# Neural Networks Basics

**Data Intelligence and Learning (DIAL) Lab**

**Prof. Jongwuk Lee**
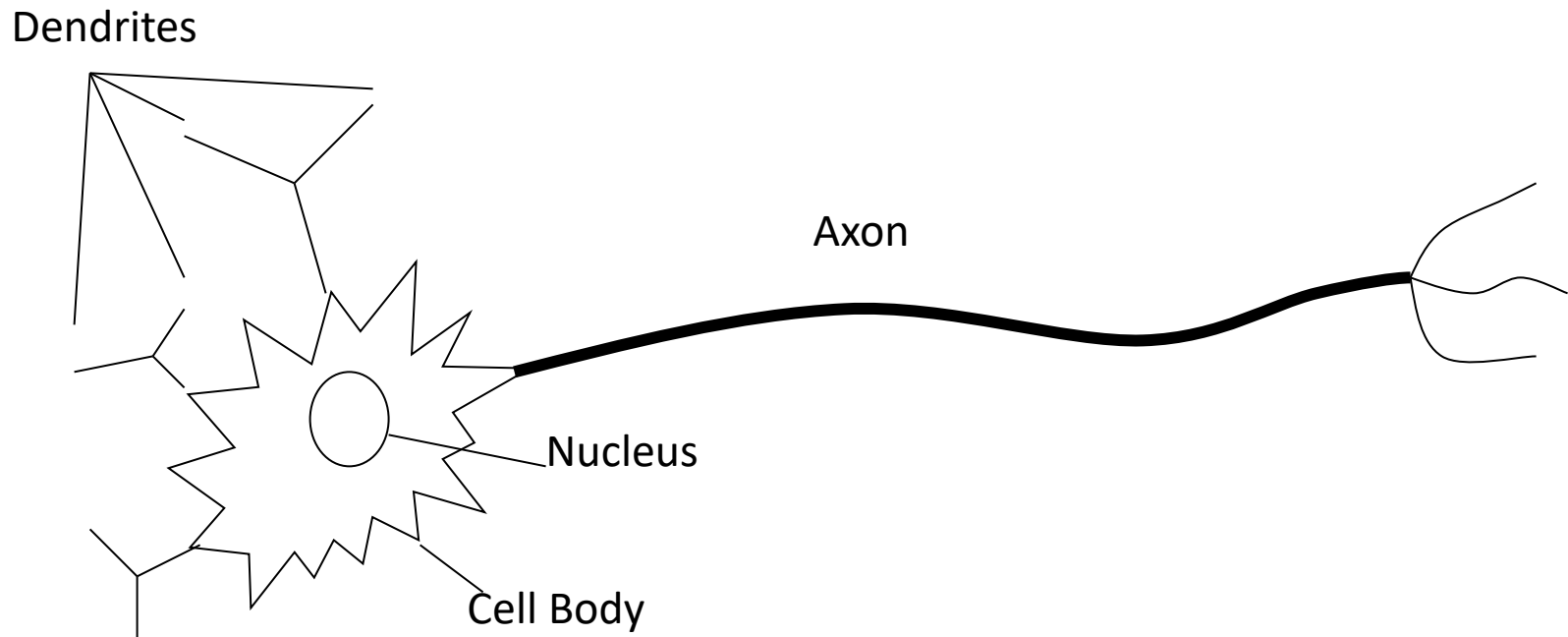
# Perceptron

# Brief History of Neural Networks



LSTM (1997) and LeNet (1998)

Deep Neural Network (Pretraining)

Multi-layered Perceptron (Backpropagation)

SVM

ADALINE

XOR Problem

AlexNet (2012)

Perceptron

Electronic Brain

Golden Age

Dark Age ("AI Winter")

| 1940 | 1950 | 1960 | 1970 | 1980 | 1990 | 2000 | 2010 |

1943    1957    1960    1969    1986    1995    2006

S. McCulloch – W. Pitts

F. Rosenblatt

B. Widrow – M. Hoff

M. Minsky – S. Papert

D. Rumelhart – G. Hinton – R. Wiliams

V. Vapnik – C. Cortes

G. Hinton – S. Ruslan

X AND Y    X OR Y    NOT X

- Adjustable Weights
- Weights are not Learned

- Learnable Weights and Threshold

- XOR Problem

Foward Activity

Backward Error

- Solution to nonlinearly separable problems
- Big computation, local optima and overfitting

- Limitations of learning prior knowledge
- Kernel function: Human Intervention

- Hierarchical feature Learning

# Concept of Neurons

➢ **Receive inputs from other neurons (via synapses).**

➢ **When input exceeds a threshold, "fires."**

  ◆ Sends output along axon to other neurons.
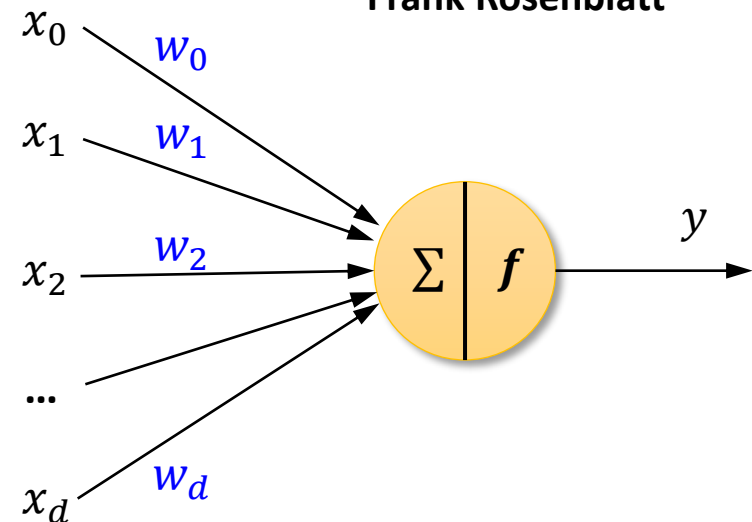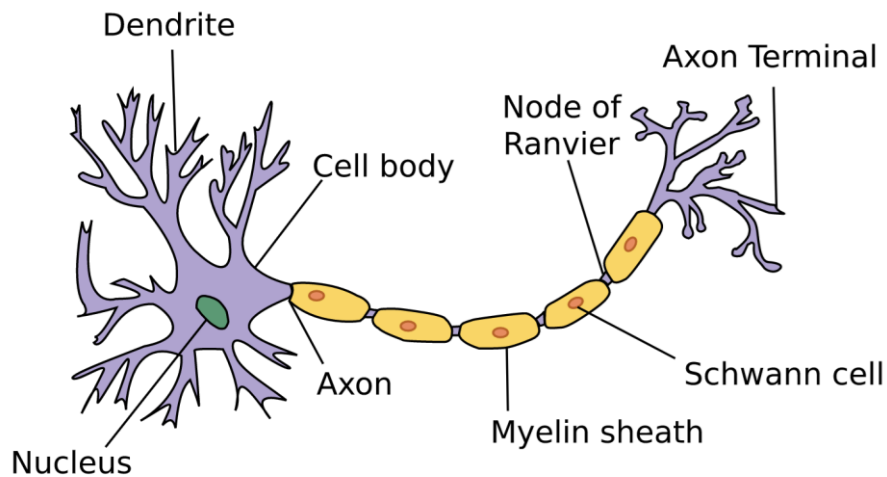
➢ **Human brain: $10^{11}$ neurons**

Dendrites

Axon

Nucleus

Cell Body

# Perceptron: Artificial Neuron (1957)

➢ A **neuron** is activated when the correlation between the input $\mathbf{x}$ and a pattern $\mathbf{w}$ exceeds a **threshold**.

