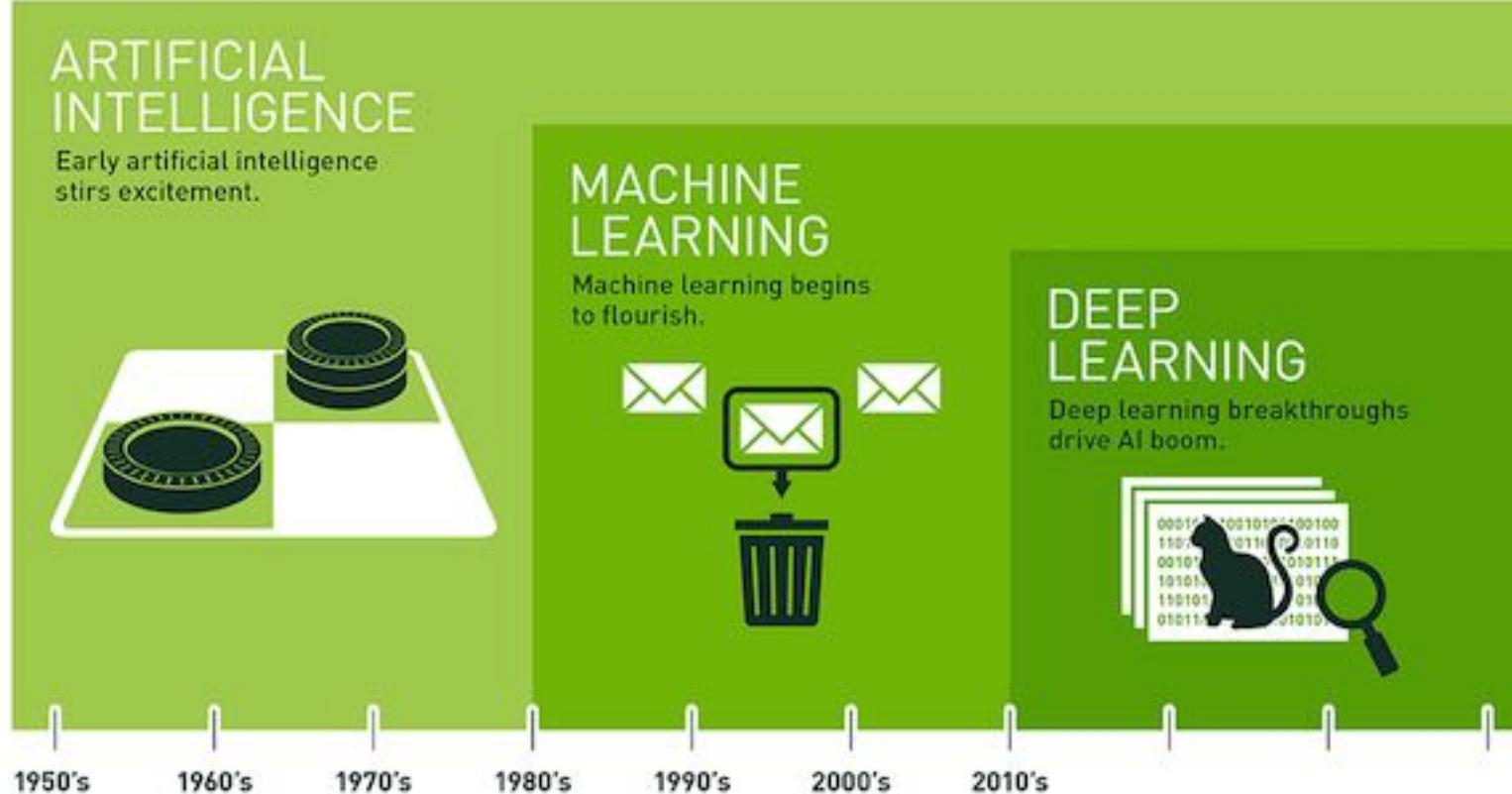


딥러닝

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



Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

# Why is Deep Learning Hot Now?

## Big Data Availability

**facebook**

350 millions  
images uploaded  
per day

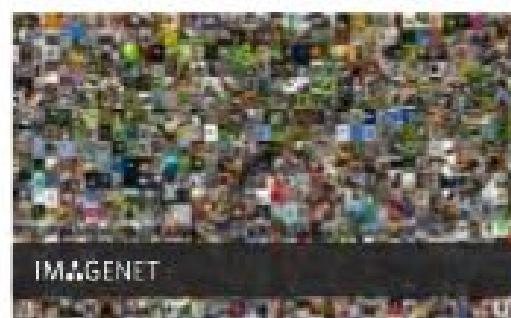
**Walmart**

2.5 Petabytes of  
customer data  
hourly

**YouTube**

300 hours of video  
uploaded every  
minute

## New ML Techniques



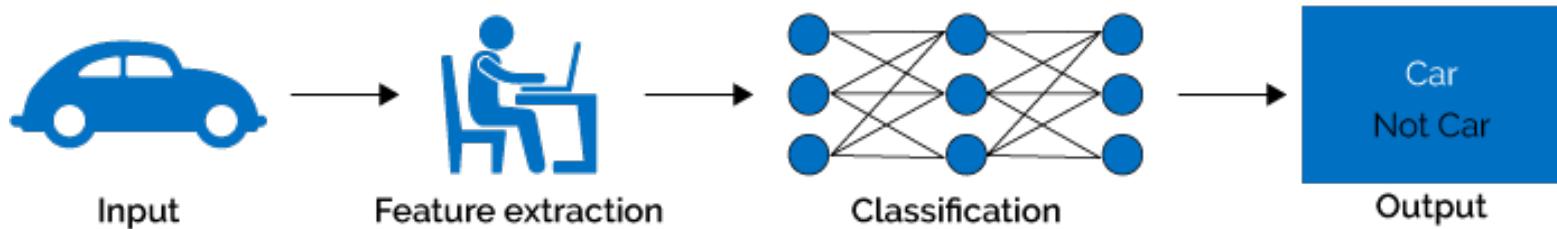
## GPU Acceleration



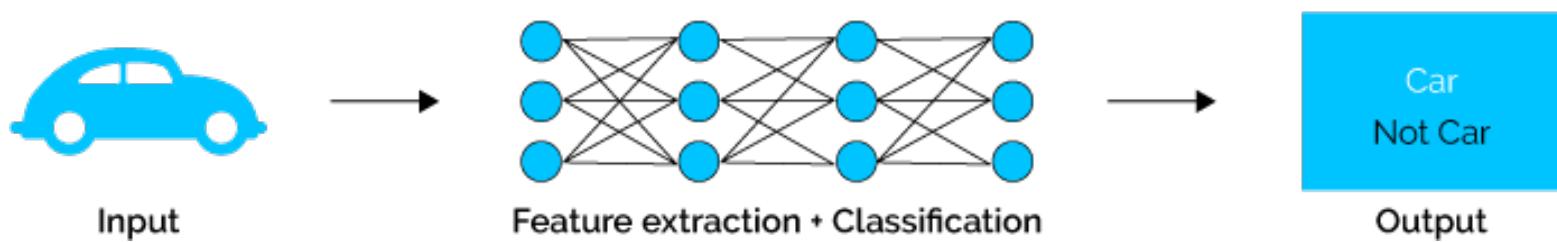
# Why is deep learning a growing trend?

- **few feature engineering**
- state-of-the-art performance

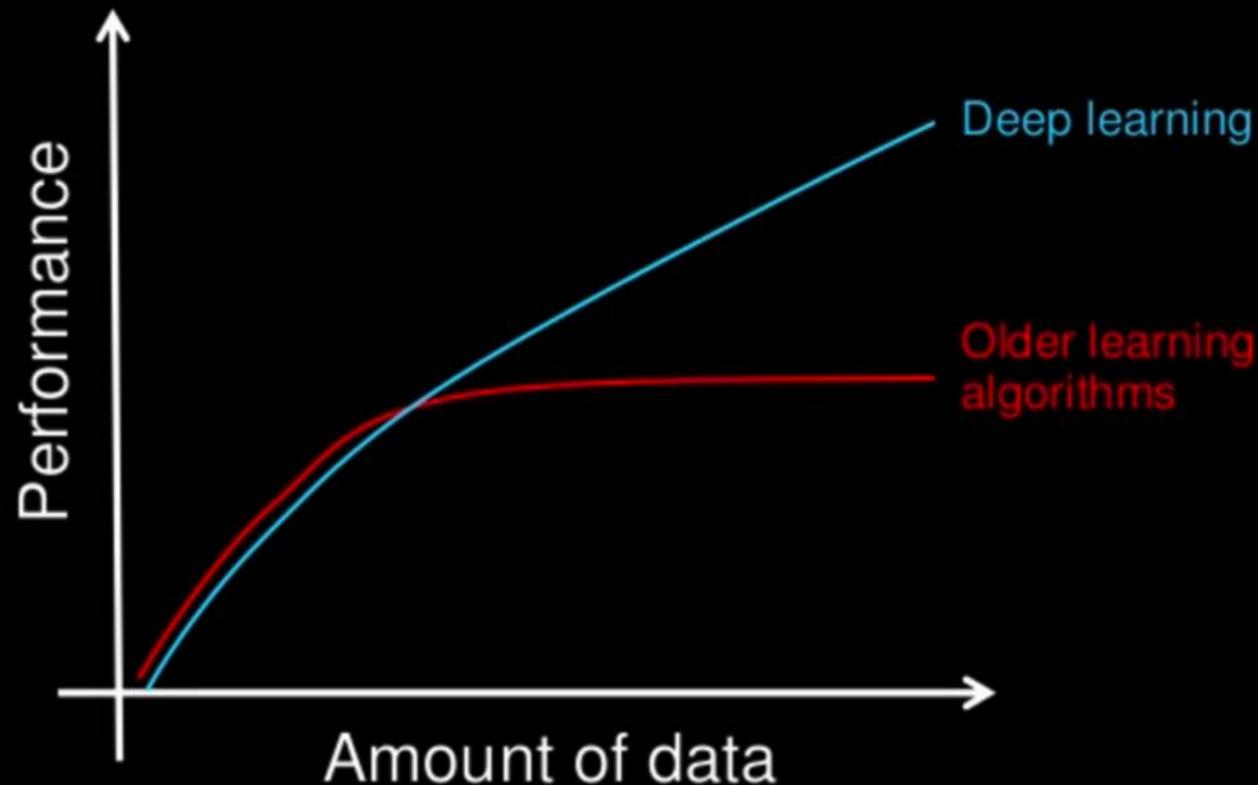
## Machine Learning



## Deep Learning



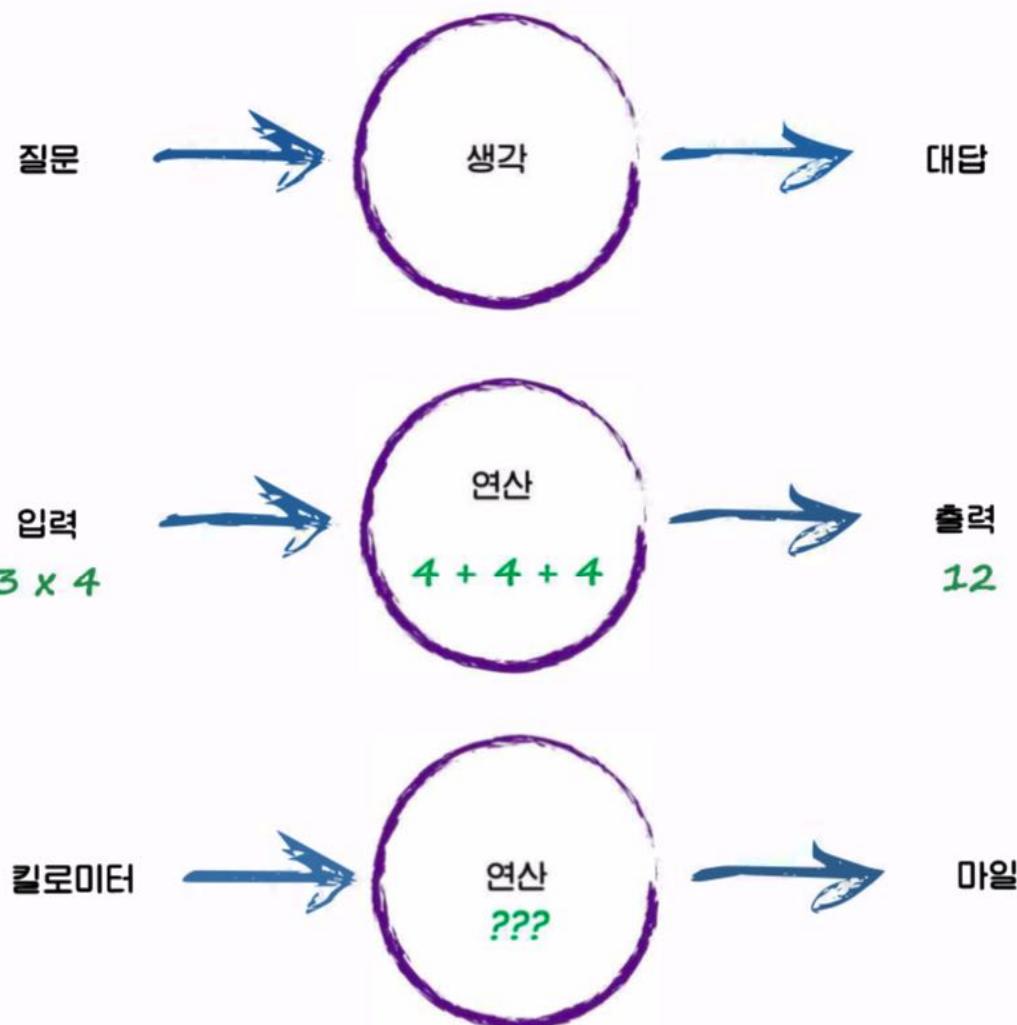
# Why deep learning

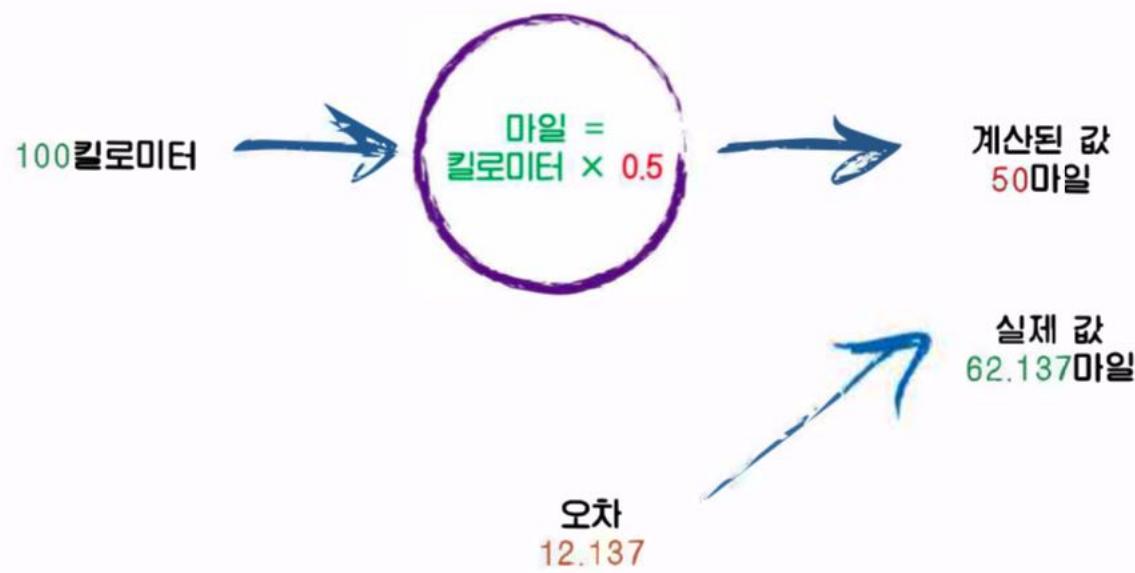


How do data science techniques scale with amount of data?

# 1. Machine Learning

---







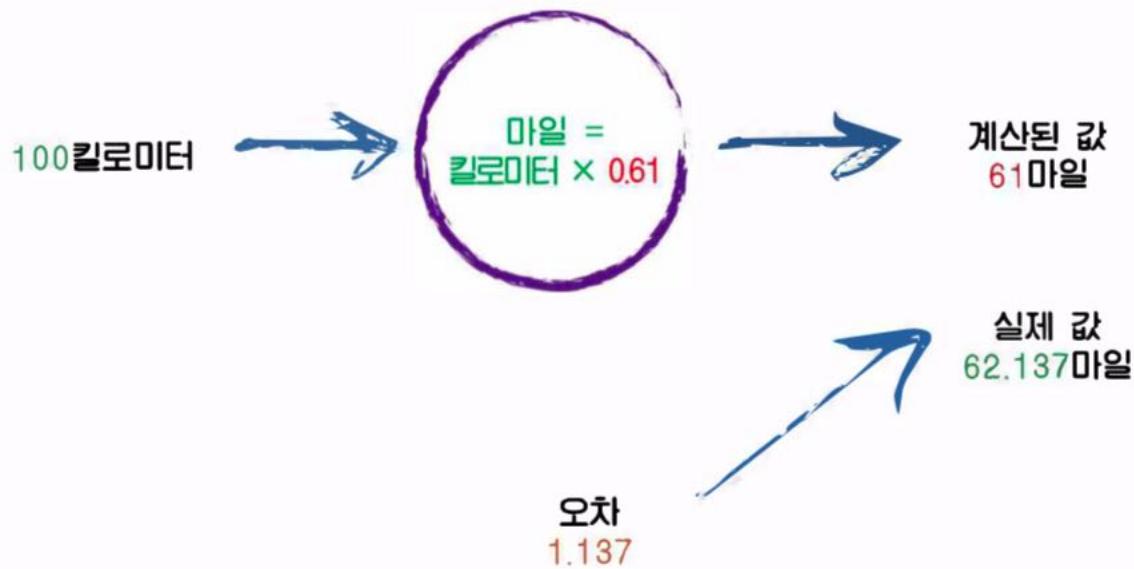
실제 값  
62.137마일

오차  
2.137



실제 값  
62.137마일

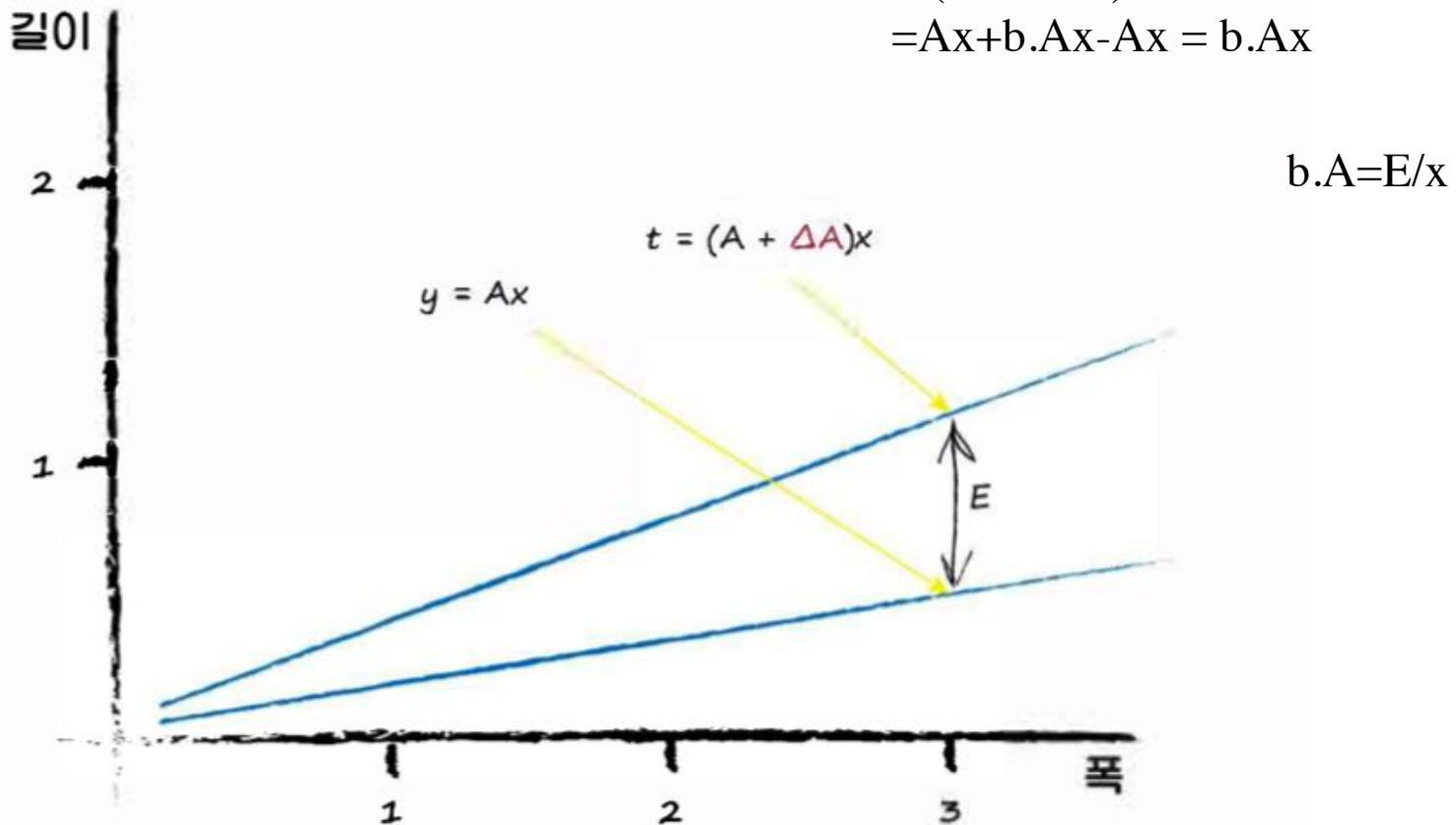
오차  
-7.863



어떤것의 동작원리를 정확히 파악할수 없을때 취할수있는한방

- 조정할 수 있는 매개변수 값을 포함하는 모델을 만들어보는 것
- 모델을 정교회해나가는 좋은 방법은 오차에 기초해 매개변수 값을 조정해나가는 것.

$$\begin{aligned}E &= t - y \\&= (A + t \cdot A)x - Ax \\&= Ax + b \cdot Ax - Ax = b \cdot Ax\end{aligned}$$

