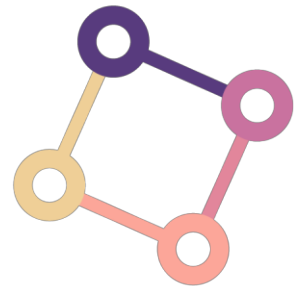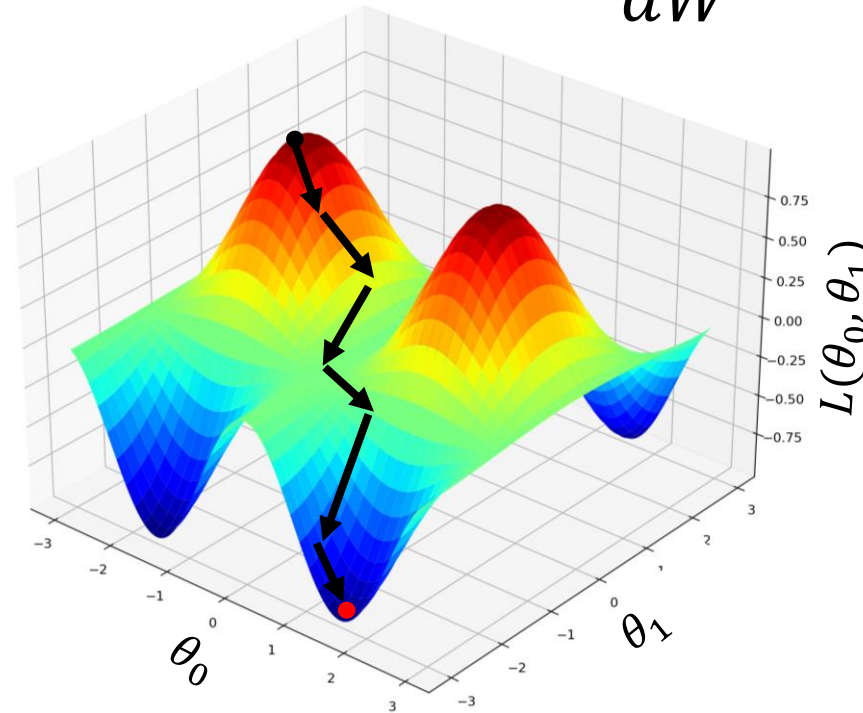# Training Neural Networks

## 주재걸 교수

KAIST 김재철AI대학원

DAVIAN
Data and Visual Analytics Lab

# Training Neural Networks via Gradient Descent

Given the optimization problem, $\min_{W} L(W)$ , where $W$ is the neural

network parameters, we optimize $W$ using gradient descent approach:

