



IF3270 Pembelajaran Mesin

Artificial Neural Network (ANN)

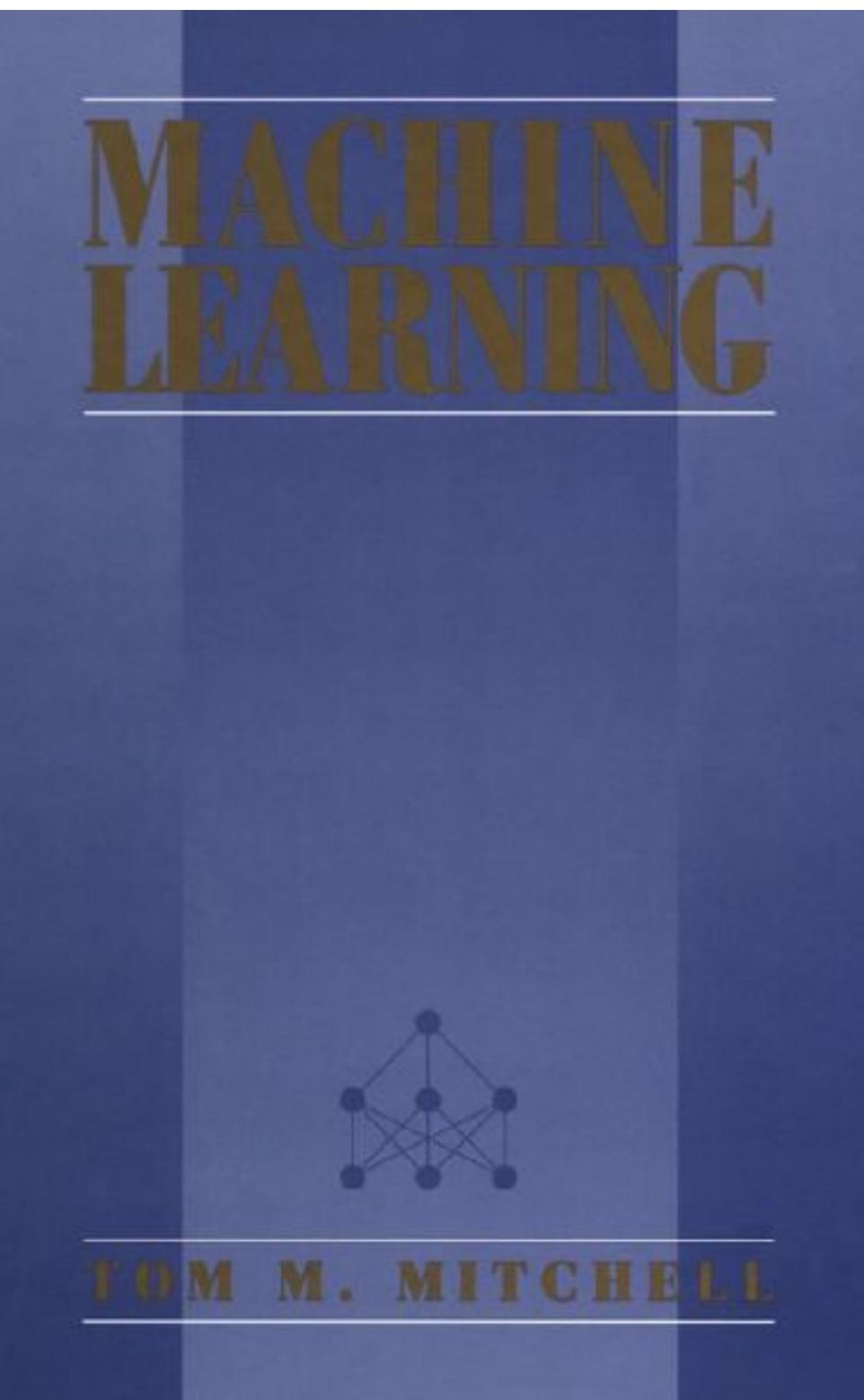
Tim Pengajar IF3270



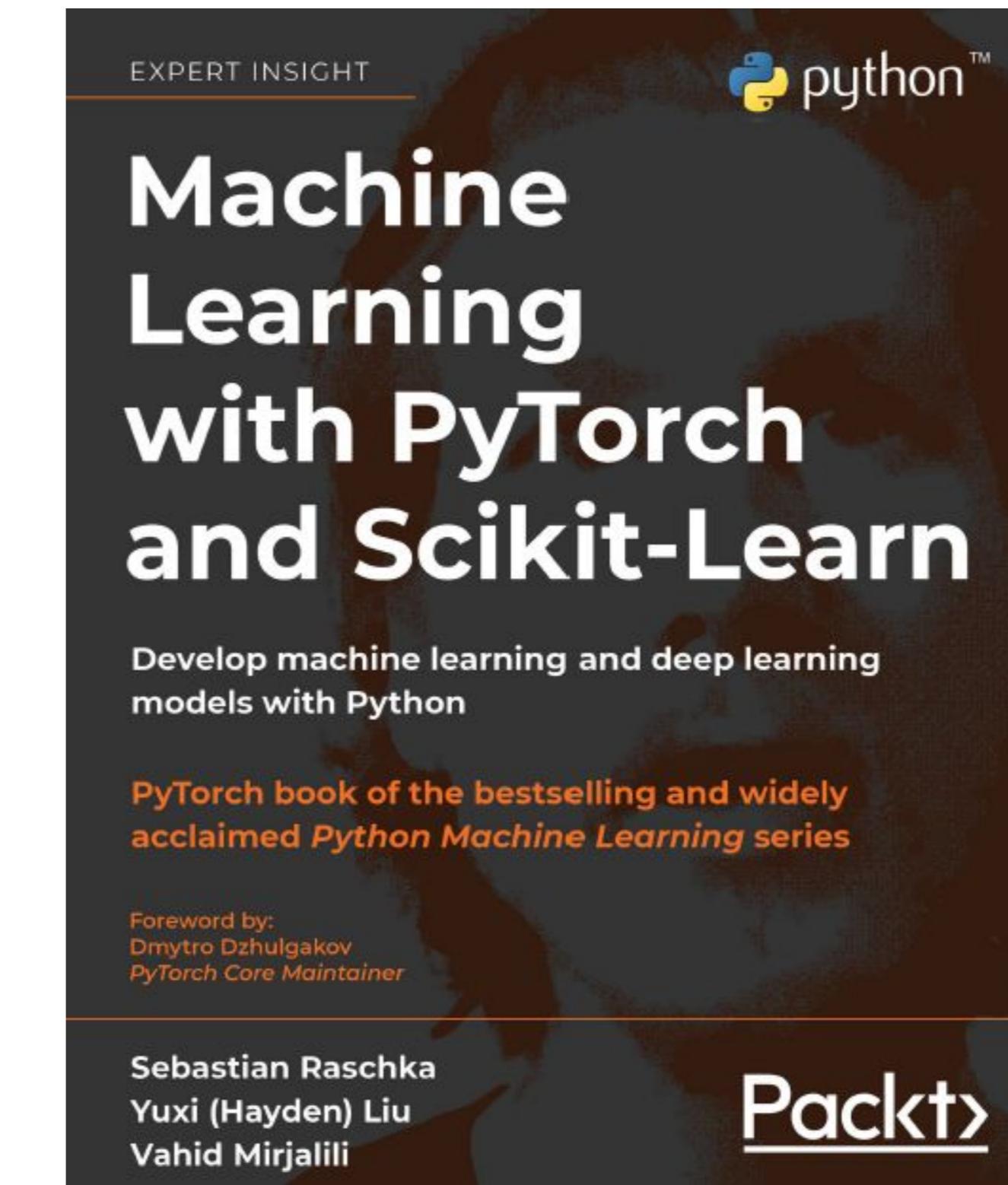
Review

- Machine learning (ML) overview
 - AI vs ML
 - Learning Type:
 - Supervised Learning: Classification, Regression
 - Unsupervised Learning
 - Reinforcement Learning
- Ensemble Methods □ Supervised Learning
 - Homogeneous Parallel Ensemble: Bagging, RF
 - Heterogeneous Parallel Ensemble: Weighting, Meta Learning (Stacking)
 - Sequential Ensemble: Adaboost, Gradient Boosting
- Supervised Learning: Perceptron
- Supervised Learning: ANN: Feed Forward Neural Network □ Today

Reference



Tom Mitchell,
Machine Learning,
McGraw Hill, 1997



Raschka, et.al., Machine Learning
with Pytorch and Scikit-Learn,
Packt Publishing Ltd., 2022
(Chapter 11)

Outline

Introduction
to ANN

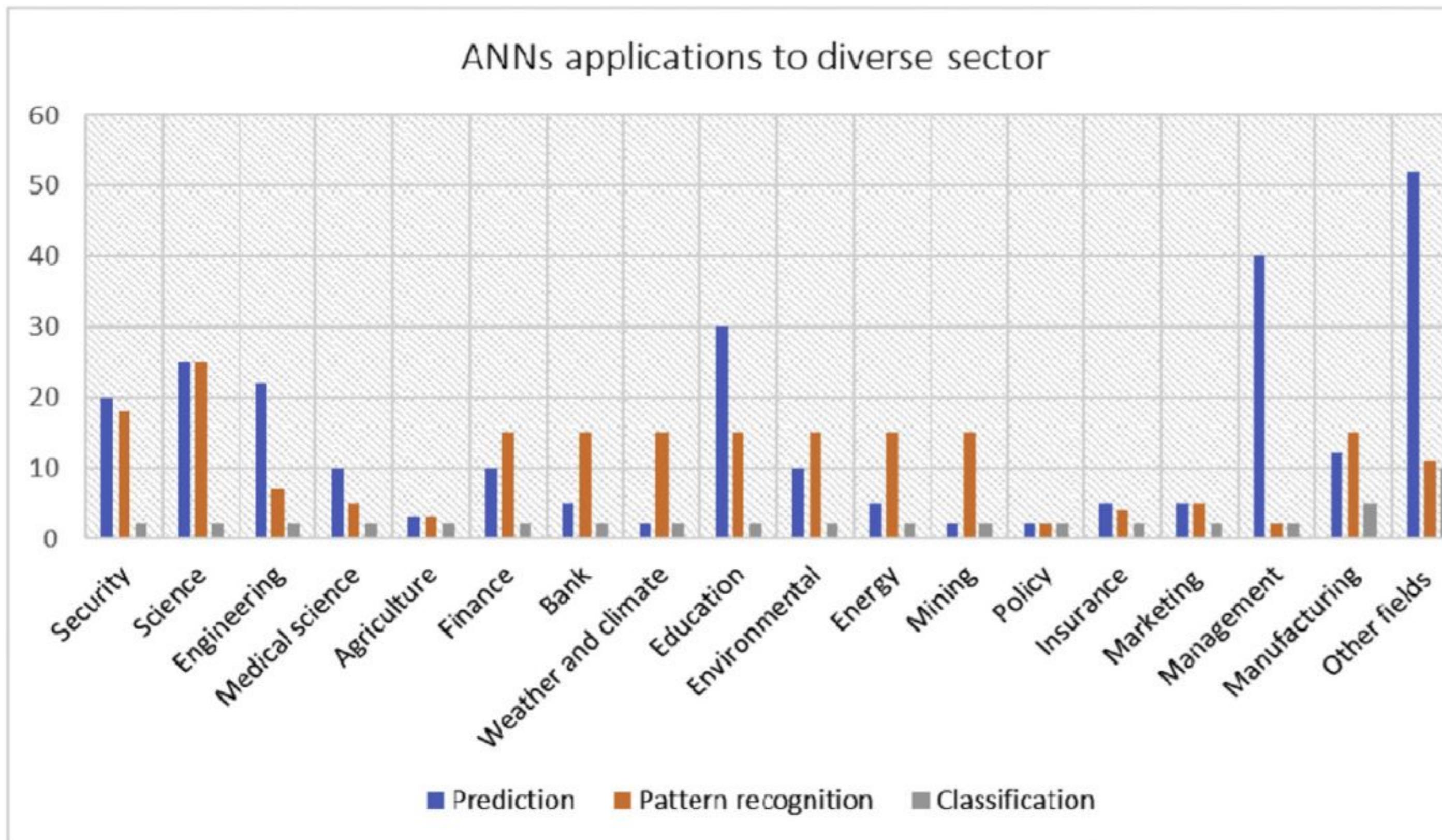
Feed Forward
NN (Forward
Propagation)

Feed Forward NN
(Backpropagation)

Introduction to ANN

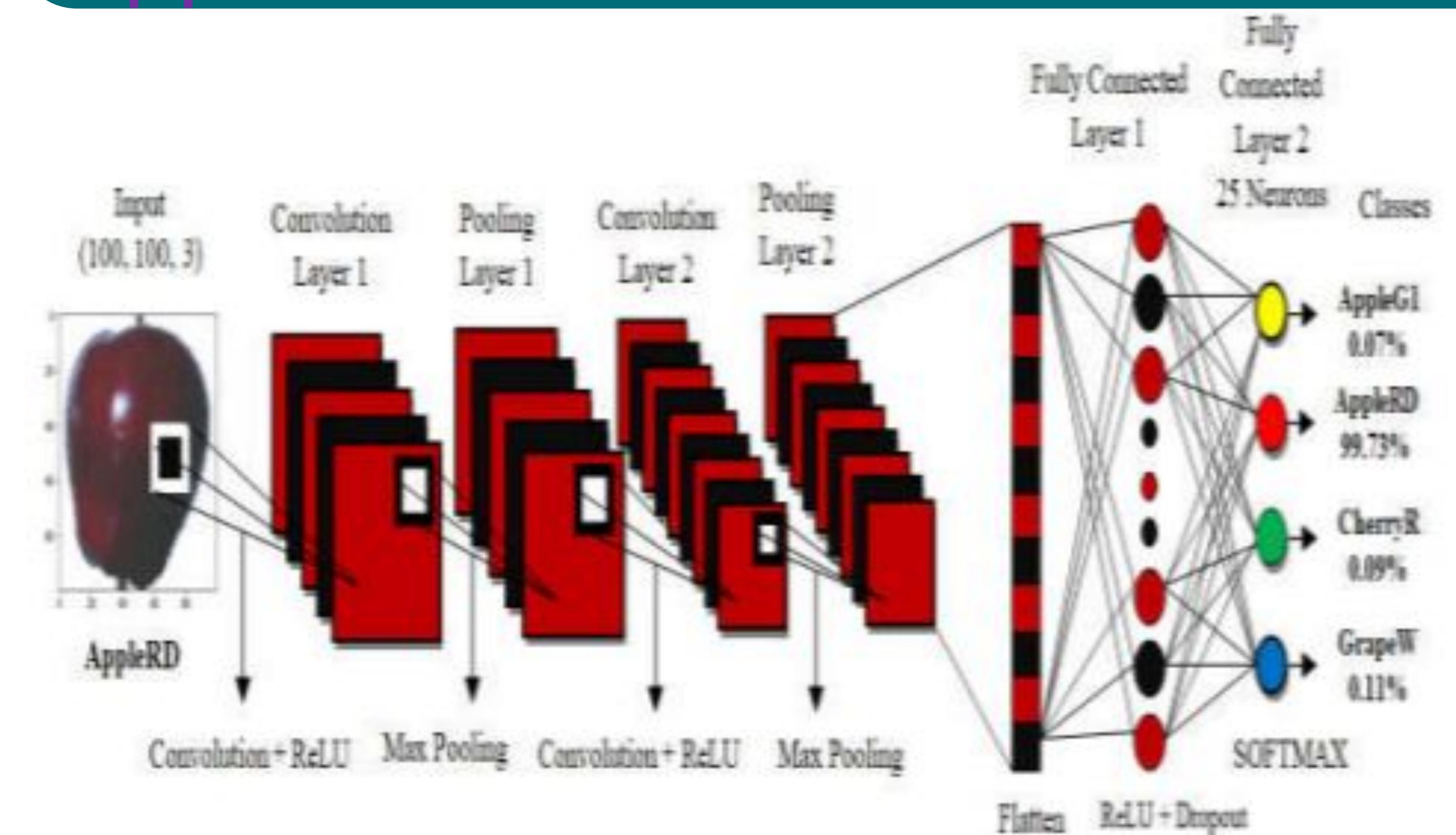
Artificial neural network: Why

ANN apps to diverse sectors,
and perform well



Abiodun, O. I., Jantan, A., Omolara, A. E., Dada, K. V., Mohamed, N. A., & Arshad, H. (2018). State-of-the-art in artificial neural network applications: A survey. *Helicon*, 4(11), e00938.

Rapid increase in ANN
architectural innovations and
application



Sakib, S., Ashrafi, Z., Siddique, M., & Bakr, A. (2019). Implementation of Fruits Recognition Classifier using Convolutional Neural Network Algorithm for Observation of Accuracies for Various Hidden Layers. <https://arxiv.org/ftp/arxiv/papers/1904/1904.00783.pdf>

Appropriate Problems for ANN Learning

- Instances are represented by many attribute-value pairs.
- Target function output may be discrete-valued, real-valued, or a vector of several real- or discrete-valued attributes
- **Training examples may contain errors**
- Long training times are acceptable
- **Fast evaluation of the learned target function may be required**
- Ability of humans to understand the learned target function is not important

Artificial neural network: What

Perceptron (incl. Adaline)

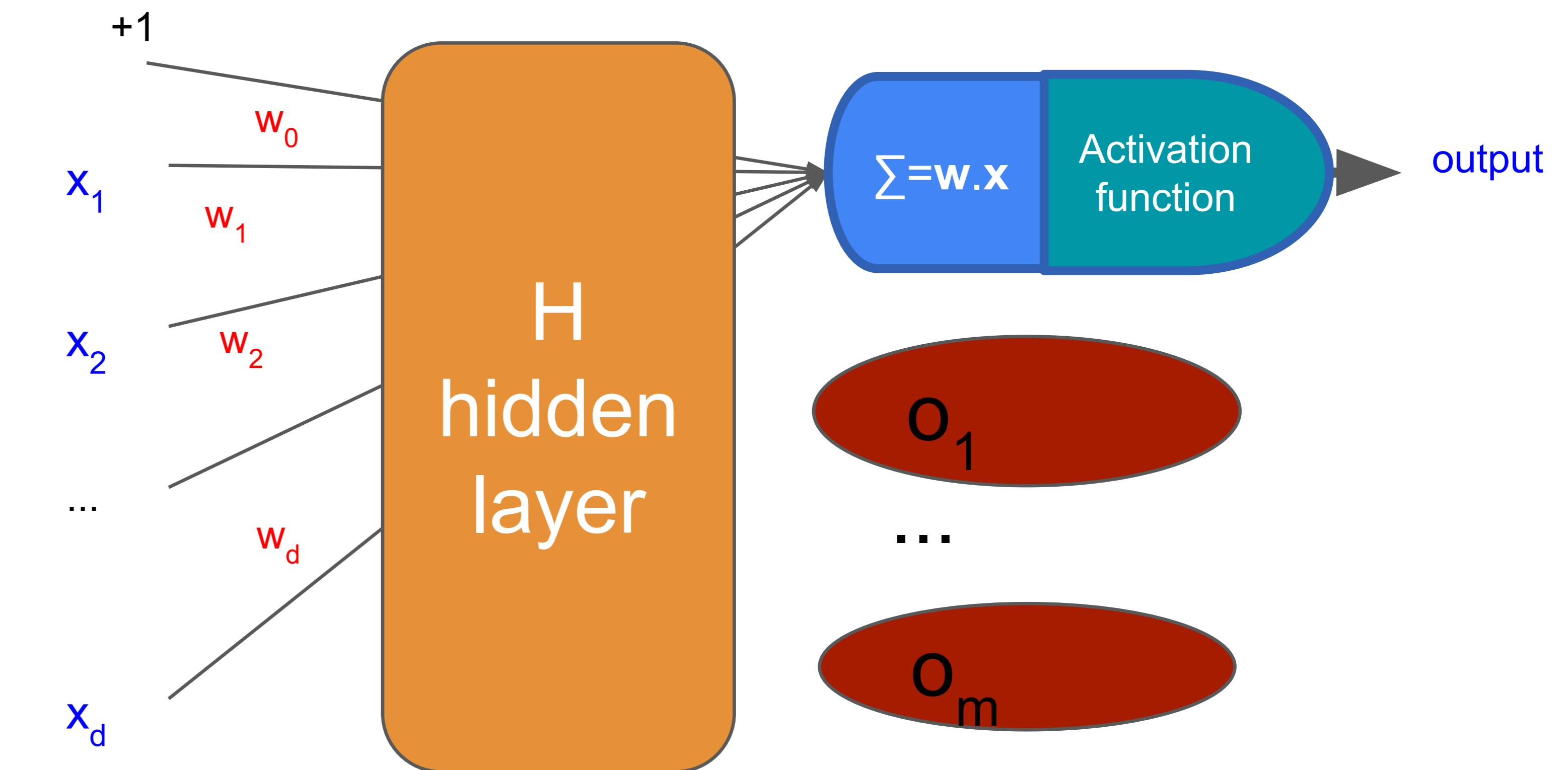
single unit/neuron, estimation function (hypothesis)

Artificial Neural

many computational units become intelligent only via their interactions with each other is inspired by the brain (neuroscience)