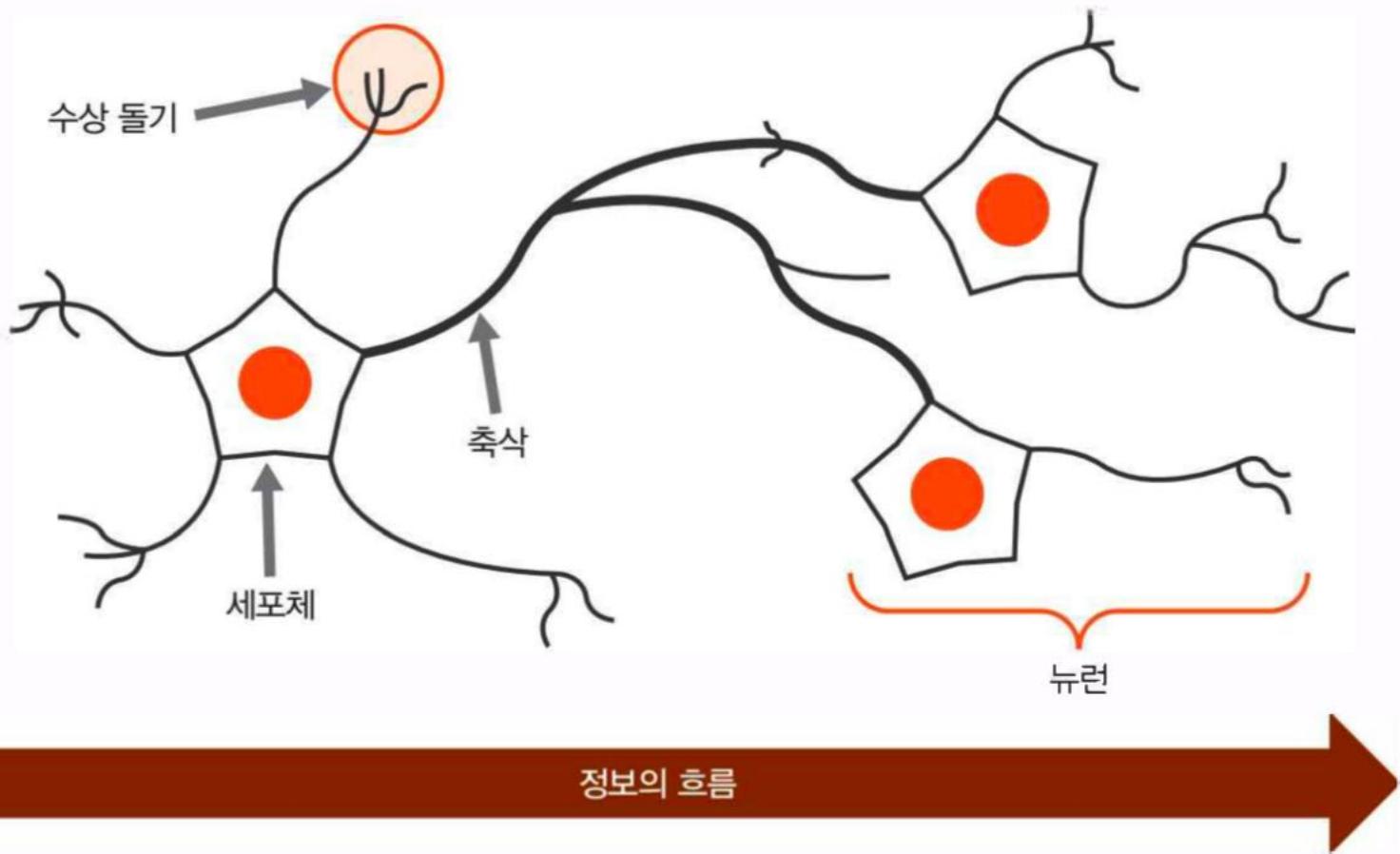


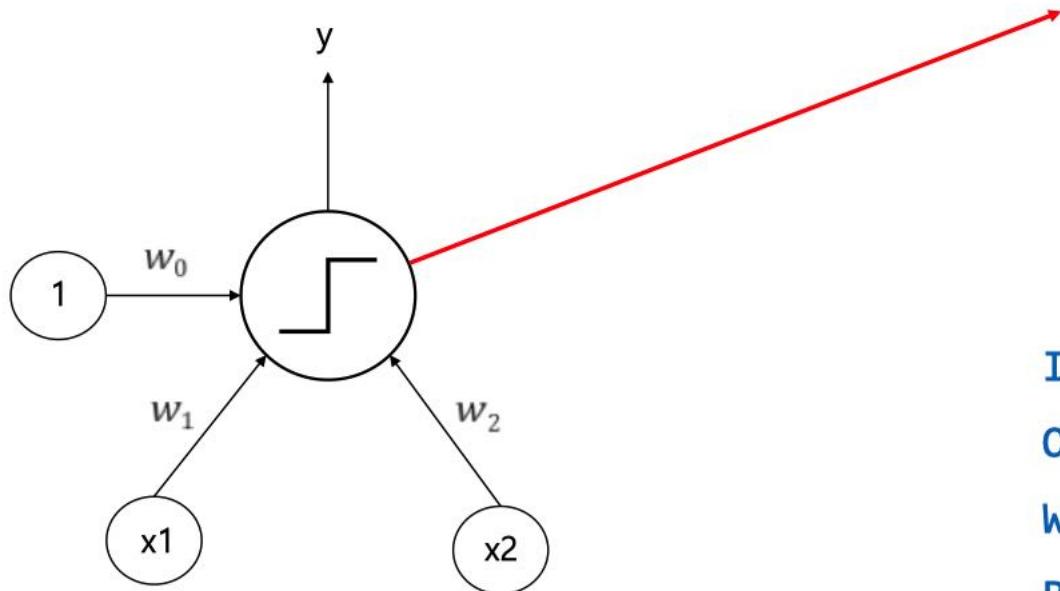
조정 과정의 문제점은 이전의 학습 데이터는 무시하고 최종 학습 데이터에만 맞춰 업데이트된다는 것  
이를 해결하기 위해 학습률을 도입해 업데이트의 정도를 조정 해줍니다.  
이를 통해 단일 학습 데이터가 학습에 지배적인 영향을 주는 것을 방지할 수 있습니다.  
현실에서 학습 데이터는 잡음이 섞여 있거나 오차를 가집니다. 학습률을 이용한 업데이트는 이러한 데이터의 오류의 영향을 제한하는 효과

# 1. Perceptron

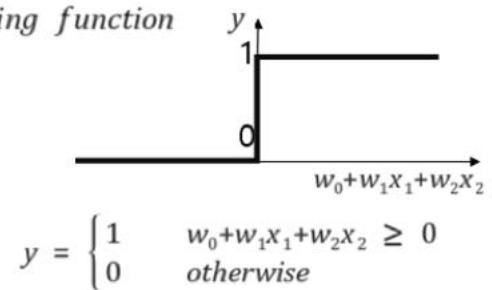
---

지능 = 뇌





*hard thresholding function*

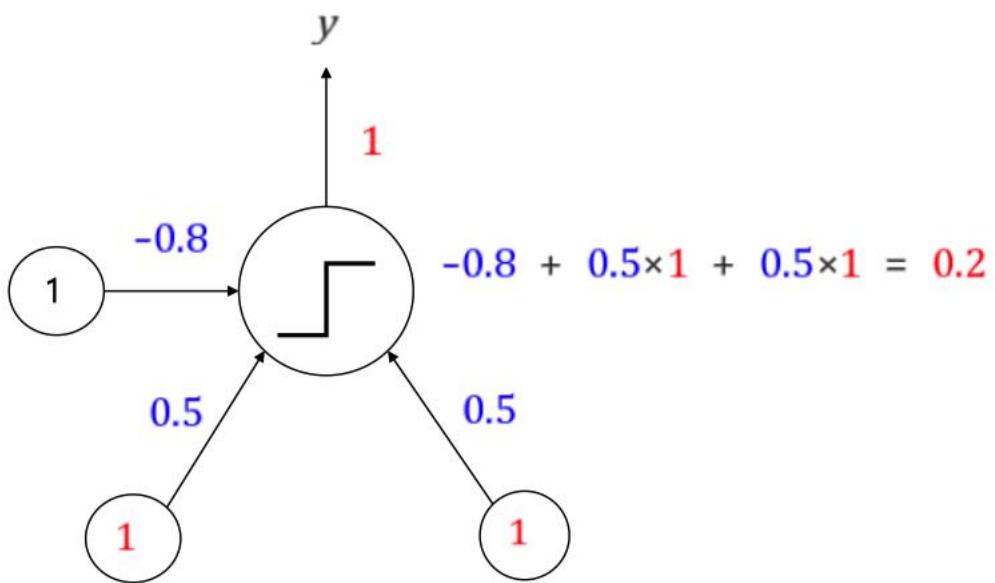


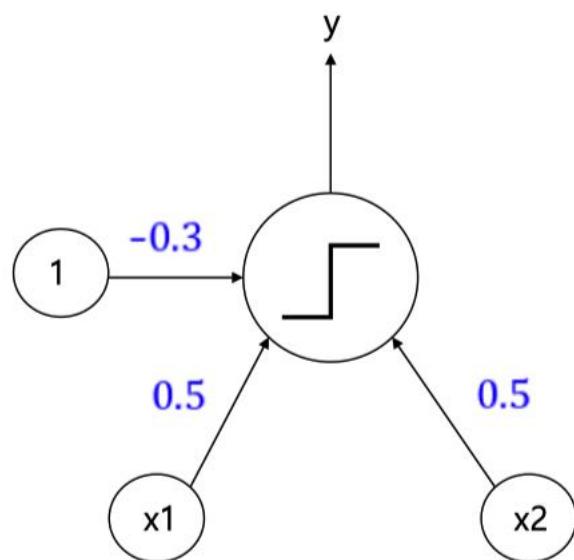
Inputs:  $x_1, x_2$

Output:  $y$

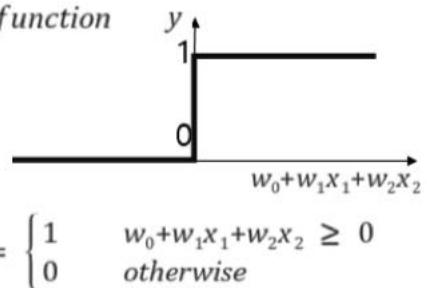
Weights:  $w_0, w_1, w_2$

Bias: 1





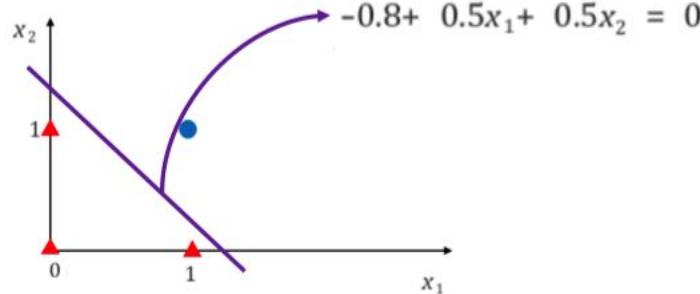
*hard thresholding function*



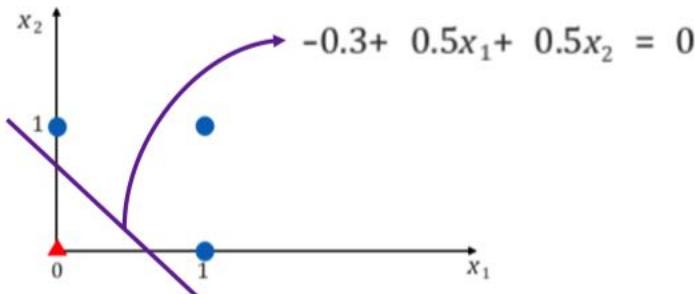
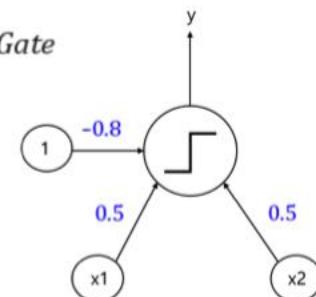
*OR Gate*

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

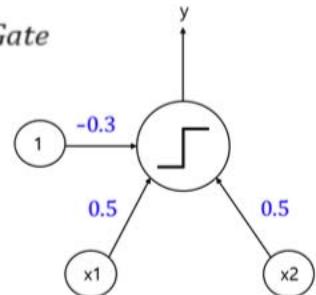
● : 1      ▲ : 0



*AND Gate*

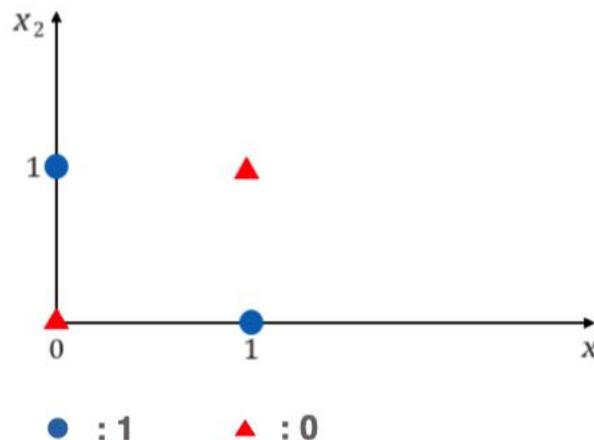


*OR Gate*



## Is it possible to solve XOR problem using a single layer perceptron?

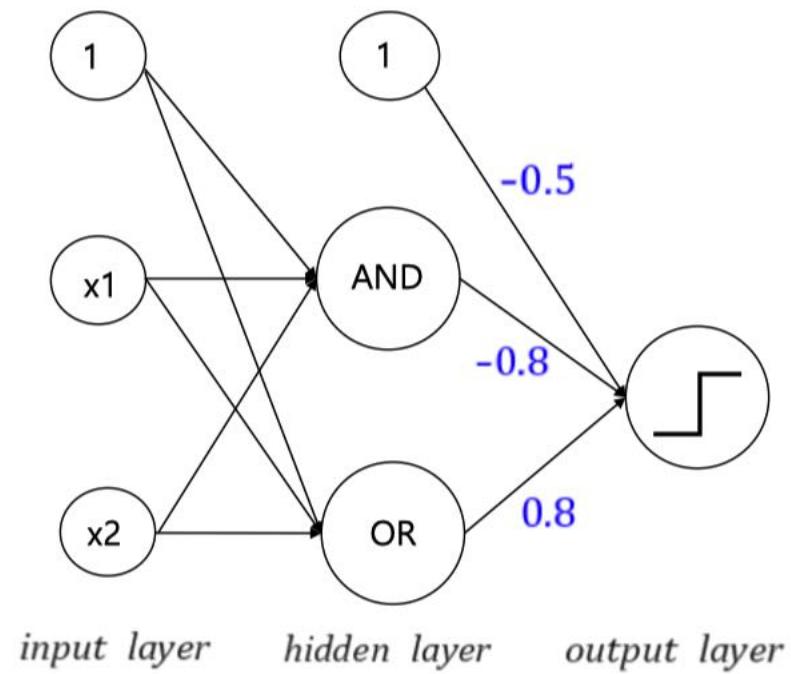
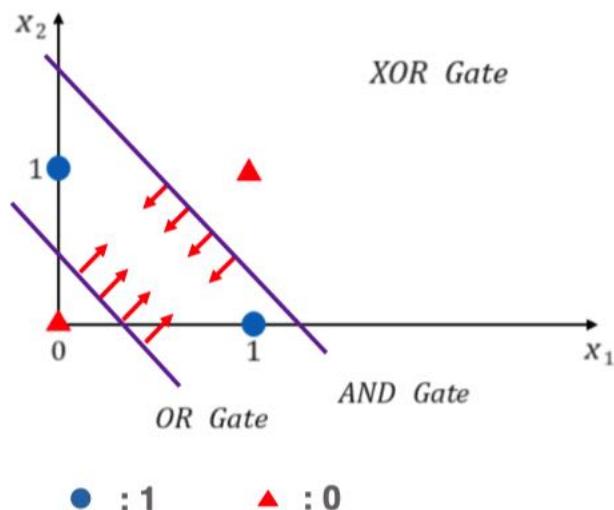
No. Single layer perceptron can only solve linear problem. XOR problem is non-linear.

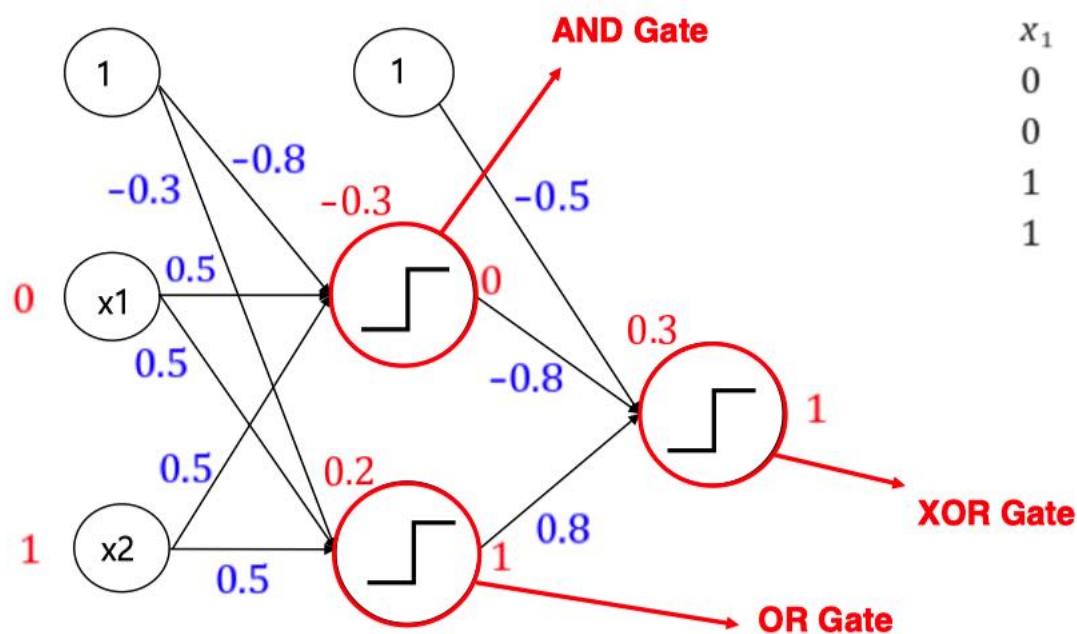


XOR Gate

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

**But if we use 2 single layer perceptron, we can solve XOR problem. This model is called multi layer perceptron.**



*OR Gate*

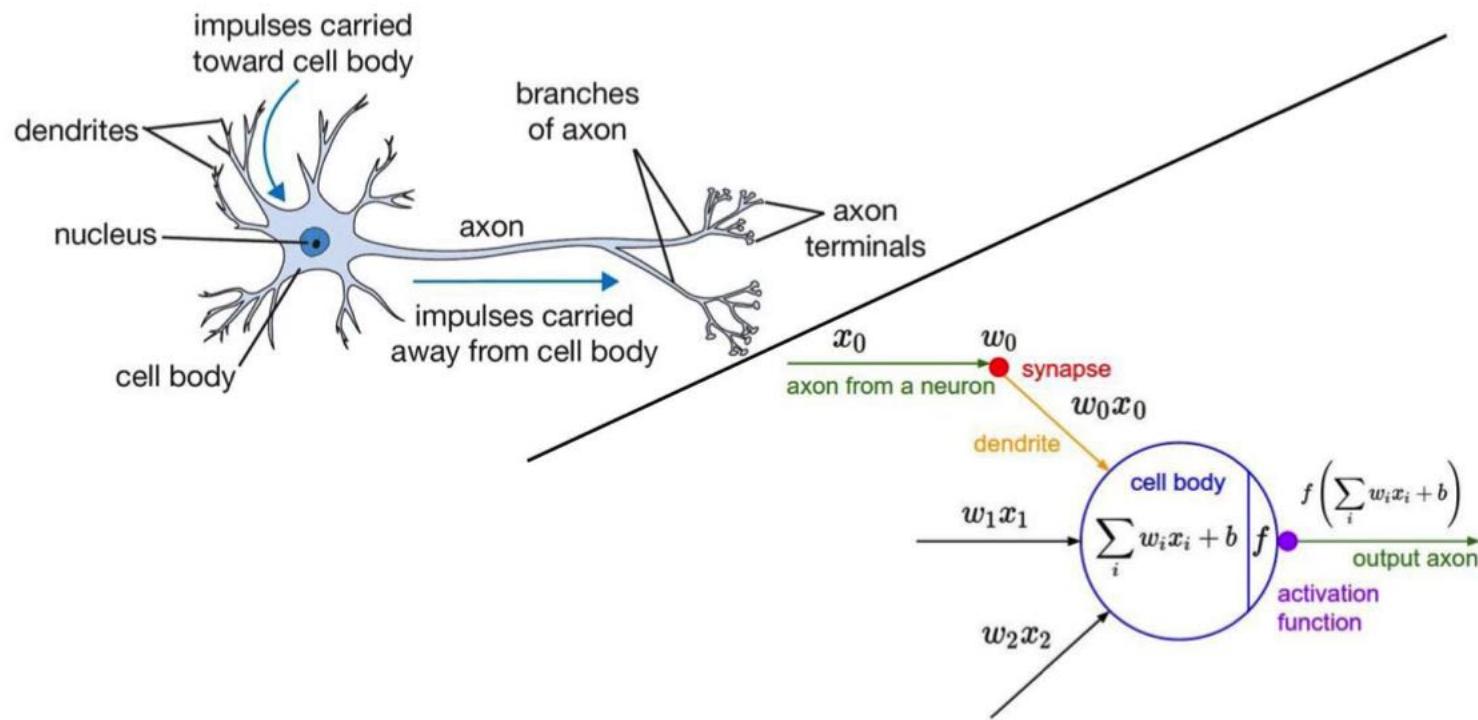
$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	1