- ◆ It is called an **artificial neuron**.
- ◆ It mimics the function of human **neurons**.

**Frank Rosenblatt**

# What is the Perceptron?

➢ **Linear combination of input x:**

$$s = w_0 + \sum_{i=1}^{d} w_i x_i$$

➢ **Nonlinear transformation of $s$:**

$$y = f(s), \text{where} \begin{cases} +1 & \text{if } s \geq 0 \\ -1 & \text{otherwise} \end{cases}$$
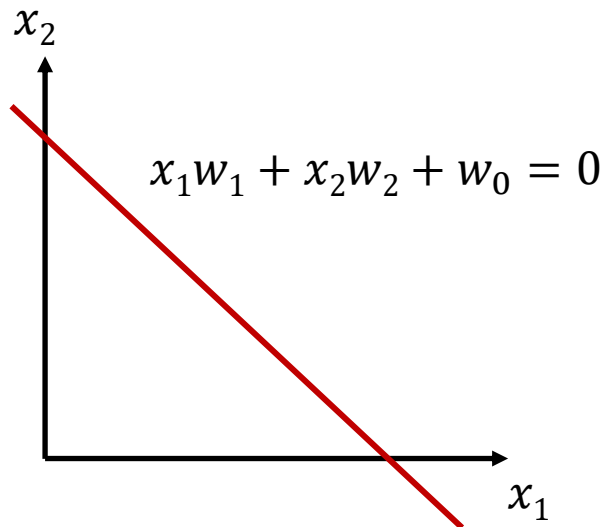


**Activation function**

$$1 \quad w_0$$
$$x_1 \quad w_1$$
$$x_2 \quad w_2$$
$$\cdots$$
$$x_d \quad w_d$$

$$\Sigma \longrightarrow \quad \longrightarrow y$$

$$s = w_0 + \sum_{i=1}^{d} w_i x_i \qquad f(s) = \begin{cases} +1 & if\ s \geq 0 \\ -1 & otherwise \end{cases}$$

# What is the Perceptron?

➢ **Linear combination of input x:**

$$s = \mathbf{w^T x} = \sum_{i=0}^{d} w_i x_i$$

➢ **Nonlinear transformation of $s$:**

$$x_1 w_1 + x_2 w_2 + w_0 = 0$$

**Hard limit**

+1

−1

$$f(s) = \begin{cases} +1 & if\ s \geq 0 \\ -1 & otherwise \end{cases}$$

# Formulating a Learning Classifier

➢ **We regard the output $y$ into $\{+1, -1\}$.**

➢ **Find $f(\mathbf{x}) = \mathbf{w}^\mathbf{T}\mathbf{x}$ that satisfies the condition.**
- ◆ $y = +1$ if $\mathbf{w}^\mathbf{T}\mathbf{x} \geq 0$
- ◆ $y = -1$ if $\mathbf{w}^\mathbf{T}\mathbf{x} < 0$

➢ **Prediction: $\hat{y} = \text{sign}(f(\mathbf{x})) = \text{sign}(\mathbf{w}^\mathbf{T}\mathbf{x})$**
- ◆ Note: $\text{sign}(\cdot)$ returns $+1$ or $-1$.

# Formulating a Learning Classifier

➢ **Given training data** $\left\{\left(\mathbf{x}^{(i)}, y^{(i)}\right): 1 \leq i \leq n\right\}$

➢ **Hypothesis function:** $\hat{y} = \text{sign}(\mathbf{w^T x})$

➢ **0-1 loss function: # of inputs such that** $y^{(i)} \neq \text{sign}(\mathbf{w^T x}^{(i)})$

$$E(\mathbf{w}) = -\sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} y^{(i)} \mathbf{w^T x}^{(i)} \mathbb{I}[\text{mistake on } \mathbf{x}^{(i)}]$$

$$\mathbb{I}[\text{mistake on } \mathbf{x}^{(i)}] = \begin{cases} 1 & \text{if } y^{(i)} \neq \text{sign}(\mathbf{w^T x}^{(i)}) \\ 0 & \text{otherwise} \end{cases}$$

# Formulating a Learning Classifier

➢ **Representing a simplified a 0-1 loss function**

$$E(\mathbf{w}) = -\sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} y^{(i)} \mathbf{w}^{\mathbf{T}} \mathbf{x}^{(i)} \mathbb{I}[mistake\ on\ \mathbf{x}^{(i)}]$$

$$E(\mathbf{w}) = \sum_{(\mathbf{x}^{(k)}, y^{(k)}) \in \mathcal{S}} -y^{(k)} \left( \mathbf{w}^{\mathbf{T}} \mathbf{x}^{(k)} \right)$$

$\mathcal{S}$ is a set of **incorrect samples**.

# Property of the 0-1 Loss Function

➢ **For any sample** $(\mathbf{x}^{(k)}, y^{(k)}) \in \mathcal{S}, -y^{(k)}(\mathbf{w}^{\mathbf{T}}\mathbf{x}^{(k)}) > 0.$

➢ $E(\mathbf{w}) \geq 0$

$$E(\mathbf{w}) = \sum_{(\mathbf{x}^{(k)}, y^{(k)}) \in \mathcal{S}} -y^{(k)}(\mathbf{w}^{\mathbf{T}}\mathbf{x}^{(k)})$$

➢ **If w is optimal, then** $E(\mathbf{w}) = 0.$
   - ◆ $E(\mathbf{w})$ tends to increase with the number of **incorrect samples**.

# Computing the Derivative of $\mathbf{w}$

➢ **How to compute the derivative of $w_j$**

$$\frac{\partial E(\mathbf{w})}{\partial w_j} = \sum_{(\mathbf{x}^{(k)}, y^{(k)}) \in \mathcal{S}} \frac{\partial}{\partial w_j} \left( -y^{(k)} \left( w_0 x_0^{(k)} + \cdots + w_j x_j^{(k)} + \cdots + w_d x_d^{(k)} \right) \right)$$

**Other values are regarded as constants.**

$$\frac{\partial E(\mathbf{w})}{\partial w_j} = \sum_{(\mathbf{x}^{(k)}, y^{(k)}) \in \mathcal{S}} -y^{(k)} x_j^{(k)}, \qquad \forall j = 0, 1, \dots, d$$

# Perceptron Learning Algorithm (PLA)

> **Batch version**

Initialize a random weight $\mathbf{w}$.
Repeat
    $\mathcal{S} \leftarrow \emptyset$

**Step 1: finding a set of incorrect samples**

for $i = 1$ to $n$
    $\hat{y} \leftarrow \text{sign}(\mathbf{w}^\mathbf{T}\mathbf{x}^{(i)})$
    $if\ (\hat{y} \neq y^{(i)})\ \mathcal{S} \leftarrow \mathcal{S} \cup (\mathbf{x}^{(i)}, y^{(i)})$

**Step 2: updating the weight w**

if $(\mathcal{S} \neq \emptyset)$
    $\Delta\mathbf{w} = -\sum_{(\mathbf{x}^{(k)}, y^{(k)}) \in \mathcal{S}} y^{(k)}\mathbf{x}^{(k)}$
    $\mathbf{w} \leftarrow \mathbf{w} - \eta\Delta\mathbf{w}$

Until $(\mathcal{S} = \emptyset)$

# Perceptron Learning Algorithm (PLA)

➢ **Stochastic version**

Initialize a random weight $\mathbf{w}$.
Repeat
  quit ← true
  for $i \leftarrow 1$ to $n$