loss function
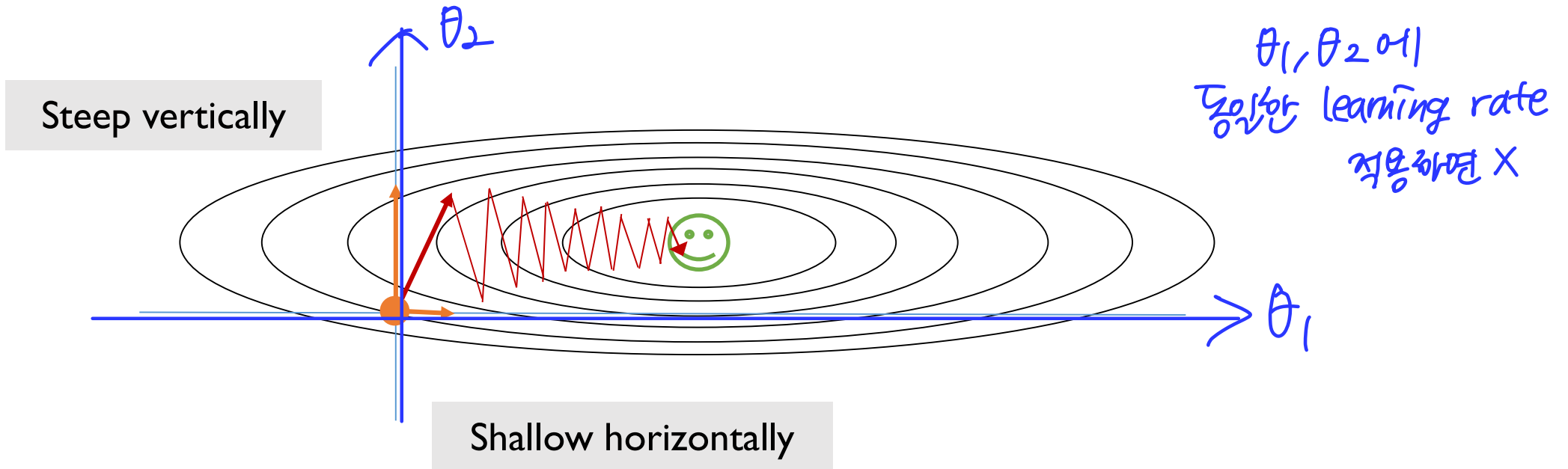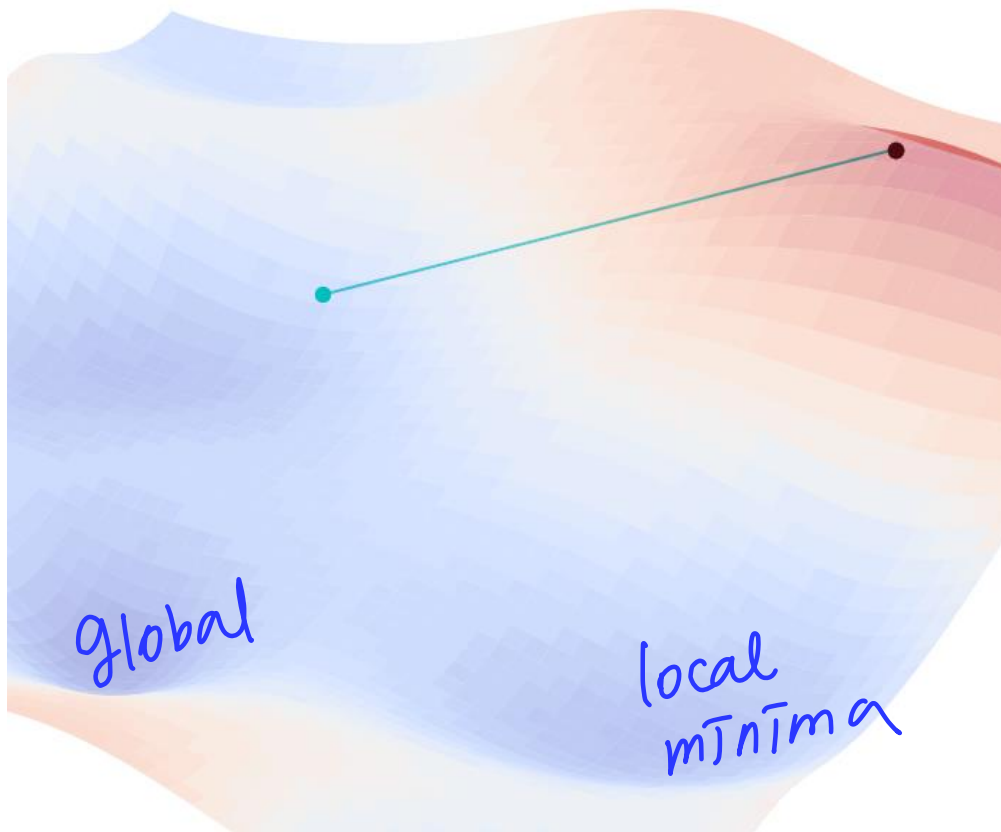
$$W := W - \alpha \frac{dL(W)}{dW}$$

Suppose loss function is steep vertically but shallow horizontally:

Q: What is the trajectory along which we converge towards the minimum with SGD?

