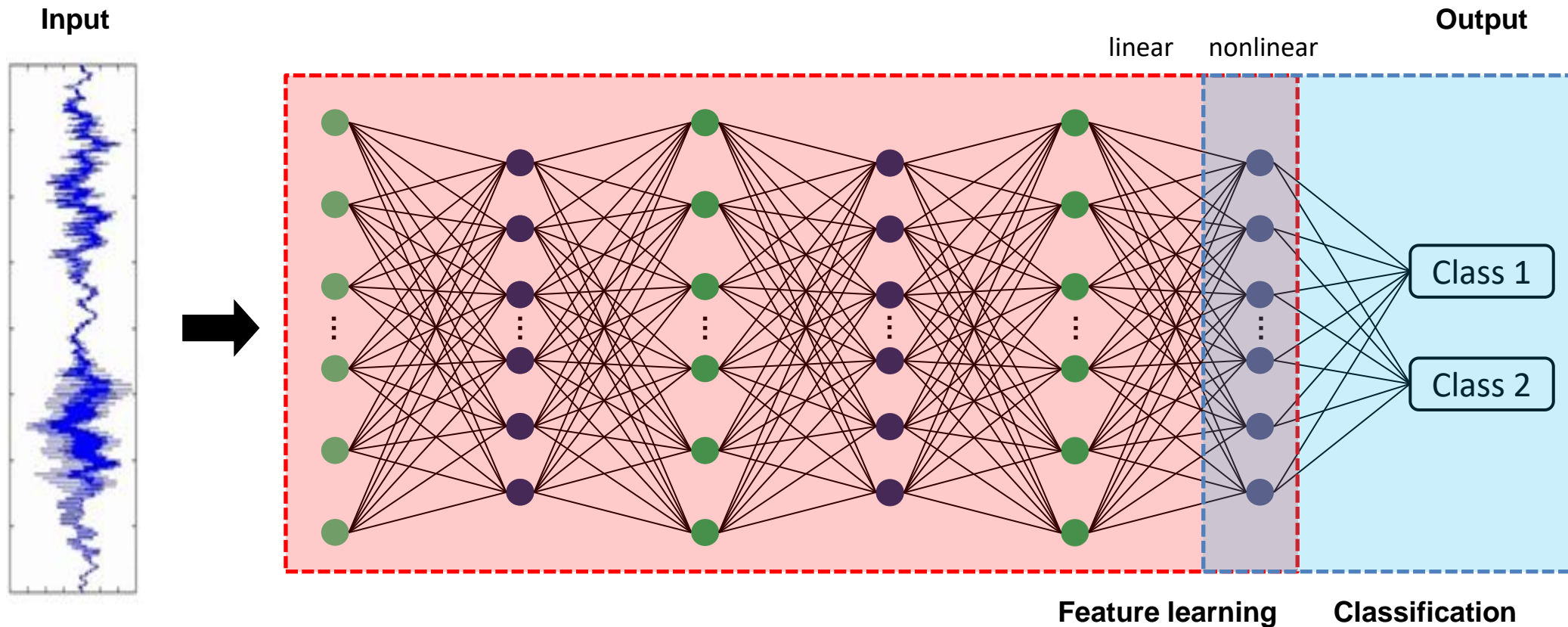
Artificial Neural Networks

- Complex/Nonlinear universal function approximator
 - Linearly connected networks
 - Simple nonlinear neurons



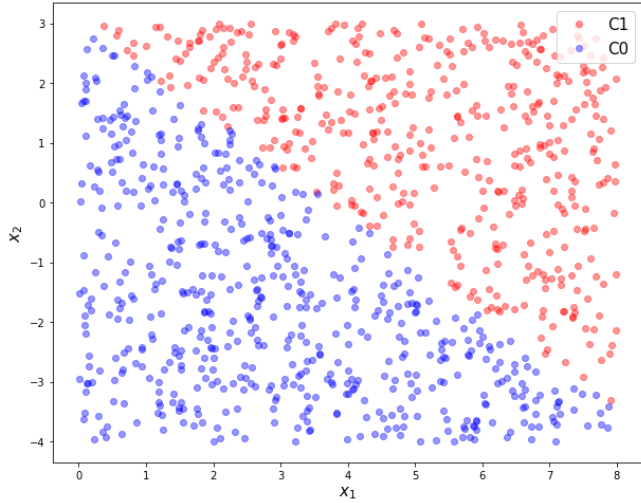
Deep Artificial Neural Networks

- Complex/Nonlinear universal function approximator
 - Linearly connected networks
 - Simple nonlinear neurons