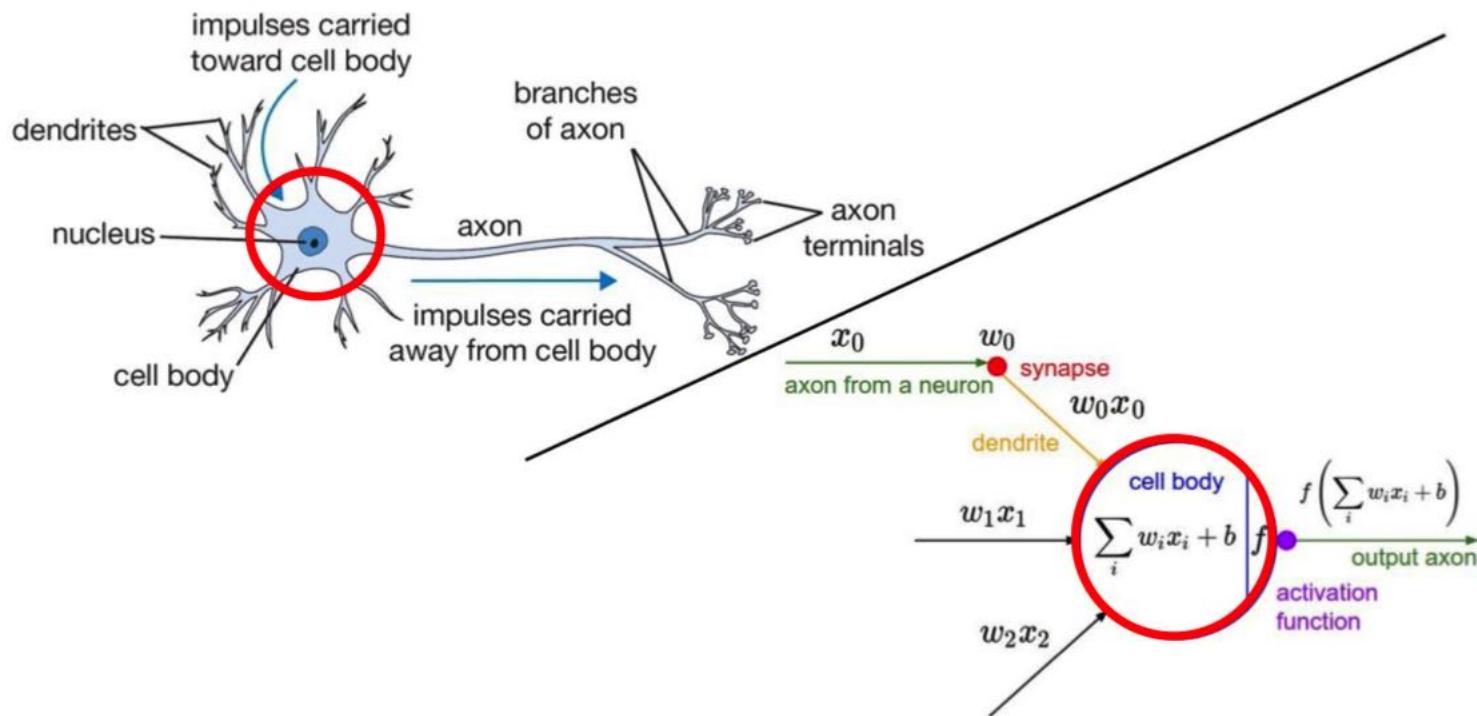
*AND Gate*

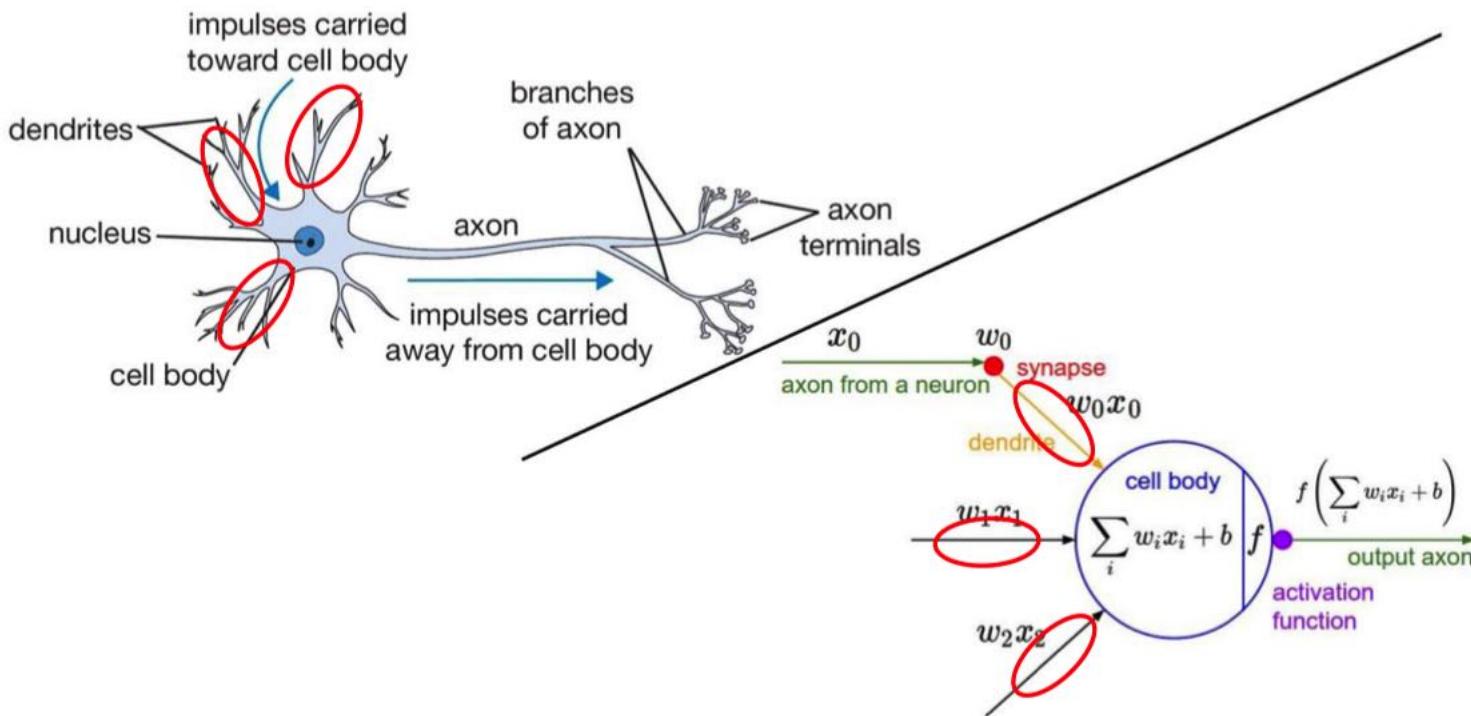
$x_1$	$x_2$	$y$
0	0	0
0	1	0
1	0	0
1	1	1

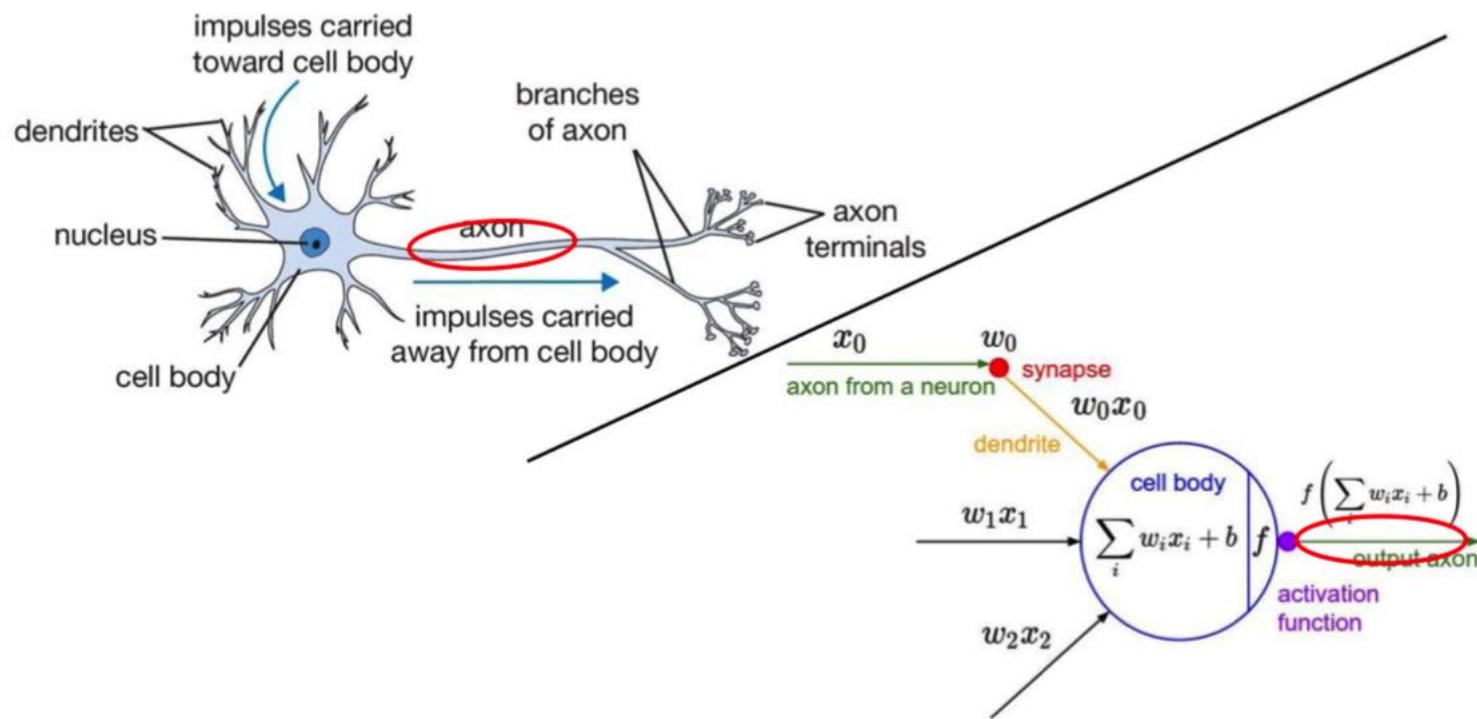
*XOR Gate*

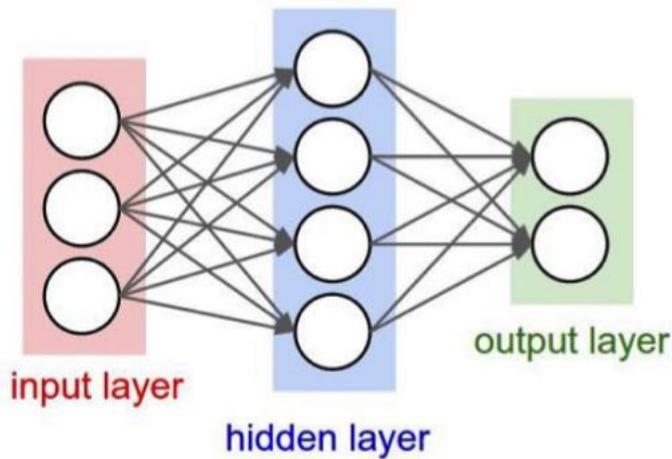
$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	0



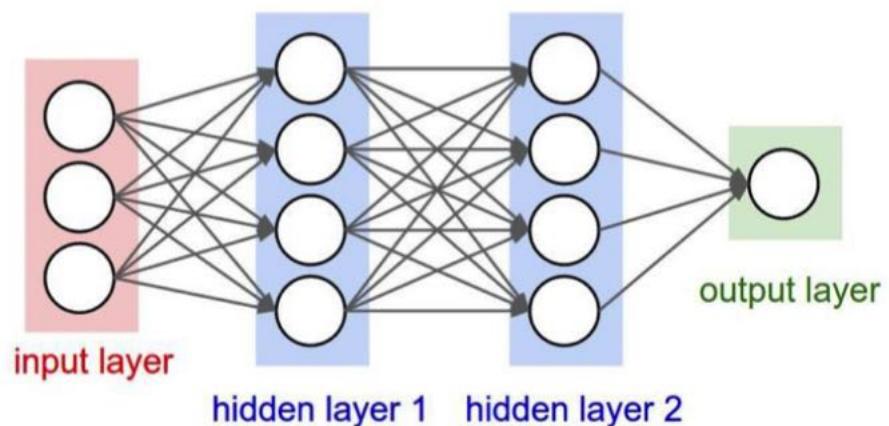




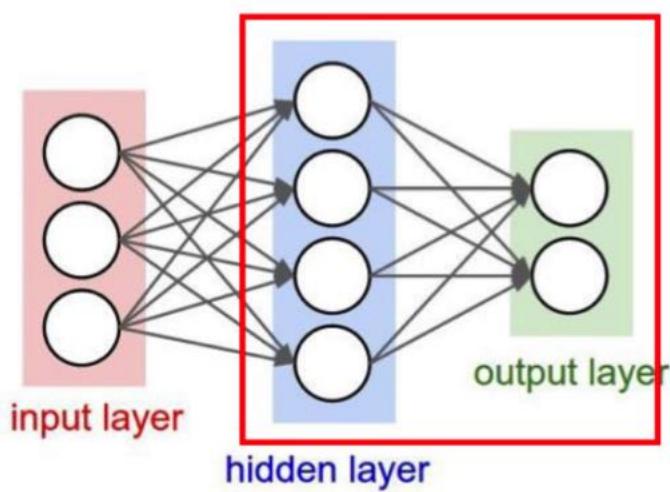




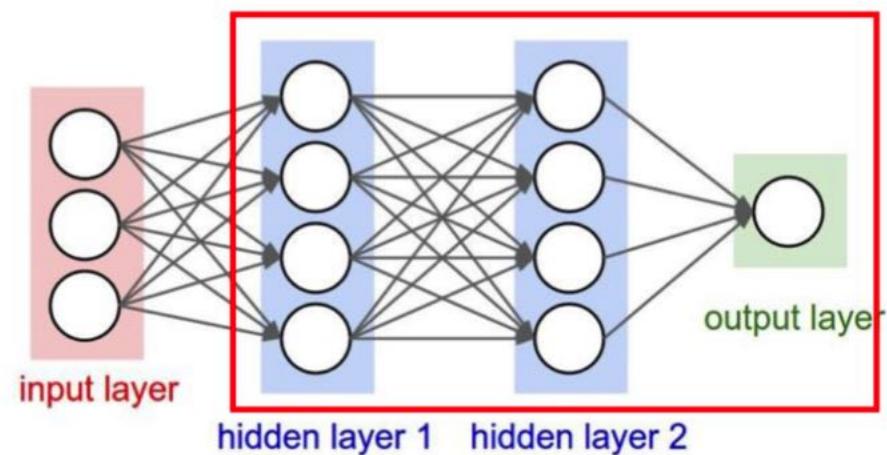
**"2-layer Neural Network"**  
or  
**"1-hidden-layer Neural Network"**



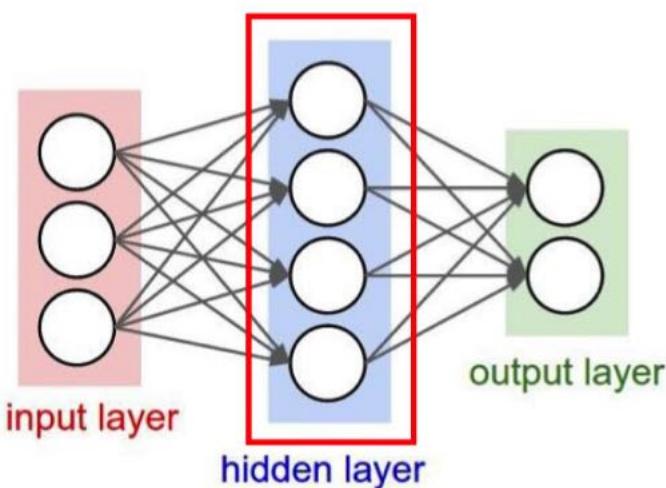
**"3-layer Neural Network"**  
or  
**"2-hidden-layer Neural Network"**



**"2-layer Neural Network"**  
or  
**"1-hidden-layer Neural Network"**

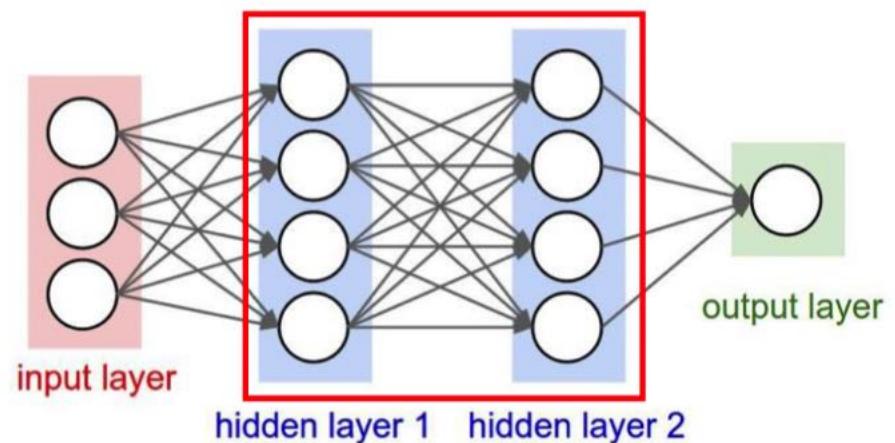


**"3-layer Neural Network"**  
or  
**"2-hidden-layer Neural Network"**

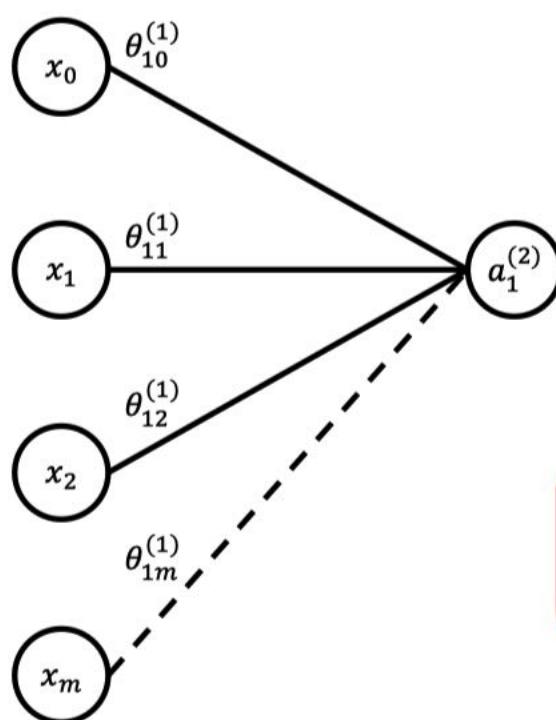


“2-layer Neural Network”  
or

“1-hidden-layer Neural Network”



“3-layer Neural Network”  
or  
“2-hidden-layer Neural Network”

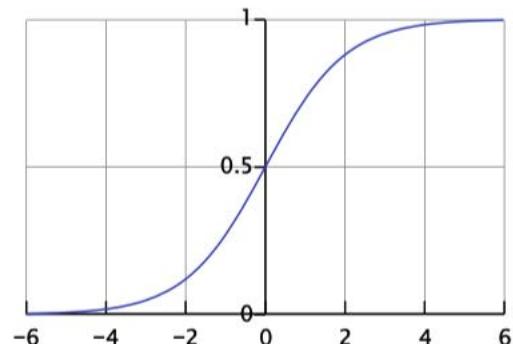


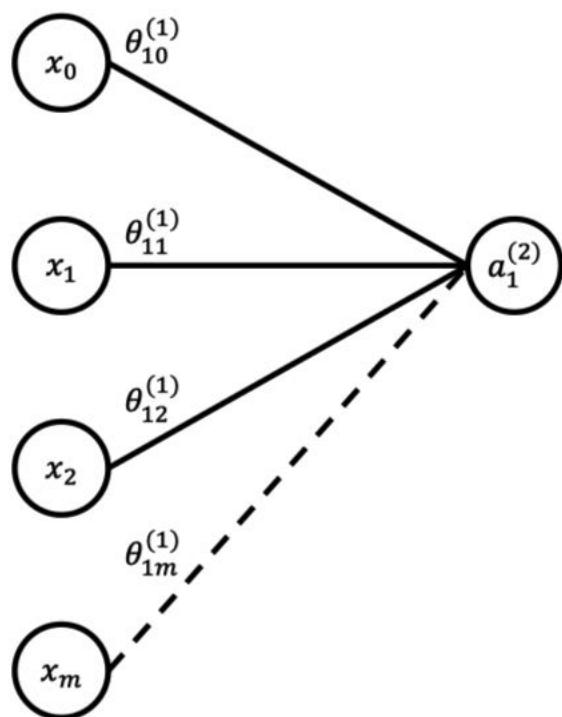
$$\begin{aligned}z_1^{(2)} &= \theta_{10}^{(1)}x_0 + \theta_{11}^{(1)}x_1 + \theta_{12}^{(1)}x_2 + \dots + \theta_{1m}^{(1)}x_m \\&= \sum_{i=0}^m \theta_{1i}^{(1)}x_i\end{aligned}$$

$$a_1^{(2)} = \sigma(z_1^{(2)})$$

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

*Logistic sigmoid* →





$$z_1^{(2)} = \theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \dots + \theta_{1m}^{(1)} x_m$$

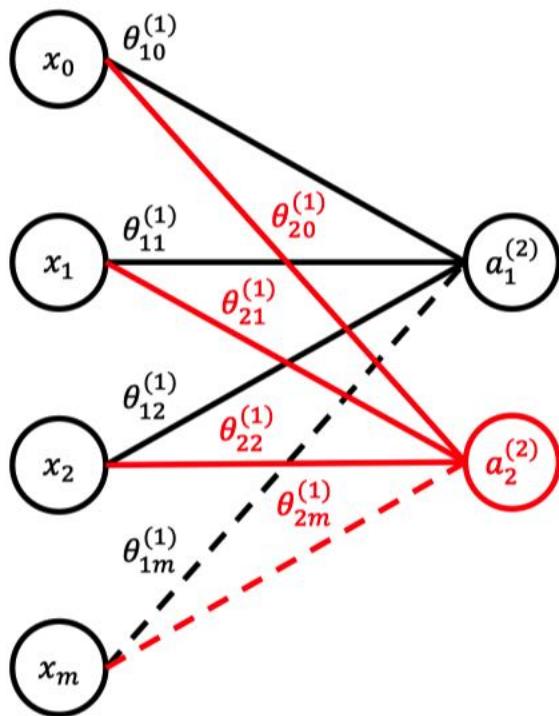
$$= \sum_{i=0}^m \theta_{1i}^{(1)} x_i$$

*Linear combination*

$$a_1^{(2)} = \sigma(z_1^{(2)})$$

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

$$[\theta_{10}^{(1)} \quad \theta_{11}^{(1)} \quad \theta_{12}^{(1)} \quad \dots \quad \theta_{1m}^{(1)}] \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