Network

composing together many different functions

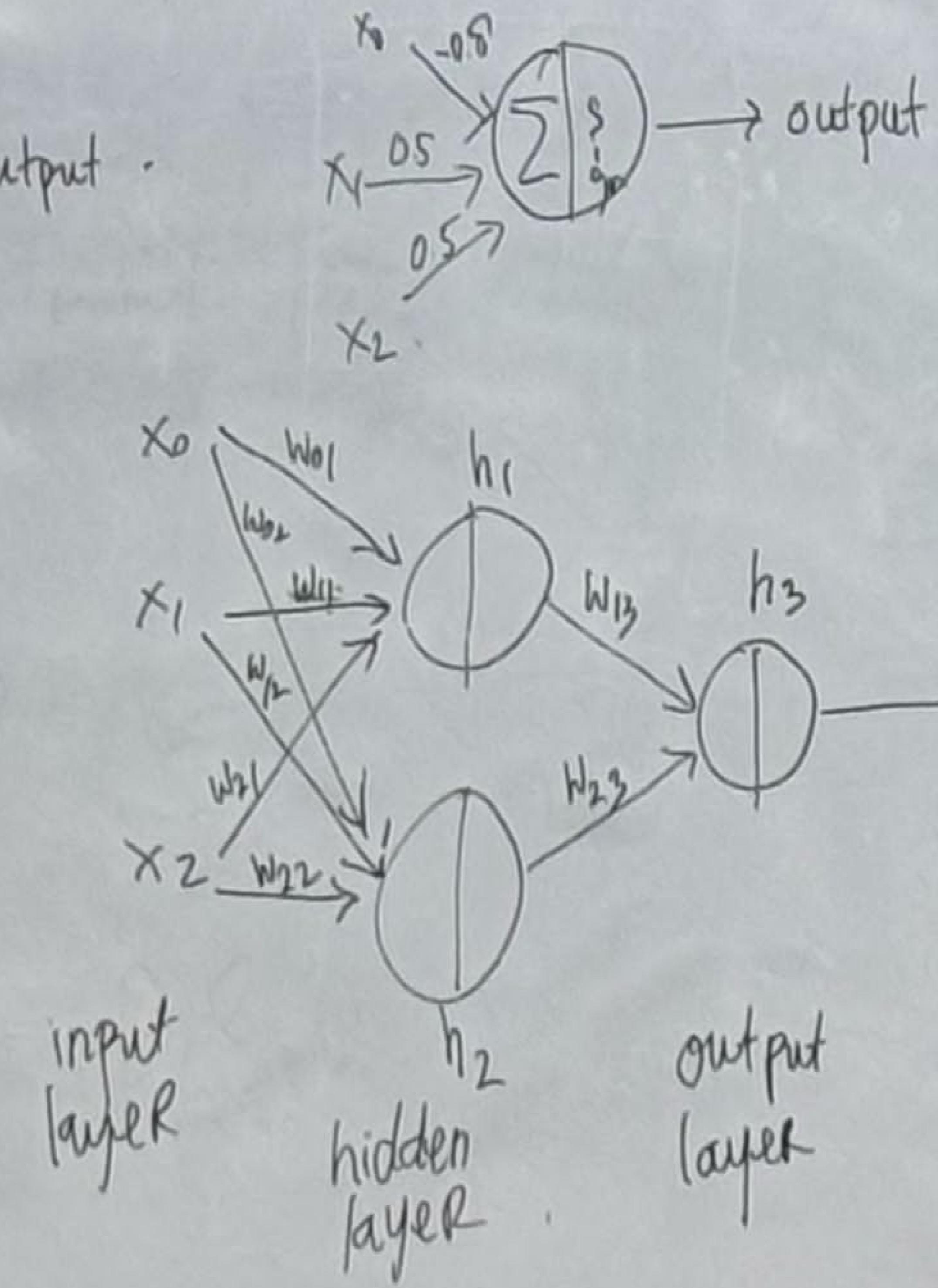
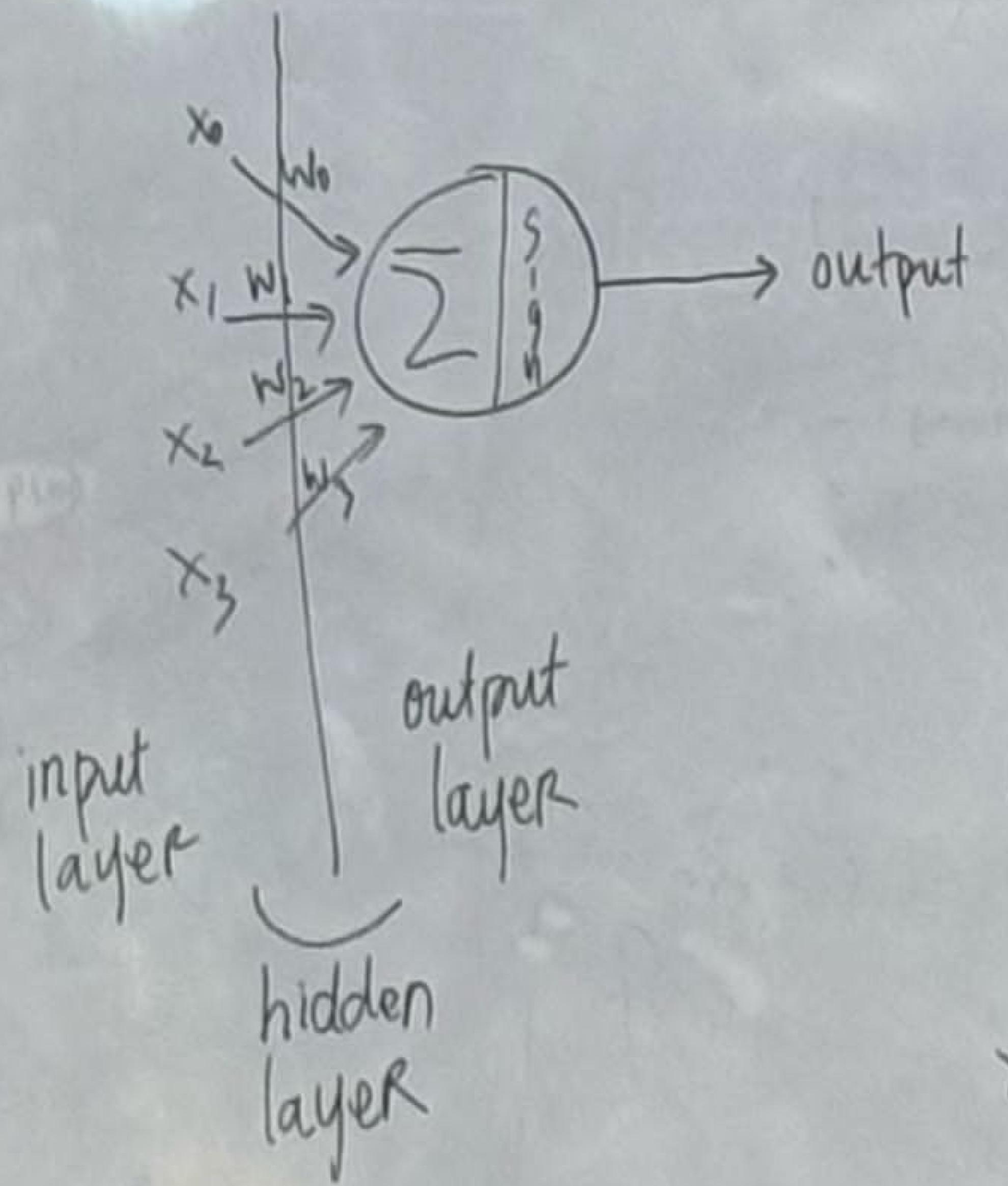


Input layer has d input neurons

Hidden layer

Output layer has m neurons

$$\underline{\text{XOR}}(A, B) = \underline{A} \wedge \overline{B} \vee \overline{A} \wedge \underline{B}$$



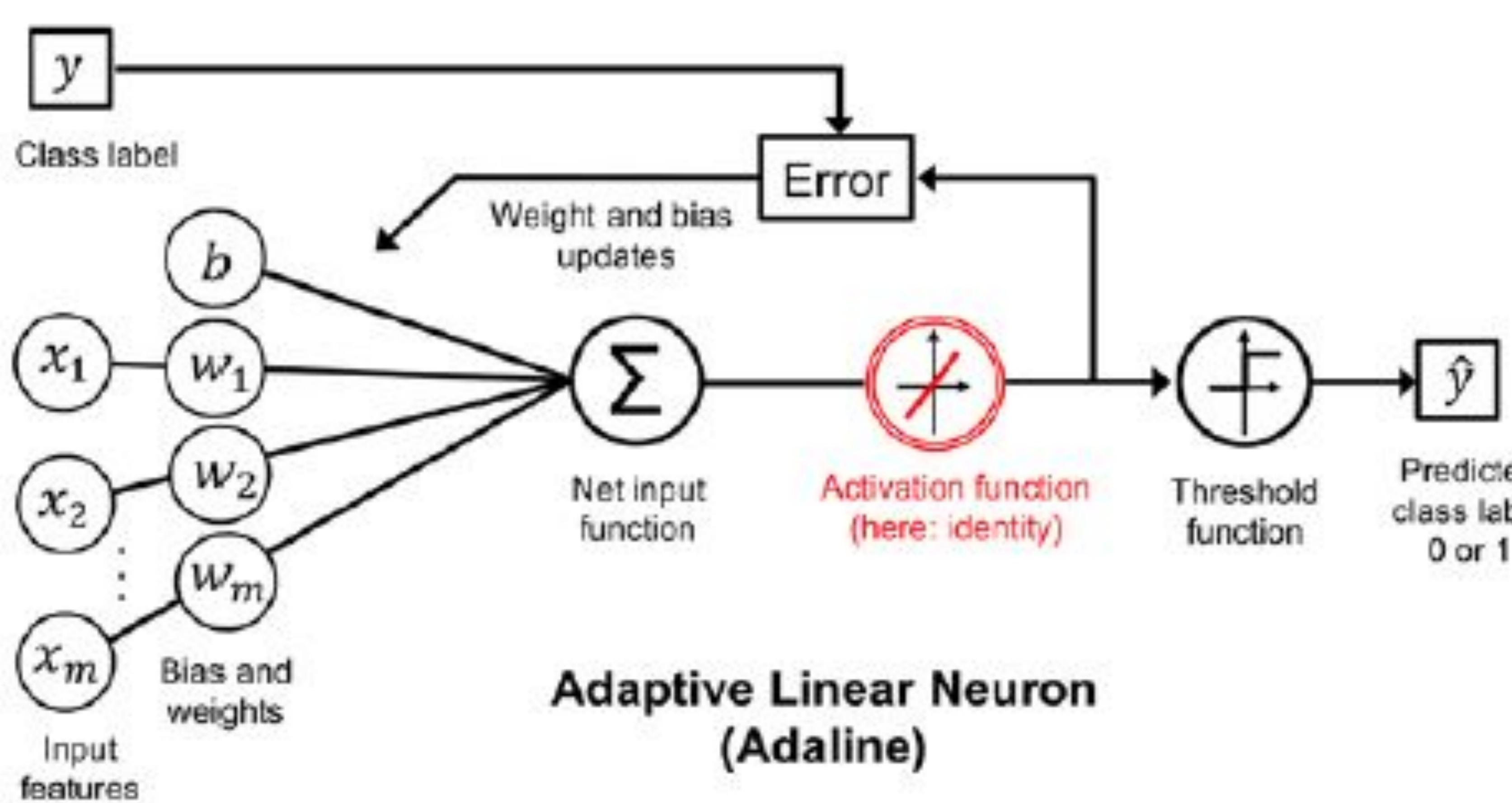
Recap “Artificial Neuron”

Single Layer Neural Network → ADaptive LInear NEuron (Adaline) Algorithm

Update Rule: $w = w + \Delta w; b = b + \Delta b;$ where Δ

$$\Delta w_j = -\eta \frac{\partial L}{\partial w_i}$$

$$\Delta b = -\eta \frac{\partial L}{\partial b}$$



$L(w)$: Squared Error lost function =

$$\frac{1}{2} (y^{(i)} - a^{(i)})^2$$

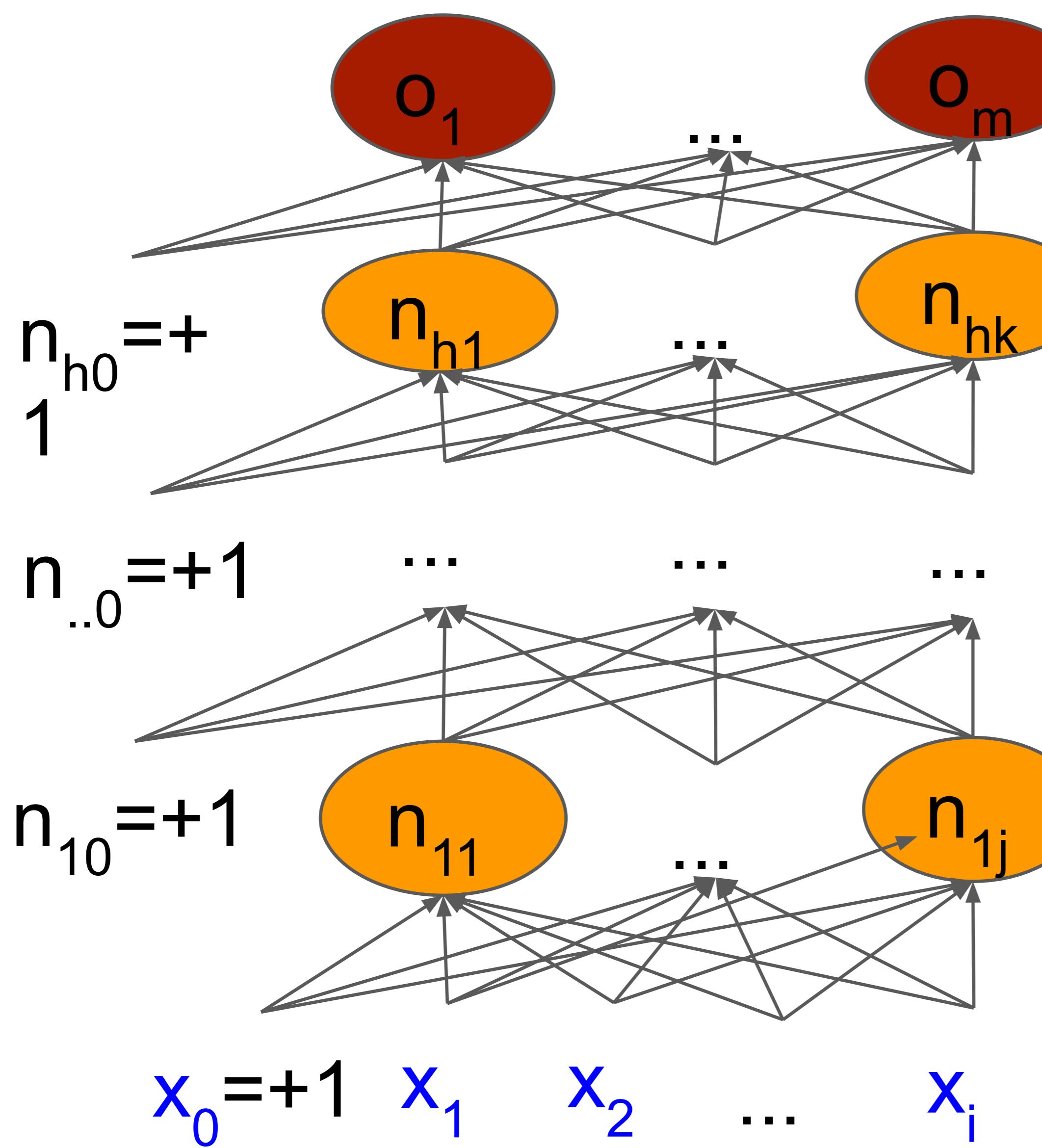
$$a = \sigma \left(\sum_j w_j \cdot x_j + b \right)$$

$$\hat{y} = \begin{cases} 1, & \text{if } a \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

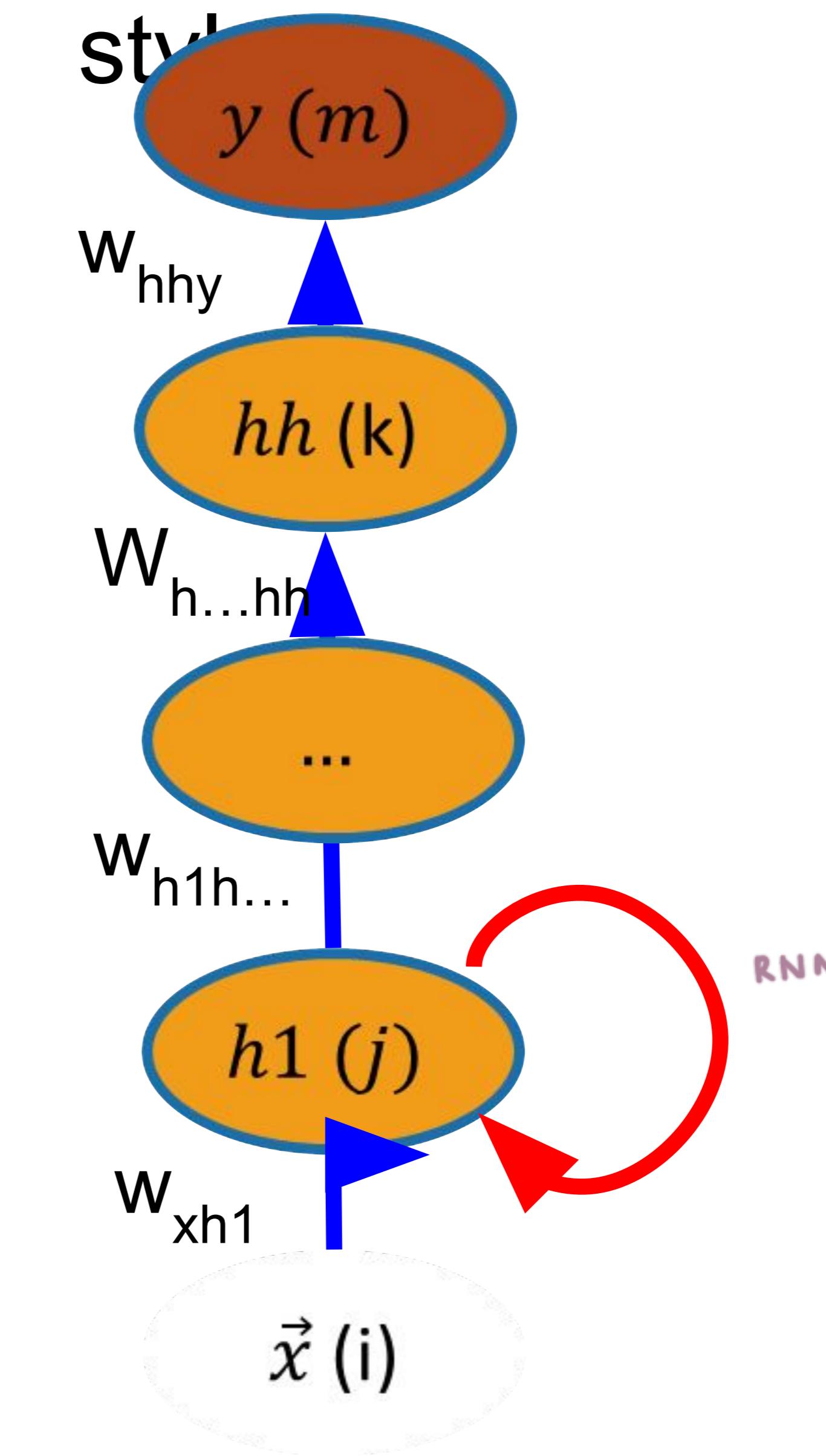
Raschka, Figure 11.1: The Adaline algorithm

Artificial Neural Network Architecture

Explicit style



Compact style



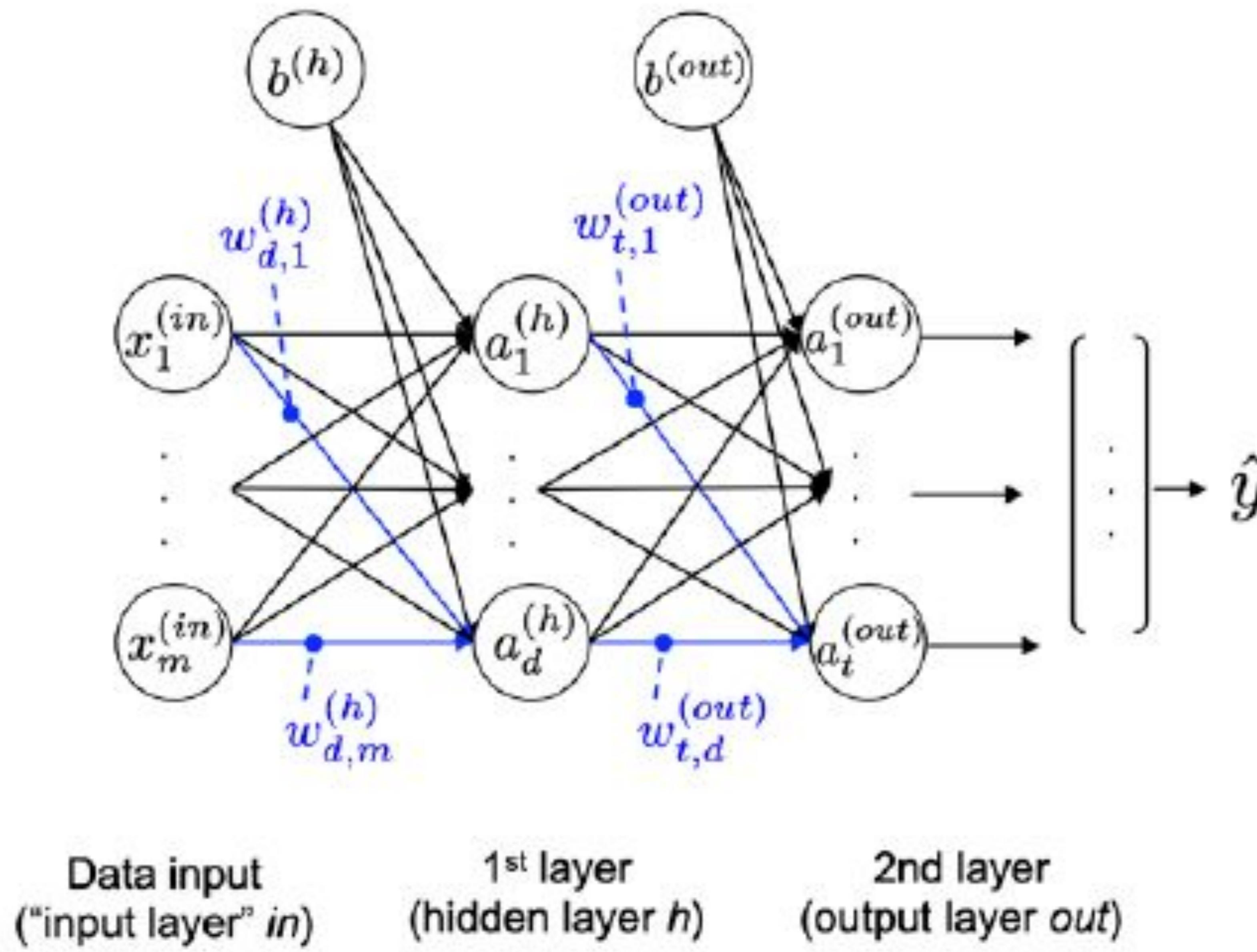
depth

ANN Architecture (Russell & Norvig, AIMA 4th ed, 2022)