**Step 1: finding an incorrect sample**

$$\hat{y} \leftarrow \text{sign}(\mathbf{w^T x}^{(i)})$$
$$\boldsymbol{if}\ (\hat{y} \neq y^{(i)})$$
   quit ← false

$$\Delta \mathbf{w} = -y^{(i)} \mathbf{x}^{(i)}$$
$$\mathbf{w} \leftarrow \mathbf{w} - \eta \Delta \mathbf{w}$$

**Step 2: updating the weight w**

Until (quit = true)

# Perceptron Learning Algorithm (PLA)

➢ **Start with a random weight vector.**

➢ **Given an input, predict its class.**

  ◆ For the mistake on a **positive sample**,

$$\mathbf{w}_{new} = \mathbf{w} - \eta \nabla E\big(\mathbf{w}; (\mathbf{x}^{(i)}, y^{(i)})\big) = \mathbf{w} + \eta \mathbf{x}^{(i)}$$

  ◆ For the mistake on a **negative sample**,

$$\mathbf{w}_{new} = \mathbf{w} - \eta \nabla E\big(\mathbf{w}; (\mathbf{x}^{(i)}, y^{(i)})\big) = \mathbf{w} - \eta \mathbf{x}^{(i)}$$

➢ **Repeat it until there is no mistake.**

# Perceptron Learning Algorithm (PLA)

➢ **Start with a random weight vector.**

➢ **Given an input, predict its class.**

  ◆ For the mistake on a **positive sample**,

$$\mathbf{w}_{new} = \mathbf{w} - \eta \nabla E\big(\mathbf{w}; (\mathbf{x}^{(i)}, y^{(i)})\big) = \mathbf{w} + \mathbf{1}\mathbf{x}^{(i)}$$

  ◆ For the mistake on a **negative sample**,

$$\mathbf{w}_{new} = \mathbf{w} - \eta \nabla E\big(\mathbf{w}; (\mathbf{x}^{(i)}, y^{(i)})\big) = \mathbf{w} - \mathbf{1}\mathbf{x}^{(i)}$$

When $\eta = 1$

➢ **Repeat it until there is no mistake.**

# Recap: Learning a Linear Classifier

➢ **Execute the algorithm until no mistakes are encountered.**

Randomly choose an initial solution $\mathbf{w^0}$.

For $t = 0, 1, \dots$

    Find a mistake sample $(\mathbf{x}^{(i)}, y^{(i)})$ of $\mathbf{w^t}$
        $\text{sign}(\mathbf{w^T x}^{(i)}) \neq y^{(i)}$

    Correct the mistake by
        $\mathbf{w^{t+1}} = \mathbf{w^t} + y^{(i)}\mathbf{x}^{(i)}$

Until no more mistake is found

Return last $\mathbf{w^t}$ as the learned model.

If $y = +1$    $\mathbf{w} + y\mathbf{x}$   $\mathbf{x}$   $\mathbf{w}$

If $y = -1$    $\mathbf{w}$   $\mathbf{x}$   $\mathbf{w} + y\mathbf{x}$   $-\mathbf{x}$

# AND Operation

> **How to represent the AND operation**

$$w_1 = 1.0, w_2 = 1.0, w_0 = -1.5$$

| $x_1$ | $x_2$ | $\sum$ | $y$ |
|-------|-------|--------|-----|
| 0 | 0 | $-1.5$ | 0 |
| 0 | 1 | $-0.5$ | 0 |
| 1 | 0 | $-0.5$ | 0 |
| 1 | 1 | $+0.5$ | 1 |

**1**

**1**

$$1.0x_1 + 1.0x_2 - 1.5 = 0$$

# OR Operation

> **How to represent the OR operation**

$$w_1 = 1.0, w_2 = 1.0, w_0 = -0.5$$

| $x_1$ | $x_2$ | $\sum$ | $y$ |
|---|---|---|---|
| 0 | 0 | $-0.5$ | 0 |
| 0 | 1 | 0.5 | 1 |
| 1 | 0 | 0.5 | 1 |
| 1 | 1 | 1.5 | 1 |

$$1.0x_1 + 1.0x_2 - 0.5 = 0$$

# Linear Separability

➢ **If PLA halts (i.e., no more mistakes),**

◆ (**necessary condition**) $D$ allows some **w** to make no mistake.

➢ **Call such $D$ linearly separable.**

| | | |
|---|---|---|
| **Linear separable** | **Linear non-separable** | **Linear non-separable** |
| **Good!** | **Need a linear model that allows some errors.** | **Need a non-linear model.** |

# XOR Operation (1969)

➤ **How to train a linear decision boundary?**



**"No one on earth had found a viable way to train" by Marvin Minsky**

➤ **Guaranteed to converge if linearly separable.**

➤ **Otherwise, many simple functions are NOT learnable.**

# Multilayer Perceptron (MLP)

# Multilayer Perceptron (MLP)

➢ **It is a neural network of multiple artificial neurons.**



$x_1$  $x_2$  ...  $x_d$  1

$y_1$  $y_2$  ...  $y_k$

1

**Input layer**        **Hidden layer**        **Output layer**

# Multilayer Perceptron (MLP)

➢ **Structure**

◆ **Input layer**

• Simply pass the input values to the next layer.

• # of nodes = # of inputs

◆ **Hidden layer**

• There can be several hidden layers.

• # of nodes should be given.

◆ **Output layer**

• # of nodes = # of outputs

➢ **MLP is also called the feed-forward neural network.**

# Multilayer Perceptron (MLP)



$$h_1 = \frac{1}{1 + \exp\left(-(w_{11}x_1 + w_{12}x_2 + w_{10})\right)}$$

$$h_2 = \frac{1}{1 + \exp\left(-(w_{21}x_1 + w_{22}x_2 + w_{20})\right)}$$

$$y = \frac{1}{1 + \exp\left(-(w_{31}h_1 + w_{32}h_2 + w_{30})\right)}$$

$$y = \frac{1}{1 + \exp\left(-\left(w_{31}\left(\underbrace{\frac{1}{1 + \exp(-(w_{11}x_1 + w_{12}x_2 + w_{10}))}}_{h_1}\right) + w_{32}\left(\underbrace{\frac{1}{1 + \exp(-(w_{21}x_1 + w_{22}x_2 + w_{20}))}}_{h_2}\right) + w_{30}\right)\right)}$$

# How to Solve XOR Operation?

$$w_{11} = 1.0, w_{12} = 1.0, w_{10} = -1.5 \qquad w_{21} = 1.0, w_{22} = 1.0, w_{20} = -0.5 \qquad w_{31} = -1.0, w_{32} = 1.0, w_{30} = -0.5$$