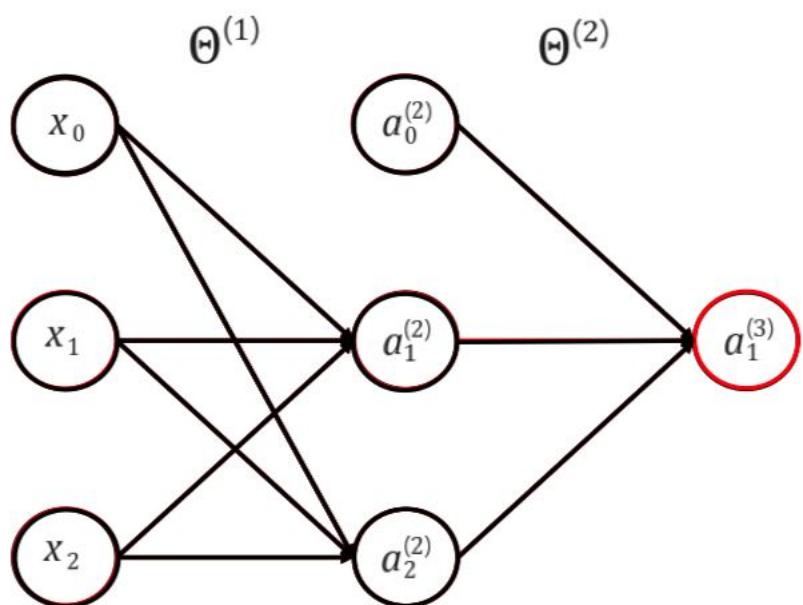


$$\begin{aligned} z_1^{(2)} &= \theta_{10}^{(1)}x_0 + \theta_{11}^{(1)}x_1 + \theta_{12}^{(1)}x_2 + \cdots + \theta_{1m}^{(1)}x_m \\ &= \sum_{i=0}^m \theta_{1i}^{(1)}x_i \end{aligned}$$

$$\begin{aligned} z_2^{(2)} &= \theta_{20}^{(1)}x_0 + \theta_{21}^{(1)}x_1 + \theta_{22}^{(1)}x_2 + \cdots + \theta_{2m}^{(1)}x_m \\ &= \sum_{i=0}^m \theta_{2i}^{(1)}x_i \end{aligned}$$

$$\begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \end{bmatrix} = \begin{bmatrix} \theta_{10}^{(1)} & \theta_{11}^{(1)} & \theta_{12}^{(1)} & \cdots & \theta_{1m}^{(1)} \\ \theta_{20}^{(1)} & \theta_{21}^{(1)} & \theta_{22}^{(1)} & \cdots & \theta_{2m}^{(1)} \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

$$(2 \times m) \times (m \times 1)$$

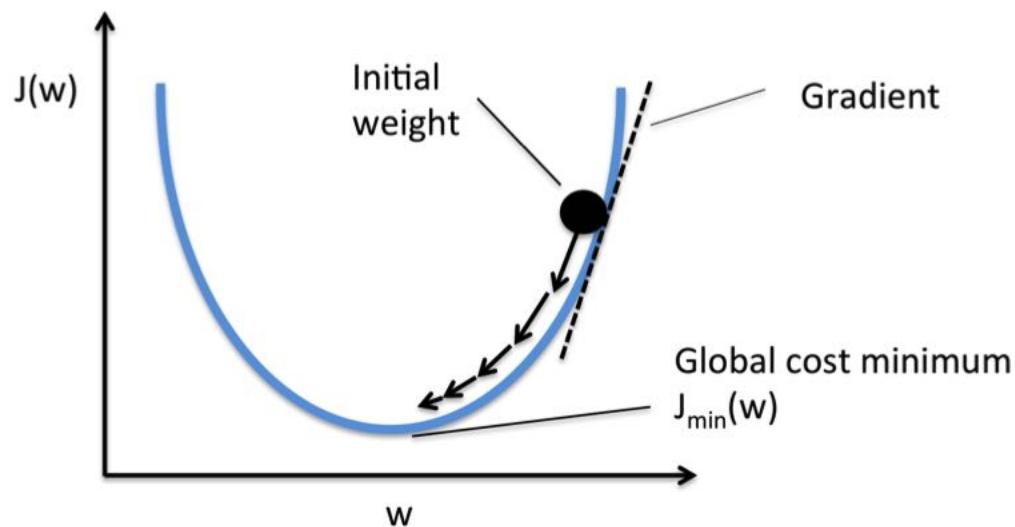


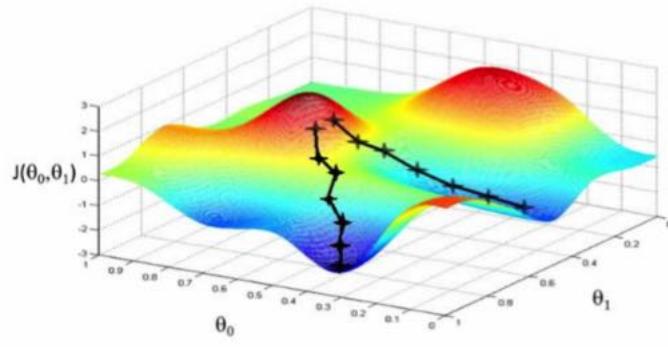
$$\begin{pmatrix} z_1^{(2)} \\ z_2^{(2)} \end{pmatrix} = \begin{pmatrix} \Theta_{10}^{(1)} & \Theta_{11}^{(1)} & \Theta_{12}^{(1)} \\ \Theta_{20}^{(1)} & \Theta_{21}^{(1)} & \Theta_{22}^{(1)} \end{pmatrix} \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix}$$

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

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

$$a_1^{(3)} = g(z_1^{(3)})$$

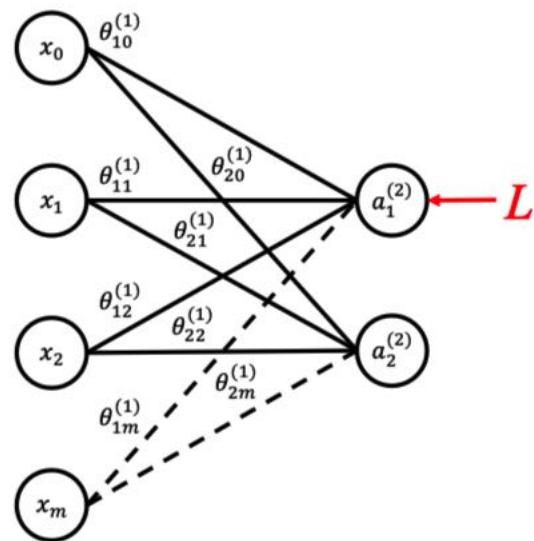




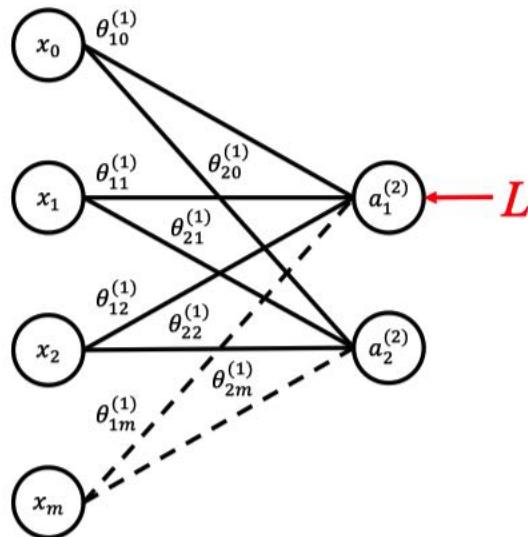
- Compute  $\frac{\partial J}{\partial \theta_0}, \frac{\partial J}{\partial \theta_1}$
- Update weights with

$$\theta_0 := \theta_0 - \alpha \frac{\partial J}{\partial \theta_0}$$

$$\theta_1 := \theta_1 - \alpha \frac{\partial J}{\partial \theta_1}$$



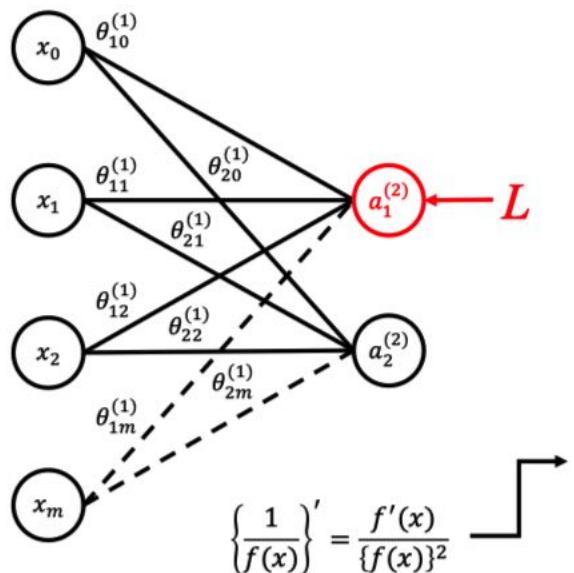
- Use chain rule  
(for MSE loss)



- Use chain rule  
(for MSE loss)

$$\frac{\partial L}{\partial a_1^{(2)}}$$

$$\frac{\partial L}{\partial a_1^{(2)}} = \frac{\partial (y_1 - a_1^{(2)})^2}{\partial a_1^{(2)}} = -2(y_1 - a_1^{(2)})$$

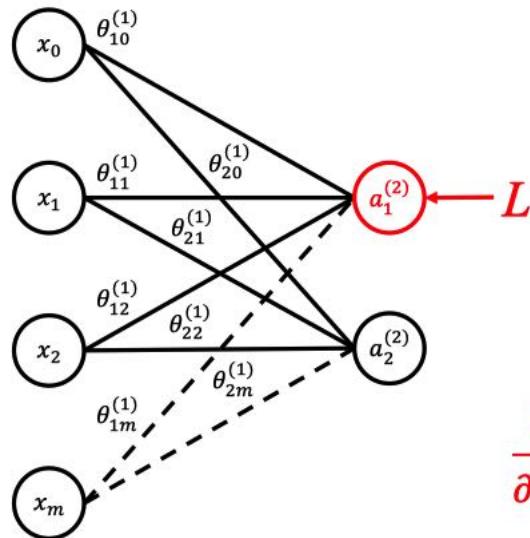


- Use chain rule (for MSE loss)

$$\frac{\partial L}{\partial z_1^{(2)}} = \frac{\partial L}{\partial a_1^{(2)}} \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}}$$

$$\frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} = \frac{\partial \sigma(z_1^{(2)})}{\partial z_1^{(2)}} = \sigma(z_1^{(2)}) (1 - \sigma(z_1^{(2)}))$$

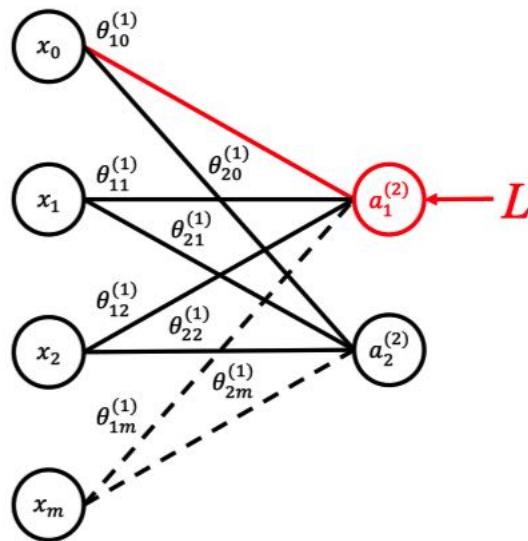
$$\left\{ \frac{1}{f(x)} \right\}' = \frac{f'(x)}{\{f(x)\}^2}$$



- Use chain rule  
(for MSE loss)

$$\frac{\partial L}{\partial z_1^{(2)}} = \frac{\partial L}{\partial a_1^{(2)}} \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}}$$

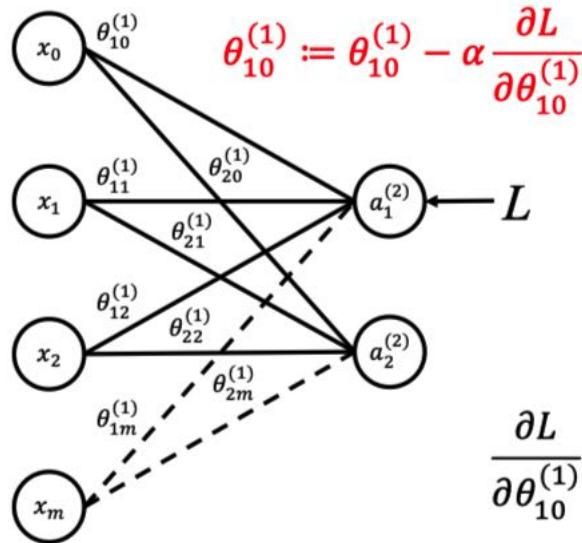
$$\frac{\partial L}{\partial z_1^{(2)}} = -2(y_1 - a_1^{(2)})\sigma(z_1^{(2)})(1 - \sigma(z_1^{(2)}))$$



- Use chain rule  
(for MSE loss)

$$\frac{\partial L}{\partial \theta_{10}^{(1)}} = \frac{\partial L}{\partial a_1^{(2)}} \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \frac{\partial z_1^{(2)}}{\partial \theta_{10}^{(1)}}$$

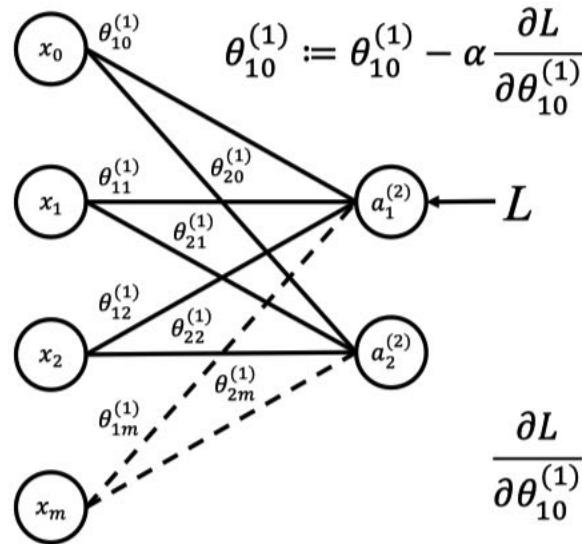
$$\frac{\partial z_1^{(2)}}{\partial \theta_{10}^{(1)}} = \frac{\partial \sum_{i=0}^m \theta_{1i}^{(1)} x_i}{\partial \theta_{10}^{(1)}} = x_0$$



- Use chain rule (for MSE loss)

$$\frac{\partial L}{\partial \theta_{10}^{(1)}} = \frac{\partial L}{\partial a_1^{(2)}} \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \frac{\partial z_1^{(2)}}{\partial \theta_{10}^{(1)}}$$

$$\frac{\partial L}{\partial \theta_{10}^{(1)}} = -2(y_1 - a_1^{(2)})\sigma(z_1^{(2)})(1 - \sigma(z_1^{(2)}))x_0$$



- Use chain rule (for MSE loss)

$$\frac{\partial L}{\partial \theta_{10}^{(1)}} = \frac{\partial L}{\partial a_1^{(2)}} \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \frac{\partial z_1^{(2)}}{\partial \theta_{10}^{(1)}}$$

$$\frac{\partial L}{\partial \theta_{10}^{(1)}} = -2(y_1 - a_1^{(2)})\sigma(z_1^{(2)})(1 - \sigma(z_1^{(2)}))x_0$$

I don't like math ππ

## 전제

신경망에는 적응 가능한 가중치와 편향이 있고, 이 가중치와 편향을 훈련 데이터에 적응하도록 조정하는 과정을 ‘학습’이라 합니다. 신경망 학습은 다음과 같이 4단계로 수행합니다.

### 1단계 – 미니배치

훈련 데이터 중 일부를 무작위로 가져옵니다. 이렇게 선별한 데이터를 미니배치라 하며, 그 미니배치의 손실 함수 값을 줄이는 것이 목표입니다.

### 2단계 – 기울기 산출

미니배치의 손실 함수 값을 줄이기 위해 각 가중치 매개변수의 기울기를 구합니다. 기울기는 손실 함수의 값을 가장 작게 하는 방향을 제시합니다.

### 3단계 – 매개변수 갱신

가중치 매개변수를 기울기 방향으로 아주 조금 갱신합니다.

### 4단계 – 반복

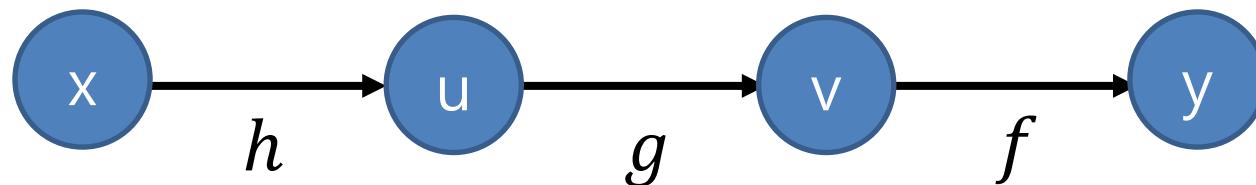
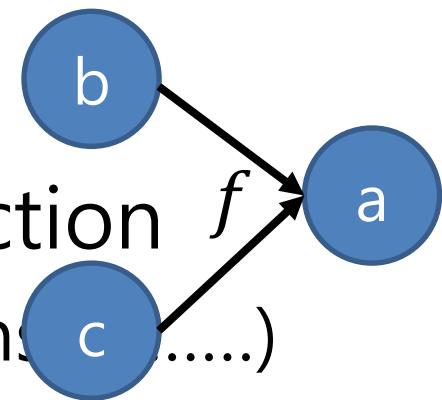
1~3단계를 반복합니다.

# Computational Graph

- A “language” describing a function
  - **Node**: variable (scalar, vector, tensor.....)
  - **Edge**: operation (simple function)