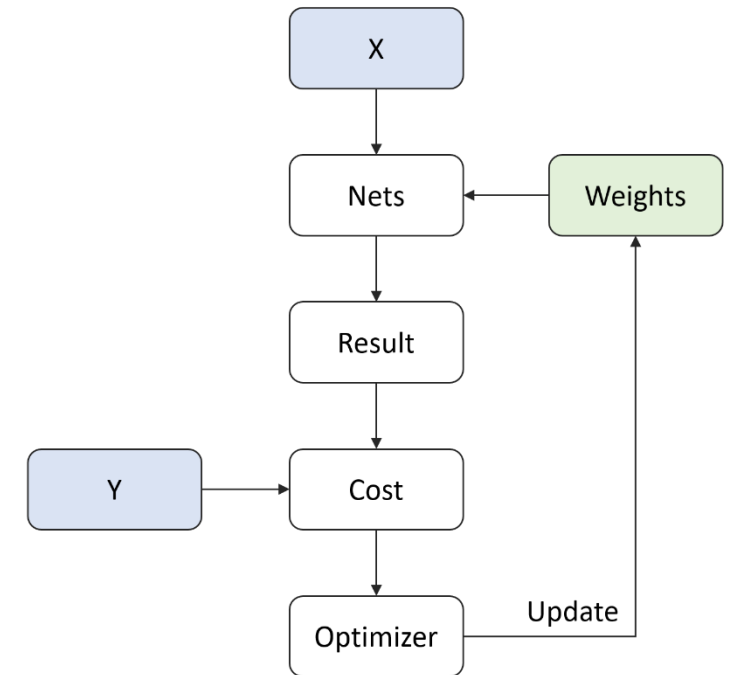
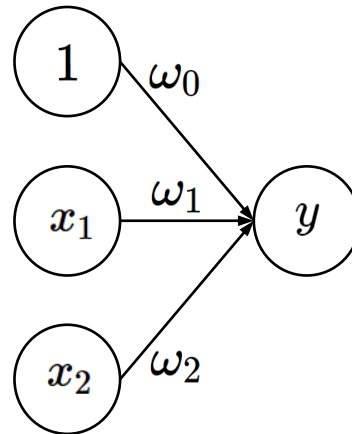


Looking at Parameters

Logistic Regression in a Form of Neural Network



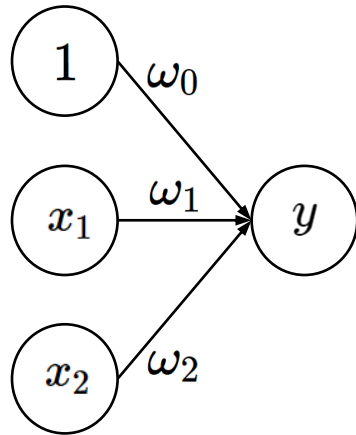
$$y = \sigma(\omega_0 + \omega_1 x_1 + \omega_2 x_2)$$



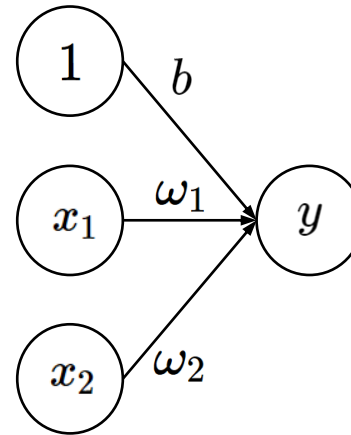
Logistic Regression in a Form of Neural Network

- Neural network convention

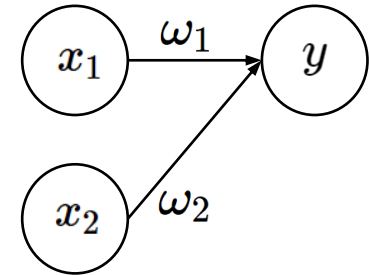
$$y = \sigma(\omega_0 + \omega_1 x_1 + \omega_2 x_2)$$