$$h_1 = \text{sigmoid}(x_1 + x_2 - 1.5) \qquad h_2 = \text{sigmoid}(x_1 + x_2 - 0.5) \qquad \hat{y} = \text{sigmoid}(-x_1 + x_2 - 0.5)$$
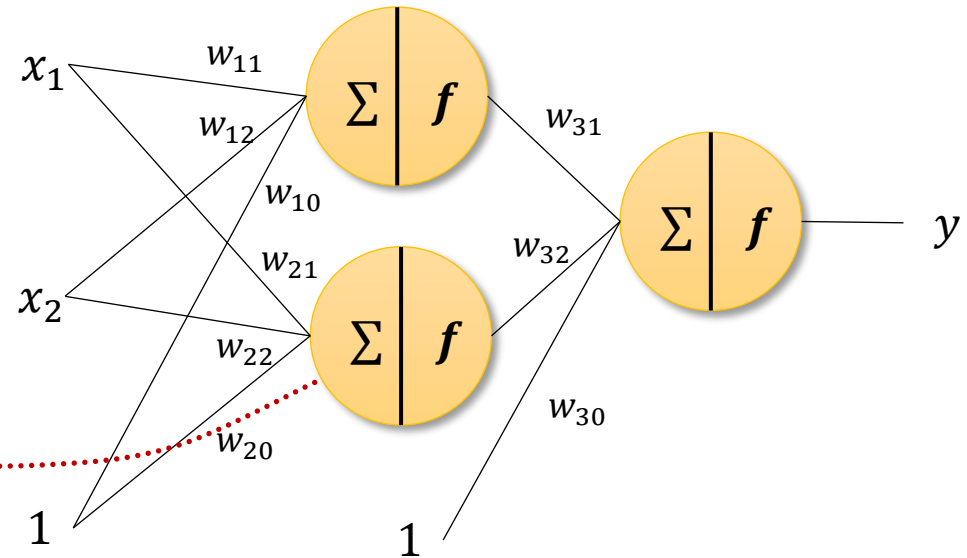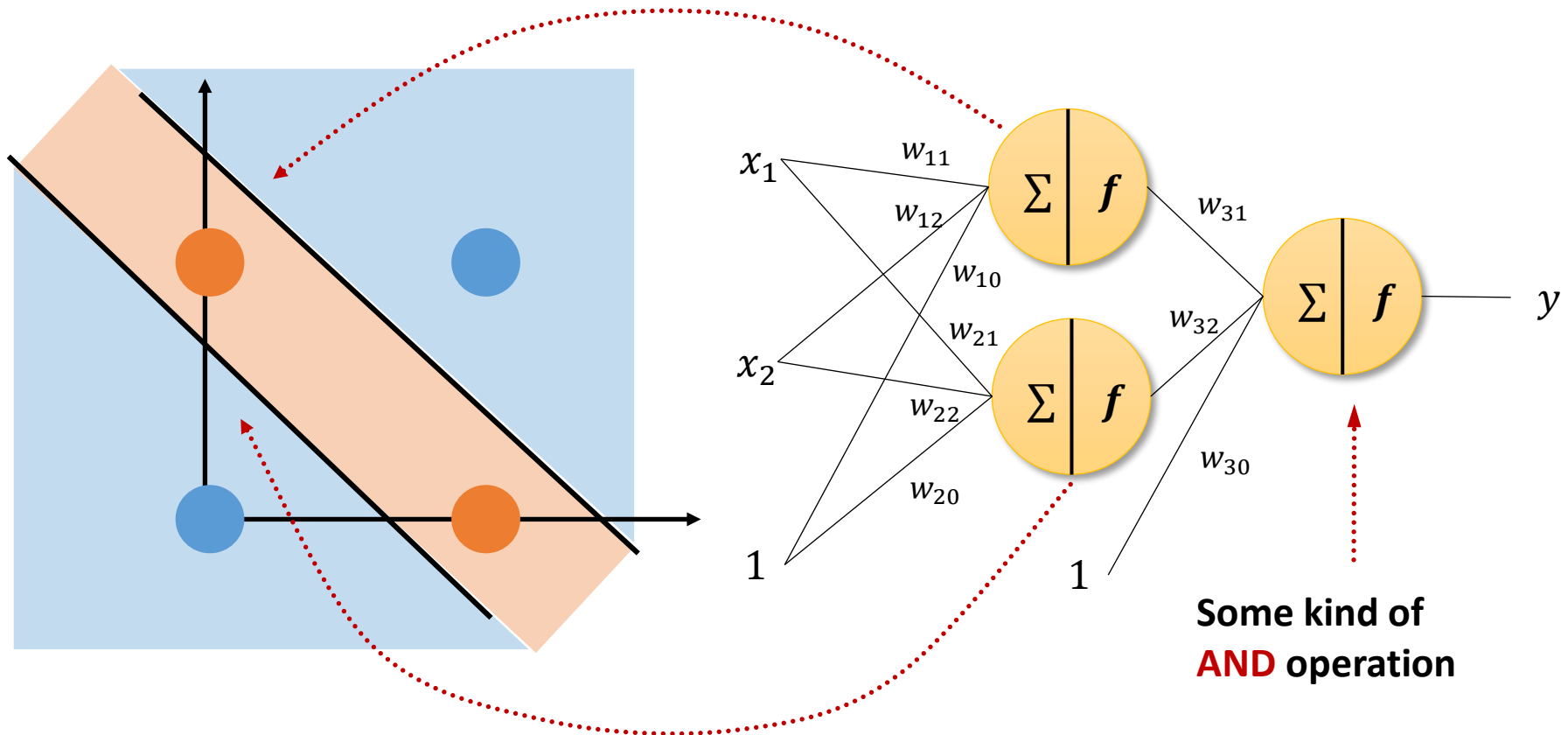
| $x_1$ | $x_2$ | $\sum$ | $h_1$ |
|-------|-------|--------|-------|
| 0 | 0 | -1.5 | 0 |
| 0 | 1 | -0.5 | 0 |
| 1 | 0 | -0.5 | 0 |
| 1 | 1 | +0.5 | 1 |

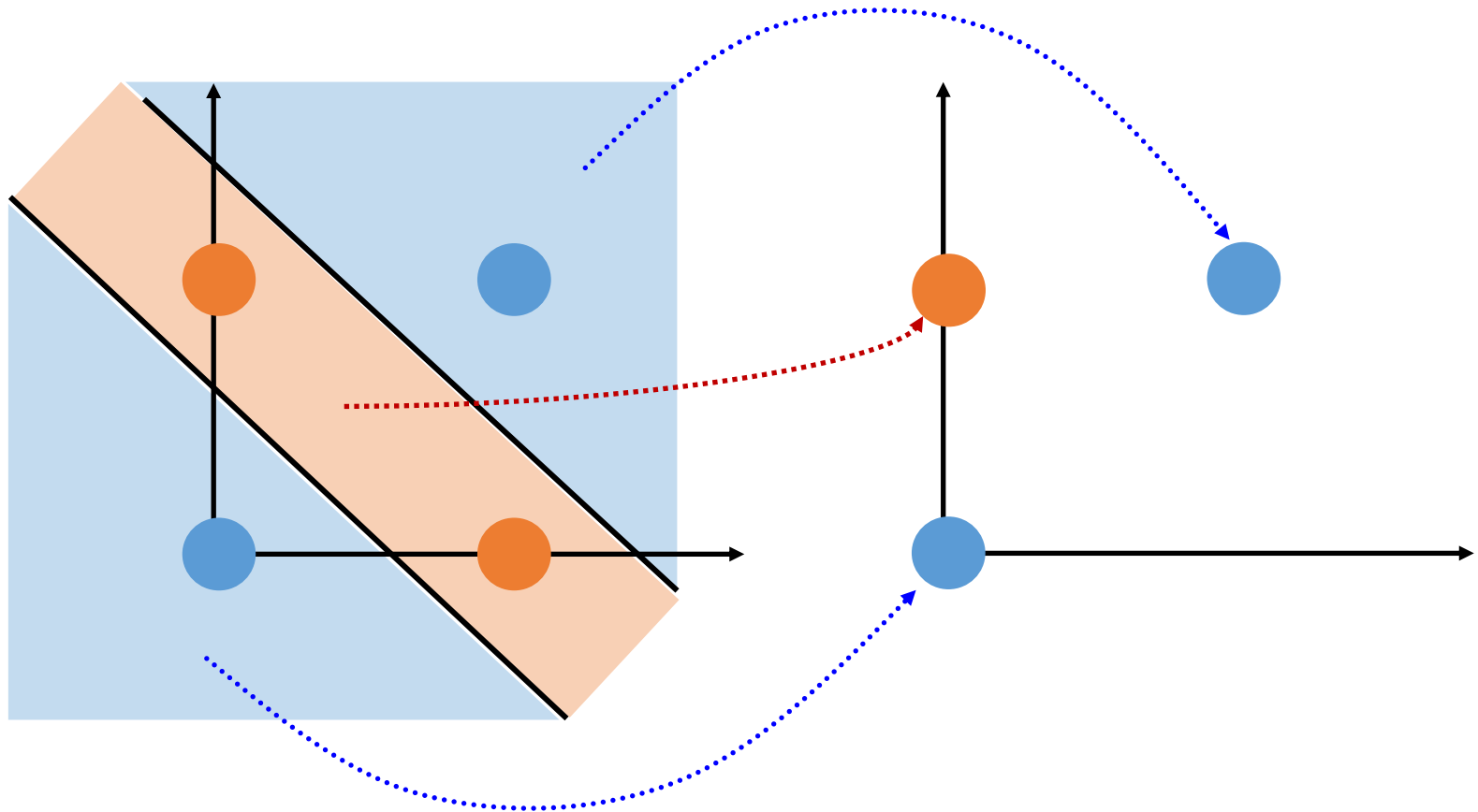| $x_1$ | $x_2$ | $\sum$ | $h_2$ |
|-------|-------|--------|-------|
| 0 | 0 | -0.5 | 0 |
| 0 | 1 | 0.5 | 1 |
| 1 | 0 | 0.5 | 1 |
| 1 | 1 | 1.5 | 1 |