- Feed Forward Neural Network (FFNN)
 - connections only in one direction □ it forms a directed acyclic graph with designated input and output nodes
 - Multi layer Perceptron: Fully Connected Feed Forward, non linear activation function
 - Information flows through the network from the input nodes to the output nodes, and there are no loops
- Recurrent Neural Network □ LATER
 - feeds its intermediate or final outputs back into its own inputs (loops)
 - the signal values within the network form a dynamical system that has internal state or memory

Feed Forward Neural Network (FFNN)

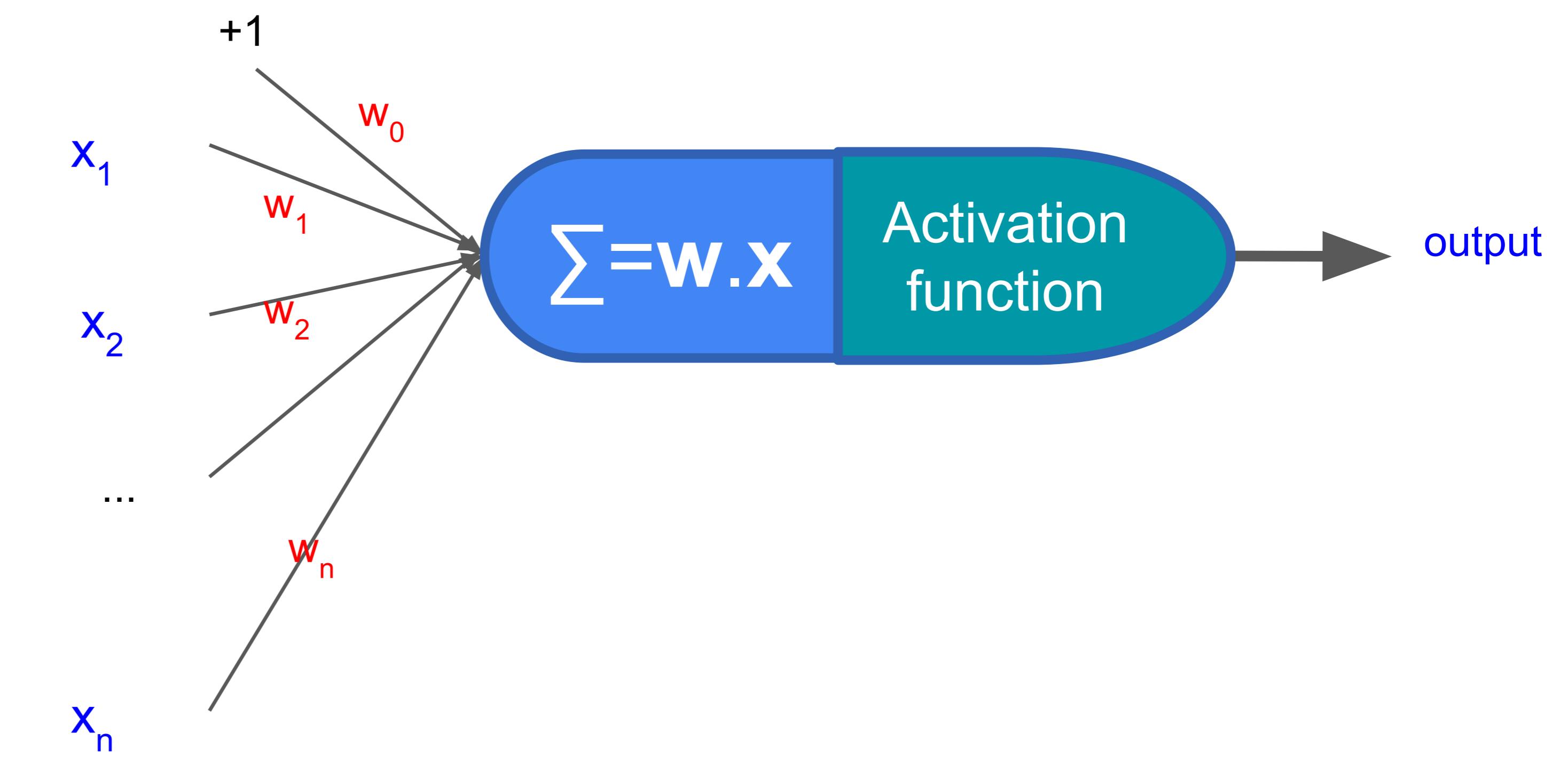
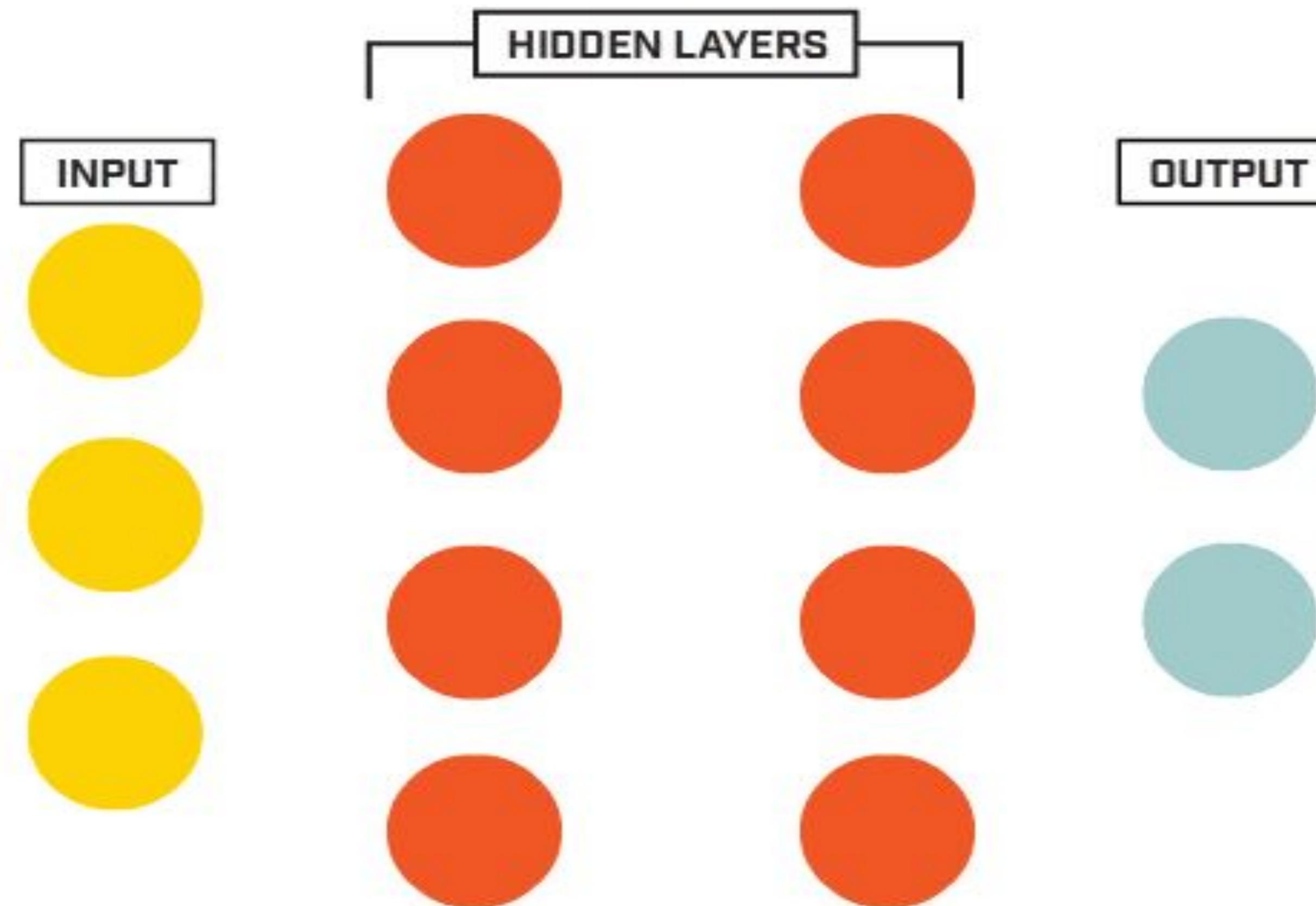
Example of Multi Layer Perceptron with 2 Layer



Raschka, Figure 11.2: A two-layer MLP

Note that in some contexts, the inputs are also regarded as a layer. However, in this case, it would make the Adaline model, which is a single-layer neural network, a two-layer neural network, which may be counterintuitive

Artificial Neural Network: Inference (Forward Propagation)

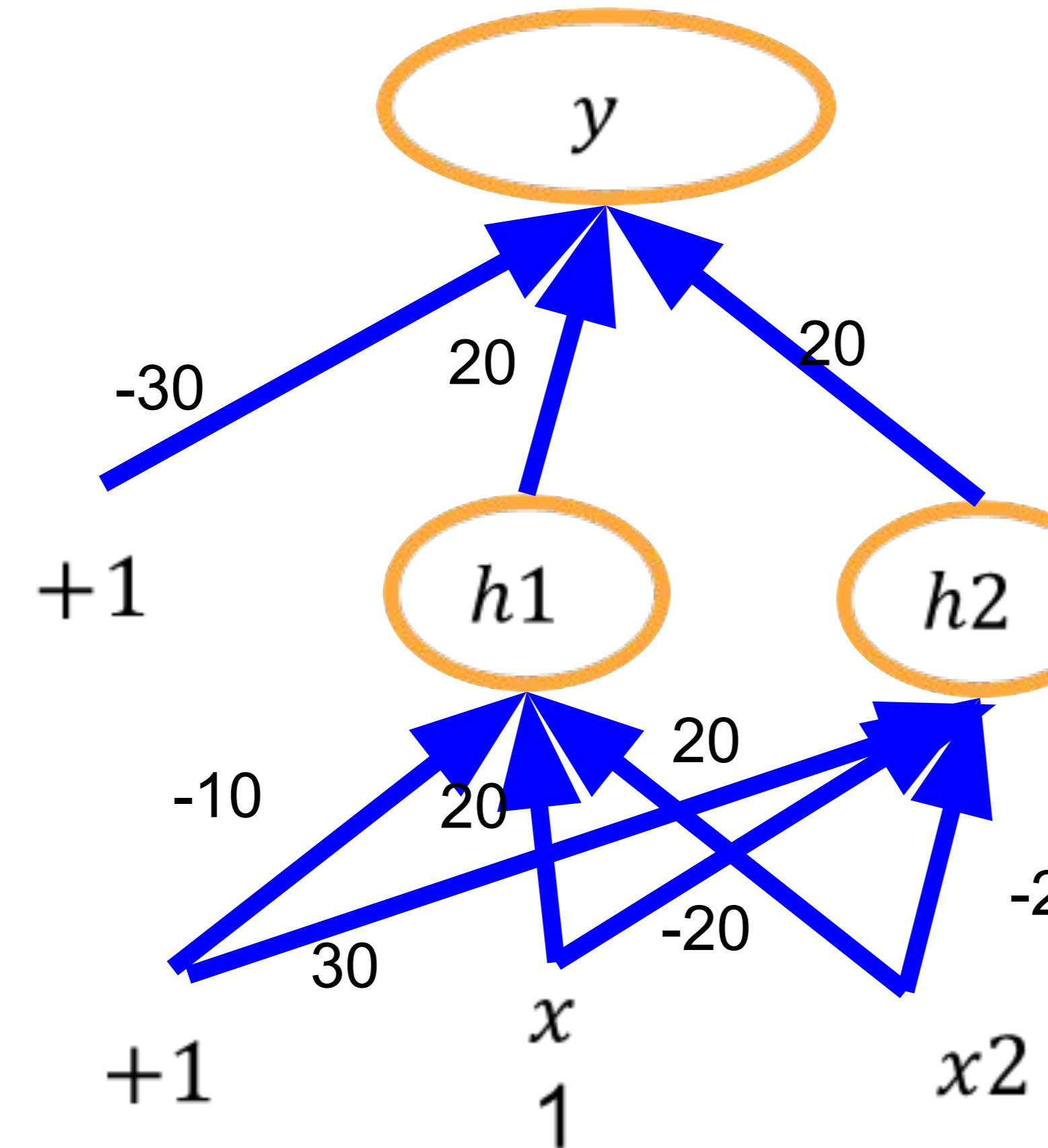
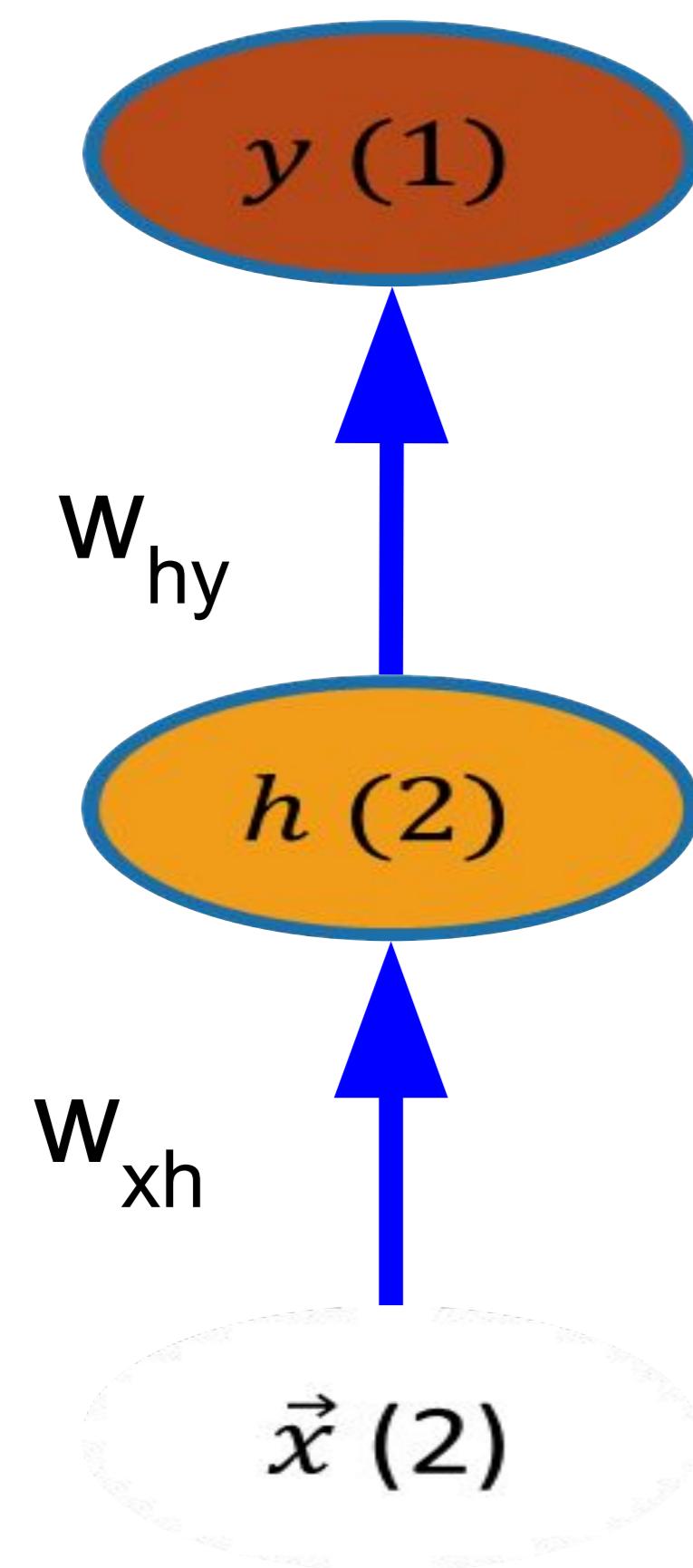


Forward propagation: Information flows through the function being evaluated from vector x , through intermediate computations used to define f , and finally to the output y □ calculate the output of FFNN

Input $\mathbf{x}=(1, x_1, x_2, \dots, x_n)$
Model $\mathbf{w}=(w_0, w_1, w_2, \dots, w_n)$
 $\sum = \mathbf{w} \cdot \mathbf{x} = w_0 + w_1 x_1 + \dots + w_n x_n$
Output = $f(\sum)$

Example: XOR (Sigmoid Model)

$$\sigma(\text{net}): 1/(1+e^{-\text{net}})$$



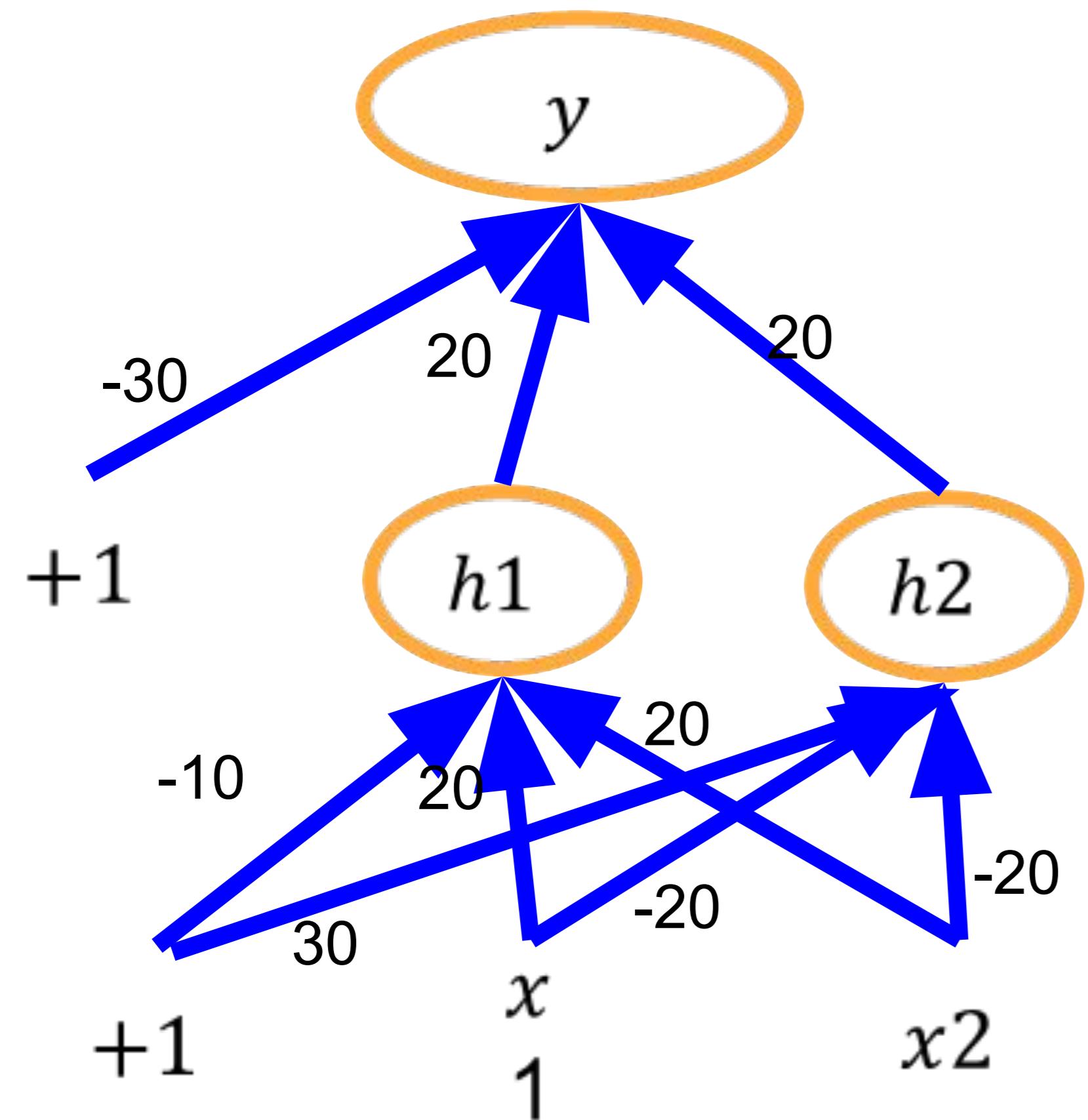
sigmoid

$$h_1 = \sigma(-10 + 20 \cdot x_1 + 20 \cdot x_2)$$
$$h_2 = \sigma(30 - 20 \cdot x_1 - 20 \cdot x_2)$$
$$y = \sigma(-30 + 20 \cdot h_1 + 20 \cdot h_2)$$

x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0

Target function classifies two input into predefined classes {0, 1}.