Very slow progress along flat direction, jitter along steep one

Steep vertically

Shallow horizontally

$\theta_2$

$\theta_1$

$\theta_1, \theta_2$ 에 동일한 learning rate 적용하면 X

# Various Gradient Descent Methods



global

local
minima

http://www.denizyuret.com/2015/03/alec-radfords-animations-for.html

**Legend:**
- GD
- Momentum
- Adagrad
- RMSProp
- Adam 고! (circled)
- Adadelta[1]
- Ftrl[2]

[1] Zeiler, M. D. (2012). ADADELTA: An Adaptive Learning Rate Method. http://arxiv.org/abs/1212.5701
[2] H. Brendan McMahan et. al., (2013). Ad Click Prediction: a View from the Trenches, KDD
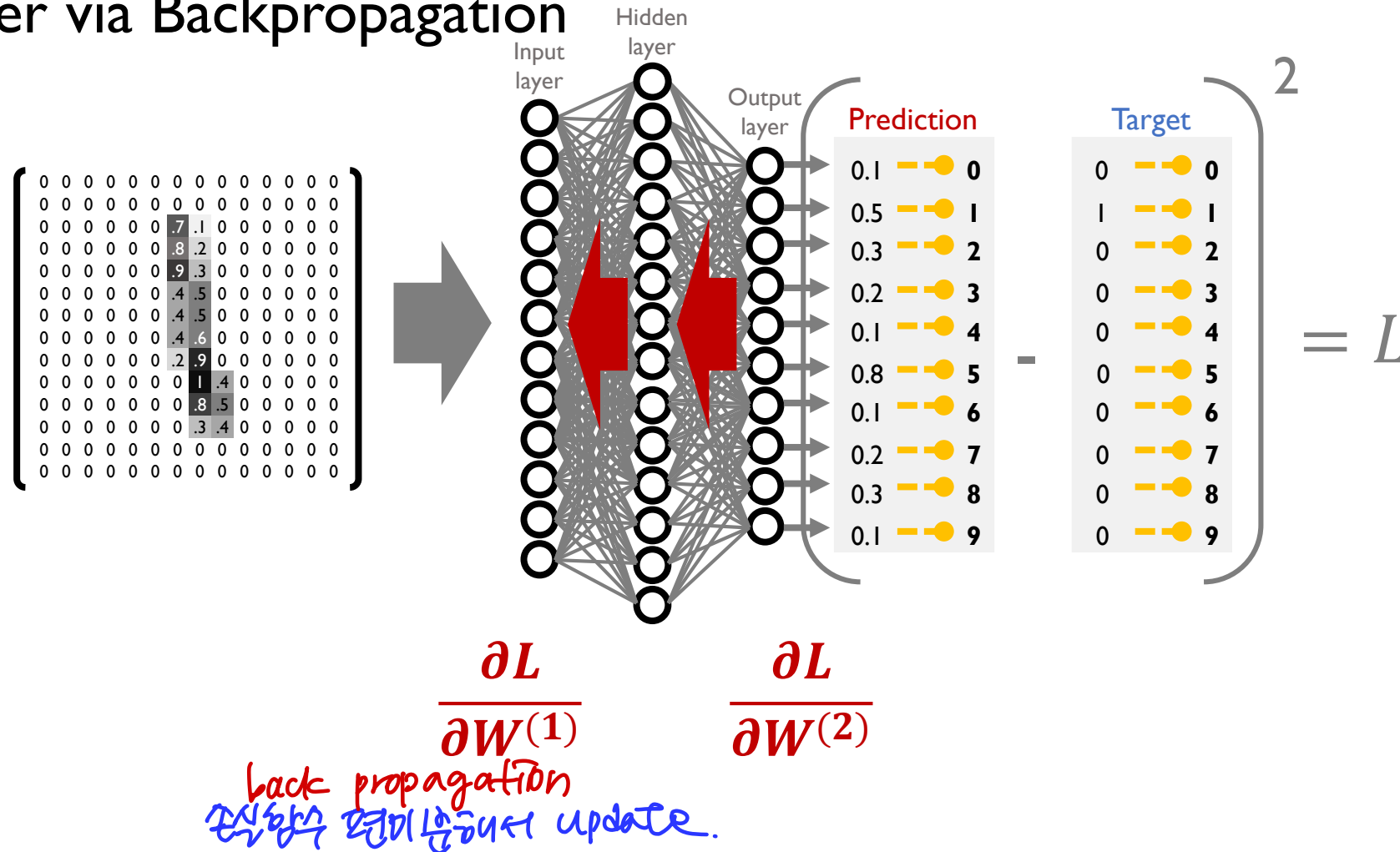
# Backpropagation to Compute Gradient in Neural Networks