| $h_1$ | $h_2$ | $\sum$ | $\hat{y}$ |
|-------|-------|--------|-----------|
| 0 | 0 | -0.5 | 0 |
| 0 | 1 | 0.5 | 1 |
| 0 | 1 | 0.5 | 1 |
| 1 | 1 | -0.5 | 0 |

# How to Solve XOR Operation?

➢ **A single neuron divides a region into two subregions.**



| $x_1$ | $x_2$ | $\Sigma$ | $h_1$ |
|:---:|:---:|:---:|:---:|
| 0 | 0 | -1.5 | 0 |
| 0 | 1 | -0.5 | 0 |
| 1 | 0 | -0.5 | 0 |
| 1 | 1 | +0.5 | 1 |

# How to Solve XOR Operation?

➢ **A single neuron divides a region into two subregions.**



| $x_1$ | $x_2$ | $\Sigma$ | $h_2$ |
|-------|-------|----------|-------|
| 0 | 0 | -1.5 | 0 |
| 0 | 1 | 0.5 | 1 |
| 1 | 0 | 0.5 | 1 |
| 1 | 1 | 1.5 | 1 |

# How to Solve XOR Operation?

> **Combining two neurons with AND operation**



Some kind of **AND** operation

# How to Solve XOR Operation?



**Mapping the regions in the original space
to the points in another space**

# Non-linear Classifier using MLP

➢ **How to classify two classes?**

# Non-linear Classifier using MLP

> **Hyperspace partitioning**

# Non-linear Classifier using MLP

➢ **Use three nodes in the hidden layer.**

# Example: Visualization of MLP

https://playground.tensorflow.org/

# Deep Neural Networks (DNNs)

# Recap: Machine Learning 1-2-3

- Collect data and extract features.

- Choose the **hypothesis function** $\mathcal{H}$ and the **loss function** $\mathcal{L}$.

- Find an optimal parameter that minimizes the empirical loss.

- Q: How to represent feature vectors for images?

- A: It is difficult to design the vector.

# Simplest Case: Linear Classifiers

➢ **Consider the original perceptron model on raw data.**

➢ **To classify the image of "dog" (+1) or "not dog" (-1), find the line that separates the +1's from the -1's.**

# Non-Linear Features, Linear Classifiers

➢ **Most problems need non-linear features.**

   ◆ Image classification, machine translation, speech recognition, etc.

➢ **How to represent non-linear features?**
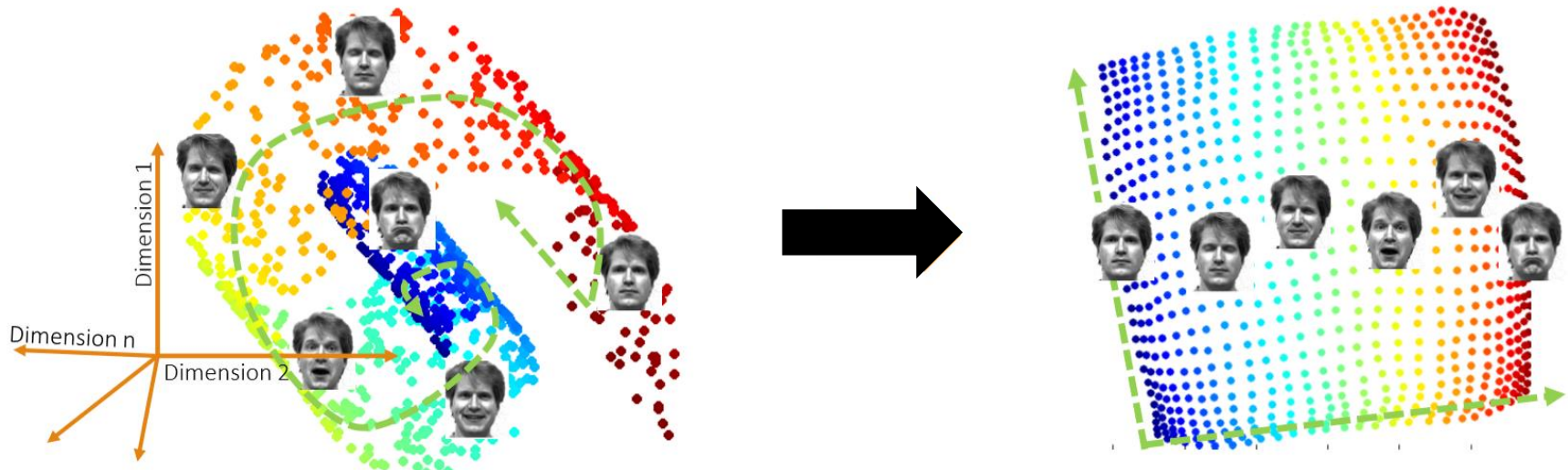
   ◆ Non-linear kernels in SVM
   ◆ Explicit design of features, e.g., SIFT, HOG



➢ **Then, linear machines are good enough.**
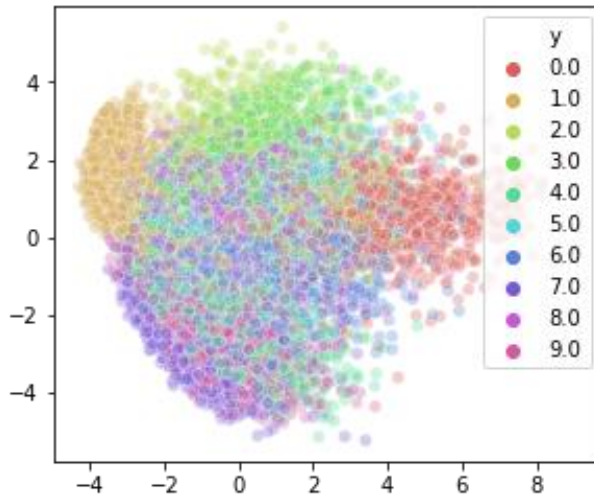
# How to Represent Good Features?