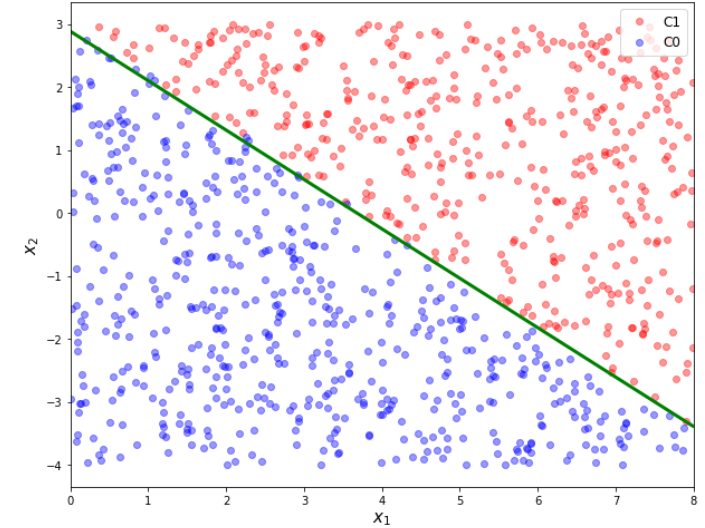
$$y = \sigma(b + \omega_1 x_1 + \omega_2 x_2)$$



Do not indicate bias units

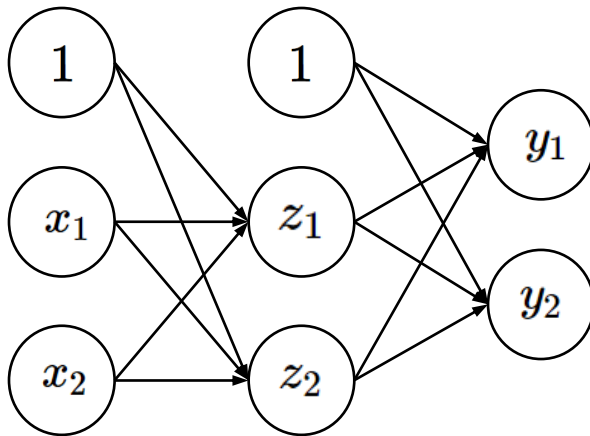


```
n_input = 2  
n_output = 1
```

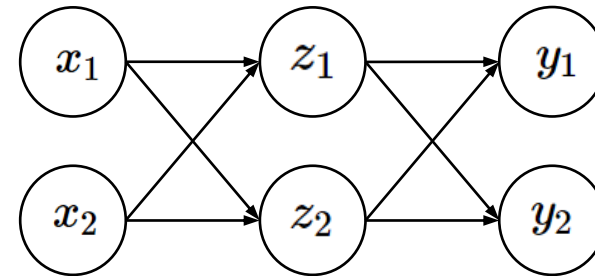


Nonlinearly Distributed Data

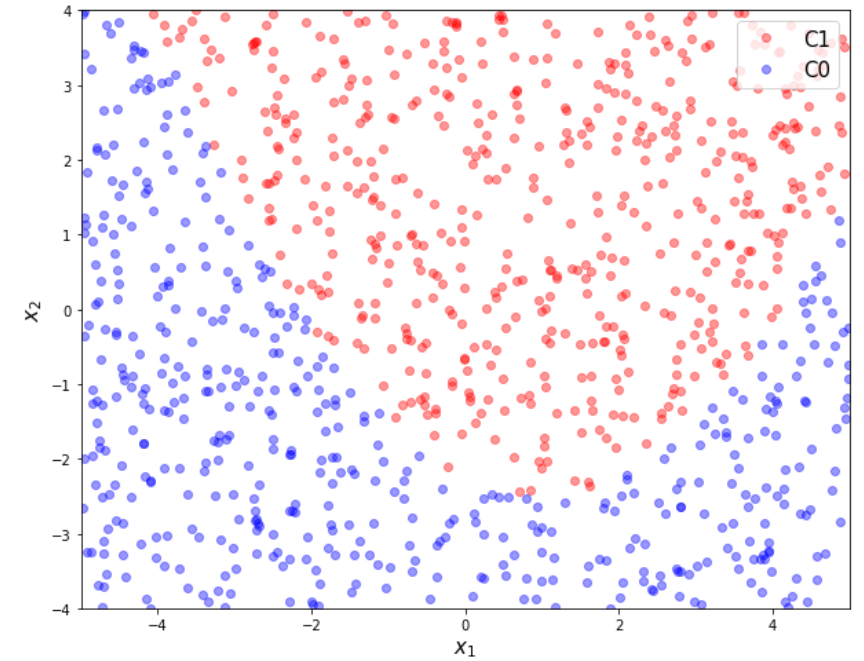
- Example to understand network's behavior
 - Include a hidden layer