FFNN: XOR (Sigmoid Model)



$$h_1 = \sigma(-10 + 20*x_1 + 20*x_2)$$
$$h_2 = \sigma(30 - 20*x_1 - 20*x_2)$$
$$y = \sigma(-30 + 20*h_1 + 20*h_2)$$

bias	x_0	x_1	x_2	y	Σh_1	h_1	Σh_2	h_2	$\Sigma \hat{y}$	\hat{y}
	1	0	0	0	-10	0.00	30	1.00	-10.00	0.00
	1	0	1	1	10	1.00	10	1.00	10.00	1.00
	1	1	0	1	10	1.00	10	1.00	10.00	1.00
	1	1	1	0	30	1.00	-10	0.00	-10.00	0.00

$$\sigma(\text{net}): 1/(1+e^{-\text{net}})$$

Latihan 1

Diketahui sebuah arsitektur jaringan Neural Network sebagai berikut yang memiliki dua unit input x , tiga unit output o dan satu unit hidden h , serta menggunakan fungsi aktivasi Sigmoid di setiap unitnya. Proses pembelajaran menerapkan algoritma Backpropagation.

Tuliskan formula untuk menghitung luaran h dan o_1

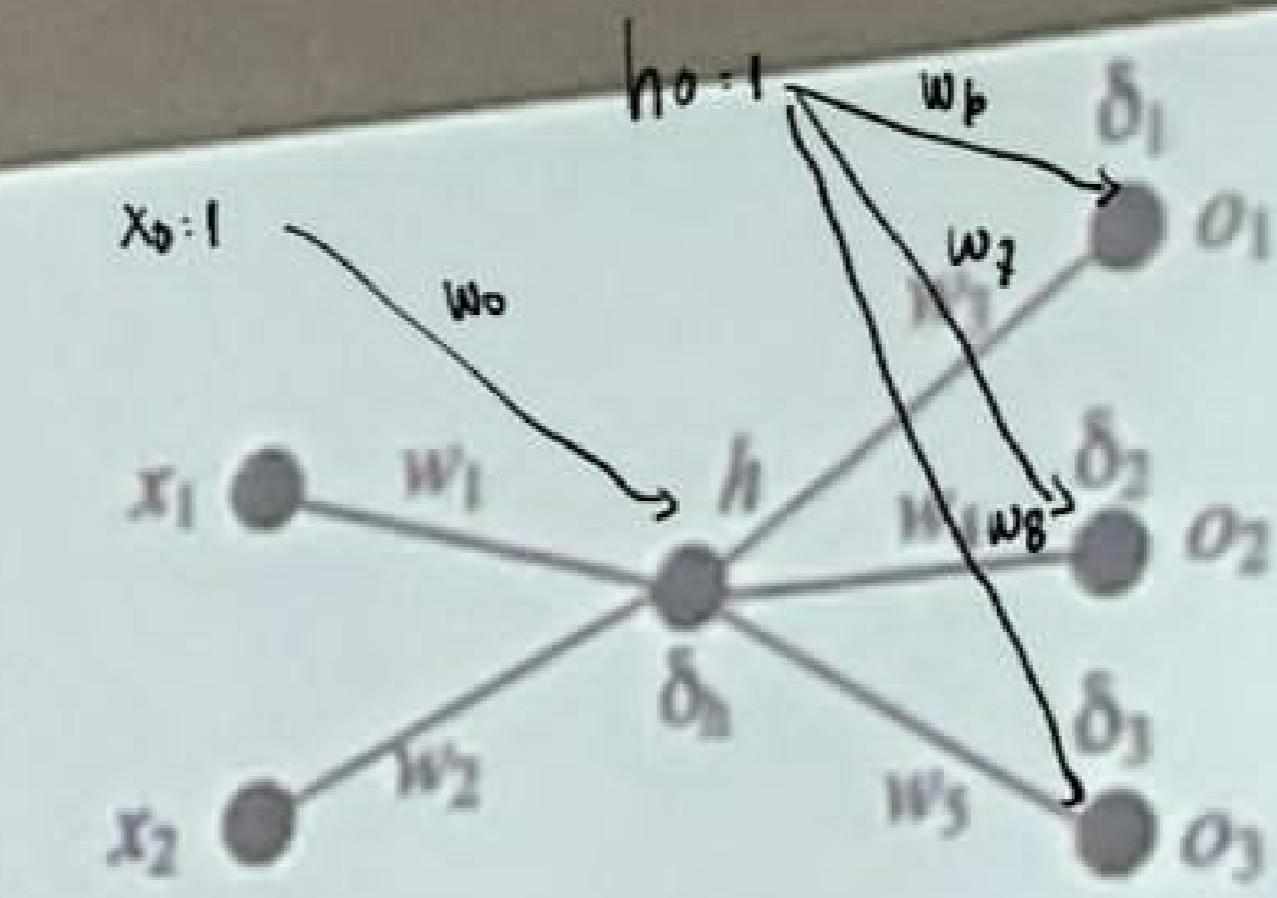
$$h = \sigma(w_0 + w_1 \cdot x_1 + w_2 \cdot x_2)$$

(yg masuk ke o_3 cuma w_3 dan w_5)

$$o_3 = \sigma(w_3 \cdot h + w_5 \cdot h)$$

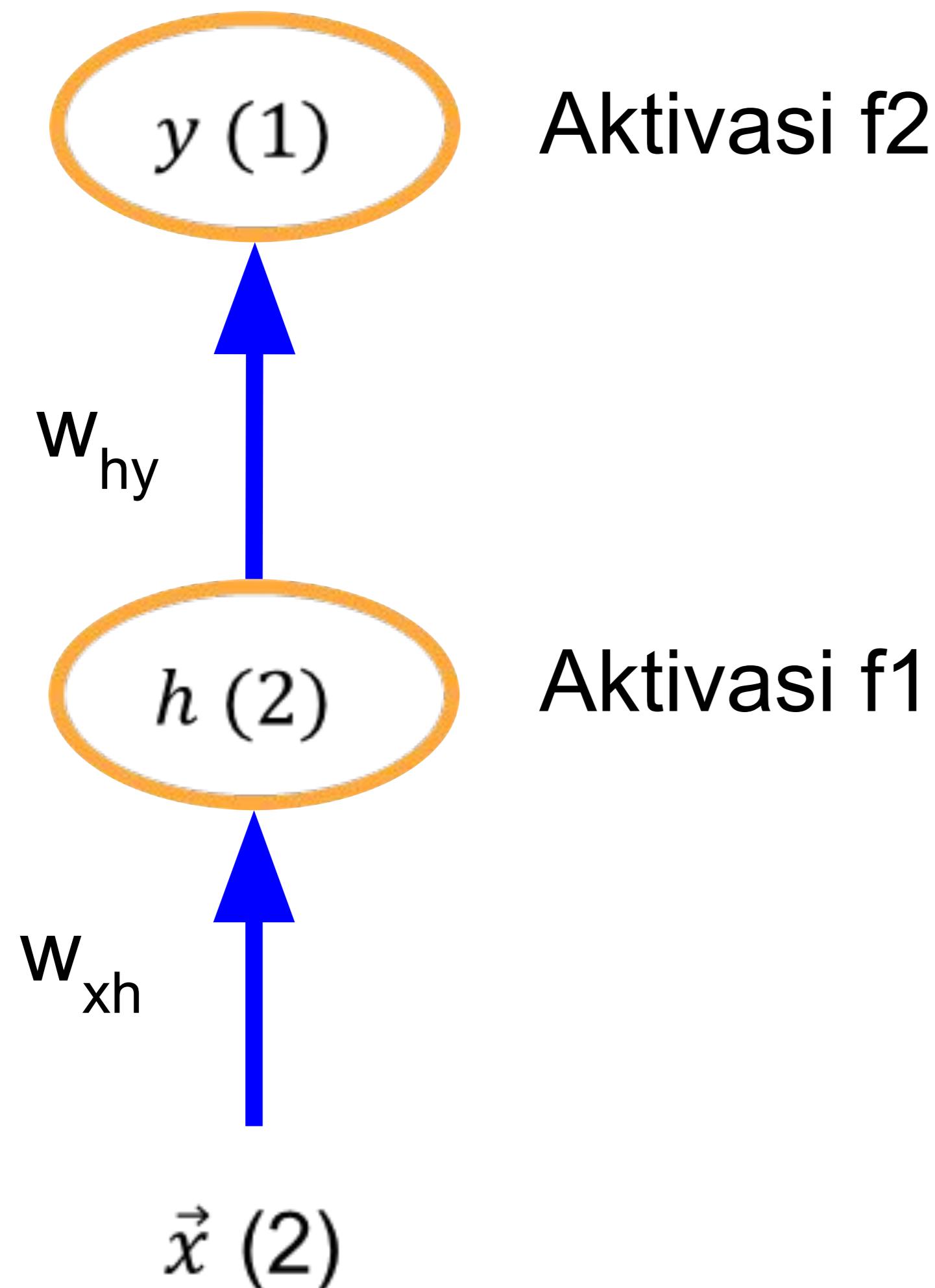
$$o_1 = \sigma(w_6 + w_7 \cdot h)$$

$$o_2 = \sigma(w_8 + w_9 \cdot h)$$



FFNN: Forward Propagation dengan Input Vector \mathbf{x}

Compact style



Input vector \mathbf{x} :

$$f(\mathbf{x}; \mathbf{w}_{xh}, c, \mathbf{w}_{hy}, b) = f_2(\mathbf{w}_{hy}^T \cdot f_1(\mathbf{w}_{xh}^T \cdot \mathbf{x} + c) + b)$$

Jika bobot bias c menjadi bagian \mathbf{w}_{xh} dan bobot bias b menjadi bagian dari \mathbf{w}_{hy} , maka:

$$f(\mathbf{x}; \mathbf{w}_{xh}, \mathbf{w}_{hy}) = f_2(\mathbf{w}_{hy}^T \cdot f_1(\mathbf{w}_{xh}^T \cdot \mathbf{x}))$$

tips soal :

kalo cm dikasih 1 instance \rightarrow vektor

$$XOR = \Gamma$$

$$XOR = RELU \times Linear$$

$$f_2 \rightarrow W_{hy} = \begin{bmatrix} e \\ f \end{bmatrix}; b = [h]$$

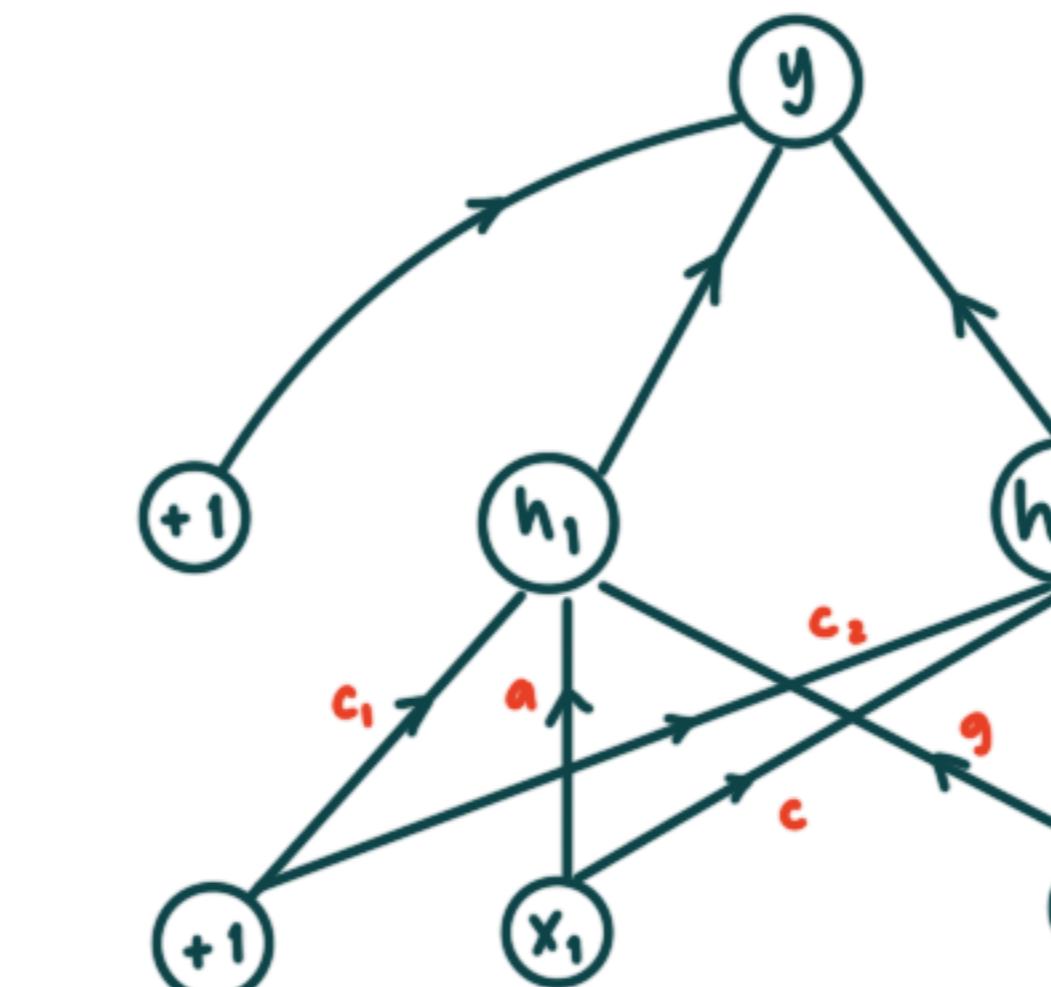
$$f_1 \rightarrow W_{xh} = \begin{bmatrix} a & c \\ g & d \end{bmatrix}; c = \begin{bmatrix} c_1 \\ c_2 \end{bmatrix}$$

formula 2

$$W_{hy} = \begin{bmatrix} b_1 \\ e \\ f \end{bmatrix}$$

$$W_{xh} = \begin{bmatrix} c_1 & a & c \\ c_2 & g & d \end{bmatrix}$$

$$f(\vec{x}) = f_2(W_{hy} \cdot f_1(W_{xh} \cdot x))$$



kalo datanya cuma 1,
jgn lupa di transpose

formula 1

$$h_1 = f_1(x_1 \cdot a + x_2 \cdot g + c_1)$$

$$h_2 = f_1(c_2 + x_1 \cdot c + x_2 \cdot d)$$

$$y = f_2([b_1] + h_1 \cdot e + h_2 \cdot f) \rightarrow f_2 = ([b_1] + \begin{bmatrix} e \\ f \end{bmatrix}^T \cdot \begin{bmatrix} h_1 \\ h_2 \end{bmatrix})$$

(e & f adlh W_{hy} ,
dikalikan dgn h_1, h_2 ,
ditambah dgn b)

$$\begin{bmatrix} a & g \\ c & d \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} c_1 \\ c_2 \end{bmatrix}$$

$$f(\vec{x}) = f_2(W_{hy}^T \cdot f_1(W_{xh}^T \cdot x + c) + b)$$

$$\begin{bmatrix} e \\ f \end{bmatrix}$$

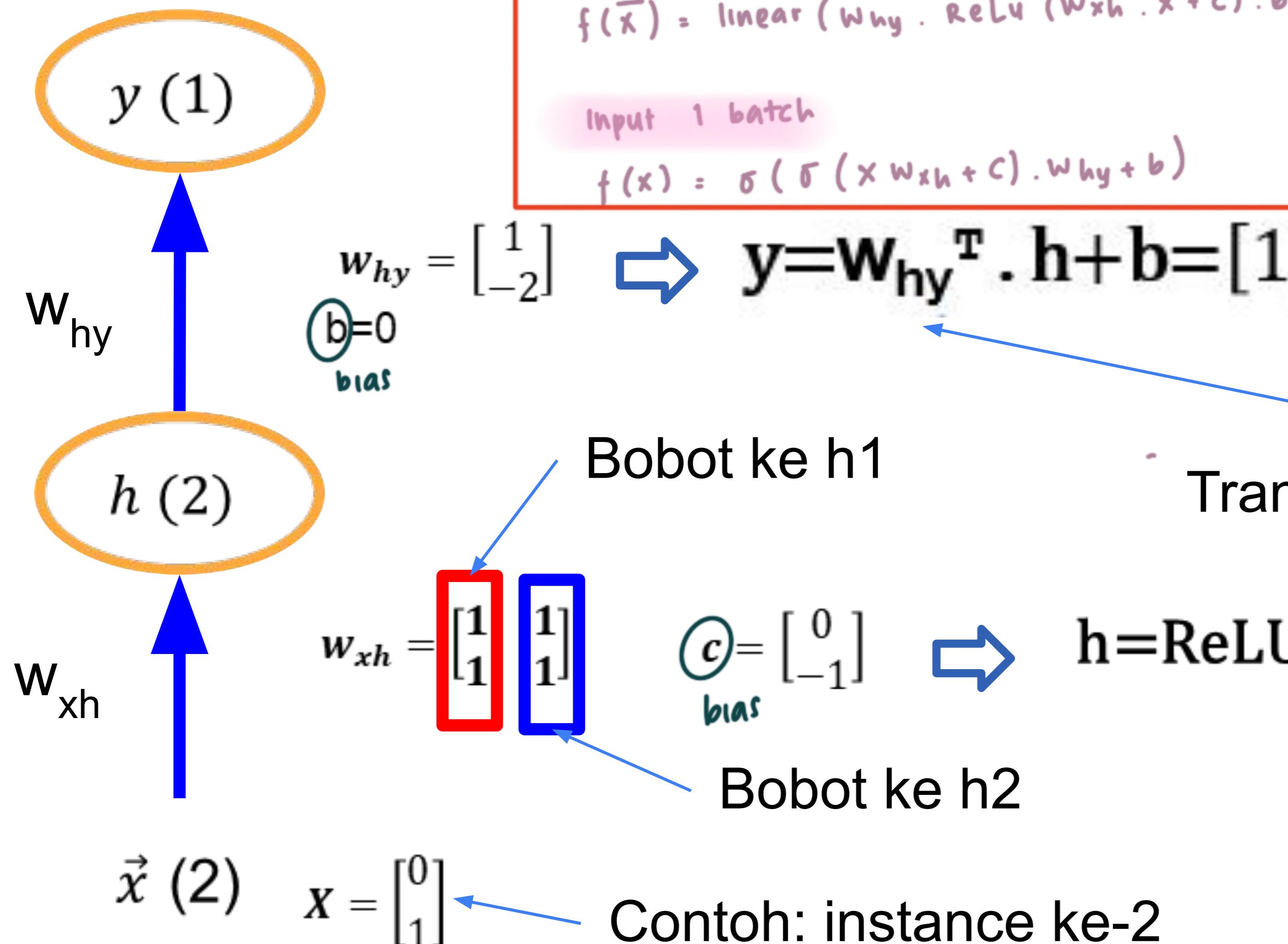
$$= [eh_1 + fh_2] + [b_1]$$

$$= [eh_1 + eh_2 + b_1]$$

FFNN: XOR (RELU & LINEAR) □ input vector x

rectified linear unit

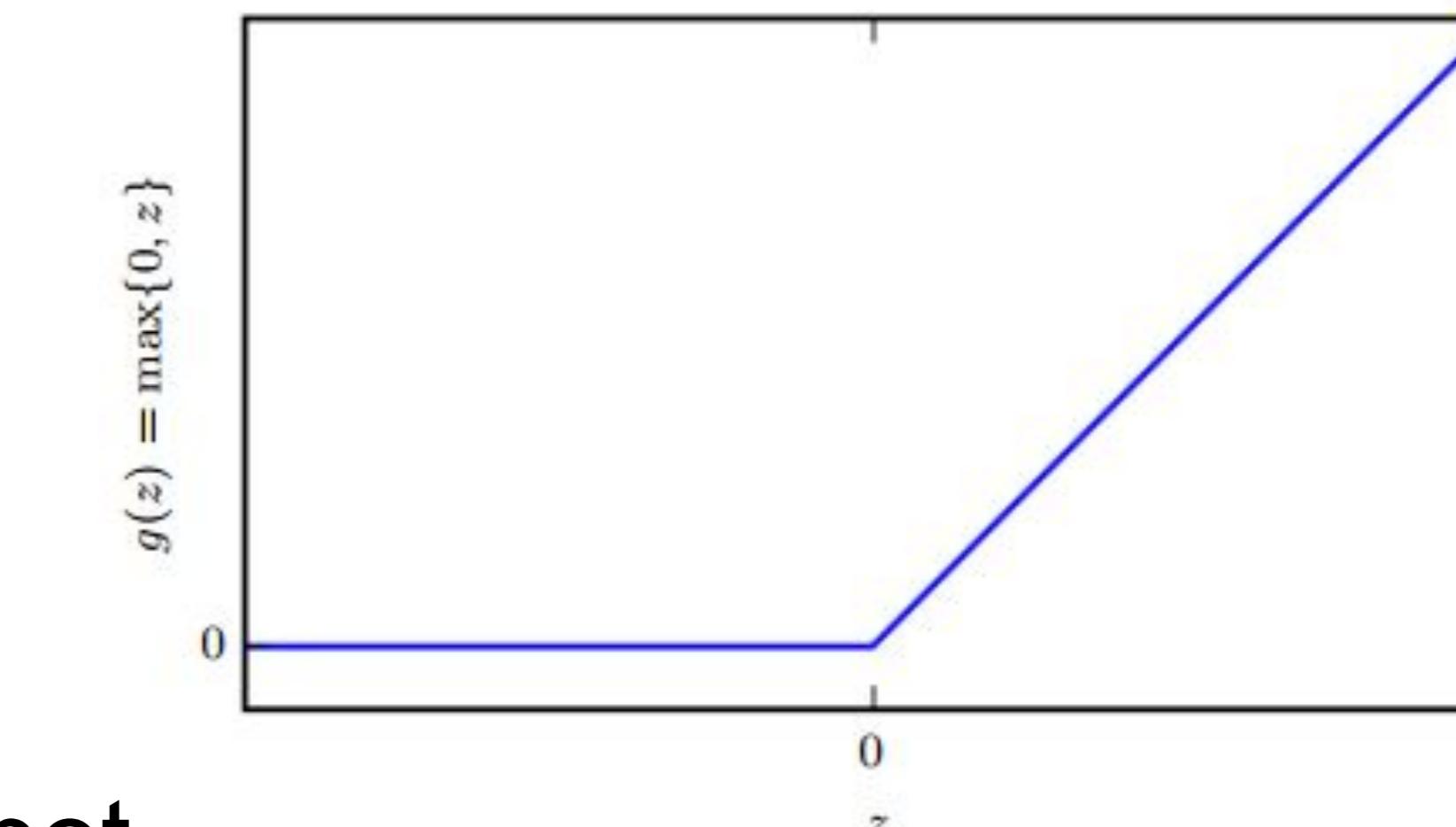
Compact style



Input vector \mathbf{x} :