**Example**  $y = f(g(h(x)))$

$$u = h(x) \quad v = g(u) \quad y = f(v)$$



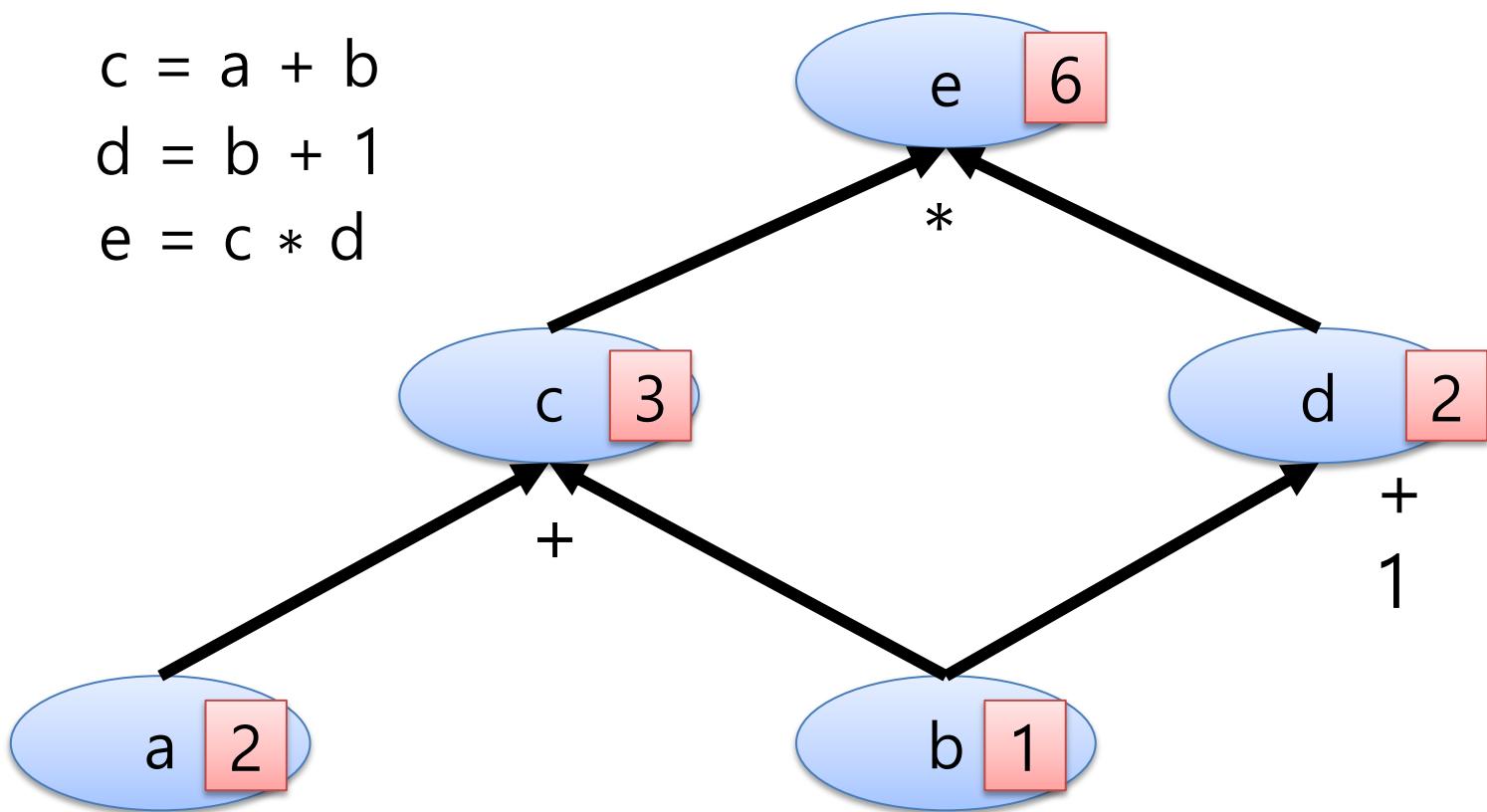
# Computational Graph

- Example:  $e = (a+b) * (b+1)$

$$c = a + b$$

$$d = b + 1$$

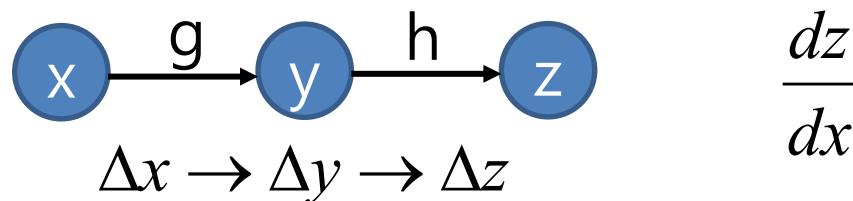
$$e = c * d$$



# Review: Chain Rule

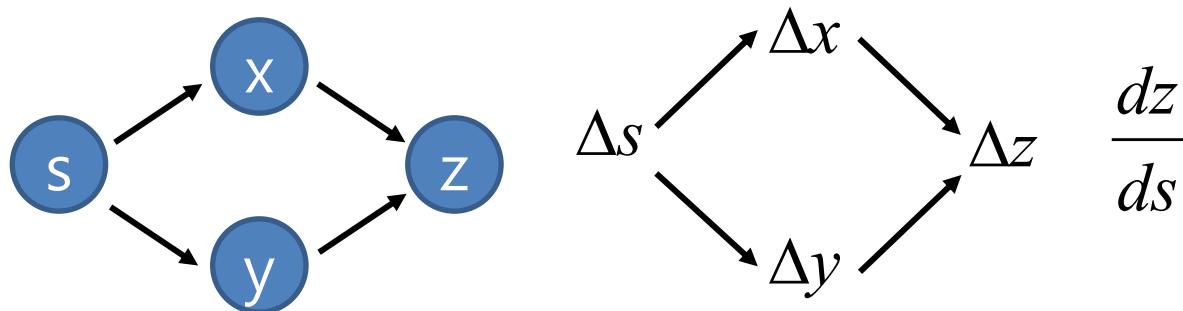
**Case**  
**1**

$$z = f(x) \rightarrow y = g(x) \quad z = h(y)$$



**Case**  
**2**

$$z = f(s) \rightarrow x = g(s) \quad y = h(s) \quad z = k(x, y)$$



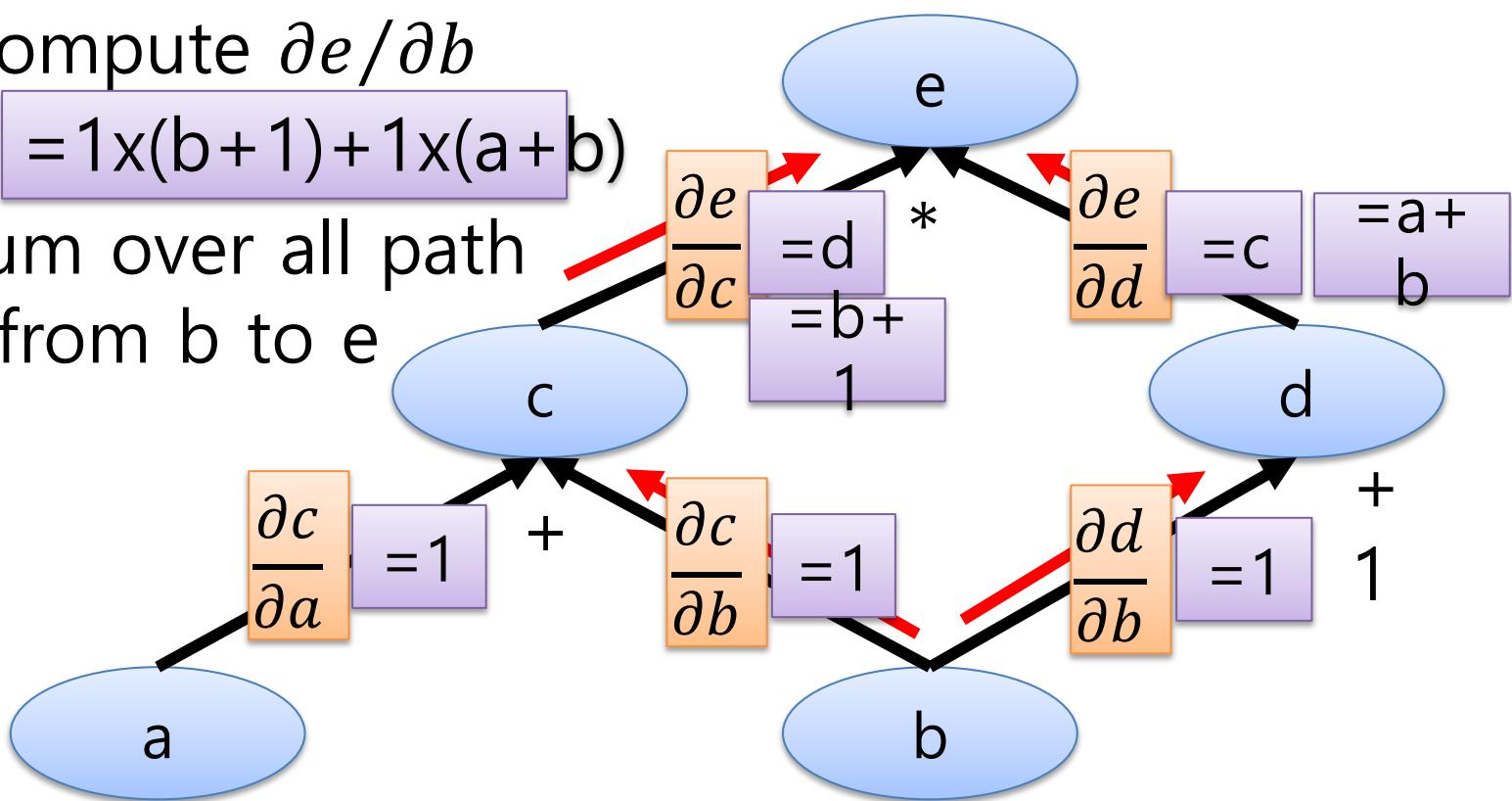
# Computational Graph

- Example:  $e = (a+b) * (b+1)$

Compute  $\partial e / \partial b$

$$= 1 \times (b+1) + 1 \times (a+b)$$

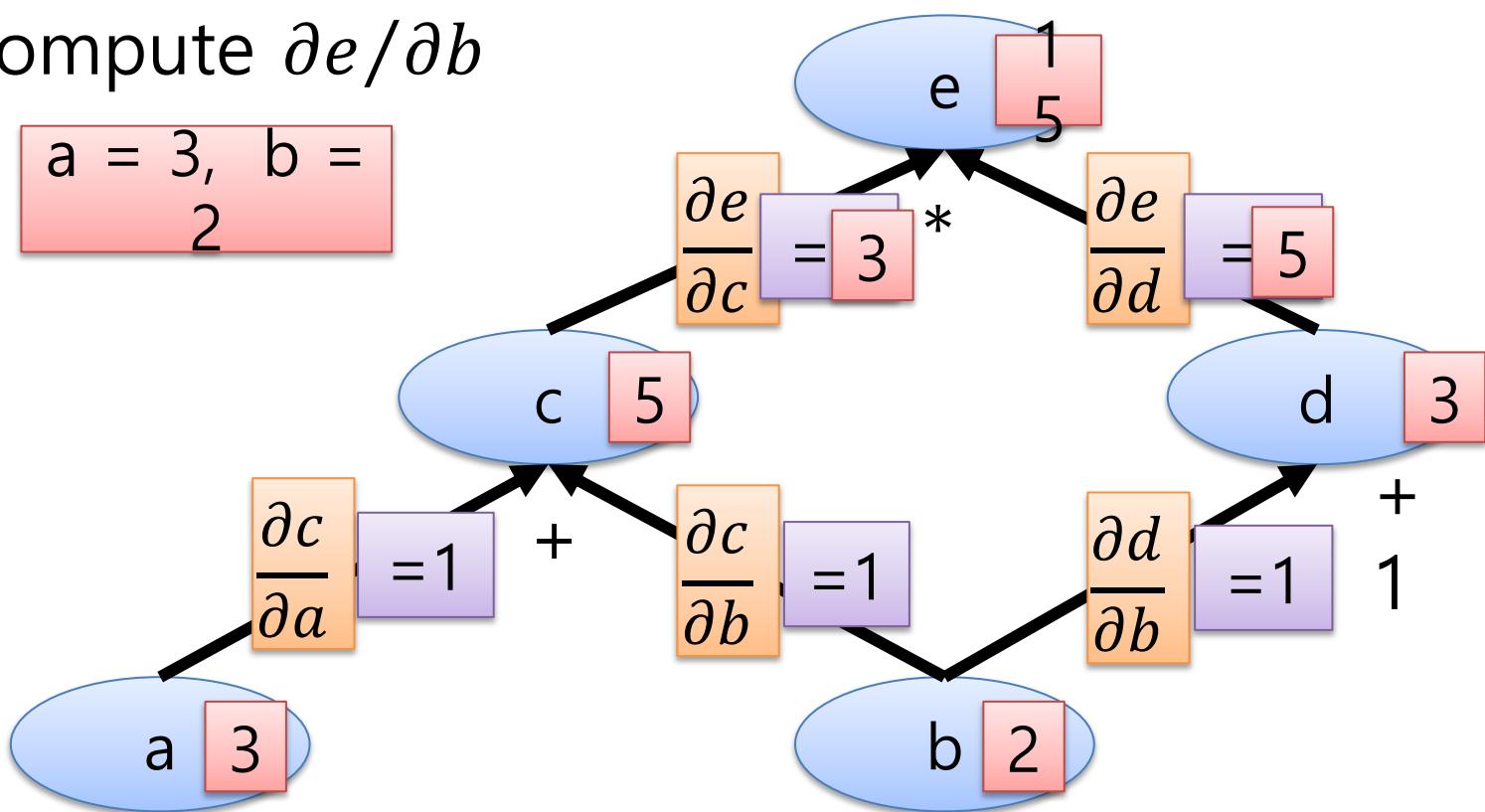
Sum over all paths  
from b to e



# Computational Graph

- Example:  $e = (a+b) * (b+1)$
- $\frac{\partial e}{\partial b} = \frac{\partial e}{\partial c} \frac{\partial c}{\partial b} + \frac{\partial e}{\partial d} \frac{\partial d}{\partial b}$
- Compute  $\frac{\partial e}{\partial b}$

$$a = 3, b = 2$$

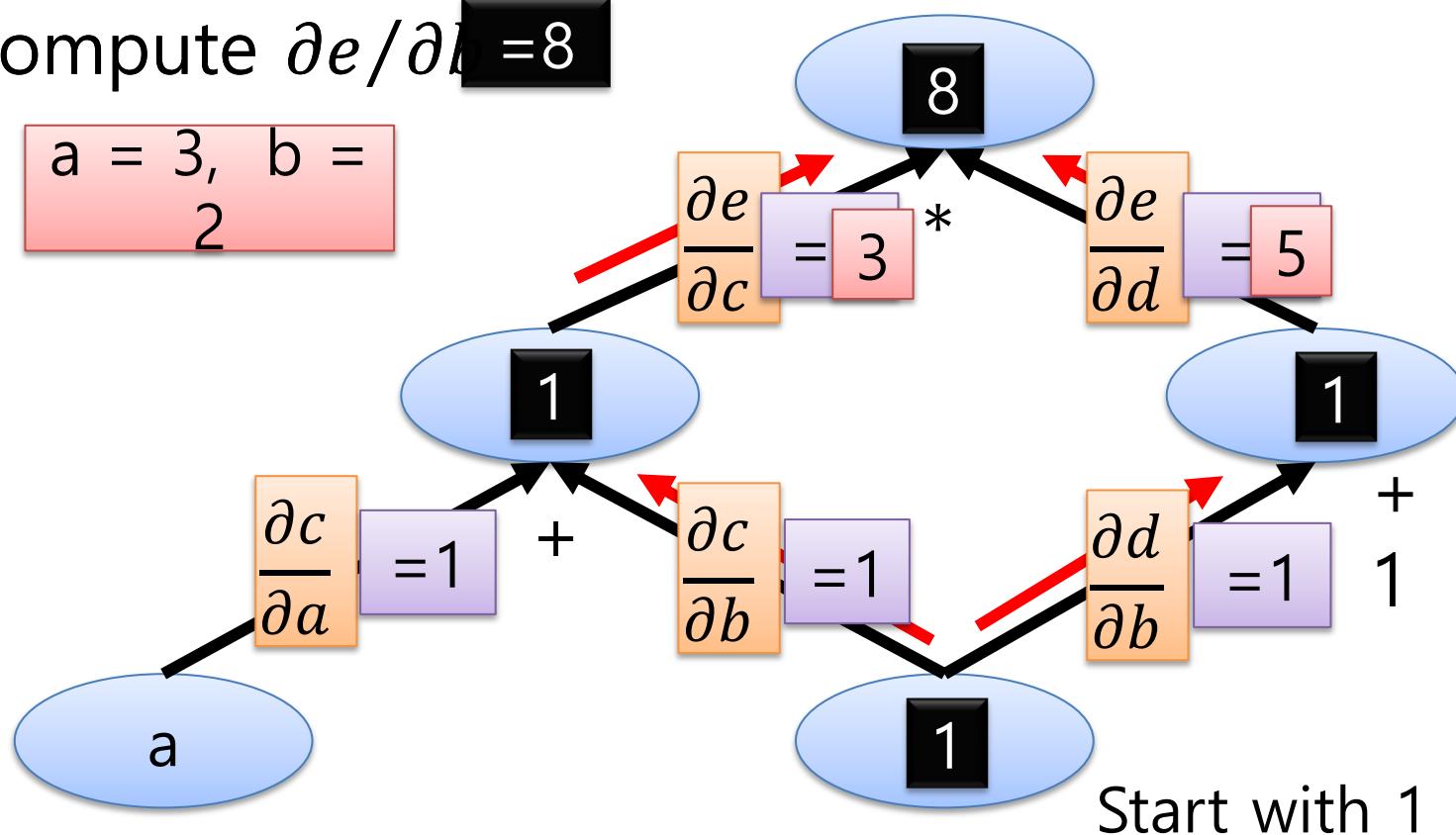


# Computational Graph

- Example:  $e = (a+b) * (b+1)$

Compute  $\partial e / \partial a = 8$

$$a = 3, b = 2$$

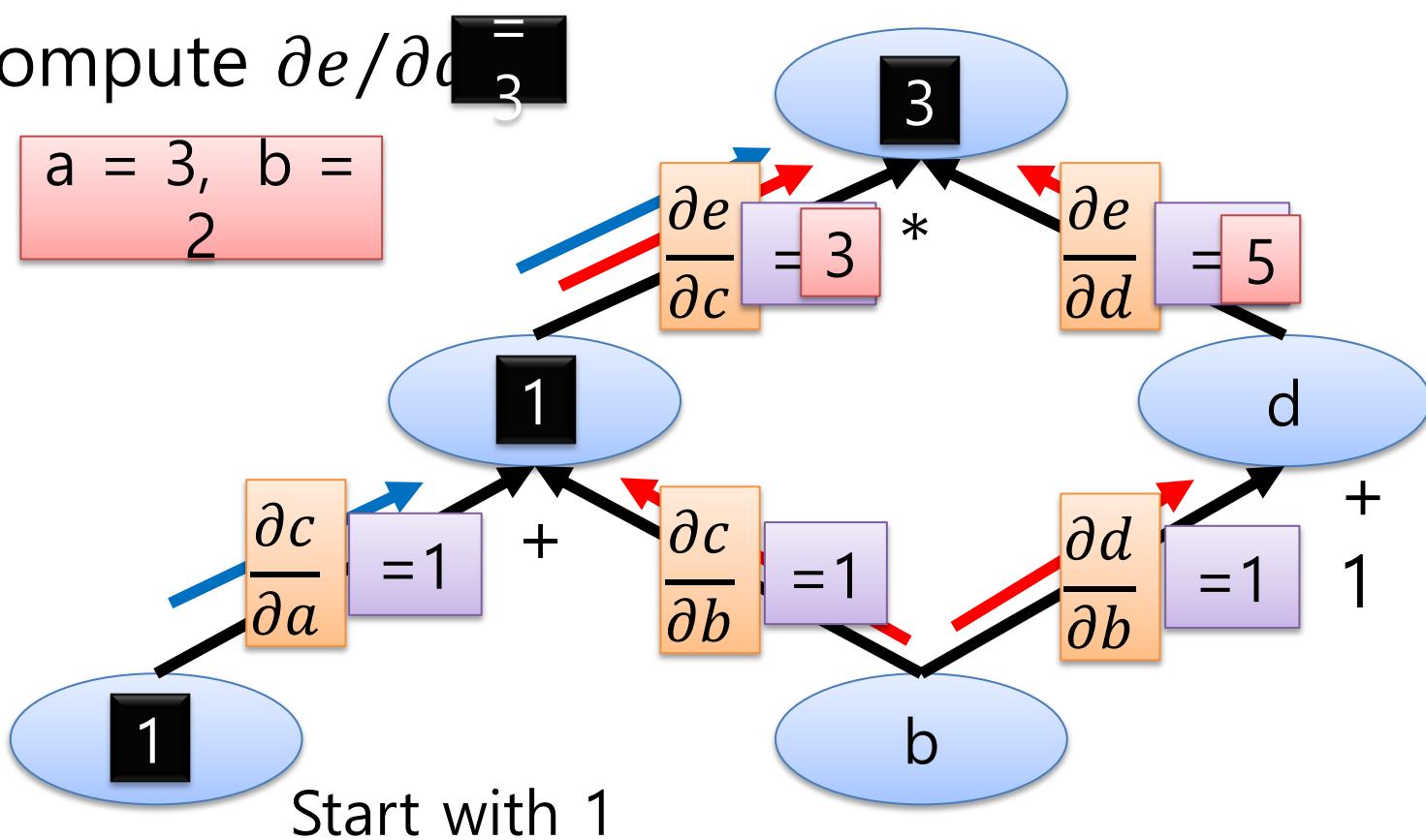


# Computational Graph

- Example:  $e = (a+b) * (b+1)$

Compute  $\partial e / \partial a$

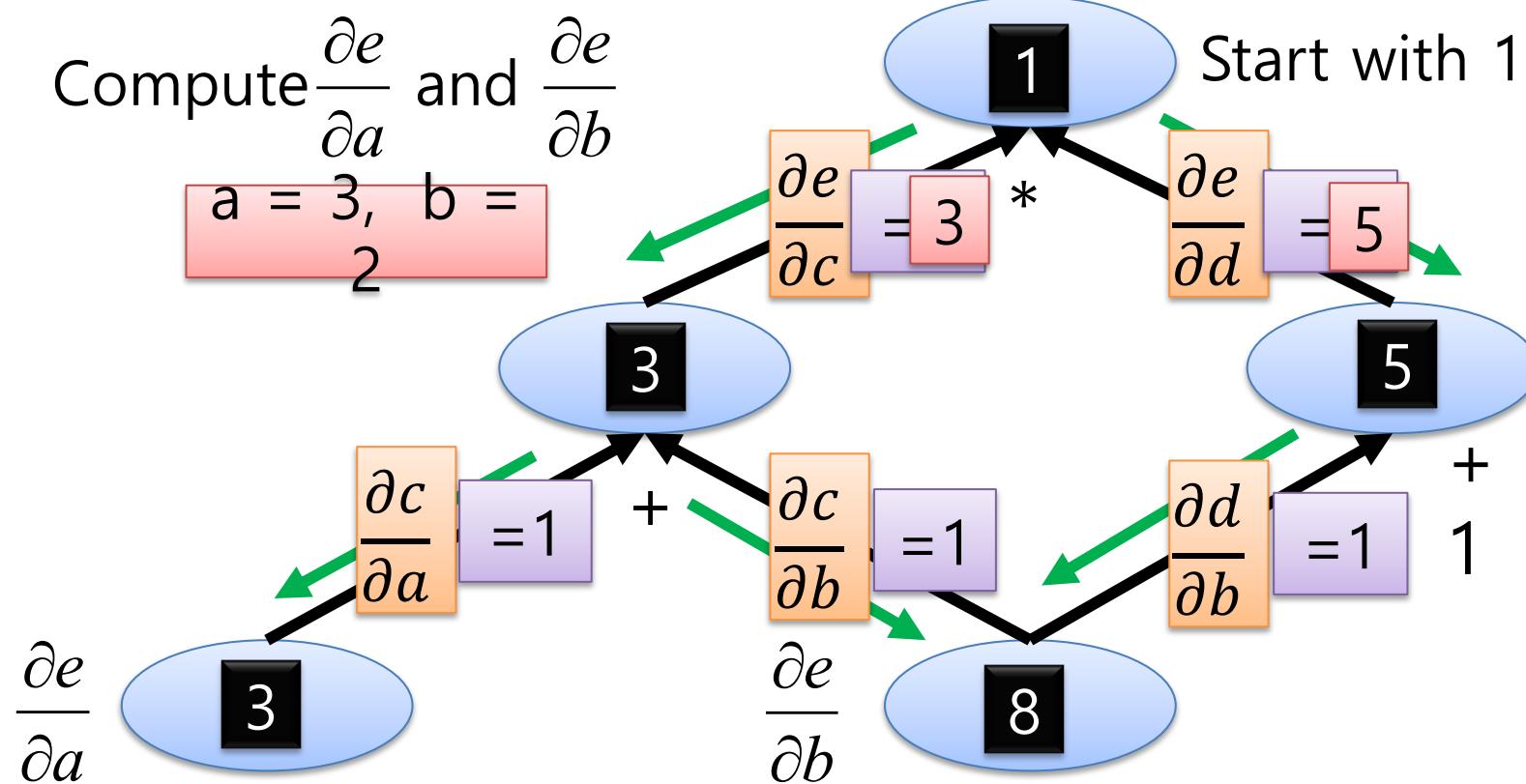
$$a = 3, b = 2$$



# Computational Graph

Reverse mode

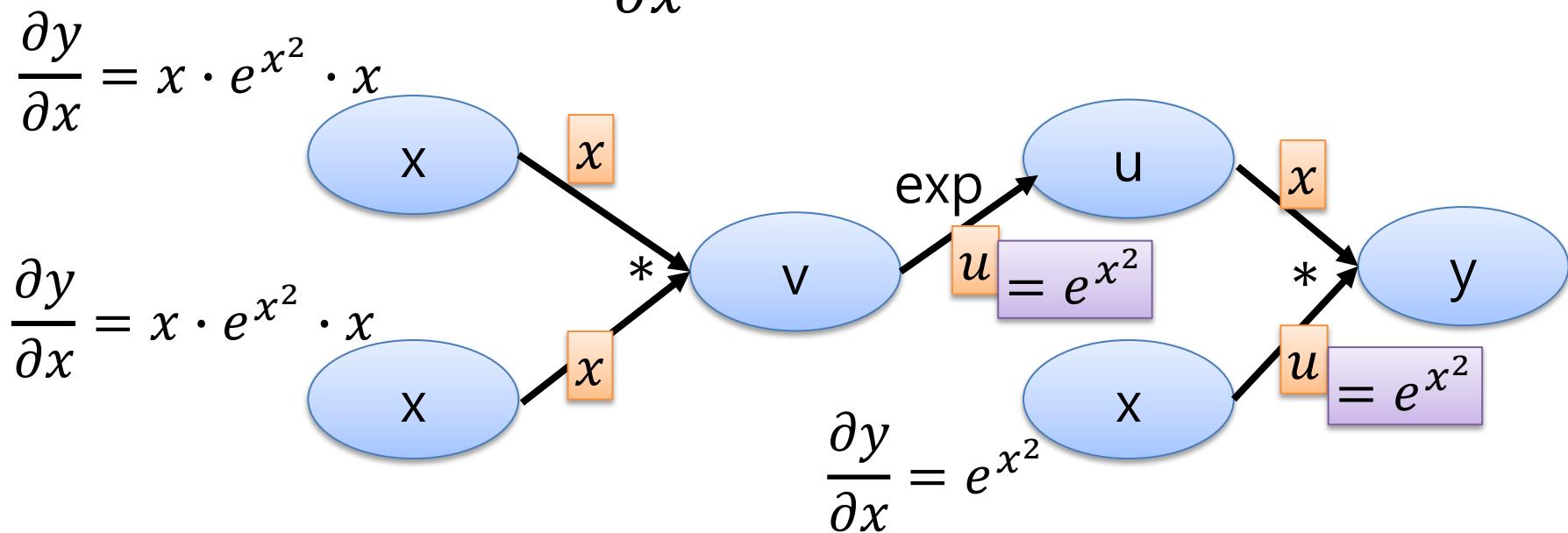
- Example:  $e = (a+b) * (b+1)$  What is the benefit?