Do not include bias units

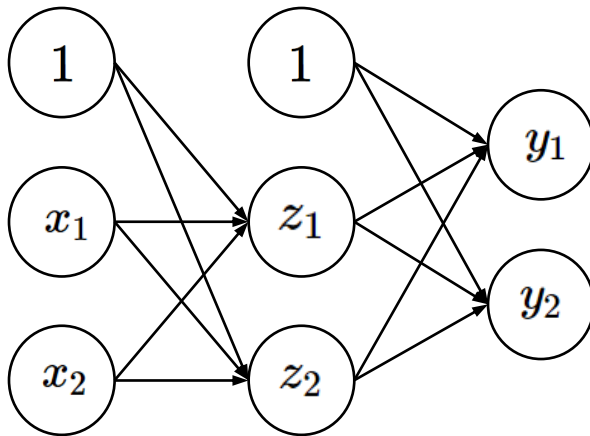


```
n_input = 2  
n_hidden = 2  
n_output = 2
```

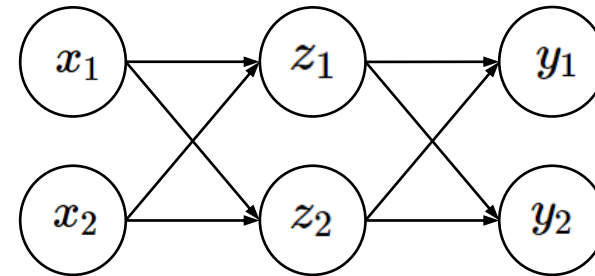


Multi Layers

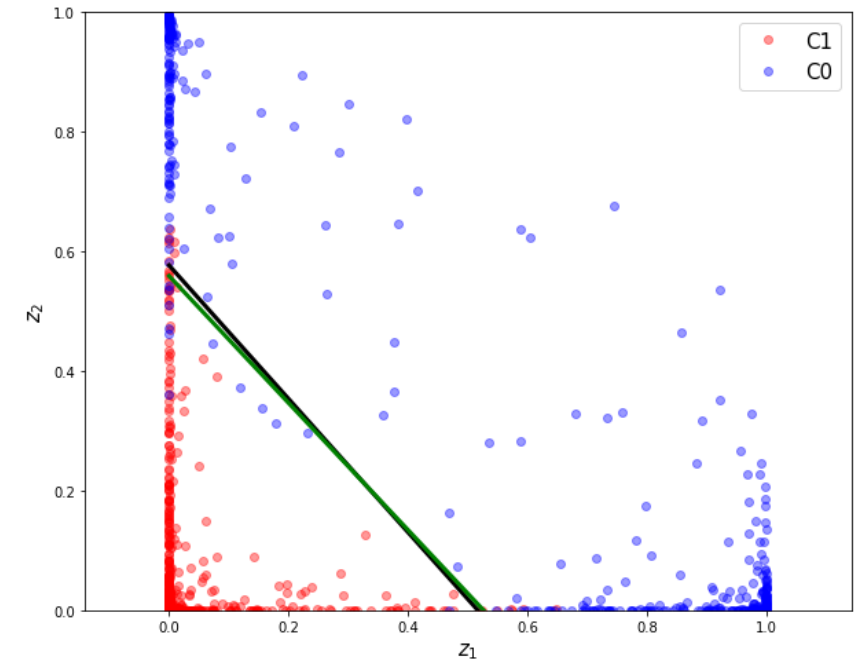
- z space



Do not include bias units

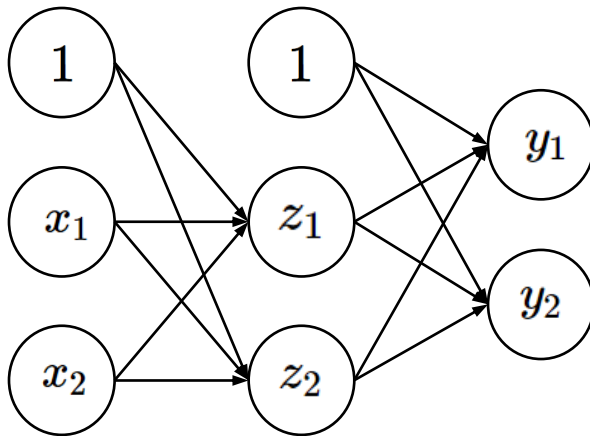


n_input = 2
n_hidden = 2
n_output = 2

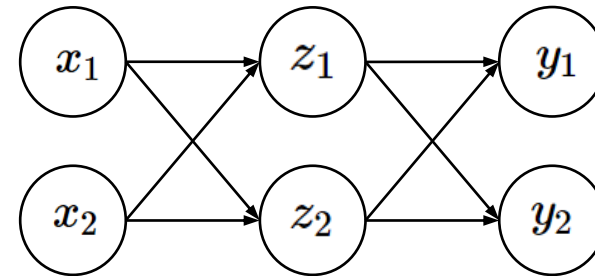


Multi Layers

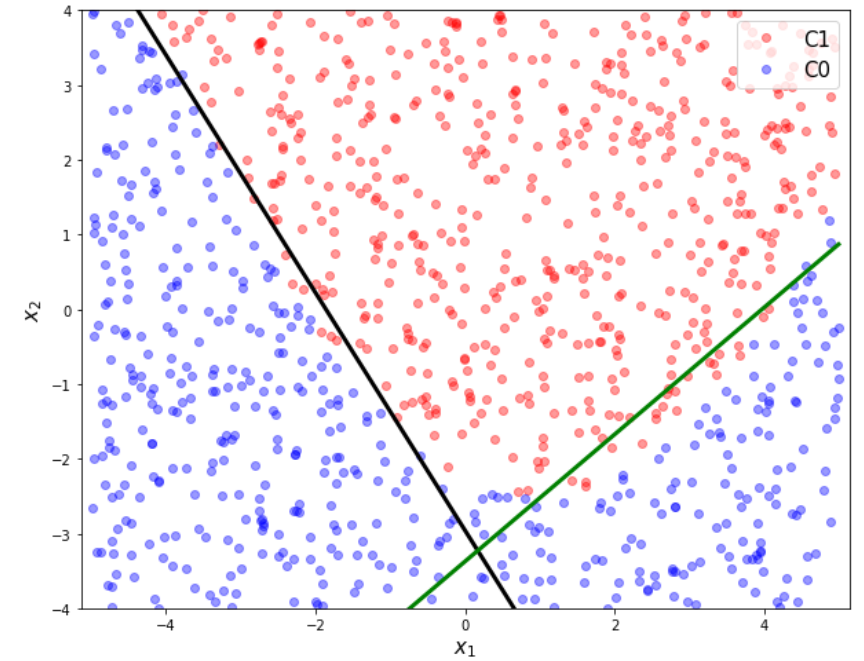
- x space



Do not include bias units



```
n_input = 2  
n_hidden = 2  
n_output = 2
```





(Artificial) Neural Networks: Training

Industrial AI
Prof. Seungchul Lee

Training Neural Networks: Loss Function

- Measures error between target values and predictions

$$\min_{\omega} \sum_{i=1}^m \ell \left(h_{\omega} \left(x^{(i)} \right), y^{(i)} \right)$$

- Example

- Squared loss (for regression):