$$f(\mathbf{x}; \mathbf{w}_{xh}, \mathbf{c}, \mathbf{w}_{hy}, \mathbf{b}) = \mathbf{w}_{hy}^T \cdot (\max\{0, \mathbf{w}_{xh}^T \cdot \mathbf{x} + \mathbf{c}\}) + \mathbf{b}$$

$$\text{ReLU}(z) = \max\{0, z\}$$

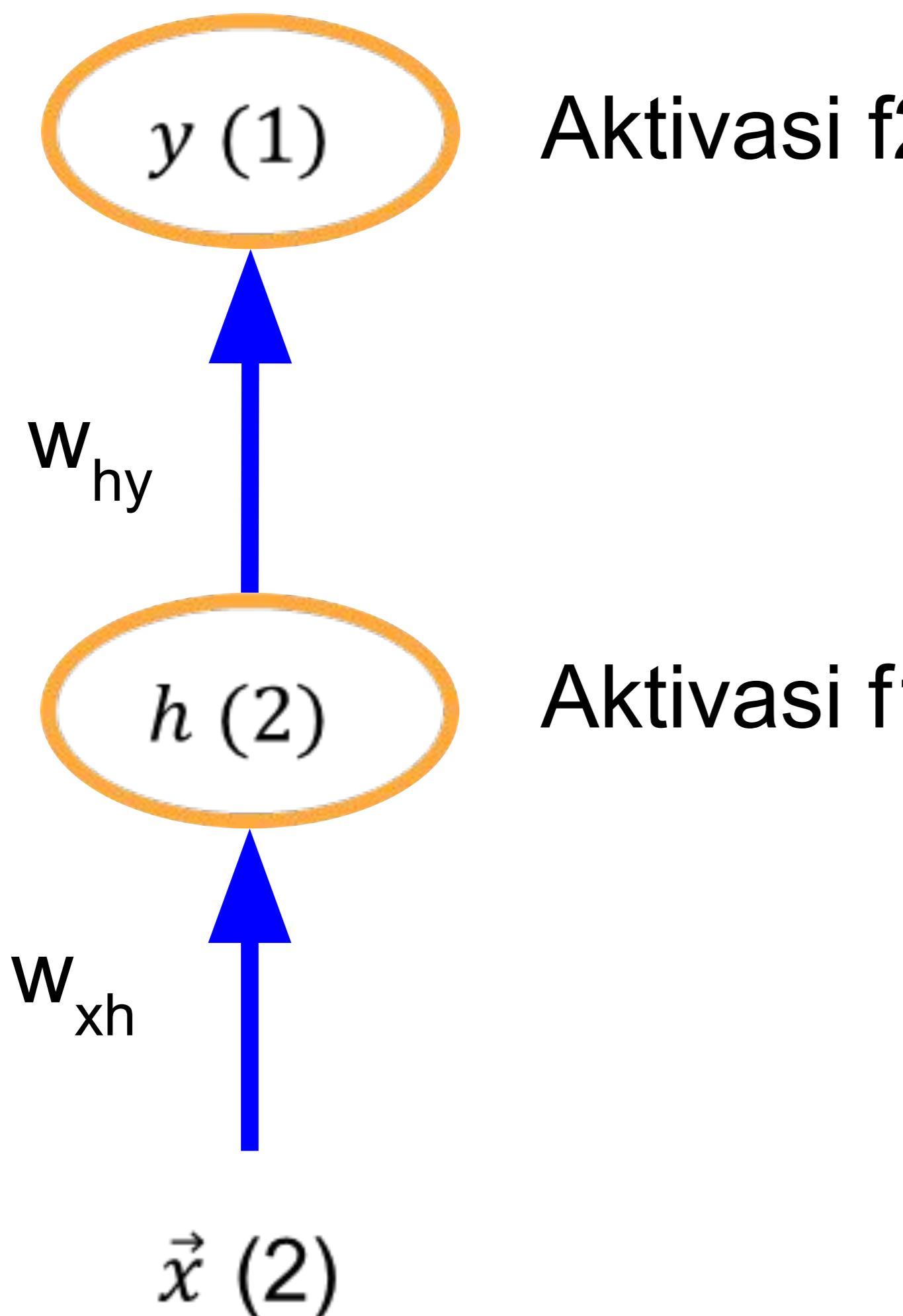


$$h = \text{ReLU}(\mathbf{w}_{xh}^T \cdot \mathbf{x} + \mathbf{c}) = \max\{0, \mathbf{w}_{xh}^T \cdot \mathbf{x} + \mathbf{c}\} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

h1 = $\text{ReLU}(0+x_1+x_2)$
 h2 = $\text{ReLU}(-1+x_1+x_2)$
 $y = (0+h1-2*h2)$

FFNN: Forward Propagation dengan Input Mini-batch \mathbf{X}

Compact style



Input vector x :

$$f(x; w_{xh}, c, w_{hy}, b) = f_2(w_{hy}^T \cdot f_1(w_{xh}^T \cdot x + c) + b)$$

Jika bobot bias c menjadi bagian w_{xh} dan bobot bias b menjadi bagian dari w_{hy} , maka:

$$f(x; w_{xh}, w_{hy}) = f_2(w_{hy}^T \cdot f_1(w_{xh}^T \cdot x))$$

Input mini-batch \mathbf{x} :

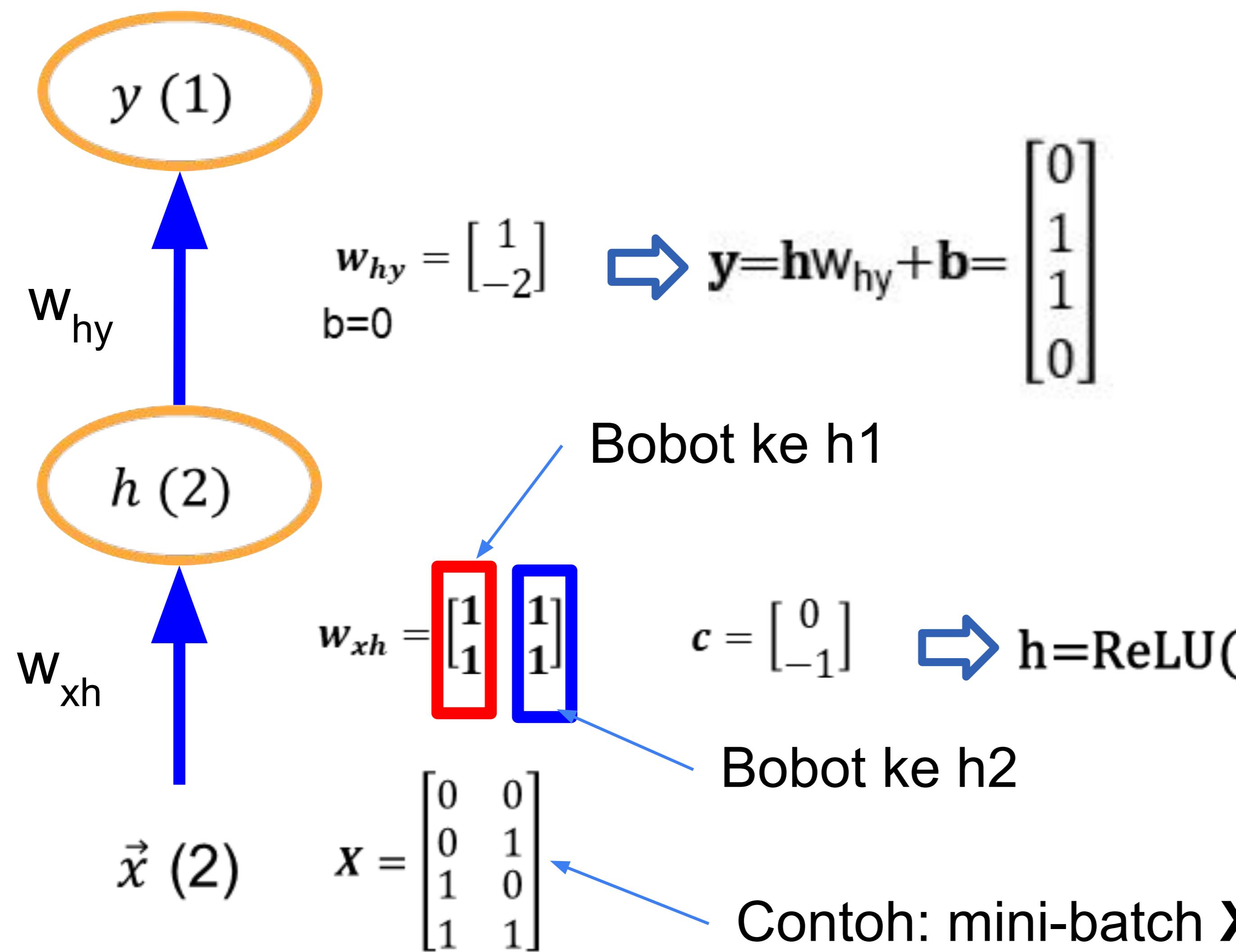
$$f(\mathbf{x}; w_{xh}, c, w_{hy}, b) = f_2(f_1(\mathbf{x}w_{xh} + c) \cdot w_{hy} + b)$$

Jika bobot bias c menjadi bagian w_{xh} dan bobot bias b menjadi bagian dari w_{hy} , maka:

$$f(\mathbf{x}; w_{xh}, w_{hy}) = f_2(f_1(\mathbf{x}w_{xh}) \cdot w_{hy})$$

FFNN: XOR (RELU & LINEAR) □ input mini-batch X

Compact style



$$h_1 = \text{ReLU}(0 + x_1 + x_2)$$

$$h_2 = \text{ReLU}(-1 + x_1 + x_2)$$

$$y = (0 + h_1 - 2 * h_2)$$

$$X = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$$

$$W_{xh} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 1 \end{bmatrix} \leftarrow \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}_{4 \times 3} + \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 1 \end{bmatrix}_{3 \times 2} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}_{4 \times 3}$$

$$h = \text{ReLU}(XW_{xh} + C)$$

$$XW_{xh} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$$

$$\text{ReLU}(XW_{xh} + C) = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 1 \end{bmatrix}$$

Forward Propagation

Require: Network depth, l

Require: $\mathbf{W}^{(i)}, i \in \{1, \dots, l\}$, the weight matrices of the model

Require: $\mathbf{b}^{(i)}, i \in \{1, \dots, l\}$, the bias parameters of the model

Require: \mathbf{x} , the input to process

Require: \mathbf{y} , the target output

$$\hat{\mathbf{h}}^{(0)} = \mathbf{x}$$

for $k = 1, \dots, l$ **do**

$$\mathbf{a}^{(k)} = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \hat{\mathbf{h}}^{(k-1)}$$

$$\hat{\mathbf{h}}^{(k)} = f(\mathbf{a}^{(k)})$$

end for

$$\hat{\mathbf{y}} = \hat{\mathbf{h}}^{(l)}$$

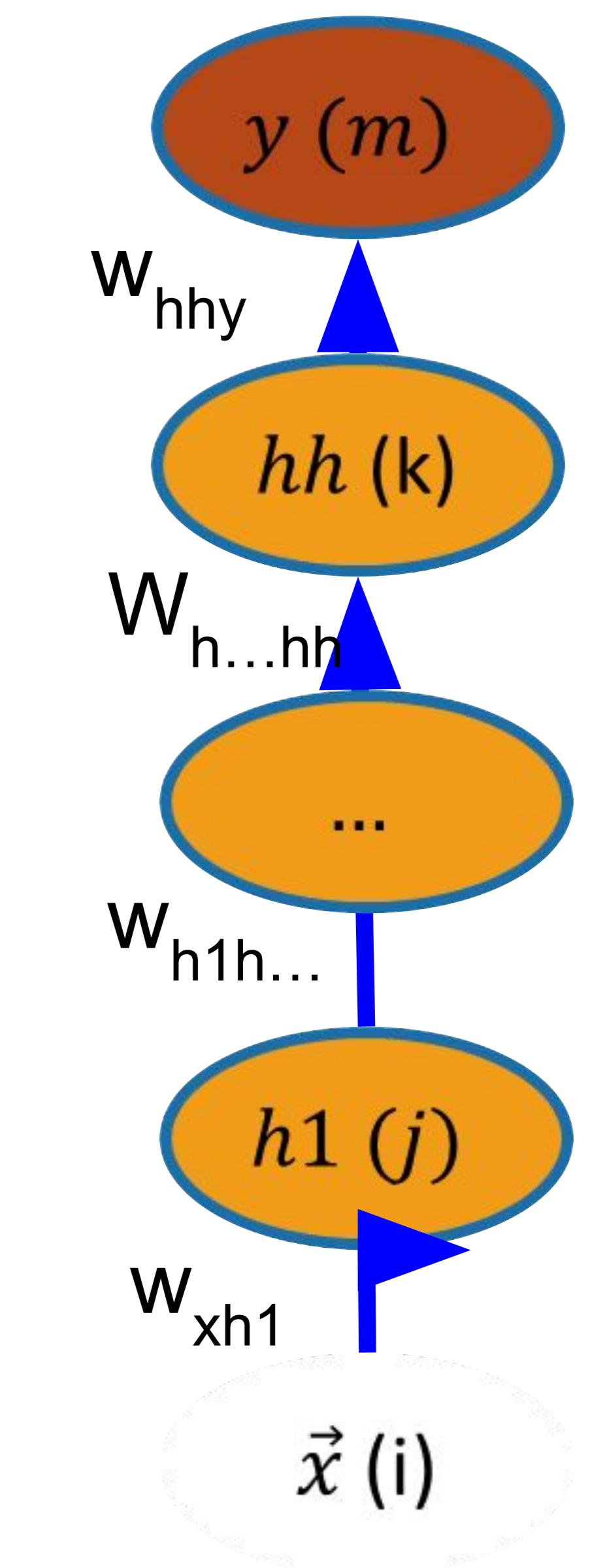
$$J = L(\hat{\mathbf{y}}, \mathbf{y}) + \lambda \Omega(\theta)$$

$f(a)$: activation function

$L(\hat{\mathbf{y}}, \mathbf{y})$: loss function

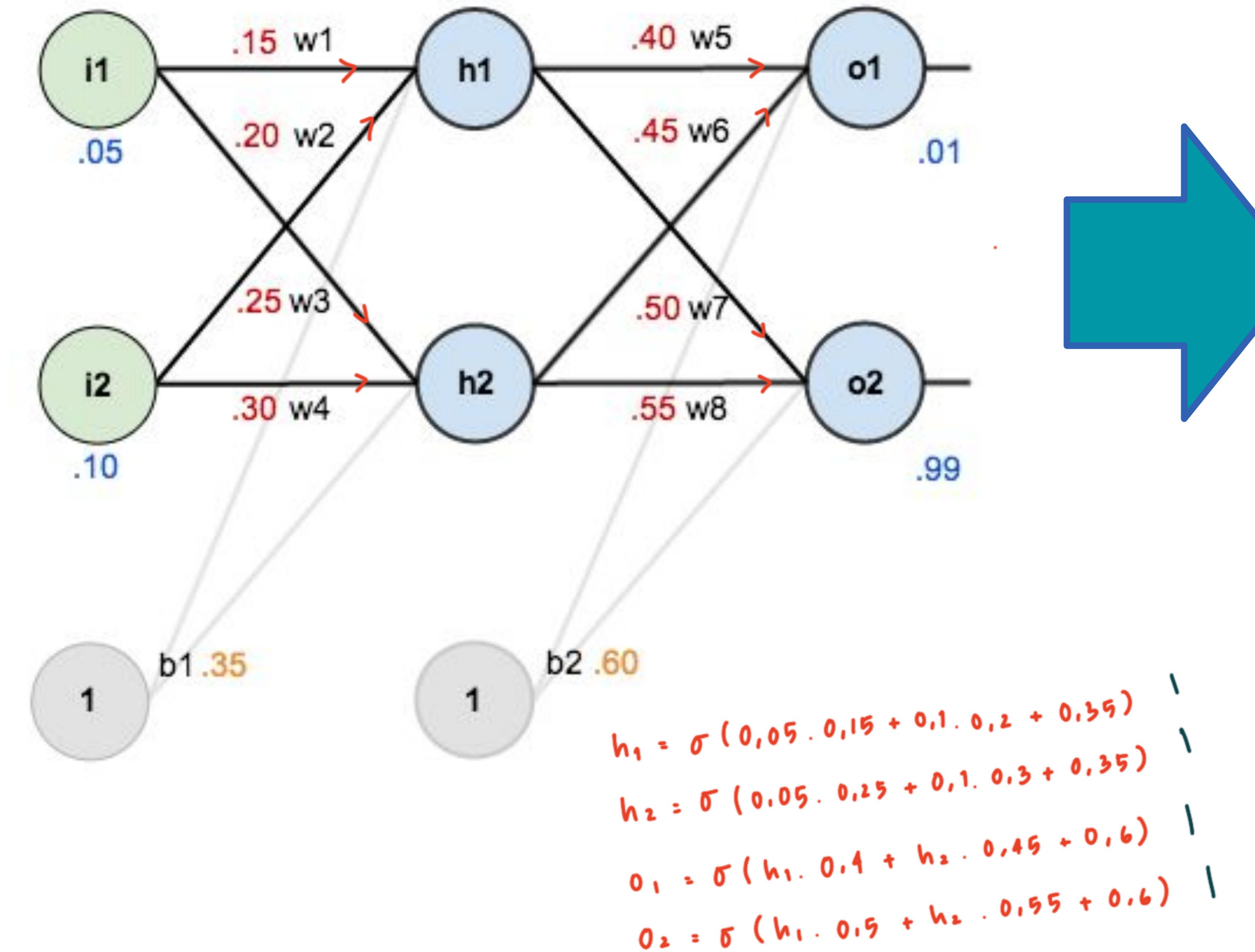
θ : all parameters (\mathbf{W}, \mathbf{b})

$\Omega(\theta)$: regularizer

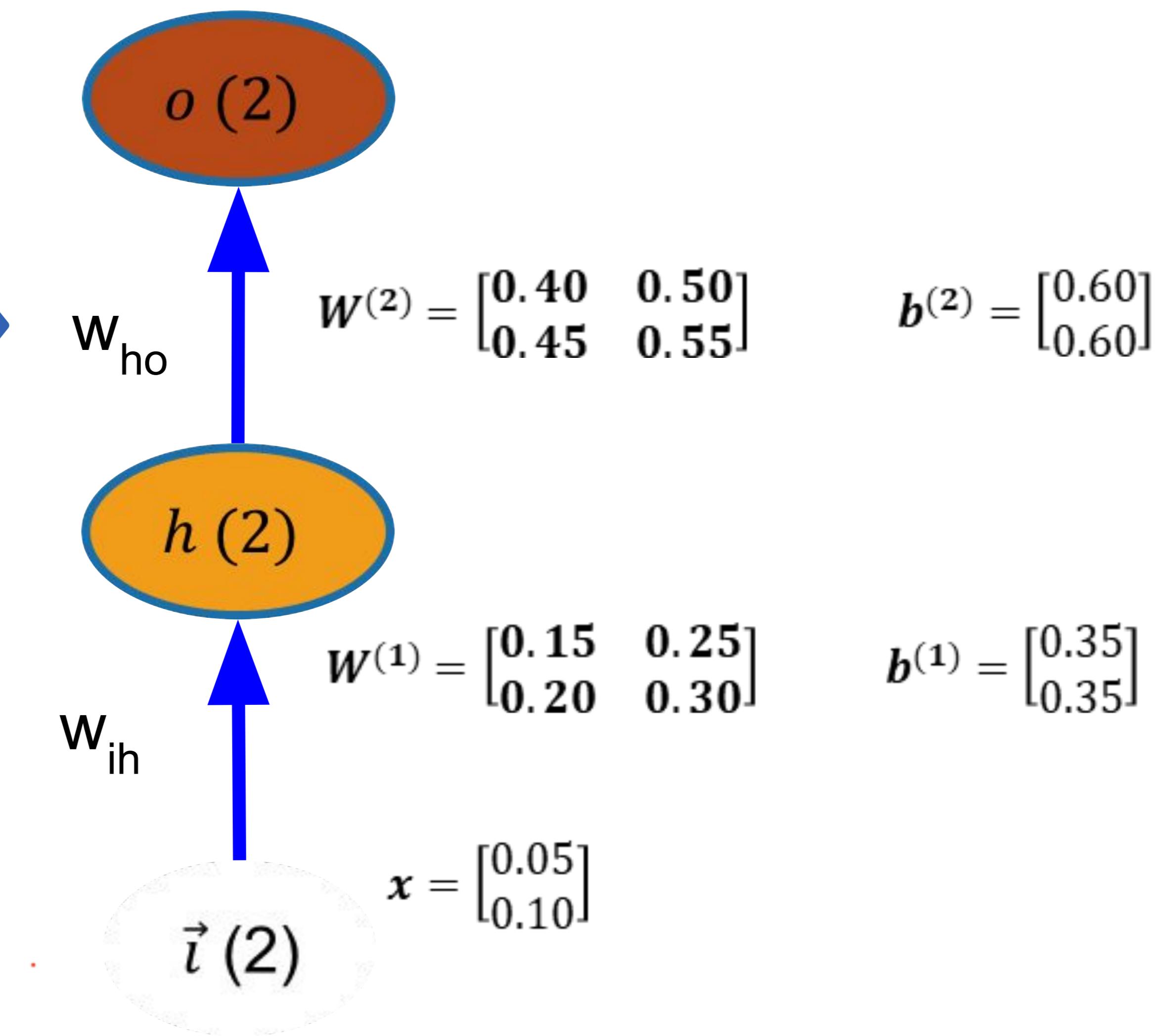


Forward Propagation: Example

(from <http://mattmazur.com/2015/03/17/a-step-by-step-backpropagation-example/>)