# Computational Graph

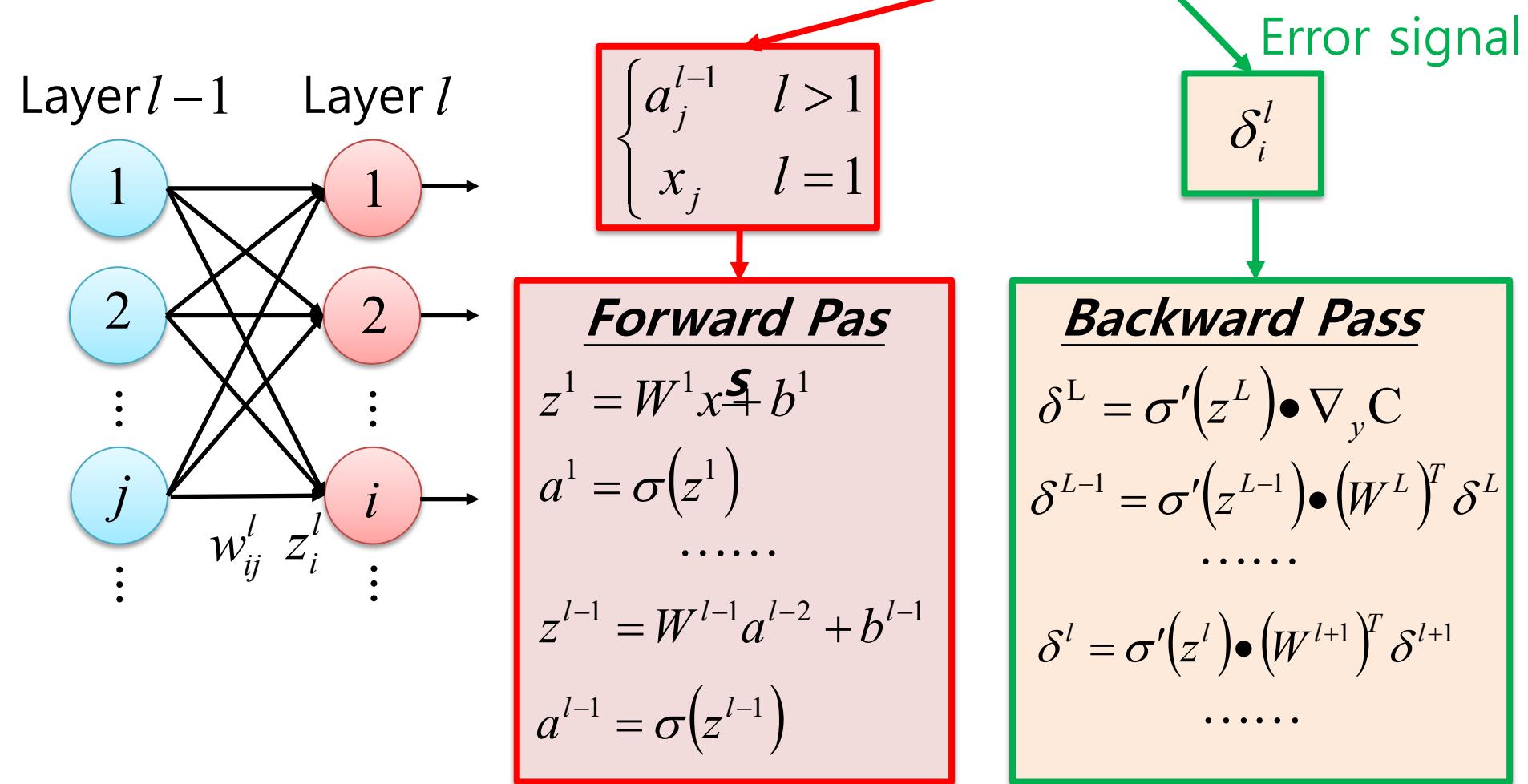
- ***Parameter sharing.*** the same parameters appearing in different nodes

$$y = xe^{x^2} \quad \frac{\partial y}{\partial x} = ? \quad e^{x^2} + x \cdot e^{x^2} \cdot 2x$$

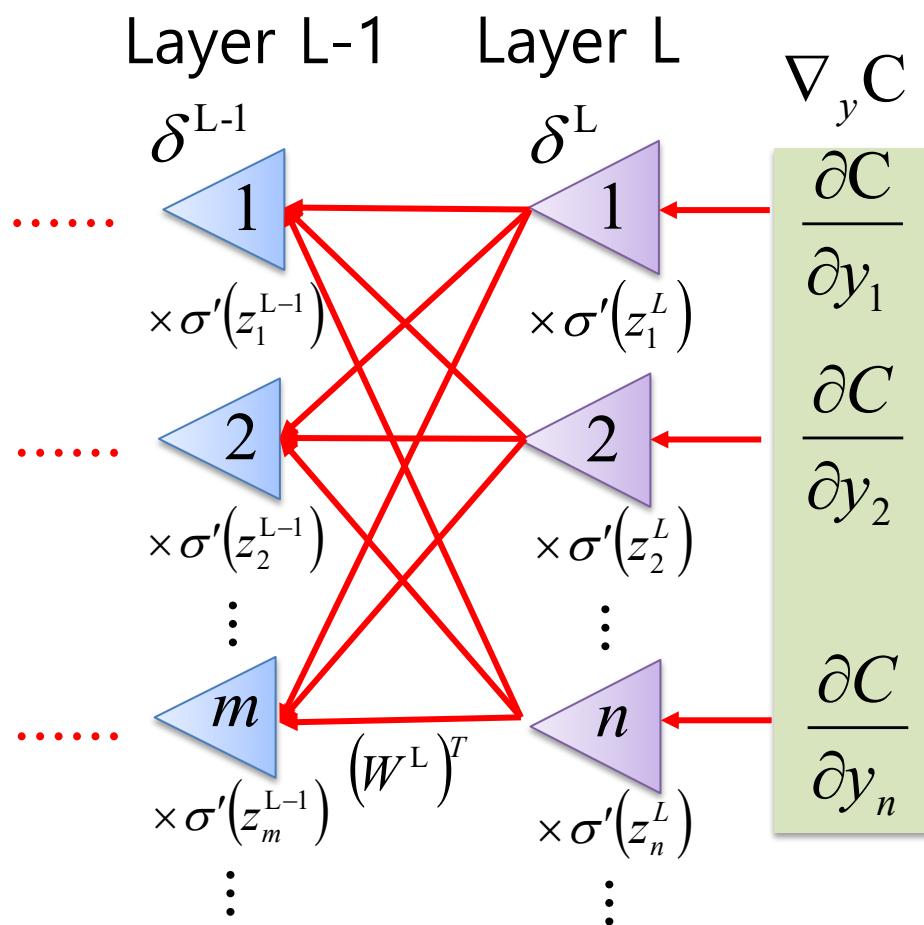


# Computational Graph for Feedforward Net

# Review: Backpropagation

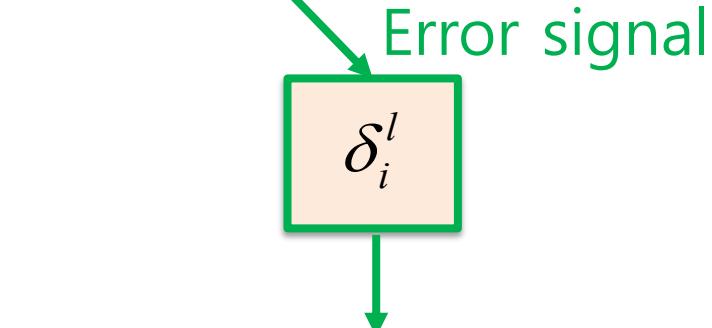


# Review: Backpropagation



(we do not use softmax here)

$$\frac{\partial C}{\partial w_{ij}} \quad \frac{\partial z_i^l}{\partial w_{ij}} \quad \frac{\partial C}{\partial z_i^l}$$



**Backward Pass**

$$\delta^L = \sigma'(z^L) \bullet \nabla_y C$$

$$\delta^{L-1} = \sigma'(z^{L-1}) \bullet (W^L)^\top \delta^L$$

$$\dots$$

$$\delta^l = \sigma'(z^l) \bullet (W^{l+1})^\top \delta^{l+1}$$

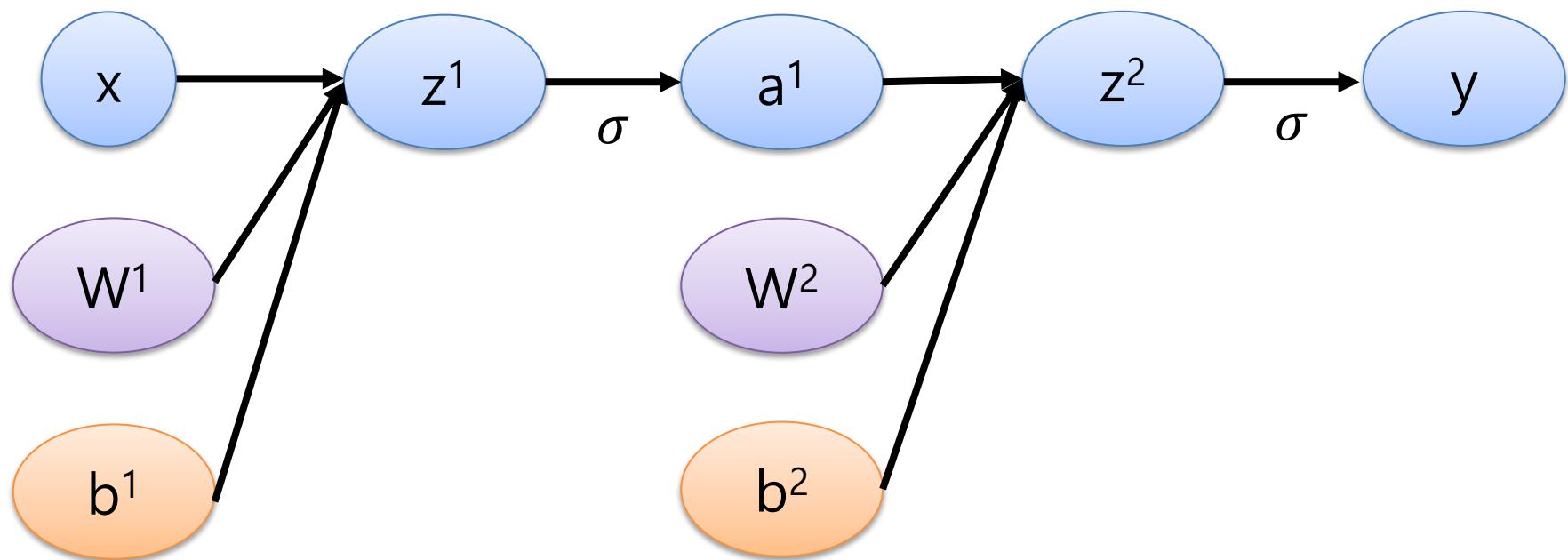
$$\dots$$

# Feedforward Network

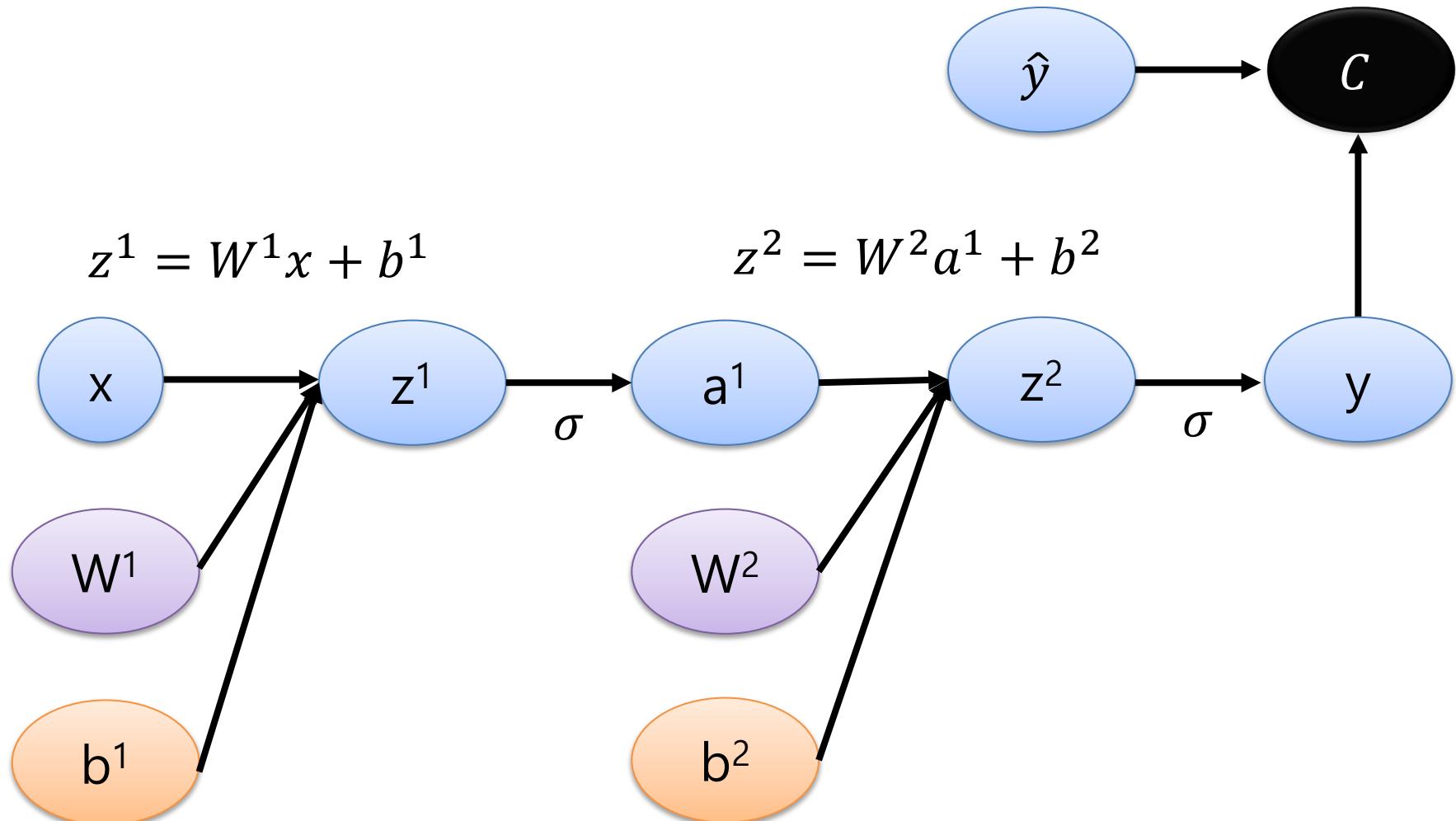
$$y = \sigma(W^L \dots \sigma(W^2 \sigma(W^1 x + b_1) + b_2) \dots + b_L)$$

$$z^1 = W^1 x + b^1$$

$$z^2 = W^2 a^1 + b^2$$



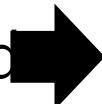
# Loss Function of Feedforward Network



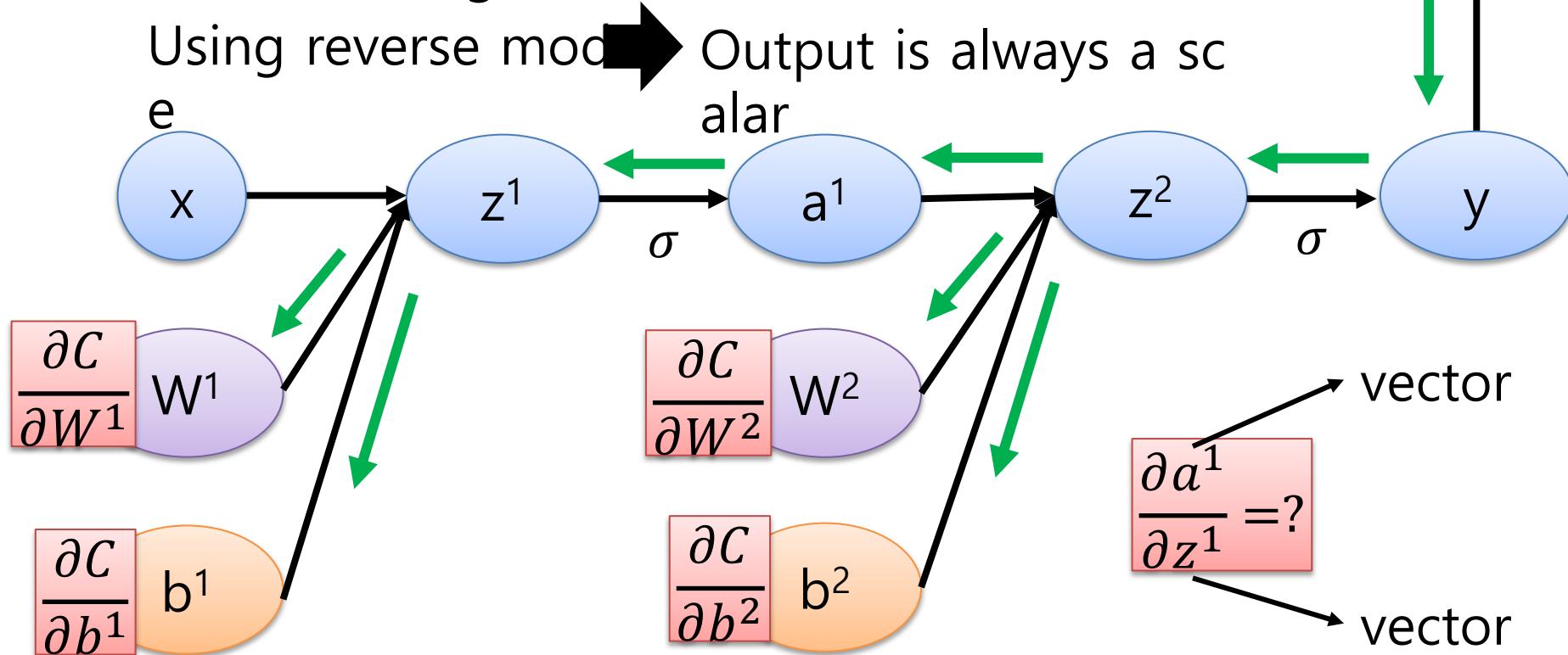
# Gradient of Cost Function

To compute the gradient

Computing the partial derivative on the edge  
Using reverse mode



Output is always a scalar



# Jacobian Matrix

$$y = f(x) \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$\frac{\partial y}{\partial x} =$$

}

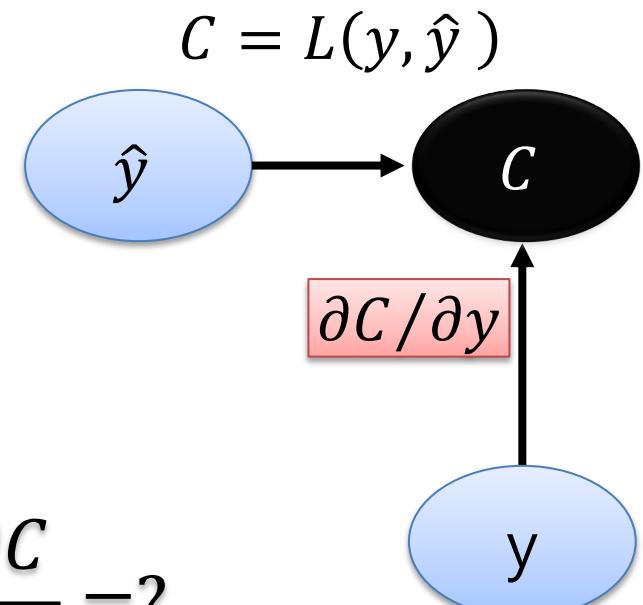
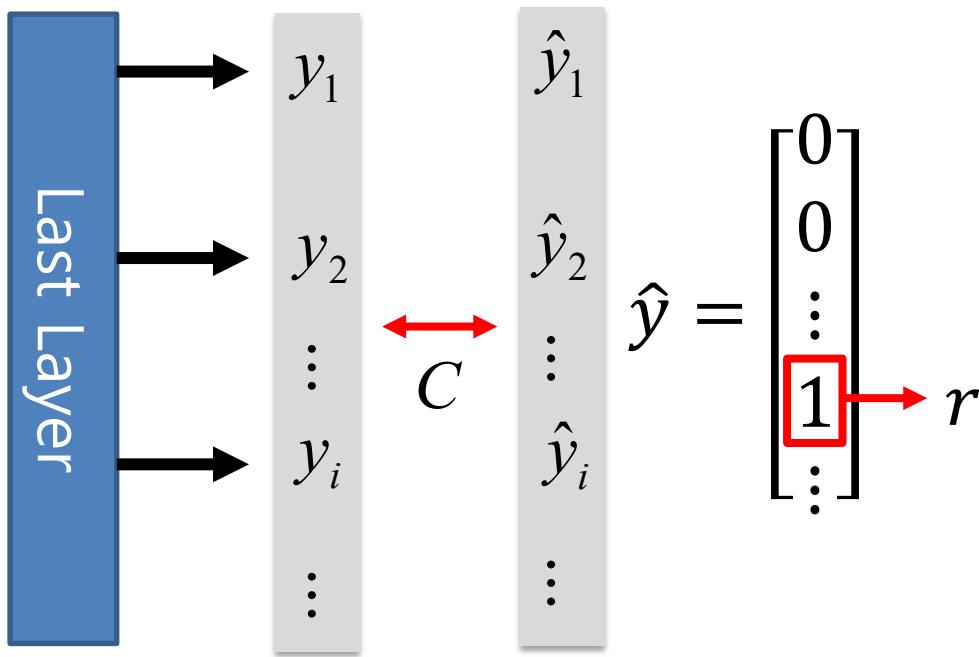
size of  
y

## Example

size of  
x

$$\begin{bmatrix} x_1 + x_2 x_3 \\ 2x_3 \end{bmatrix} = f \left( \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right) \quad \frac{\partial y}{\partial x} = [ \quad ]$$

# Gradient of Cost Function



Cross Entropy:  $C = -\log y_r$

$$\frac{\partial C}{\partial y} = [ \quad ]$$

$i = r:$

$$\frac{\partial C}{\partial y_r} = -1/y_r$$

$$i \neq r: \quad \frac{\partial C}{\partial y_i} = 0$$

# Gradient of Cost Function

$\frac{\partial y}{\partial z^2}$  is a Jacobian matrix

$$\begin{matrix} \text{squa} \\ \text{r} \\ \text{e} \\ z^2 \end{matrix} \quad \begin{matrix} y \\ \hat{y} \end{matrix}$$