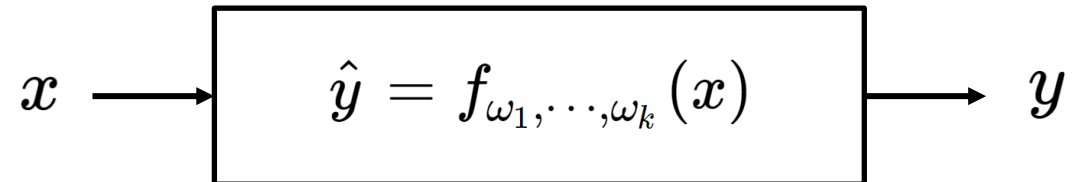
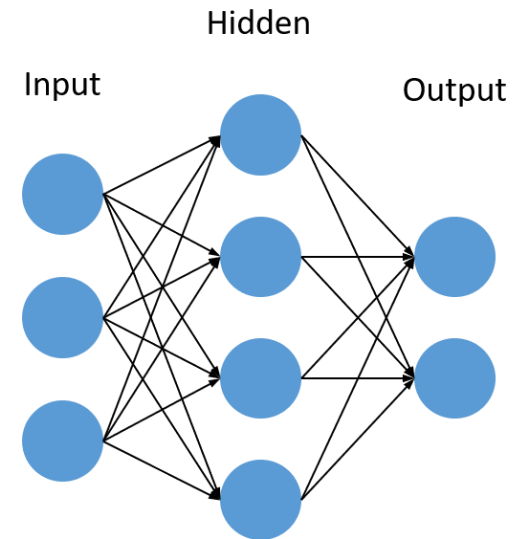
$$\frac{1}{m} \sum_{i=1}^m \left(h_{\omega} \left(x^{(i)} \right) - y^{(i)} \right)^2$$

- Cross entropy (for classification):

$$-\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \left(h_{\omega} \left(x^{(i)} \right) \right) + \left(1 - y^{(i)} \right) \log \left(1 - h_{\omega} \left(x^{(i)} \right) \right)$$

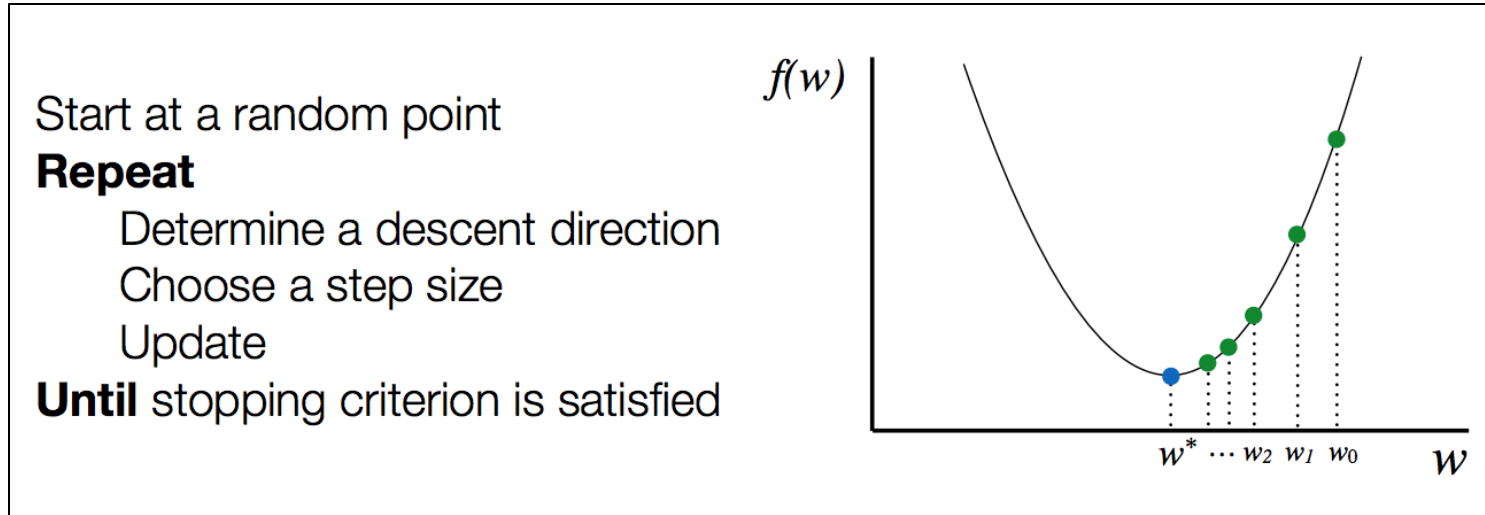
Gradients in ANN

- Learning weights and biases from data using gradient descent
- $\frac{\partial \ell}{\partial \omega}$: too many computations are required for all ω
- Structural constraint of NN:
 - Composition of functions
 - Chain rule
 - Dynamic programming