➢ **Raw data live high dimensionalities. However, data lie in low-dimensional manifolds.**

➢ **Can we discover a manifold for our data?**

  ◆ Hypothesis: **Semantically similar things** lie **closer** together than **semantically dissimilar things**.
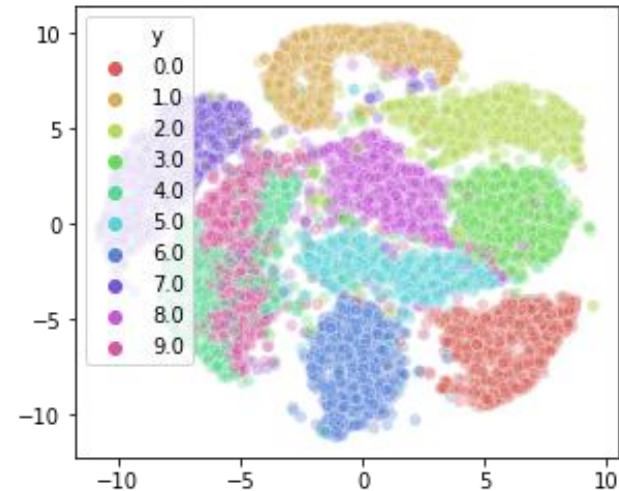
# Example: Manifolds of Digits

➢ **There are good features (manifolds).**

➢ **Each image has 784 dimensions on the MNIST dataset.**
   ◆ 28 pixels x 28 pixels = 784 dimensions



**Entangled features
(PCA)**

**Disentangled features
(t-SNE)**

# Why Learn Features?

➢ **Manually designed features**

- ◆ Often, take **a lot of time to** implement
- ◆ Often, take a lot of time to validate
- ◆ Often, they are **incomplete**, as one cannot know if they are optimal for the task

➢ **Learned features**

- ◆ If data is enough, easy to learn
- ◆ Compact and specific **to the task**

➢ **Time spent for designing features vs. Time spent on designing architectures**

# Shallow Learning vs. Deep Learning

**Artificial intelligence**

**Hard-coded knowledge & the rule-driven approach is used.**

**Machine learning**

**Extract the patterns & co-relation between features.**
It is difficult to know whether all features are effective and necessary.

**Representation learning**

**Learns the representation instead of the features.**
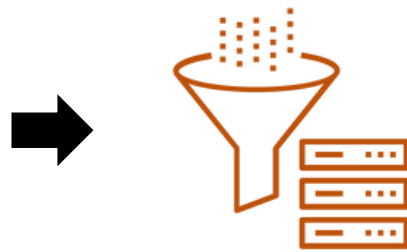There are many factors of variations for a given representation.

**Complex representations are expressed in terms of simple representations.**

**Deep learning**

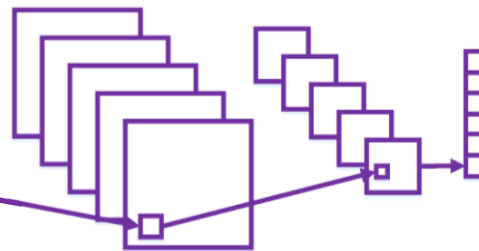# Shallow Learning vs. Deep Learning

> **Traditional machine learning**



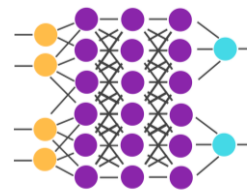**Hand-crafted feature extractor**     **Trainable classifier**

Dog
Not Dog

> **End-to-end learning → features are also learned from data!**



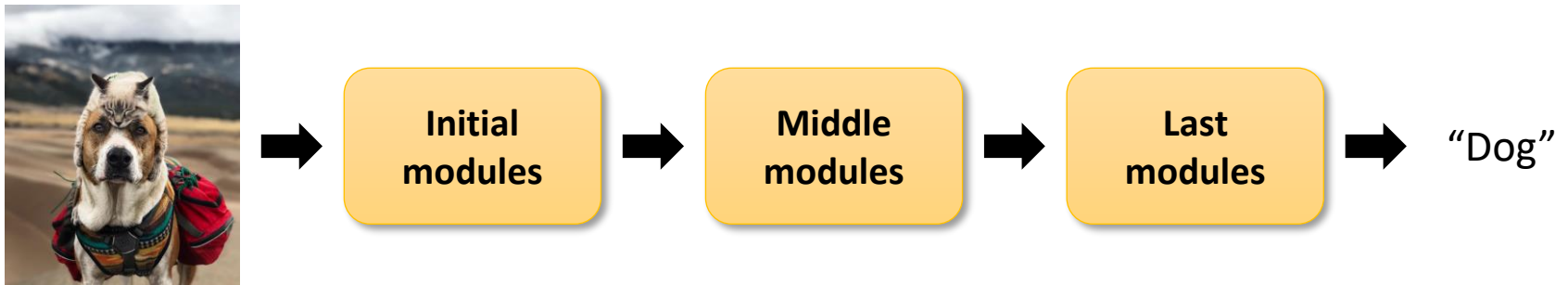**Trainable feature extractor**     **Trainable classifier**

Dog
Not Dog

# Deep Neural Networks (DNNs)

- **It consists of a pipeline of successive, differentiable modules (transformations).**
    - Each module's output is used as the input for the next module.

- **Learns hierarchical representation from data.**
- **Each sequential module produces a higher abstraction feature.**



Image → **Initial modules** → **Middle modules** → **Last modules** → "Dog"

# Deep Neural Networks (DNNs)

> **A family of parametric, non-linear, and hierarchical representation learning functions**

  - ◆ They are optimized with **stochastic gradient descent** to **encode domain knowledge**, i.e., domain invariances and stationarity.

$$a_L(x; \theta_{1,\ldots,L}) = h_L(h_{L-1}(\ldots h_1(x, \theta_1), \theta_{L-1}), \theta_L)$$
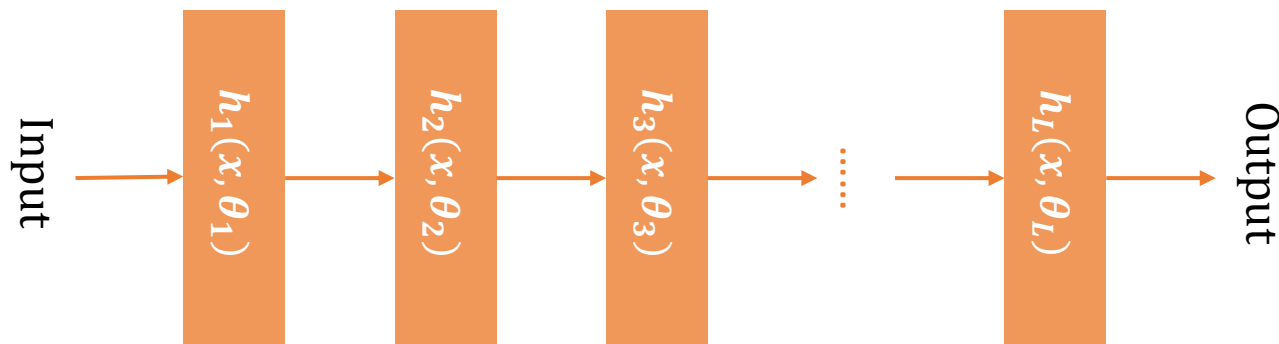
> **Given training corpus $\{X, Y\}$, find optimal parameters.**

$$\theta^* \leftarrow \underset{\theta}{\operatorname{argmin}} \sum_{(x,y) \subseteq (X,Y)} \mathcal{L}\left(y, a_L(x; \theta_{1,\ldots,L})\right)$$

# Deep Neural Networks (DNNs)

➢ A series of **hierarchically** connected functions

➢ This hierarchy can be very **complex**!