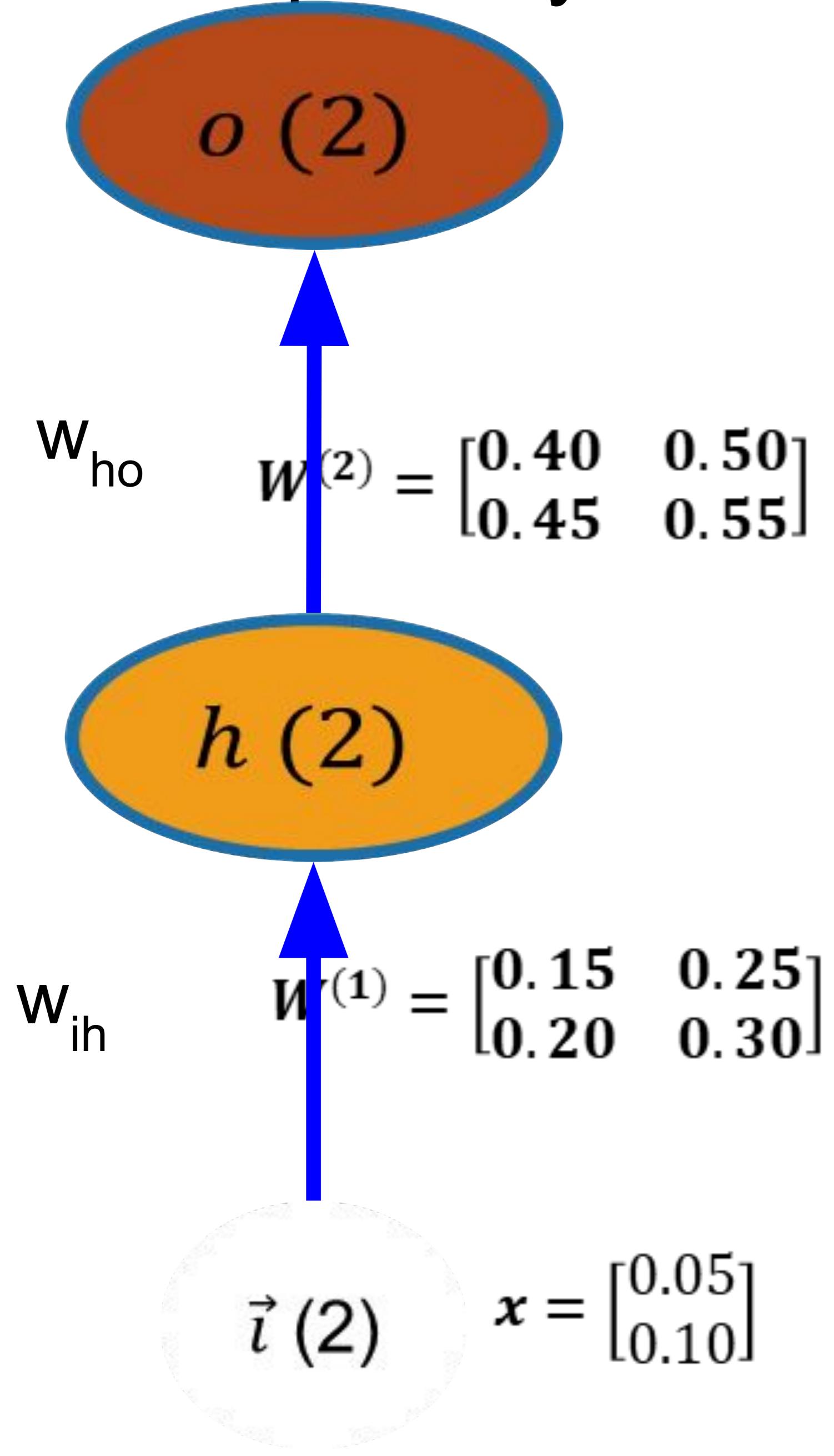


Compact style: sigmoid model



Forward Propagation

Compact style



Target:

$$\begin{aligned} o_1 &= 0.10 \\ o_2 &= 0.99 \end{aligned}$$

$$b^{(2)} = \begin{bmatrix} 0.60 \\ 0.60 \end{bmatrix}$$

$$b^{(1)} = \begin{bmatrix} 0.35 \\ 0.35 \end{bmatrix}$$

k	a	h
0	-	$\begin{bmatrix} 0.05 \\ 0.10 \end{bmatrix}$
1	$\begin{bmatrix} 0.15 * 0.05 + 0.20 * 0.10 + 0.35 \\ 0.25 * 0.05 + 0.30 * 0.10 + 0.35 \end{bmatrix} = \begin{bmatrix} 0.5933 \\ 0.5969 \end{bmatrix}$	$\begin{bmatrix} 0.5933 \\ 0.5969 \end{bmatrix}$
2	$\begin{bmatrix} 0.40 * 0.5933 + 0.45 * 0.5969 + 0.60 \\ 0.50 * 0.5933 + 0.55 * 0.5969 + 0.60 \end{bmatrix} = \begin{bmatrix} 1.1059 \\ 1.2249 \end{bmatrix}$	$\begin{bmatrix} 0.7514 \\ 0.7729 \end{bmatrix}$

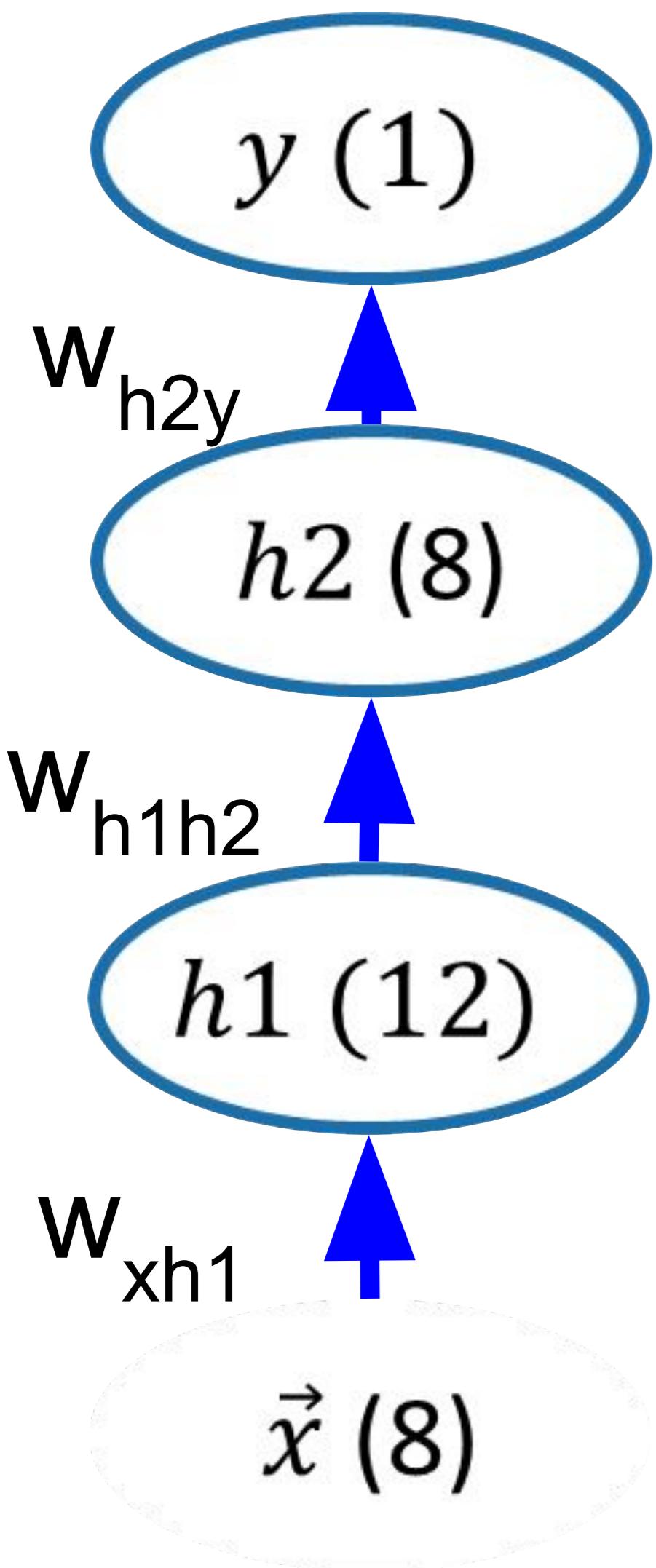
$$\hat{y} = \begin{bmatrix} 0.7514 \\ 0.7729 \end{bmatrix}$$

$$J = \frac{1}{2} [(0.01 - 0.7514)^2 + (0.99 - 0.7729)^2]$$

$$J = \mathbf{0.2984}$$

Implementing FFNN on Keras

Compact style



```
from keras.models import Sequential
from keras.layers import Dense

# define the keras model
model = Sequential()
model.add(Dense(12, input_dim=8, activation='relu'))
model.add(Dense(8, activation='relu'))
model.add(Dense(1, activation='sigmoid'))

# compile the keras model
model.compile(loss='binary_crossentropy', optimizer='adam',
metrics=['accuracy'])

# fit the keras model on the dataset
model.fit(X, y, epochs=150, batch_size=10)
```

Number of Parameter

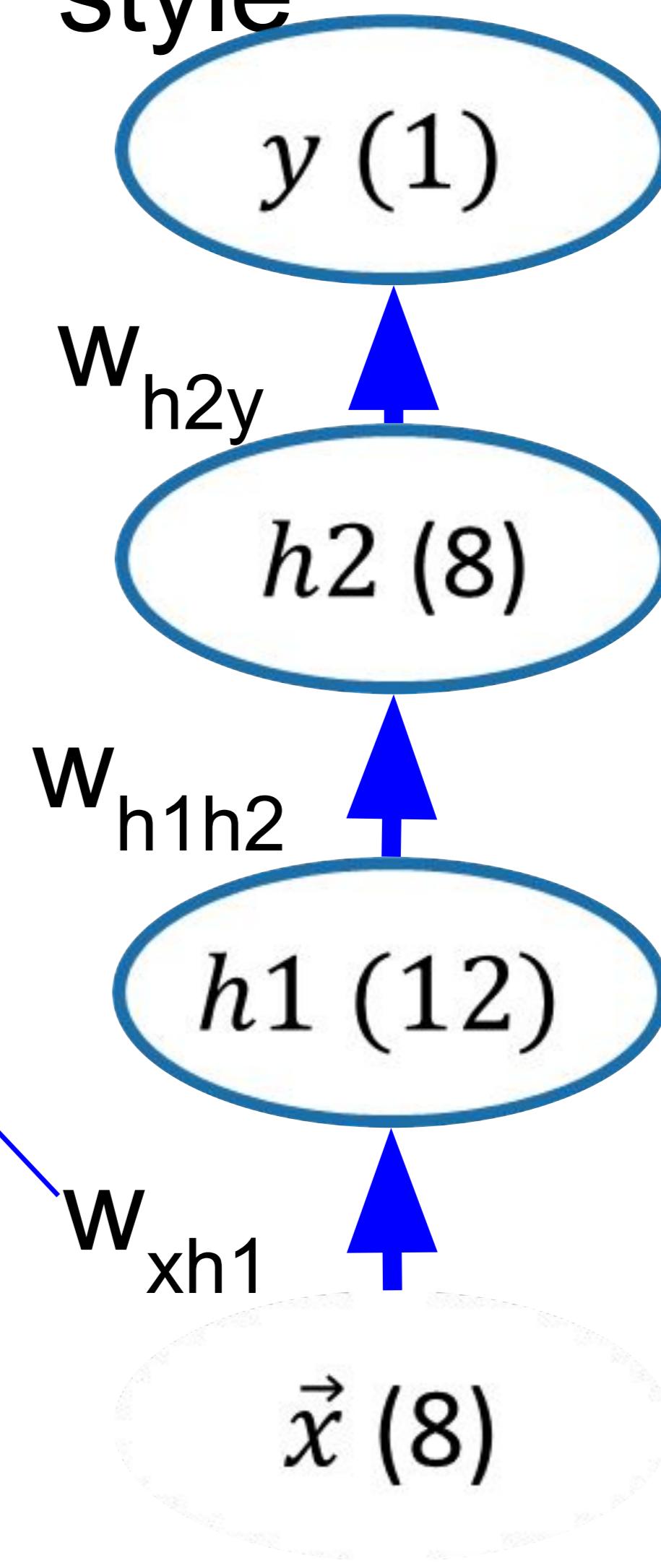
```
model.summary()
```

Model: "sequential"

Layer (type)	Output Shape	Param #
dense (Dense)	(None, 12)	108
dense_1 (Dense)	(None, 8)	104
dense_2 (Dense)	(None, 1)	9
Total params:	221	
Trainable params:	221	

$$\text{Total parameter} = (8+1)*12 + (12+1)*8 + (8+1)*1 = 221$$

Compact
style



Backpropagation (Learning)

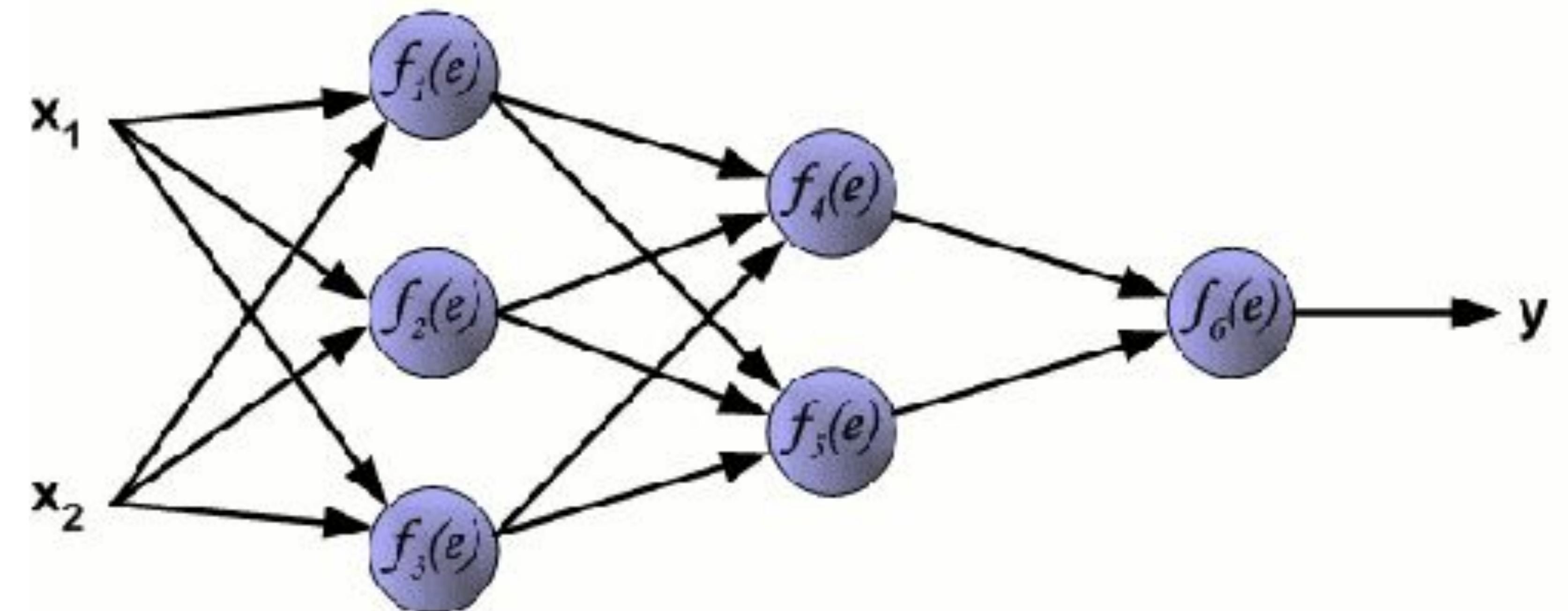
Learning Procedure

Goal: set weights to minimize error

Forward propagation: hitung output dr setiap neuron
compute output y for input instance x

Backward propagation: hitung error term utk setiap neuron
compute error term δ for each neuron (target hidden neuron is Not Available)

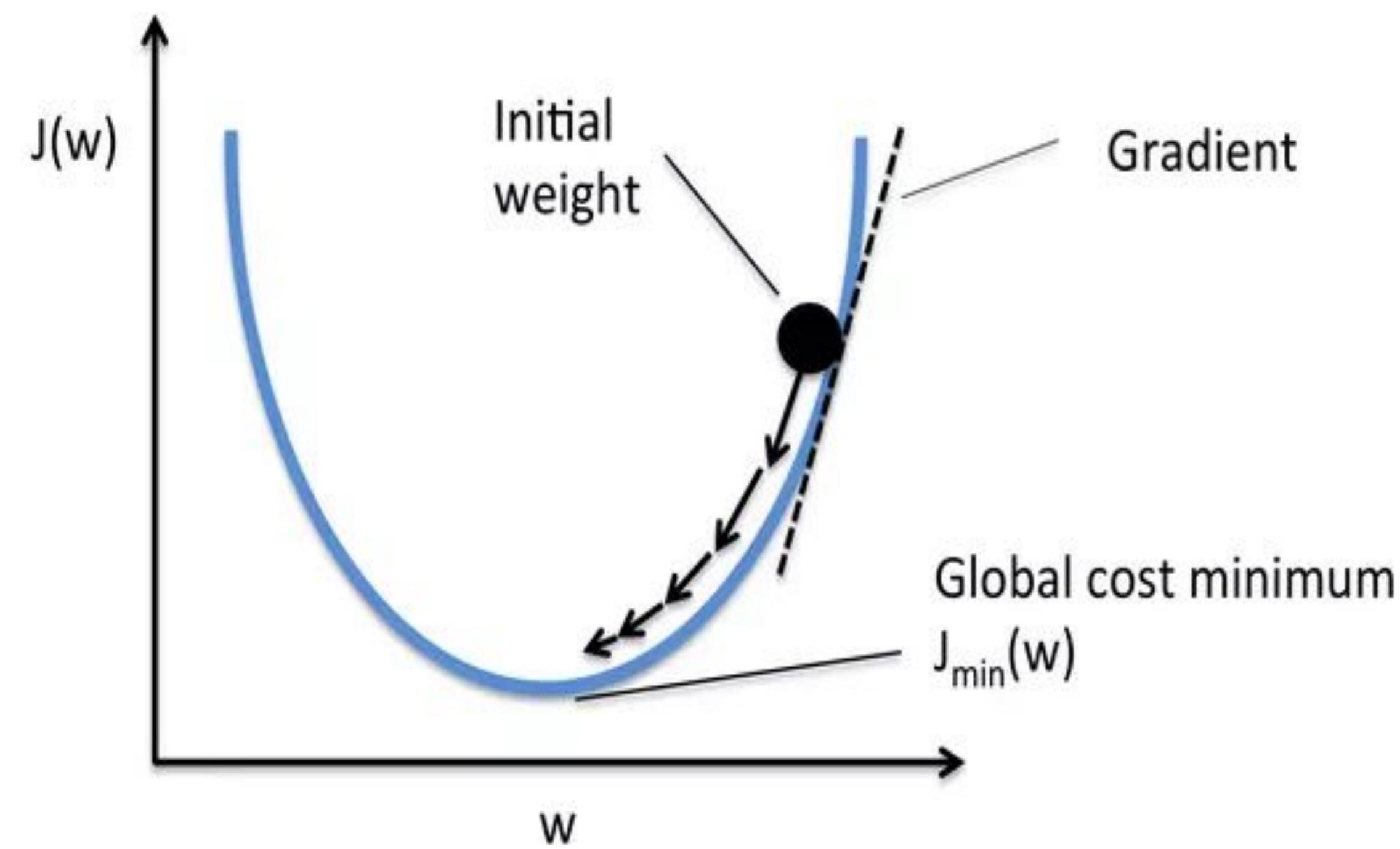
Update every weights update trsp bobot yg menghubungkan neuron



Review Gradient Descent:

Steepest Descent along Error Surface

Gradient descent: a strategy for searching through a large or **infinite hypothesis space** that can be applied whenever (1) the hypothesis space contains **continuously parameterized hypotheses** (e.g., the weights in a linear unit), and (2) the error can be **differentiated** with respect to these hypothesis parameters.