First, given an input data item, compute the loss function value via Forward Propagation

Afterwards, compute the gradient with respect to each neural network parameter via Backpropagation



$$\frac{\partial L}{\partial W^{(1)}} \qquad \frac{\partial L}{\partial W^{(2)}}$$

back propagation
뒤에서부터 편미분해서 update.

# Backpropagation to Compute Gradient in Neural Networks

Finally, update the parameters using gradient descent algorithm



$$W^{(1)} := W^{(1)} - \alpha \frac{\partial L}{\partial W^{(1)}} \qquad W^{(2)} := W^{(2)} - \alpha \frac{\partial L}{\partial W^{(2)}}$$

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

# Forward Propagation of Logistic Regression

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

Forward pass:  Computing Output

# Backpropagation of Logistic Regression

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

떨어봄 값x 구함.

Backward pass: Computing gradients



$w_0$  2.00

$x_0$  -1.00

$*$  -2.00

$+$  4.00

$w_1$  -3.00

$x_1$  -2.00

$*$  6.00

$+$  1.00  $*-1$  -1.00  exp  0.37  $+1$  1.37  $\frac{1}{x}$  0.73

1.00

$\frac{dL}{dL}$

Base Case

$w_2$  -3.00

# Backpropagation of Logistic Regression

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

Backward pass: Computing gradients



$w_0$   2.00

$x_0$   -1.00

$w_1$   -3.00

$x_1$   -2.00

$w_2$   -3.00

-2.00

4.00

6.00

1.00   -1.00   0.37   1.37   0.73

-0.53   1.00

Local Gradient

$$\frac{\partial}{\partial x}\left[\frac{1}{x}\right] = -\frac{1}{x^2}$$

×1.00

$y = \frac{1}{x}$

$x$

Downstream Gradient

Upstream Gradient

11
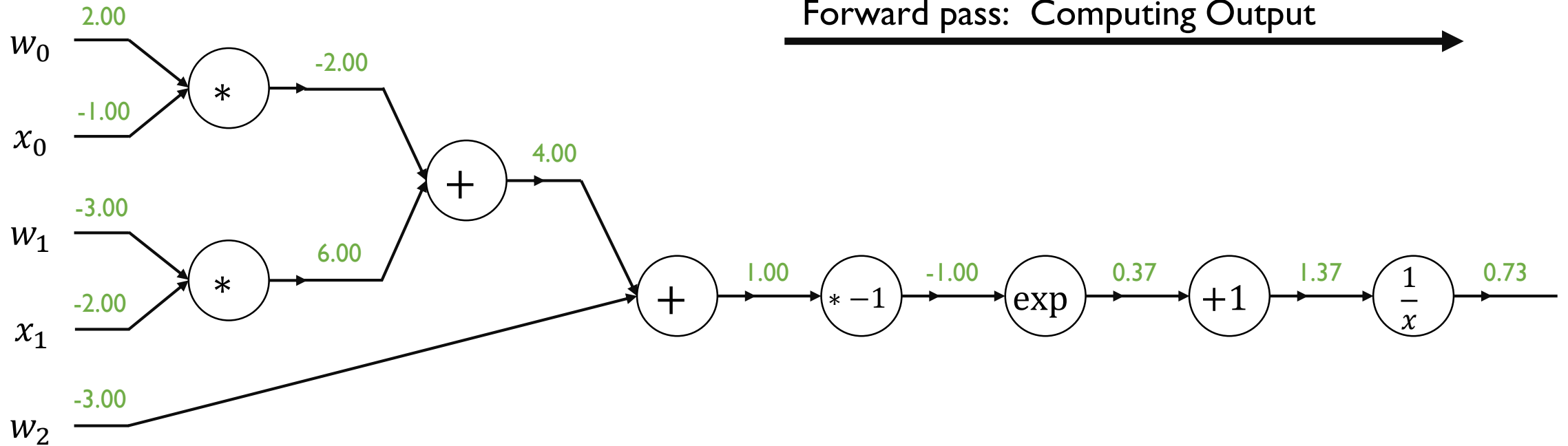
# Backpropagation of Logistic Regression