Training Neural Networks with TensorFlow

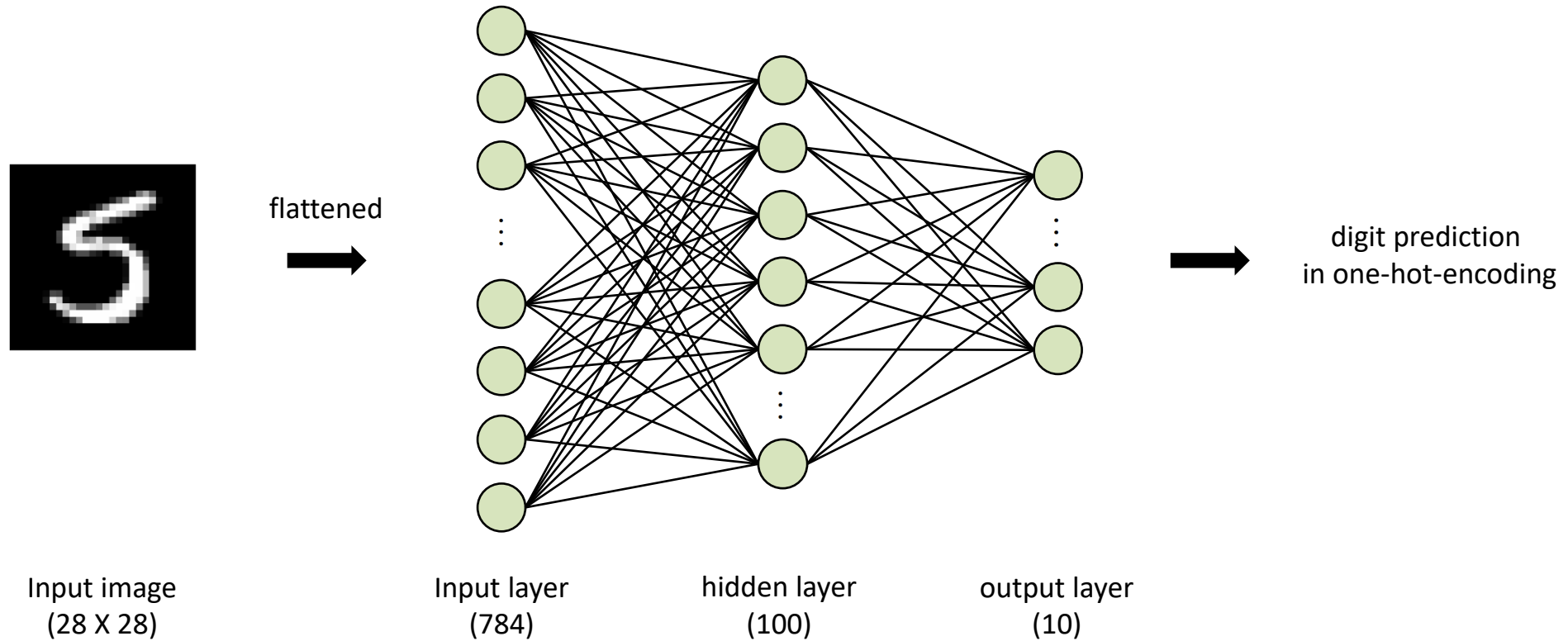
- Optimization procedure



- It is not easy to numerically compute gradients in network in general.
 - The good news: people have already done all the "hard work" of developing numerical solvers (or libraries)
 - There are a wide range of tools → We will use the TensorFlow

ANN in TensorFlow: MNIST

Our Network Model





Machine Learning and Deep Learning

Prof. Seungchul Lee
Industrial AI Lab.