<https://medium.com/data-science-group-iitr/loss-functions-and-optimization-algorithms-demystified-bb92daff331c>

Gradient

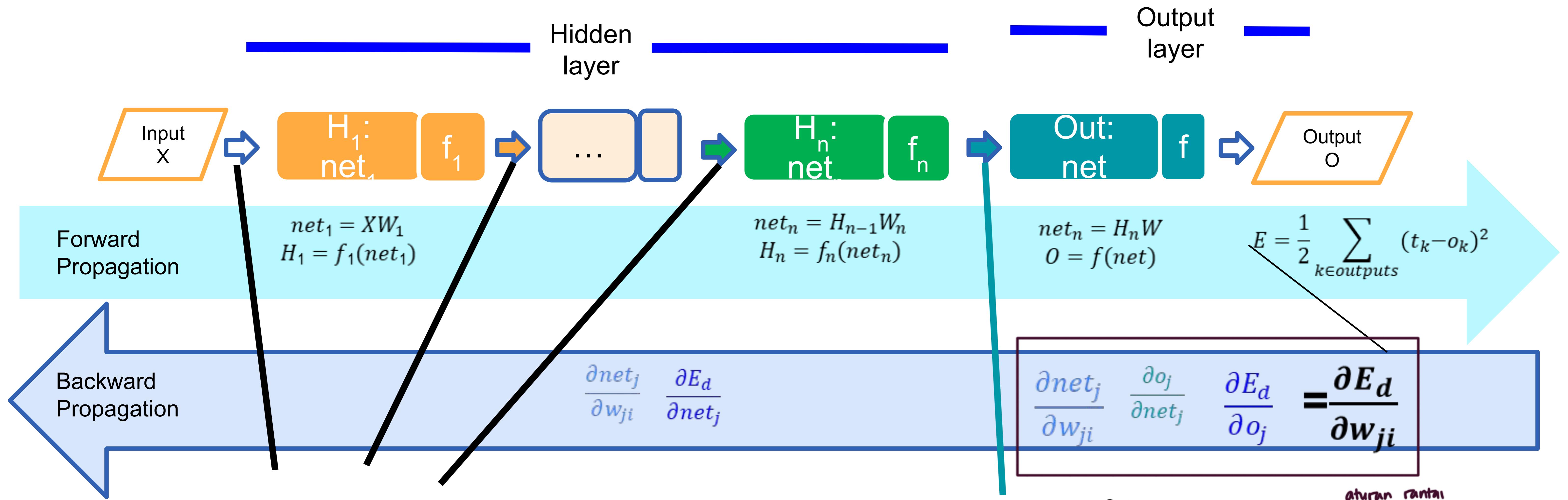
$$\nabla E[\vec{w}] \equiv \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

$$\vec{w} \leftarrow \vec{w} + \Delta \vec{w}$$

$$\Delta \vec{w} = -\eta \nabla E(\vec{w})$$

Backpropagation for MLP

w_{ji} : weight from unit i to unit j



$$w_{ji} = w_{ji} - \eta \frac{\partial E_d}{\partial w_{ji}}$$

$$= w_{ji} - \eta [-o_j(1 - o_j) \sum_{k \in outputs(j)} -\delta_k w_{kj} x_{ji}] = w_{ji} + \eta \delta_j x_{ji}$$

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial net_j} \frac{\partial net_j}{\partial w_{ji}} = -o_j(1 - o_j) \sum_{k \in outputs(j)} -\delta_k w_{kj} x_{ji}$$

aturan rantai

$$w_{ji} = w_{ji} - \eta \frac{\partial E_d}{\partial w_{ji}}$$

$$= w_{ji} - \eta [-(t_j - o_j) o_j (1 - o_j) x_{ji}]$$

$$= w_{ji} + \eta \delta_j x_{ji}$$

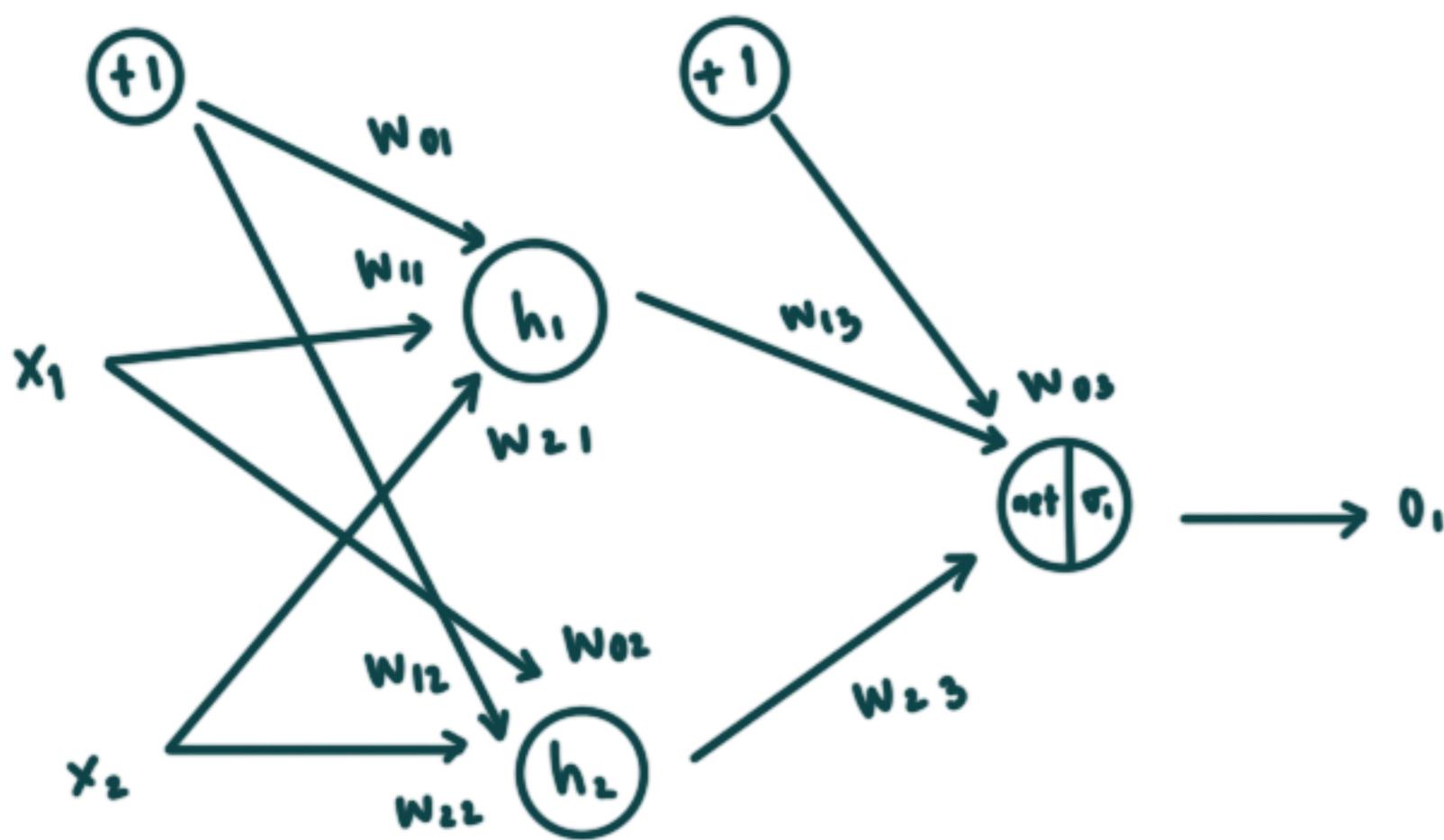
buat mau ke output doang

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial o_j} \frac{\partial o_j}{\partial net_j} \frac{\partial net_j}{\partial w_{ji}} = -(t_j - o_j) o_j (1 - o_j) x_{ji}$$

$$x \rightarrow h_1 \rightarrow h_2 \rightarrow h_3 \rightarrow \dots \rightarrow h_n \rightarrow 0$$

$$E = \frac{1}{2} \sum (t_k - o_k)^2$$

$$\sigma'(\text{net}) = \frac{\sigma(\text{net})}{o_1} \cdot \frac{(1-\sigma(\text{net}))}{(1-o_1)}$$



$$h_1 = \sigma \left(\underbrace{w_{01} + w_{11} \cdot x_1 + w_{21} \cdot x_2}_{\text{net}_{11}} \right)$$

$$h_2 = \sigma \left(\underbrace{w_{02} + w_{12} \cdot x_1 + w_{22} \cdot x_2}_{\text{net}_{12}} \right)$$

$$o_1 = \sigma \left(\underbrace{w_{03} + w_{13} \cdot h_1 + w_{23} \cdot h_2}_{\text{net}} \right)$$

$$\frac{\partial E}{\partial w_{23}} = \frac{\partial E}{\partial o_1} \cdot \frac{\partial o_1}{\partial \text{net}} \cdot \frac{\partial \text{net}}{\partial w_{23}} = -(t_1 - o_1) \cdot o_1 (1 - o_1) \cdot h_2$$

$$\frac{\partial E}{\partial o_1} = \frac{1}{2} \cdot 2 \cdot (t_1 - o_1) \cdot -1 = -(t_1 - o_1)$$

$$\frac{\partial o_1}{\partial \text{net}} = o_1 (1 - o_1)$$

$$\frac{\partial E}{\partial w_{13}} = -(t_1 - o_1) \cdot o_1 (1 - o_1) \cdot h_1$$

$$\frac{\partial \text{net}}{\partial w_{23}} = h_2$$

$$\frac{\partial E}{\partial w_{23}} = \underbrace{\frac{\partial E}{\partial o_1} \cdot \frac{\partial o_1}{\partial \text{net}}}_{\delta_1} \cdot \frac{\partial \text{net}}{\partial w_{23}} = - (t_1 - o_1) \cdot o_1 (1 - o_1) \cdot h_2$$

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial \text{net}_k} \cdot \frac{\partial \text{net}_k}{\partial o_j} \cdot \frac{\partial o_j}{\partial \text{net}_j} \cdot \underbrace{\frac{\partial \text{net}_j}{\partial w_{ji}}}_{x_{ji}}$$

$$\text{net}_0 = \sum w_{oi} \cdot x_{ji}$$

$$w_{ji} \quad \left| \begin{array}{c} \text{net}_j \\ \text{o}_j \rightarrow \text{net}_k \\ \text{o}_j \end{array} \right|$$

Backpropagation for MLP with Sigmoid Unit

Steepest gradient descent proposes new weight with negative gradient and Sum of Squared Error (SSE).

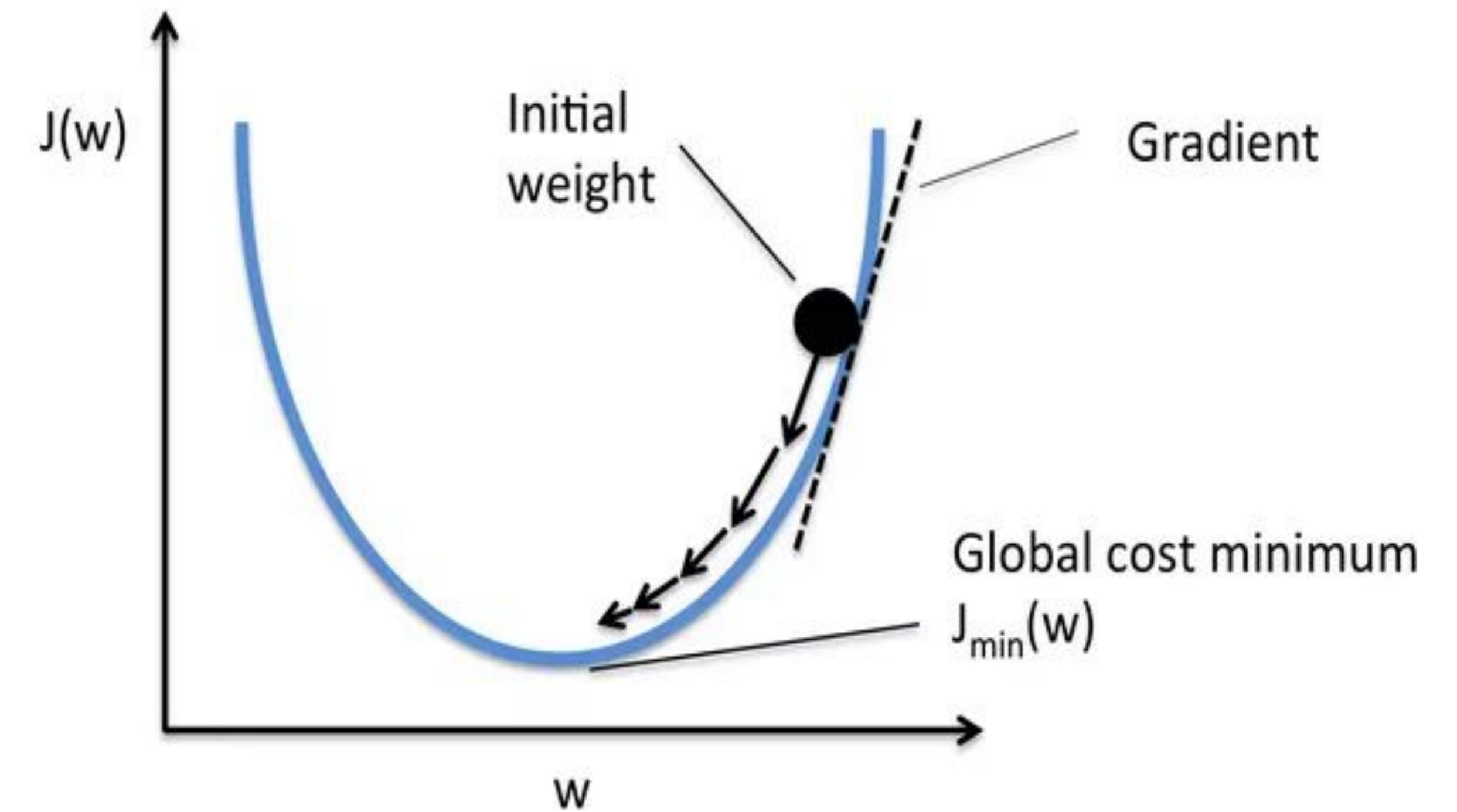
$$w_{ji} = w_{ji} + \Delta w_{ji} = w_{ji} - \eta \frac{\partial E_d}{\partial w_{ji}} = w_{ji} + \eta \delta_j x_{ji}$$

Error term for each output unit j :

$$\delta_j = o_j(1 - o_j)(t_j - o_j)$$

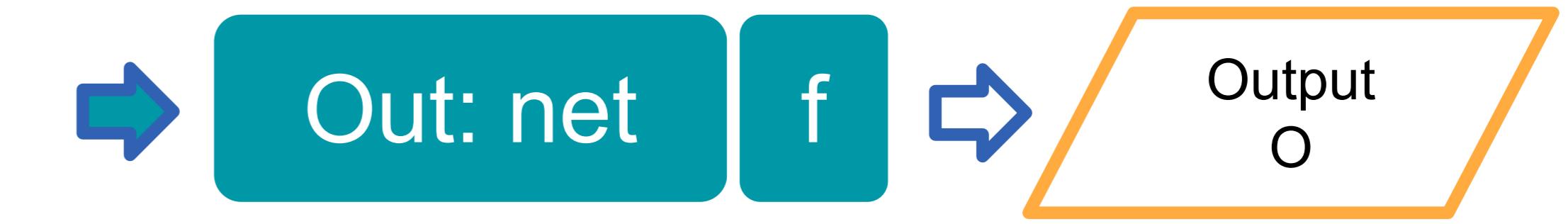
Error term for each hidden unit k :

$$\delta_h = o_h(1 - o_h) \sum_{k \in outputs} w_{kh} \delta_k$$



<https://medium.com/data-science-group-iitr/loss-functions-and-optimization-algorithms-demystified-bb92daff331c>

Computing Gradient for Output Unit Weights



Chain rule to compute gradient:

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial o_j} \frac{\partial o_j}{\partial net_j} \frac{\partial net_j}{\partial w_{ji}} = -(t_j - o_j) o_j (1 - o_j) x_{ji}$$

$$\frac{\partial net_j}{\partial w_{ji}} = x_{ji}; net_j = \sum_{i \in [0..nj]} x_{ji} w_{ji}$$

$$\frac{\partial o_j}{\partial net_j} = \frac{\partial \sigma(net_j)}{\partial net_j} = o_j (1 - o_j)$$

$$\frac{\partial E_d}{\partial o_j} = \frac{\partial}{\partial o_j} \frac{1}{2} \sum_{k \in outputs} (t_k - o_k)^2 = \frac{1}{2} \cdot 2 \cdot (t_j - o_j) \cdot -1 = -(t_j - o_j)$$

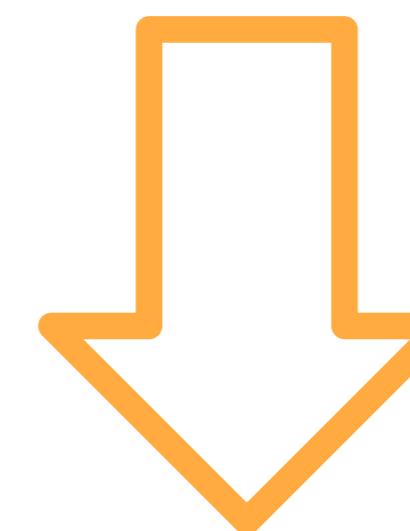
Computing Gradient for Output Unit Weights

Chain rule to compute gradient:



$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial o_j} \frac{\partial o_j}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ji}} = -(t_j - o_j) o_j(1 - o_j) x_{ji}$$

gradient



$$\delta_j = o_j(1 - o_j)(t_j - o_j)$$

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ji}} = -\delta_j x_{ji}$$

Computing Gradient for Hidden Unit Weights



Chain rule to compute gradient:

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ji}} = o_j(1 - o_j) \sum_{k \in \text{outputs}(j)} -\delta_k w_{kj} x_{ji}$$

$\frac{\partial \text{net}_j}{\partial w_{ji}} = x_{ji}; \text{net}_j = \sum_{i \in [0..n_j]} x_{ji} w_{ji}$

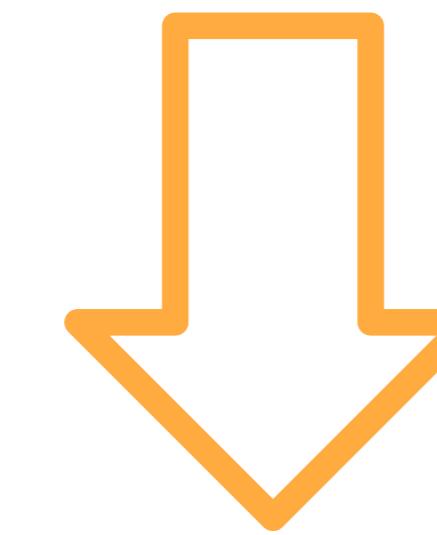
$$\frac{\partial E_d}{\partial \text{net}_j} = \sum_{k \in \text{outputs}(j)} \frac{\partial E_d}{\partial \text{net}_k} \frac{\partial \text{net}_k}{\partial o_j} \frac{\partial o_j}{\partial \text{net}_j} = \sum_{k \in \text{outputs}(j)} -\delta_k w_{kj} o_j(1 - o_j)$$
$$\frac{\partial E_d}{\partial \text{net}_j} = -\delta_j$$

Computing Gradient for Hidden Unit Weights



Chain rule to compute gradient:

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ji}} = -o_j(1 - o_j) \sum_{k \in \text{outputs}(j)} \delta_k w_{kj} x_{ji}$$



$$\delta_h = o_h(1 - o_h) \sum_{k \in \text{outputs}(h)} \delta_k w_{kh}$$

$$\frac{\partial E_d}{\partial w_{ji}} = -\delta_h x_{ji}$$

BACKPROPAGATION(*training_examples*, η , n_{in} , n_{out} , n_{hidden})

Each training example is a pair of the form (\vec{x}, \vec{t}) , where \vec{x} is the vector of network input values, and \vec{t} is the vector of target network output values.

η is the learning rate (e.g., .05). n_{in} is the number of network inputs, n_{hidden} the number of units in the hidden layer, and n_{out} the number of output units.

The input from unit i into unit j is denoted x_{ji} , and the weight from unit i to unit j is denoted w_{ji} .