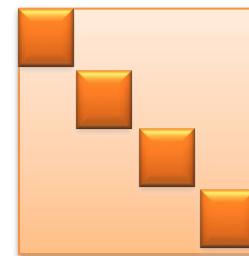
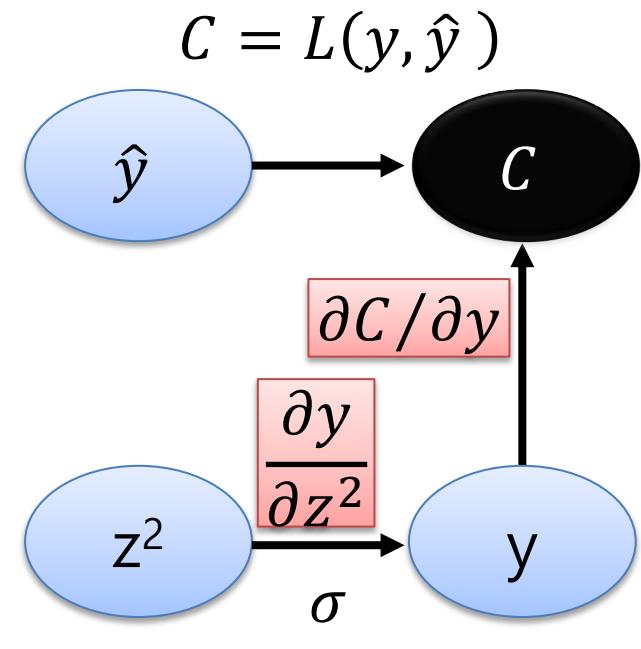
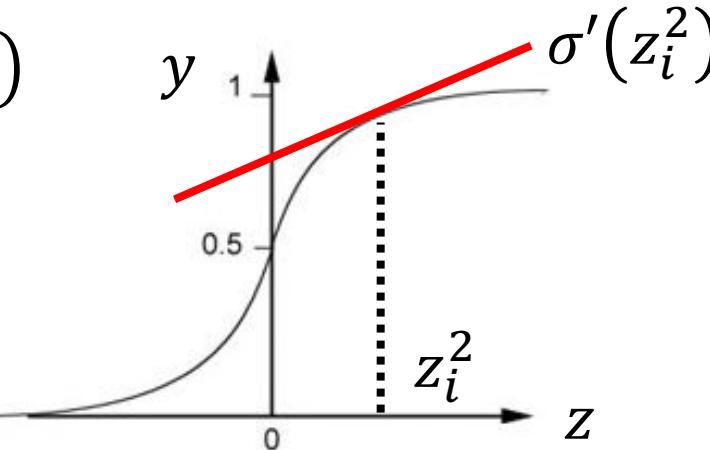
i-th row, j-th column:  $\partial y_i / \partial z_j^2$

$$i \neq j: \quad \partial y_i / \partial z_j^2 = 0$$

$$i = j: \quad \partial y_i / \partial z_i^2 = \sigma'(z_i^2)$$

$$y_i = \sigma(z_i^2)$$

How about softmax?  
😊



Diagon  
al  
Matrix

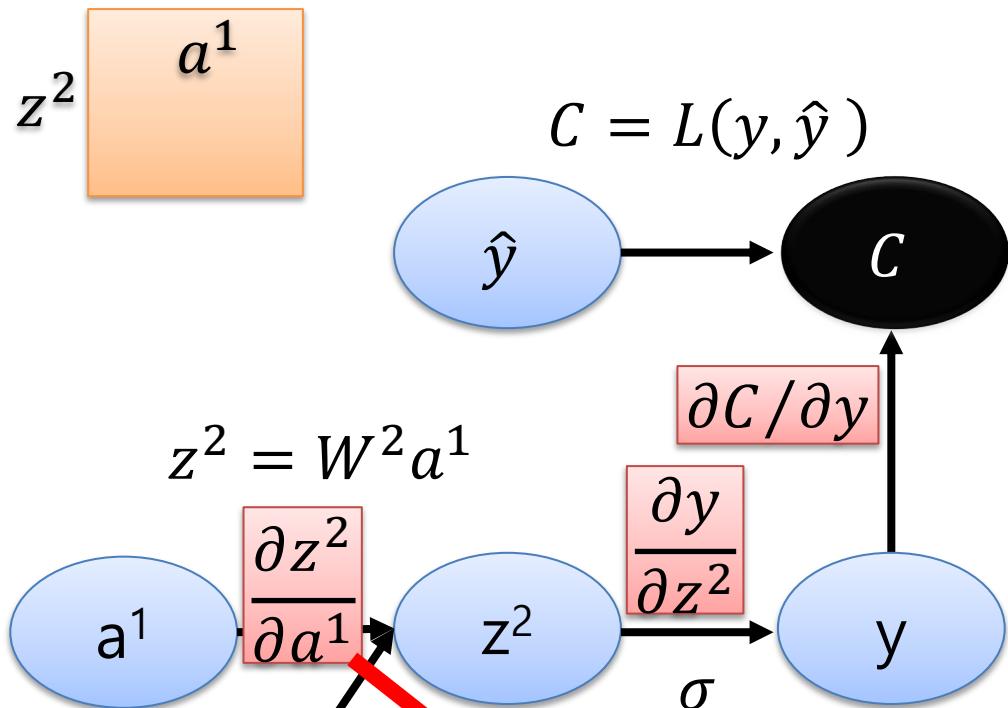
$\frac{\partial z^2}{\partial a^1}$  is a Jacobian matrix

i-th row, j-th column:

$$\frac{\partial z_i^2}{\partial a_j^1} =$$

$$z_i^2 = w_{i1}^2 a_1^1 + w_{i2}^2 a_2^1 + \dots + w_{in}^2 a_n^1$$

$$\begin{bmatrix} z_1^l \\ z_2^l \\ \vdots \\ z_i^l \\ \vdots \end{bmatrix} = \begin{bmatrix} w_{11}^l & w_{12}^l & \cdots & \\ w_{21}^l & w_{22}^l & & \ddots \\ & & \ddots & \\ & & & \ddots \end{bmatrix} \begin{bmatrix} a_1^{l-1} \\ a_2^{l-1} \\ \vdots \\ a_i^{l-1} \\ \vdots \end{bmatrix} + \begin{bmatrix} b_1^l \\ b_2^l \\ \vdots \\ b_i^l \\ \vdots \end{bmatrix}$$



$$\frac{\partial z^2}{\partial W^2} = m$$

mxn

$$(j-1)xn+k$$

i

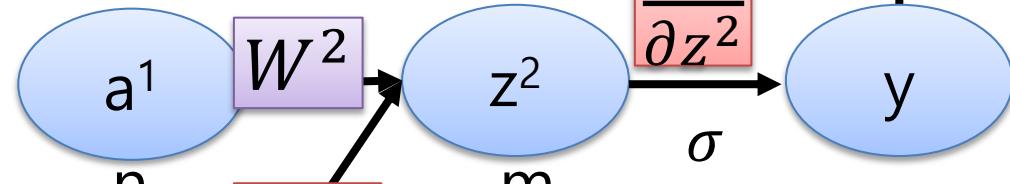
$$\frac{\partial z_i^2}{\partial W_{jk}^2}$$

$$\frac{\partial z_i^2}{\partial W_{jk}^2} = ? \quad i \neq j: \frac{\partial z_i^2}{\partial W_{jk}^2} = 0$$

$$i = j: \frac{\partial z_i^2}{\partial W_{ik}^2} = a_k^1$$

$$z^2 = W^2 a^1$$

$$z_i^2 = w_{i1}^2 a_1^1 + w_{i2}^2 a_2^1 + \dots + w_{in}^2 a_n^1$$



$$\begin{bmatrix} z_1^l \\ z_2^l \\ \vdots \\ z_i^l \\ \vdots \end{bmatrix} = \begin{bmatrix} w_{11}^l & w_{12}^l & \cdots & \\ w_{21}^l & w_{22}^l & \ddots & \\ \vdots & & & \ddots \end{bmatrix} \begin{bmatrix} a_1^{l-1} \\ a_2^{l-1} \\ \vdots \\ a_i^{l-1} \\ \vdots \end{bmatrix} + \begin{bmatrix} b_1^l \\ b_2^l \\ \vdots \\ b_i^l \\ \vdots \end{bmatrix}$$

$$W^2 \quad \frac{\partial z^2}{\partial W^2}$$

Considering  
 $W^2$  as a mxn **vector**  
 (considering  $\partial z^2 / \partial W^2$  as a tensor makes thing easier)

$$\frac{\partial z^2}{\partial W^2} = m$$

mxn

$$(j-1)xn+k$$



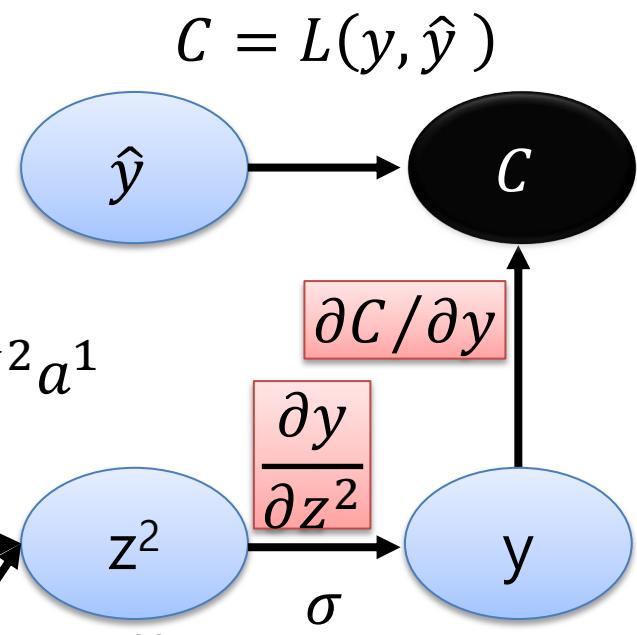
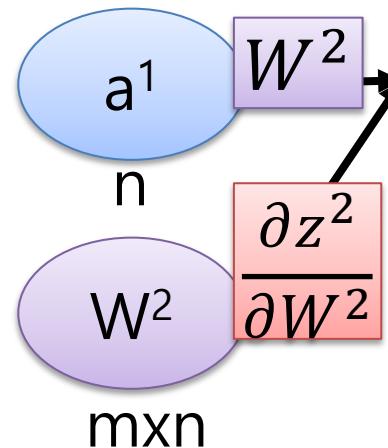
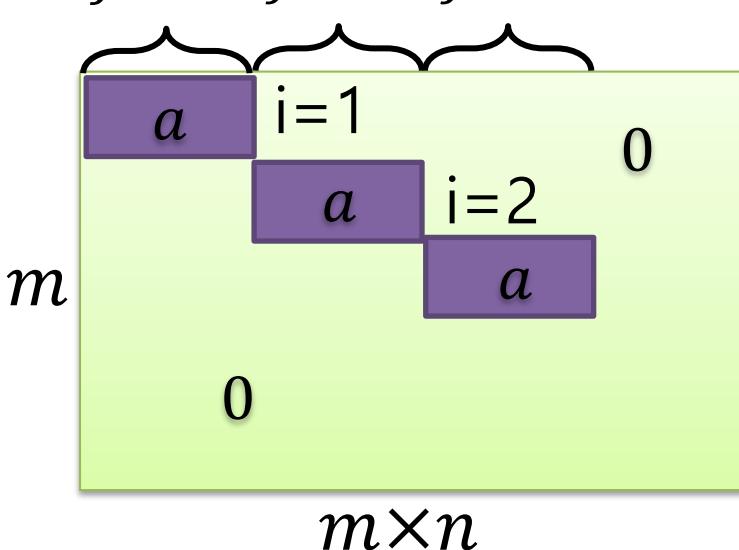
$$\frac{\partial z_i^2}{\partial W_{jk}^2}$$

$$\frac{\partial z_i^2}{\partial W_{jk}^2} = ? \quad i \neq j: \frac{\partial z_i^2}{\partial W_{jk}^2} = 0$$

$$i = j: \frac{\partial z_i^2}{\partial W_{ik}^2} = a_k^1$$

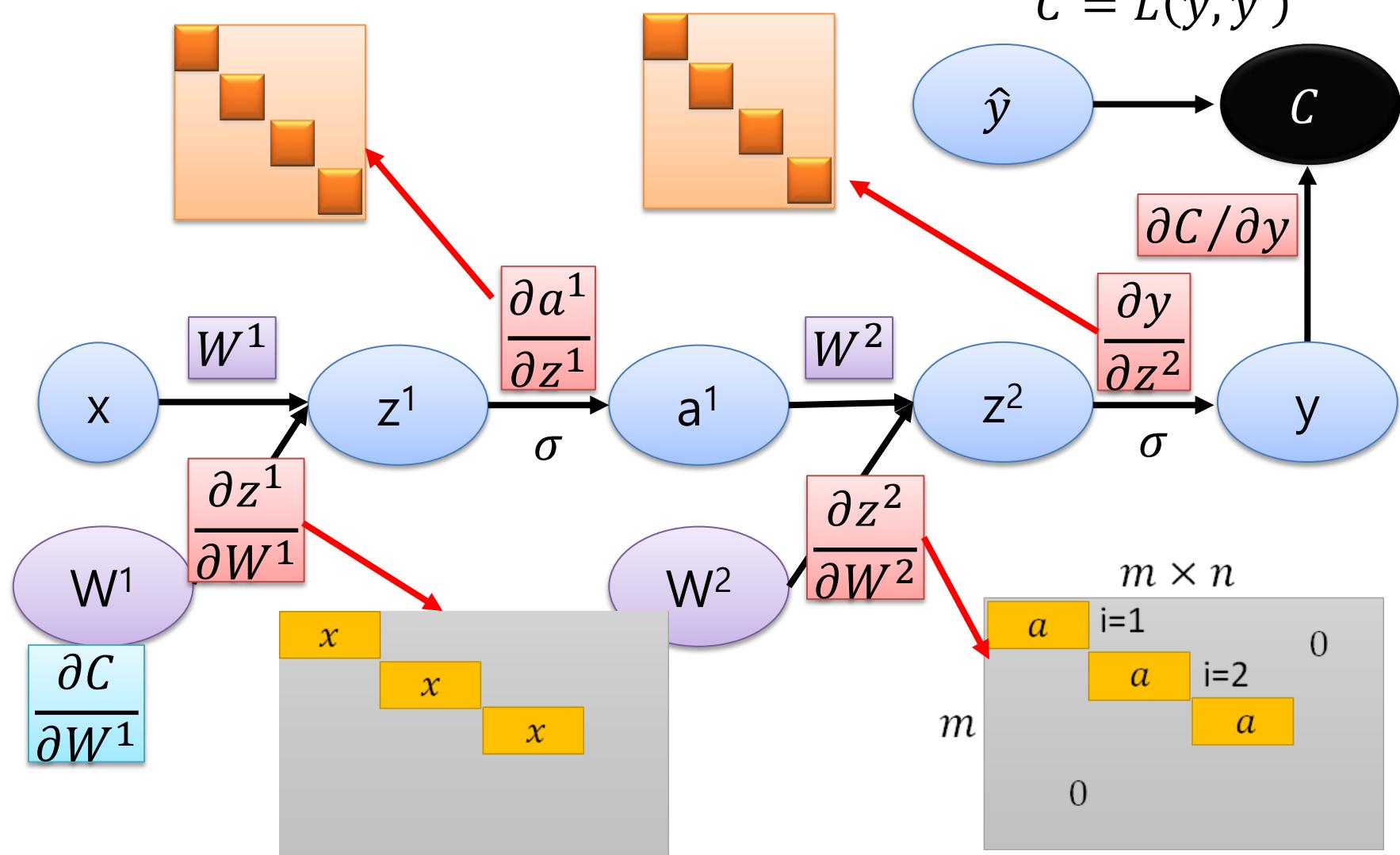
$$z^2 = W^2 a^1$$

$$j = 1 \quad j = 2 \quad j = 3$$

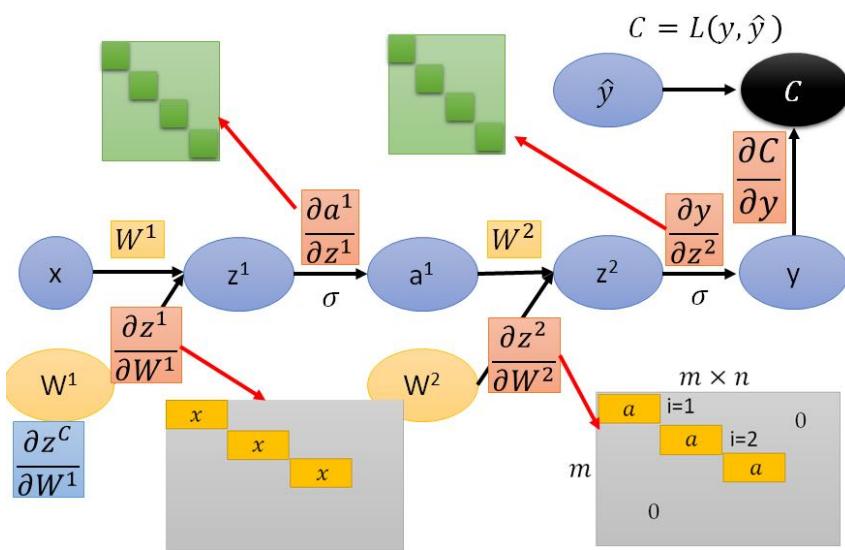


Considering  
 $W^2$  as a  $mxn$  vector

$$\frac{\partial C}{\partial W^1} = \frac{\partial C}{\partial y} \frac{\partial y}{\partial z^2} W^2 \quad \frac{\partial a^1}{\partial z^1} \frac{\partial z^1}{\partial W^1} = [\cdots \frac{\partial C}{\partial W_{ij}^1} \cdots]$$



# Question



$$\frac{\partial C}{\partial w_{ij}^l} = \boxed{\frac{\partial z_i^l}{\partial w_{ij}^l}} \boxed{\frac{\partial C}{\partial z_i^l}}$$

Error signal

### Forward Pass

$$z^1 = W^1 x + b^1$$

$$a^1 = \sigma(z^1)$$

.....

$$z^{l-1} = W^{l-1} a^{l-2} + b^{l-1}$$

$$a^{l-1} = \sigma(z^{l-1})$$

### Backward Pass

$$\delta^L = \sigma'(z^L) \bullet \nabla_y C$$

$$\delta^{L-1} = \sigma'(z^{L-1}) \bullet (W^L)^T \delta^L$$

$$\delta^l = \sigma'(z^l) \bullet (W^{l+1})^T \delta^{l+1}$$

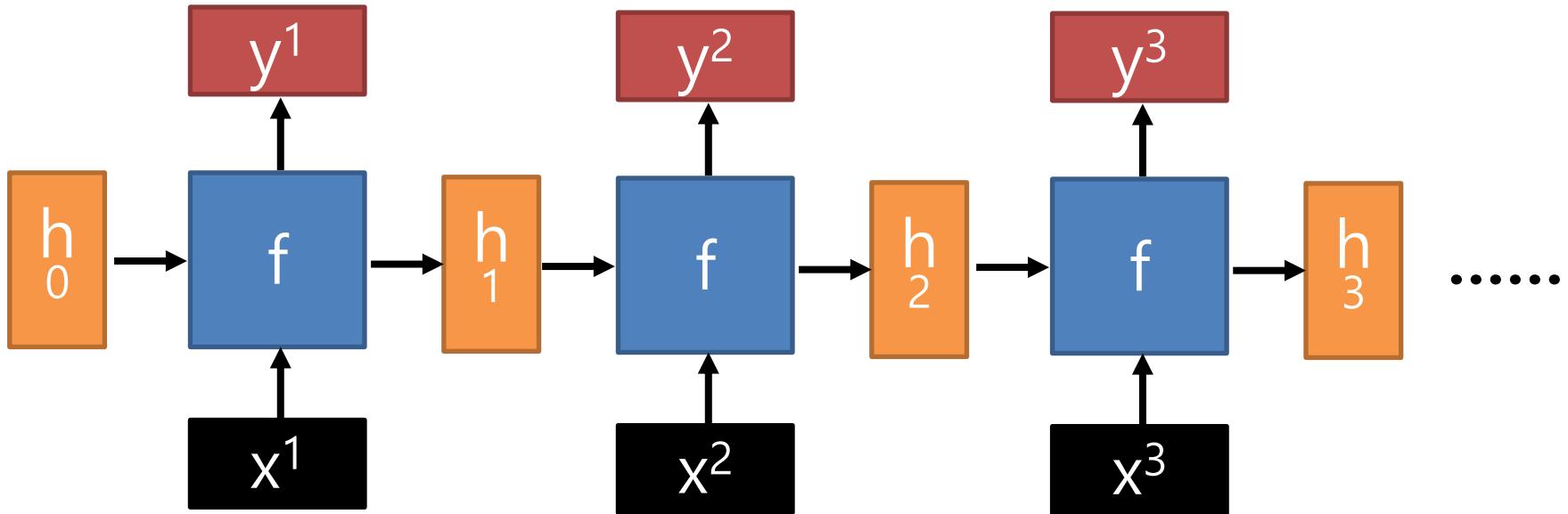
.....

Q: Only backward pass for computational graph?

Q: Do we get the same results from the two different approaches?