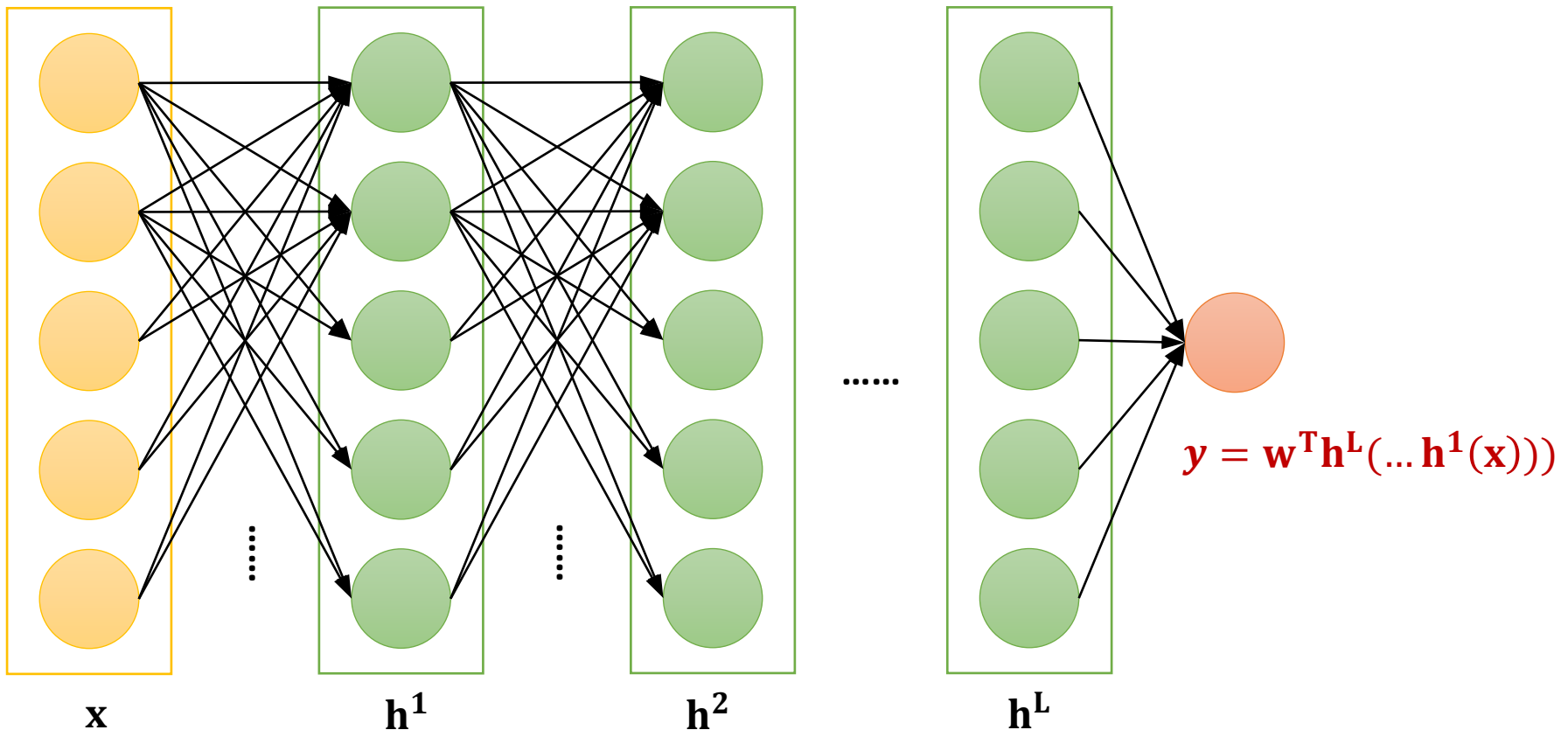$$a_L(x; \theta_{1,\dots,L}) = h_L(h_{L-1}(\dots h_1(x, \theta_1), \theta_{L-1}), \theta_L)$$



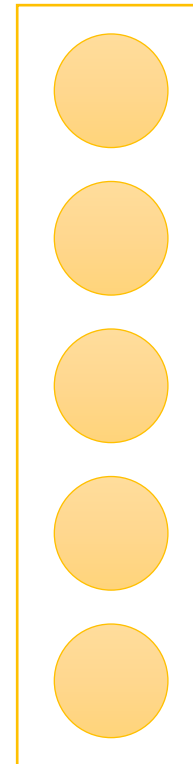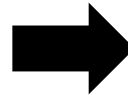*Feedforward architecture*

# Deep Neural Networks (DNNs)

➢ **What if we go deeper?**



$$y = \mathbf{w}^{\mathbf{T}}\mathbf{h}^{\mathbf{L}}(\dots \mathbf{h}^{\mathbf{1}}(\mathbf{x}))$$

$\mathbf{x}$  $\mathbf{h}^1$  $\mathbf{h}^2$  $\mathbf{h}^L$

# Input Layer of DNNs

➢ **The input is represented as a vector.**

◆ Sometimes, it requires preprocessing, e.g., **normalization**



**x**

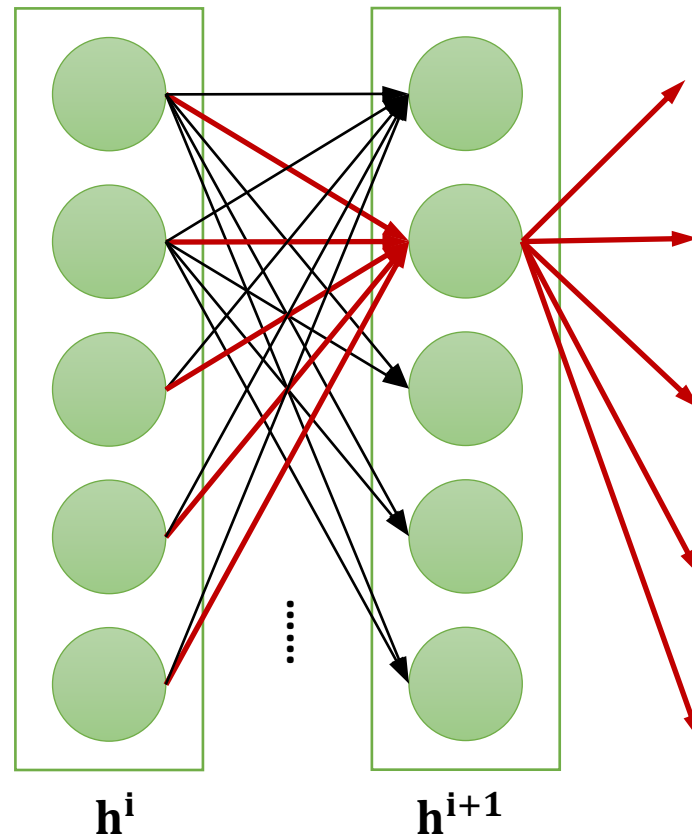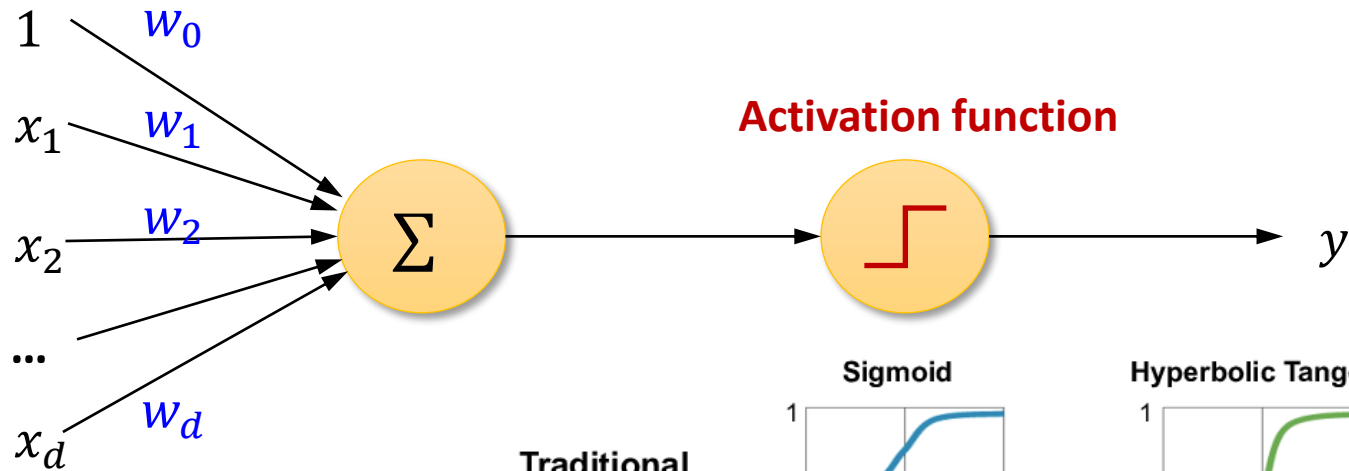# Hidden Layers of DNNs