- Create a feed-forward network with n_{in} inputs, n_{hidden} hidden units, and n_{out} output units.
- Initialize all network weights to small random numbers (e.g., between -.05 and .05).
- Until the termination condition is met, Do
 - For each (\vec{x}, \vec{t}) in *training_examples*, Do

Propagate the input forward through the network:

1. Input the instance \vec{x} to the network and compute the output o_u of every unit u in the network.

Propagate the errors backward through the network:

2. For each network output unit k , calculate its error term δ_k

$$\delta_k \leftarrow o_k(1 - o_k)(t_k - o_k) \quad (\text{T4.3})$$

3. For each hidden unit h , calculate its error term δ_h

$$\delta_h \leftarrow o_h(1 - o_h) \sum_{k \in outputs} w_{kh}\delta_k \quad (\text{T4.4})$$

4. Update each network weight w_{ji}

$$w_{ji} \leftarrow w_{ji} + \Delta w_{ji}$$

where

$$\Delta w_{ji} = \eta \delta_j x_{ji} \quad (\text{T4.5})$$

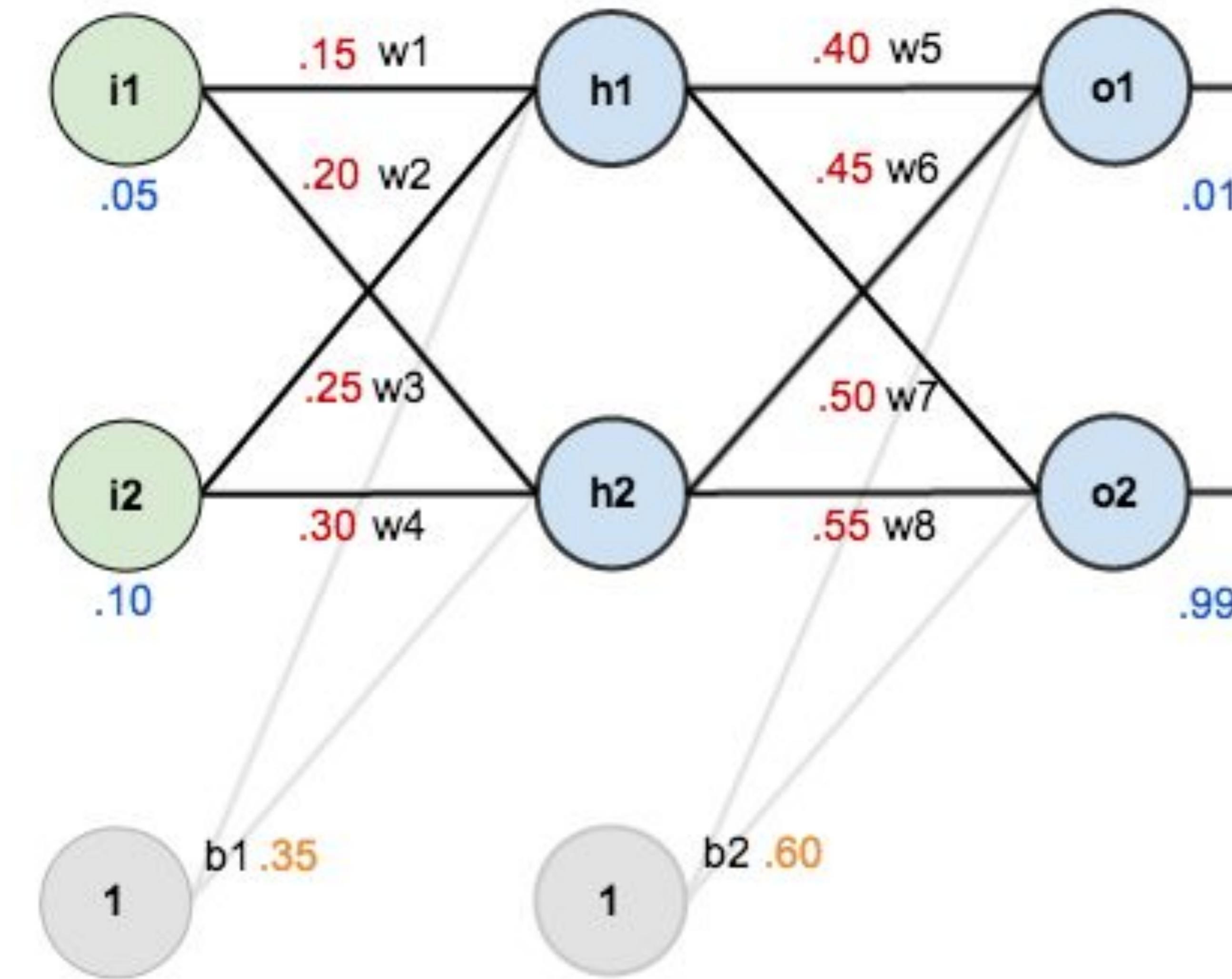
Backpropagation Algorithm (Mitchell, 1997)

Termination condition:

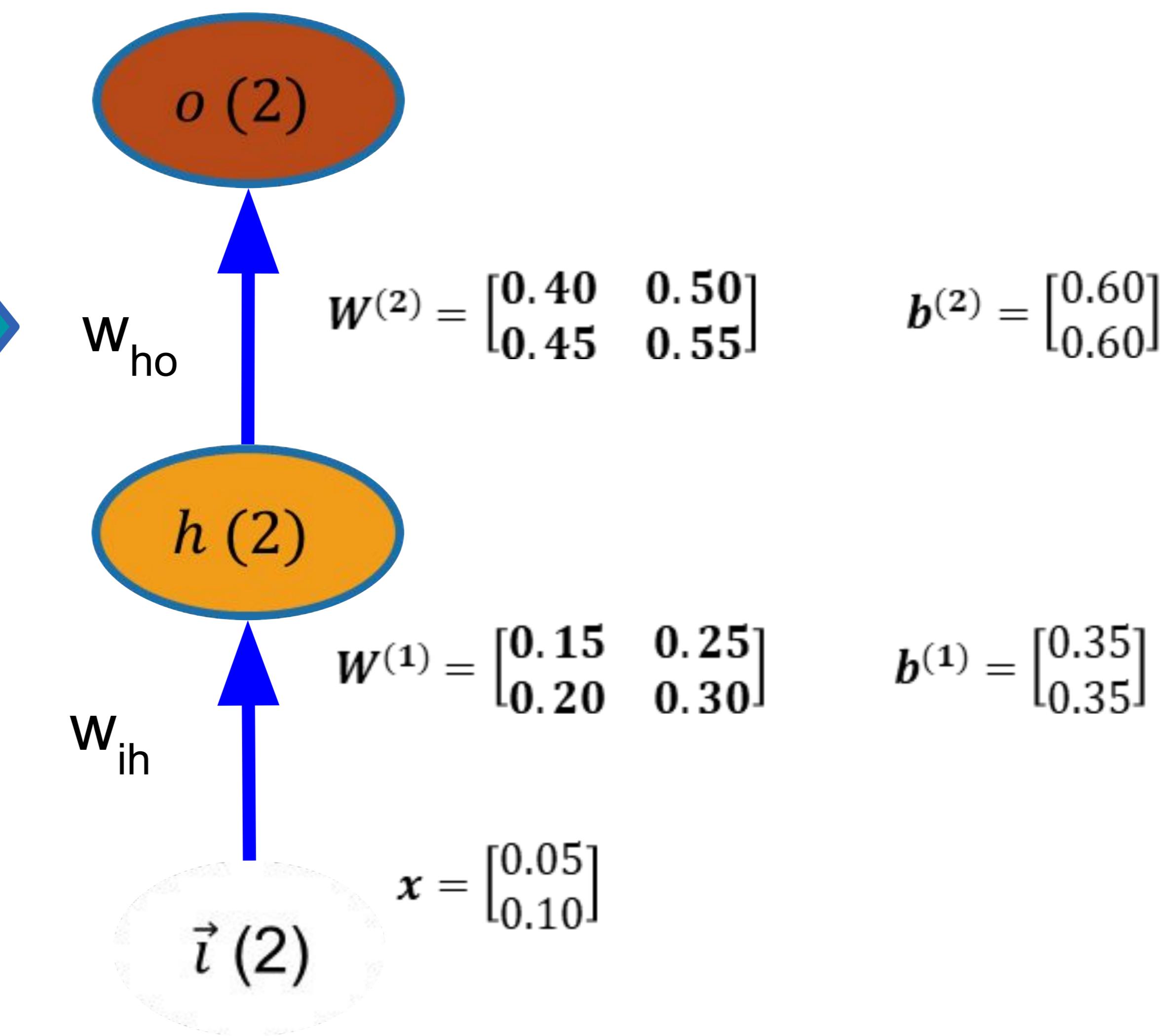
- a fixed number of iterations
- the error on the training examples < threshold
- the error on a separate validation set of examples meets some criterion

Backpropagation: Example

(from <http://mattmazur.com/2015/03/17/a-step-by-step-backpropagation-example/>)



Compact style: sigmoid model



Backpropagation: Example

Learning rate: $\eta = 0.5$

k	h
0	$\begin{bmatrix} 0.05 \\ 0.10 \end{bmatrix}$
1	$\begin{bmatrix} 0.5933 \\ 0.5969 \end{bmatrix}$
2	$\begin{bmatrix} 0.7514 \\ 0.7729 \end{bmatrix}$

$$\delta_{o_1} = 0.7514 * (1 - 0.7514) * (0.01 - 0.7514) = -0.1385$$

$$\delta_{o_2} = 0.7729 * (1 - 0.7729) * (0.99 - 0.7729) = 0.0381$$

$$\begin{aligned}\delta_{h_1} &= 0.5933 * (1 - 0.5933) * [w_5 * \delta_{o_1} + w_7 * \delta_{o_2}] \\ &= 0.2413 * [0.4 * (-0.1385) + 0.5 * 0.0381]\end{aligned}$$

Calculate also for δ_{h_2}

Example updating for w_5 ($w_{o_1 h_1}$)

$$\begin{aligned}w_5 &= 0.4 + (0.5 * \delta_{o_1} * h_1) \\ &= 0.4 + (0.5 * (-0.1385) * 0.5933) \\ &= 0.3589\end{aligned}$$

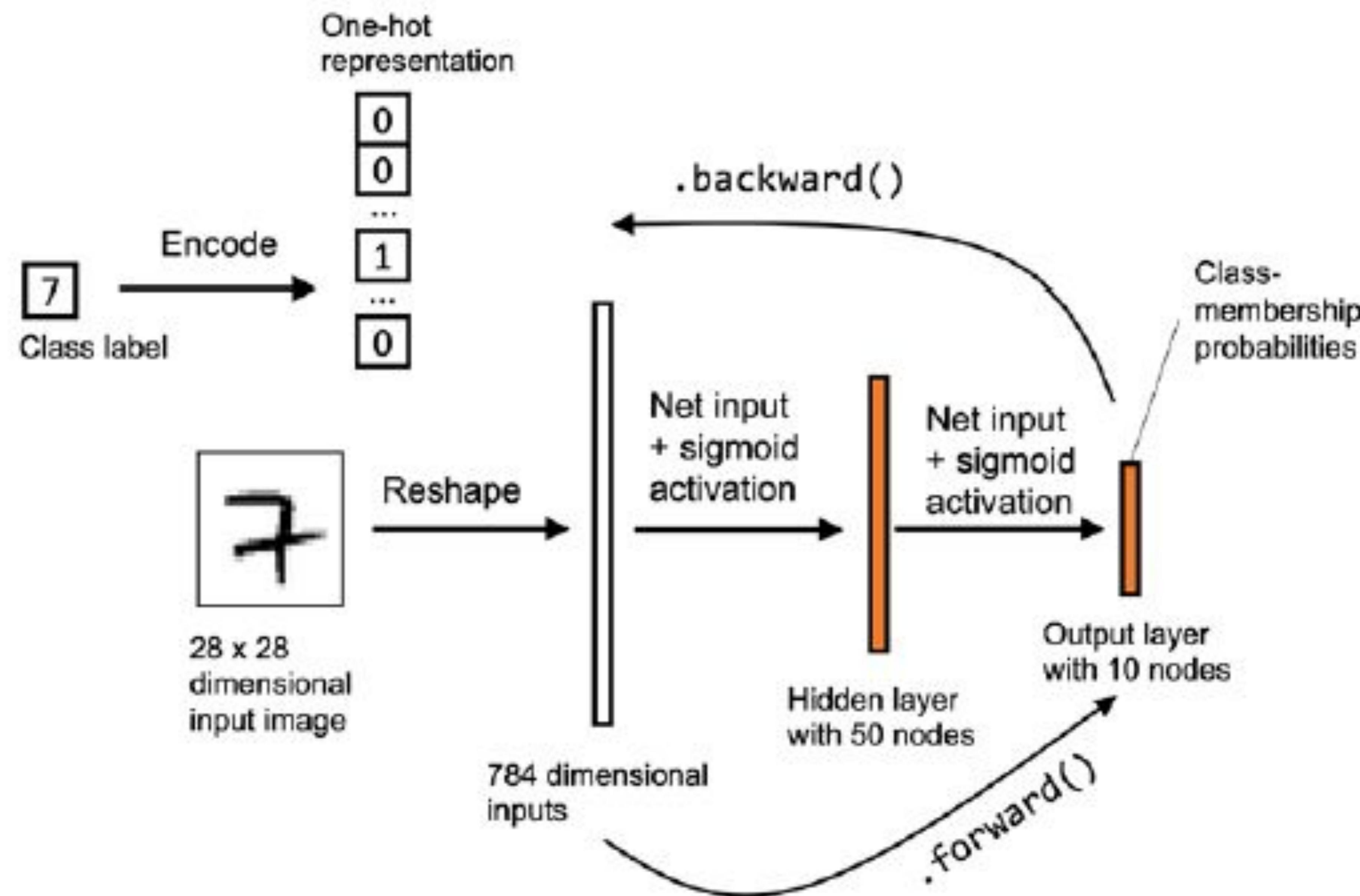
$$\delta_k \leftarrow o_k(1 - o_k)(t_k - o_k)$$

$$\delta_h \leftarrow o_h(1 - o_h) \sum_{k \in outputs} w_{kh} \delta_k$$

$$w_{ji} \leftarrow w_{ji} + \Delta w_{ji}$$

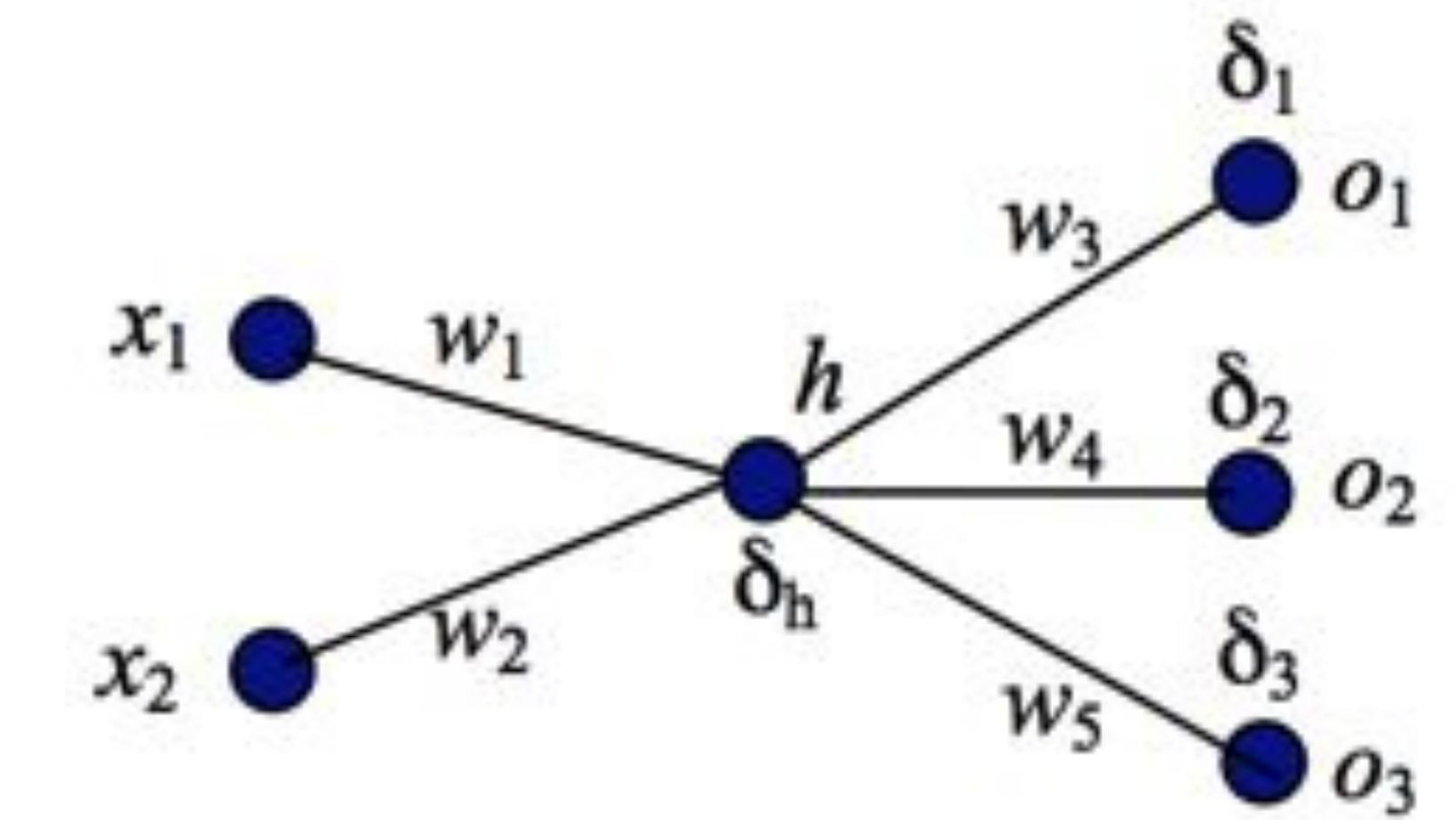
$$\Delta w_{ji} = \eta \delta_j x_{ji}$$

Example of NN Architecture for labeling handwritten digits



Raschka, *Figure 11.6: The NN architecture for labeling handwritten digits*

Latihan 1



Diketahui sebuah arsitektur jaringan Neural Network sebagai berikut yang memiliki dua unit input x , tiga unit output o dan satu unit hidden h , serta menggunakan fungsi aktivasi Sigmoid di setiap unitnya. Proses pembelajaran menerapkan algoritma Backpropagation.

- Tuliskan formula untuk menghitung luaran h dan o_1
- Jika δ_1 , δ_2 dan δ_3 adalah *error* pada ketiga unit output, δ_h adalah *error* pada hidden unit, dan t_1 adalah nilai target untuk unit output yang pertama, tuliskan formula untuk menghitung δ_1 dan δ_h .
- Jika α adalah *learning rate*, tuliskan formula untuk *update* bobot w_1 dan w_3 .

Catatan: bias selalu ada sebagai input tiap layer, simbol bebas

$$h = \sigma(w_0 + w_1 \cdot x_1 + w_2 \cdot x_2)$$

$$o_1 = \sigma(w_6 + h \cdot w_3)$$

$$\delta_k = o_k (1 - o_k) (t_{ik} - o_k)$$

$$\delta_h = o_k (1 - o_k) \cdot \sum w_{kh} \cdot \delta_k$$

$$\delta_1 = o_1 (1 - o_1) \cdot (t_1 - o_1)$$

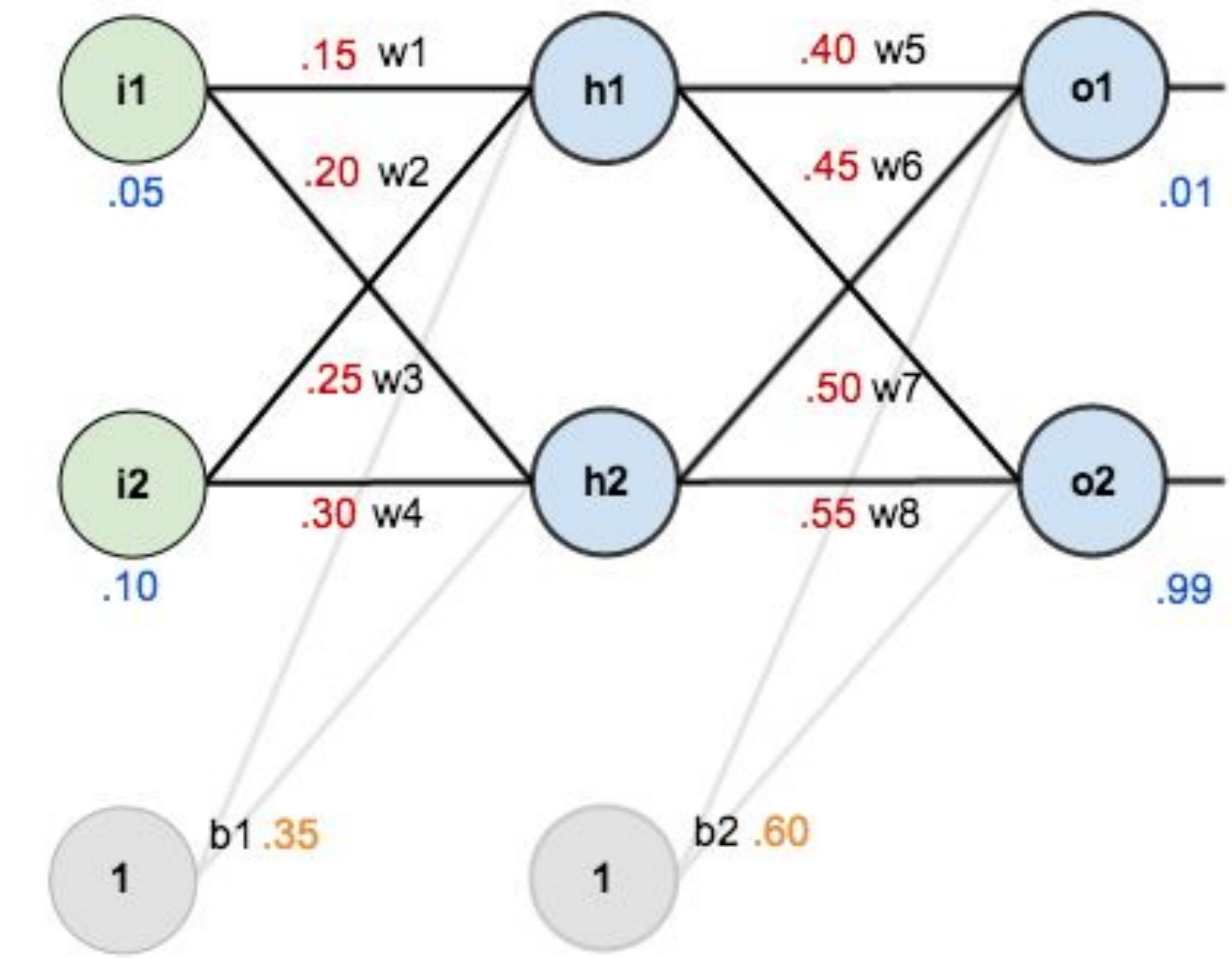
$$\delta_h = o_1 (1 - o_1) (w_3 \cdot \delta_1 + w_4 \cdot \delta_2 + w_5 \cdot \delta_3)$$

$$w_3 = w_3 + \eta \cdot \delta_1 \cdot h$$

$$w_1 = w_1 + \eta \cdot \delta_h \cdot x_1$$

Latihan 2

- Lakukan pembelajaran (backpropagation) untuk ANN gambar di samping untuk mendapatkan nilai semua bobot yang baru, lengkap dengan perhitungan.
- Terapkan bobot yang baru pada data input $i_1 = 0.05$ dan $i_2 = 0.1$
- Hitung error setelah update bobot, dan bandingkan dengan sebelum update bobot



Fungsi Aktivasi: Sigmoid
Learning Rate: 0.5

Thank you