# Computational Graph for Recurrent Network

# Recurrent Network



$$y^t, h^t = f(x^t, h^{t-1}; W^i, W^h, W^o)$$

$$h^t = \sigma(W^i x^t + W^h h^{t-1})$$

$$y^t = softmax(W^o h^t)$$

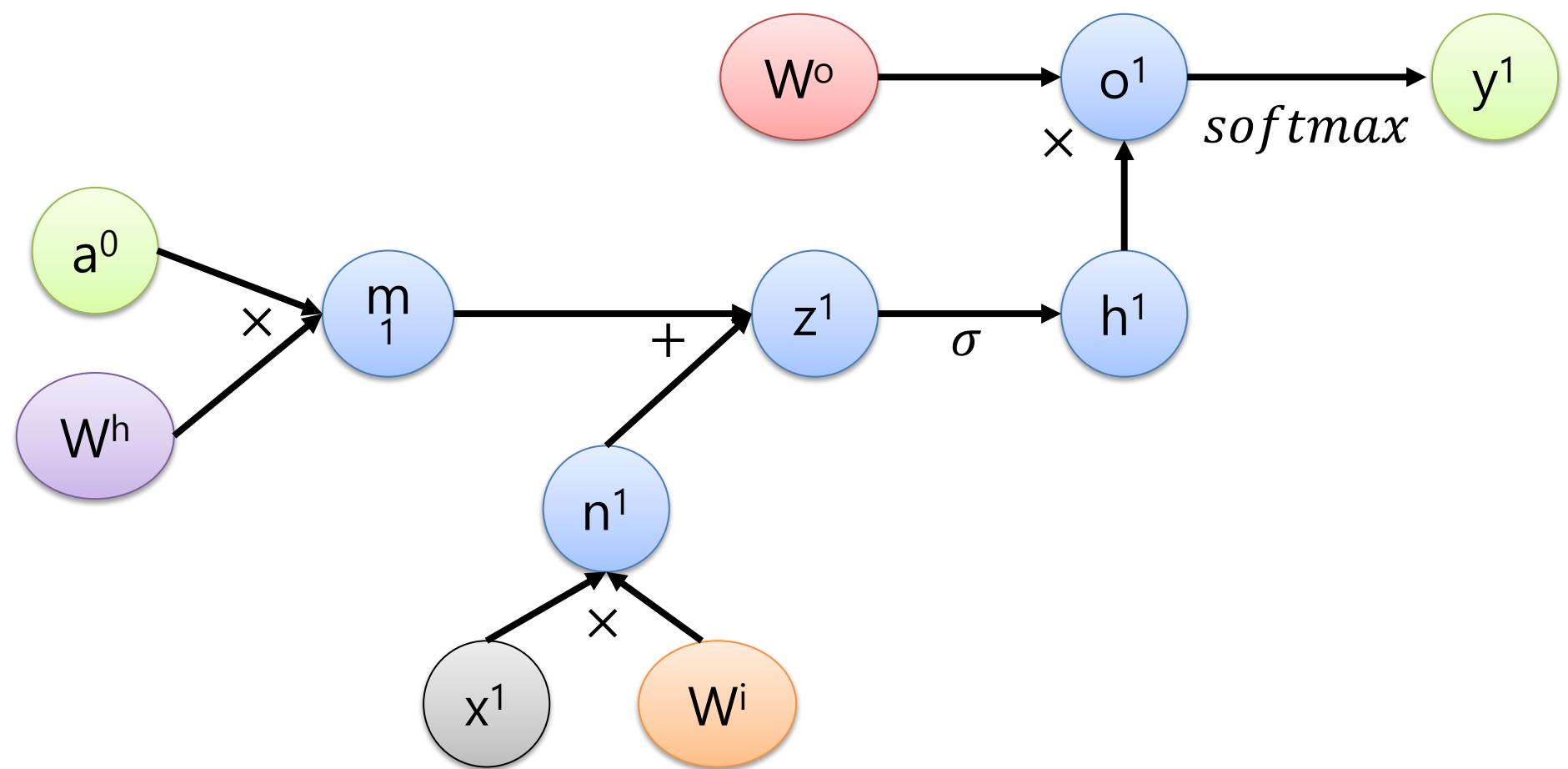
(biases are ignored here)

$$y^t, h^t = f(x^t, h^{t-1}; W^i | W^h, W^o)$$

# Recurrent Network

$$a^t = \sigma(W^i x^t + W^h h^{t-1})$$

$$y^t = \text{softmax}(W^o h^t)$$

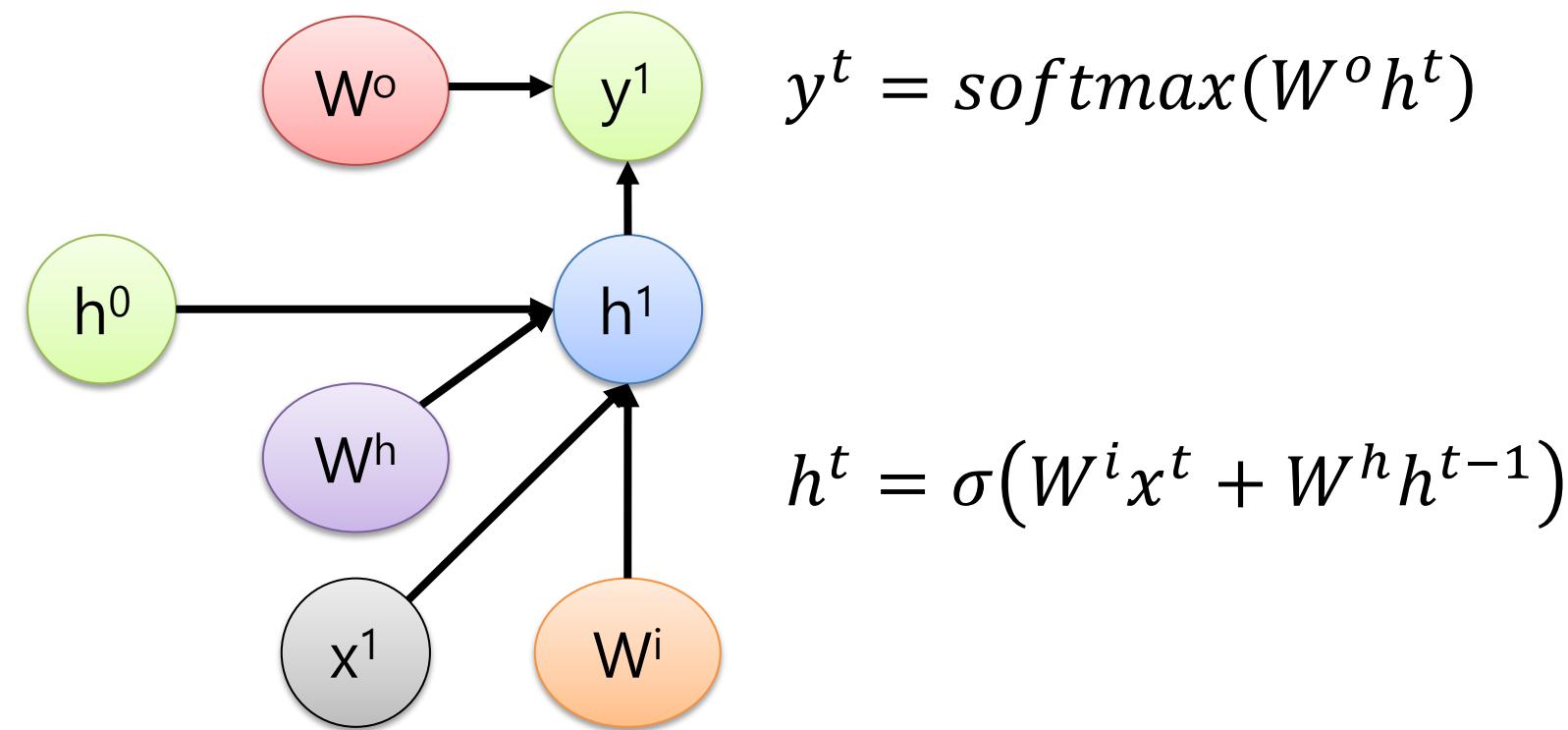


$$y^t, h^t = f(x^t, h^{t-1}; W^i, W^h, W^o)$$

# Recurrent Network

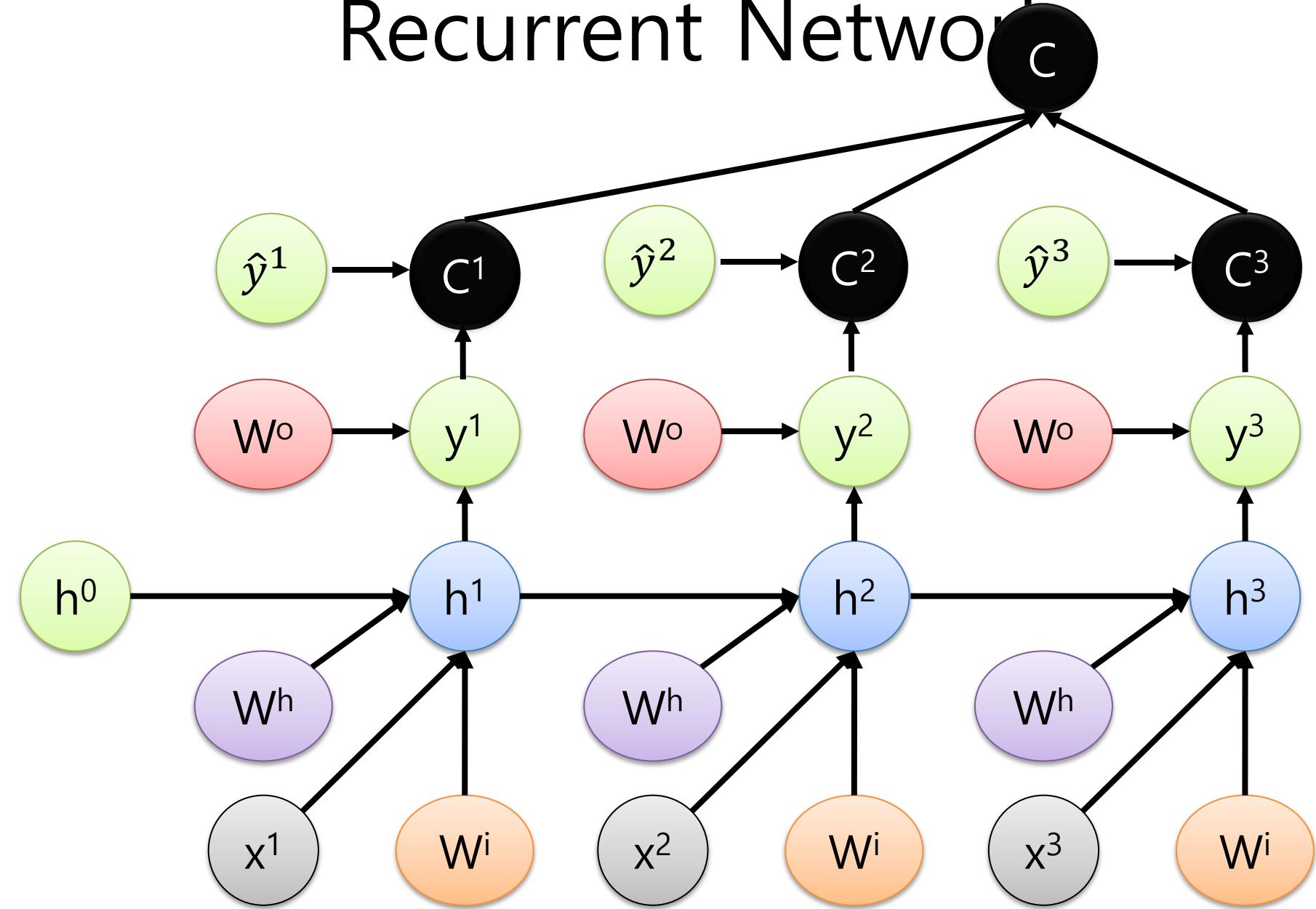
$$h^t = \sigma(W^i x^t + W^h h^{t-1})$$

$$y^t = \text{softmax}(W^o h^t)$$



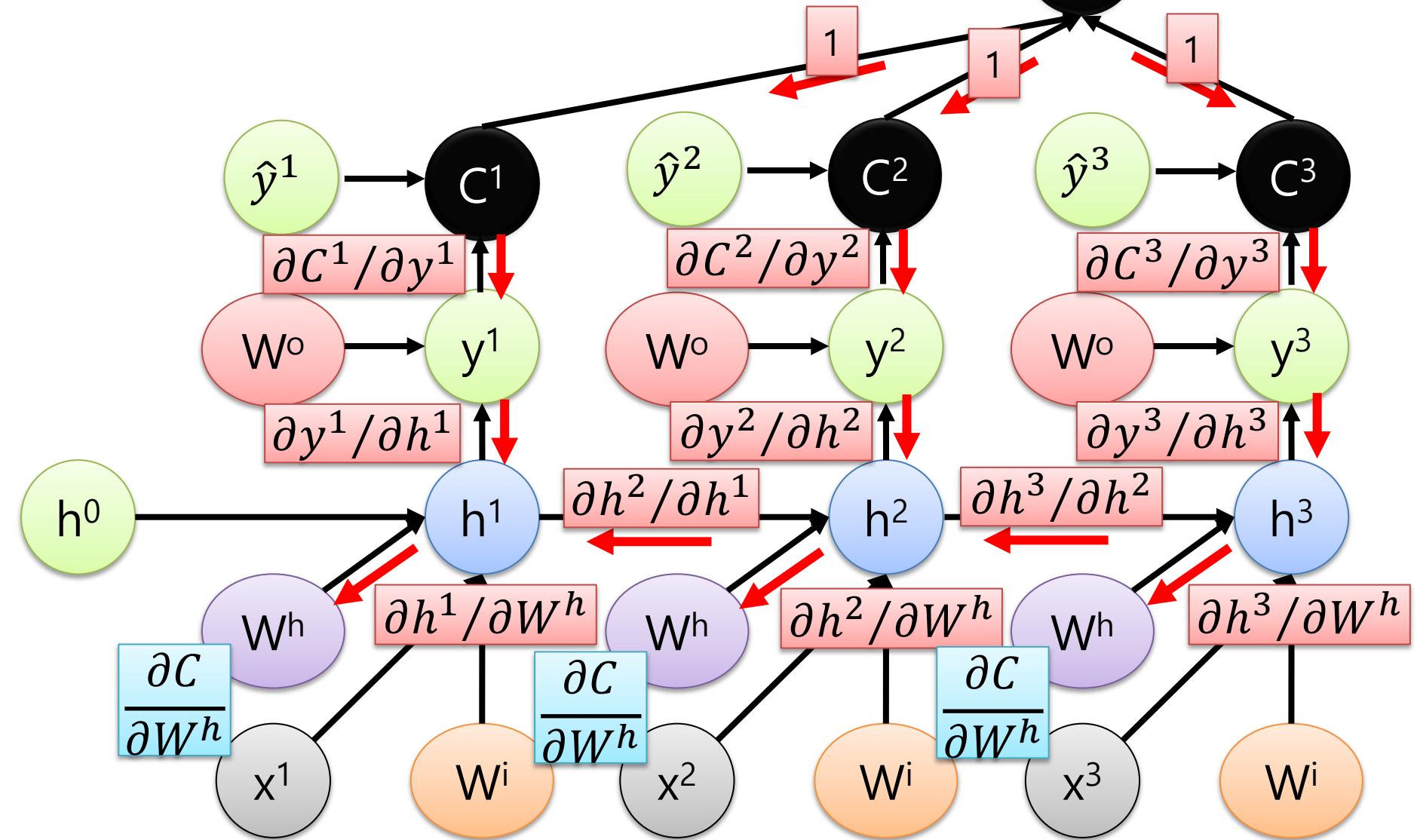
$$C = C^1 + C^{2^{73}} + C^3$$

# Recurrent Network



$$C = C^1 + C^{2^{74}} + C^3$$

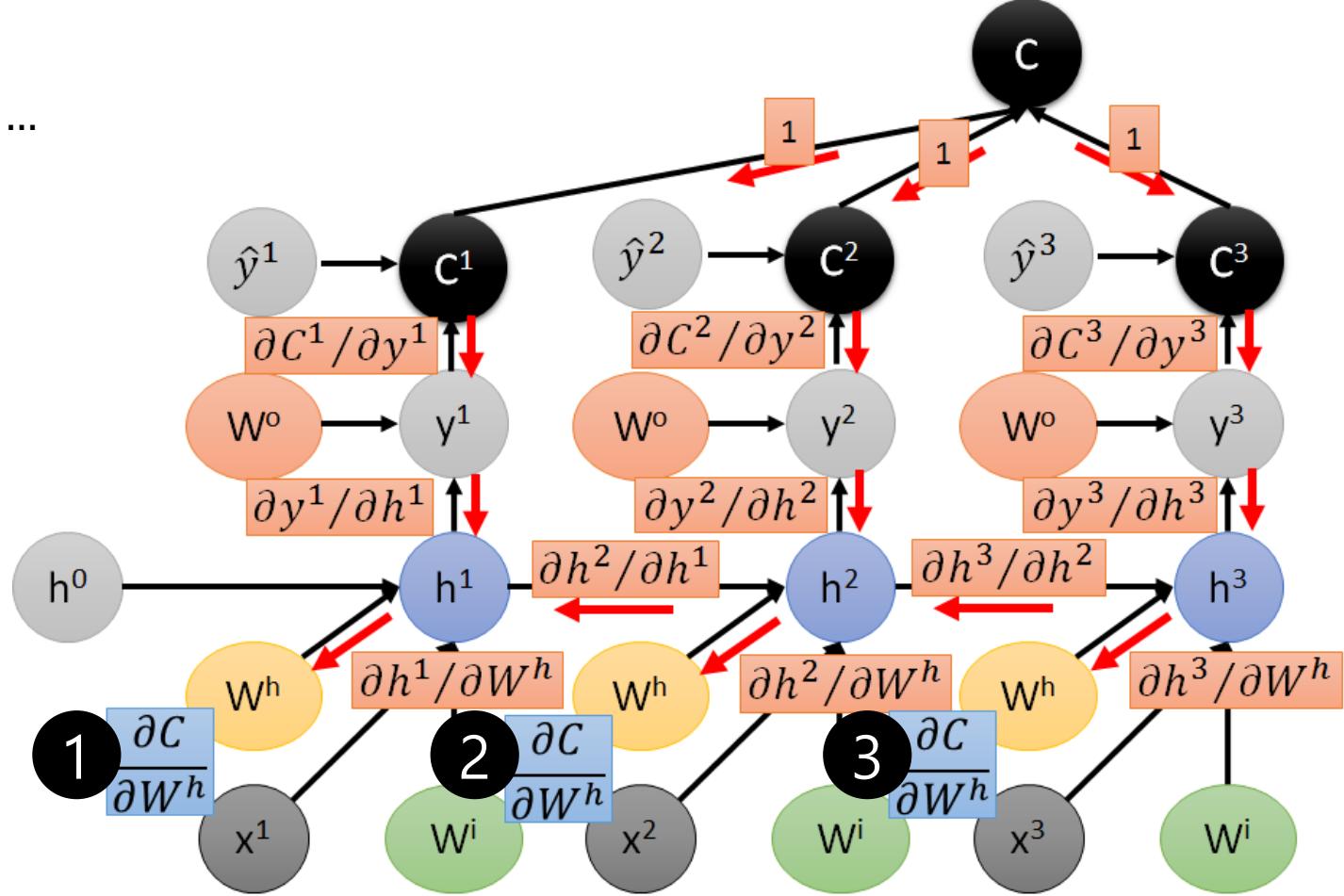
# Recurrent Network



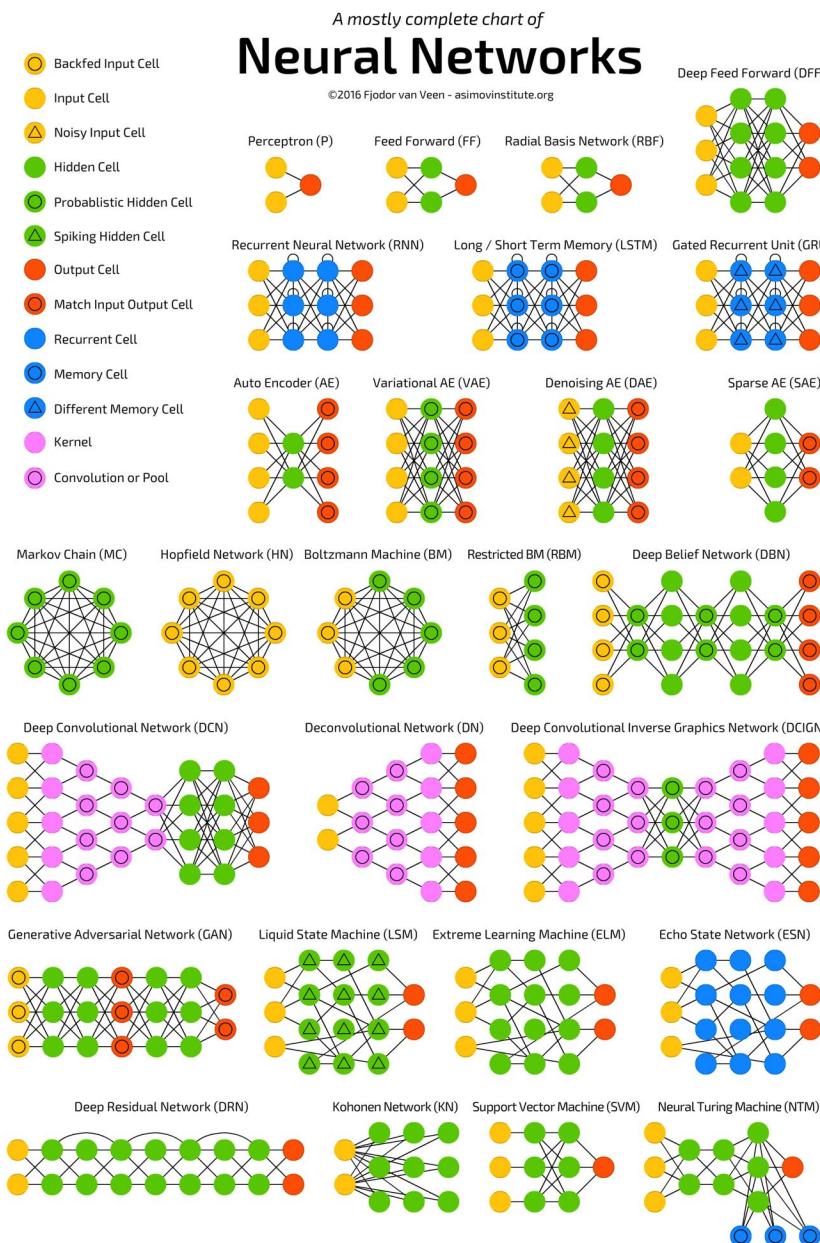
$$1 \frac{\partial C}{\partial W^h} = \left[ \frac{\partial C^1}{\partial y^1} \frac{\partial y^1}{\partial h^1} + \frac{\partial C^2}{\partial y^2} \frac{\partial y^2}{\partial h^1} + \frac{\partial C^3}{\partial y^3} \frac{\partial y^3}{\partial h^1} \frac{\partial h^2}{\partial h^1} \right] \frac{\partial h^1}{\partial W^h}$$

$$2 \frac{\partial C}{\partial W^h} = \dots \dots \quad \frac{\partial C}{\partial W^h} = 1 \frac{\partial C}{\partial W^h} + 2 \frac{\partial C}{\partial W^h} + 3 \frac{\partial C}{\partial W^h}$$

$$3 \frac{\partial C}{\partial W^h} = \dots \dots$$







Difficulties		해결 방안
학습	DNN 학습이 잘 안 됨. Vanishing Gradient	Unsupervised Pre-Training를 통한 해결 ReLU(Rectified Linear Unit Function)
계산량	학습이 많은 계산이 필요함	H/W의 발전 및 GPU 활용
성능	다른 Machine Learning Algorithm의 높은 성능	Dropout 알고리즘 등으로 Machine Learning 대비 월등한 성능

## 학습

신경망에서 원하는 결과를 얻기 위해  
뉴런 사이의 적당한 가중치를 알아내는 것

- 훈련 데이터 (Training Set) 준비 : 입력 데이터와 출력 데이터
- 신경망에 데이터 훈련 : 출력층 값 확인
- 지도학습 데이터와 차이 (오차) 계산
- 오차가 최대한 작도록 가중치를 최적화  
(Gradient Method(경사하강법)을 이용)  
(Back Propagation(오류 역전파))