> **Each neuron takes a weighted linear combination of the previous layer.**
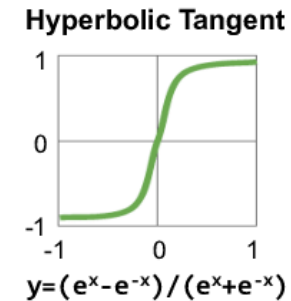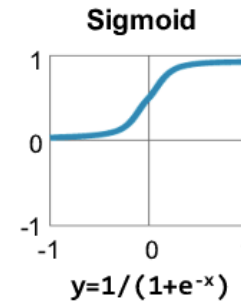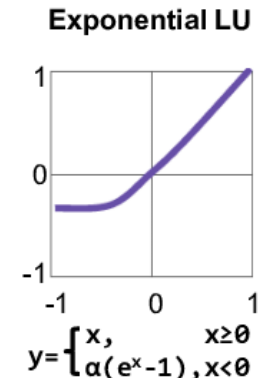>> ◆ Provide a single **aggregated** value to the next layer.



$$\mathbf{h^i} \qquad\qquad \mathbf{h^{i+1}}$$

# Nodes at the Hidden Layers

$1$

$w_0$

$x_1$

$w_1$

$x_2$

$w_2$

$\Sigma$

...

$x_d$

$w_d$

**Activation function**

$y$

There are various activation functions.

**Traditional Non-Linear Activation Functions**

**Sigmoid**

$y=1/(1+e^{-x})$

**Hyperbolic Tangent**

$y=(e^x-e^{-x})/(e^x+e^{-x})$

**Modern Non-Linear Activation Functions**

**Rectified Linear Unit (ReLU)**

$y=\max(0,x)$

**Leaky ReLU**

$y=\max(\alpha x, x)$

**Exponential LU**

$y=\begin{cases} x, & x\geq 0 \\ \alpha(e^x-1), & x<0 \end{cases}$
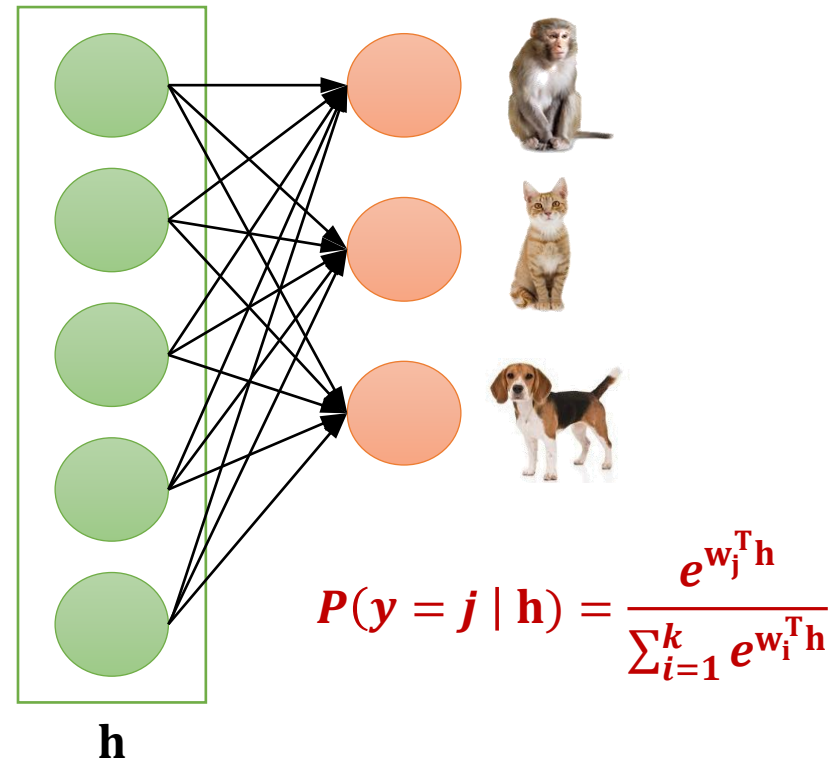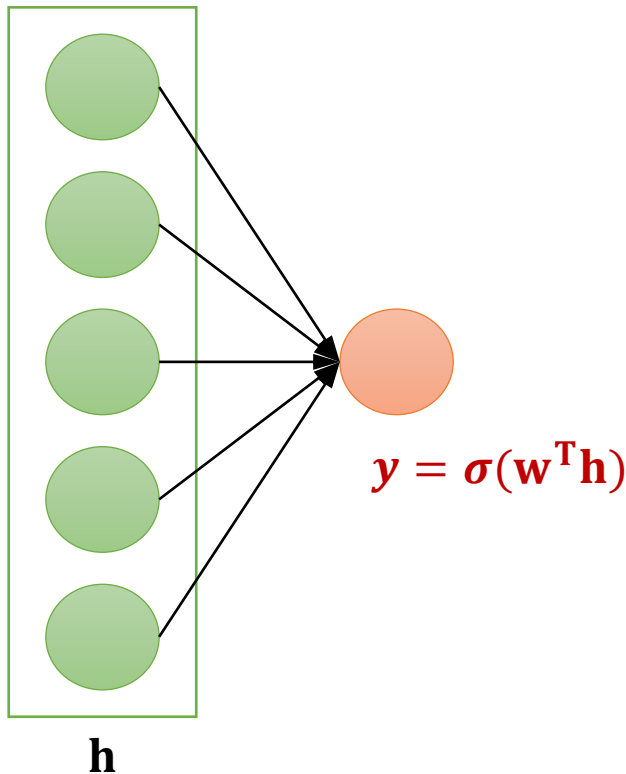
$\alpha$ = small const. (e.g. 0.1)

50

# Output Layer of DNNs

➢ **For output, there are three cases:**
- ◆ Regression: $y = \mathrm{w^T h}$
- ◆ Binary classification: $y = \sigma(\mathrm{w^T h})$
- ◆ Multi-class classification: $y = \mathrm{softmax}(\mathrm{w^T h})$



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

$$P(y = j \mid \mathbf{h}) = \frac{e^{\mathbf{w_j^T h}}}{\sum_{i=1}^{k} e^{\mathbf{w_i^T h}}}$$

$\mathbf{h}$

# Q&A