Machine Learning

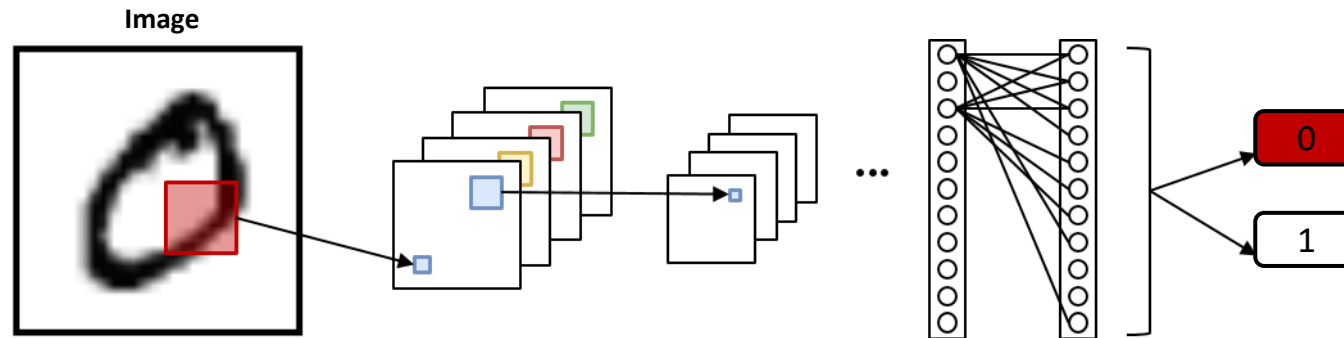
- Image of digit 0
- Image of digit 1

Image



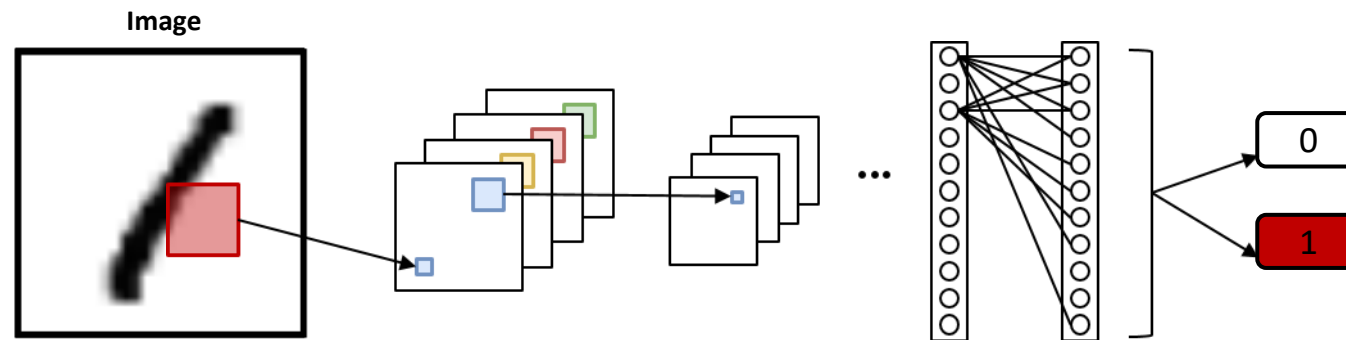
Deep Learning

- Convolutional Neural Networks (CNN)
- Image pattern recognition problems

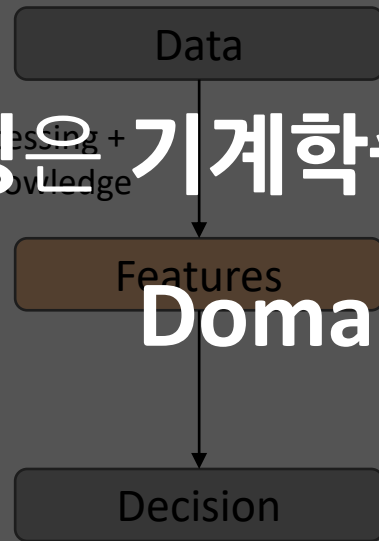


Deep Learning

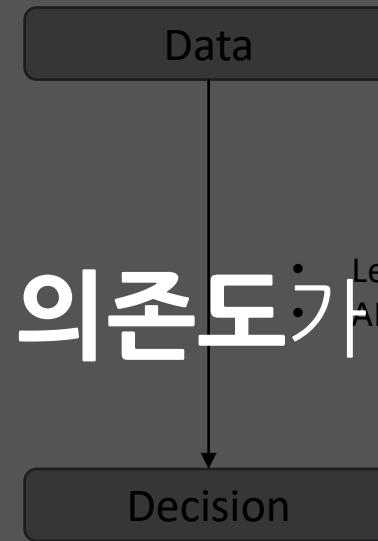
- Convolutional Neural Networks (CNN)
- Image pattern recognition problems



“딥러닝은 기계학습보다는



Domain Knowledge 의존도가 낮다.”



- Less depends on domain knowledge
- AI-discovered features