$$f(x,w) = \frac{1}{1+e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



Backward pass: Computing gradients

Local Gradient

$$\frac{\partial}{\partial x}[e^x] = e^x$$

Downstream Gradient

Upstream Gradient
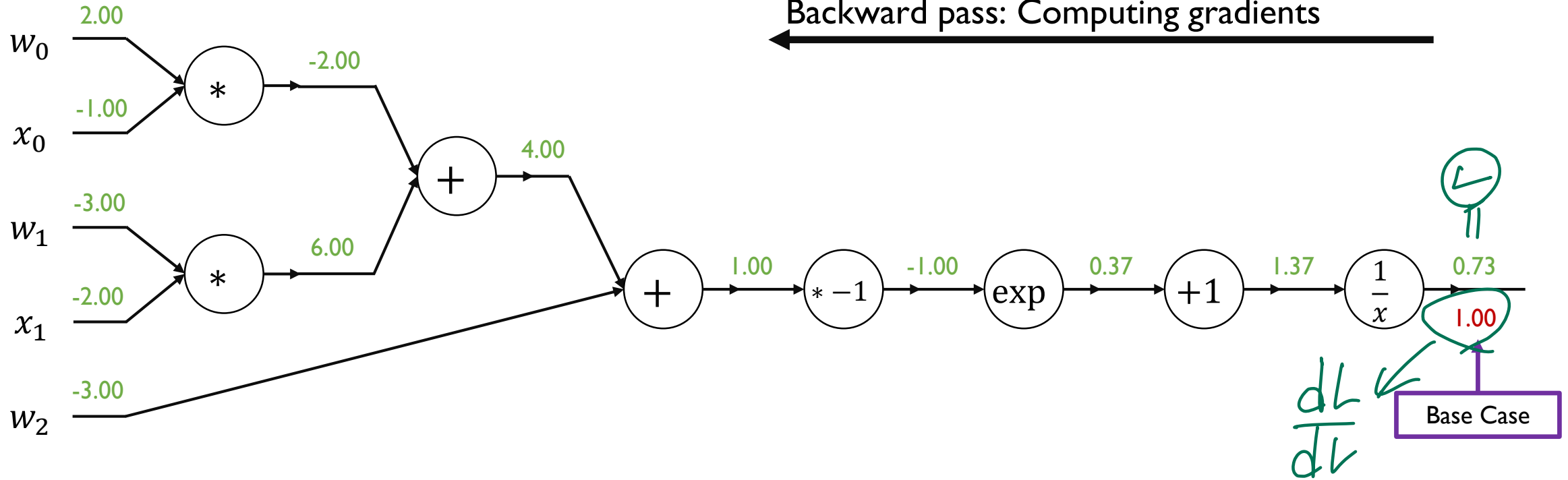
# Backpropagation of Logistic Regression

$$f(x,w) = \frac{1}{1+e^{-(w_0x_0+w_1x_1+w_2)}}$$

Backward pass: Computing gradients

Local Gradient
$$\frac{\partial}{\partial x}[-x] = -1$$

Downstream Gradient

Upstream Gradient

$w_0$   2.00

$x_0$   -1.00

-2.00

$w_1$   -3.00

$x_1$   -2.00

6.00

$w_2$   -3.00

4.00

1.00 / 0.20   *-1   -1.00 / -0.20   exp   0.37 / -0.53   +1   1.37 / -0.53   $\frac{1}{x}$   0.73 / 1.00

# Backpropagation of Logistic Regression

$$f(x,w) = \frac{1}{1+e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$
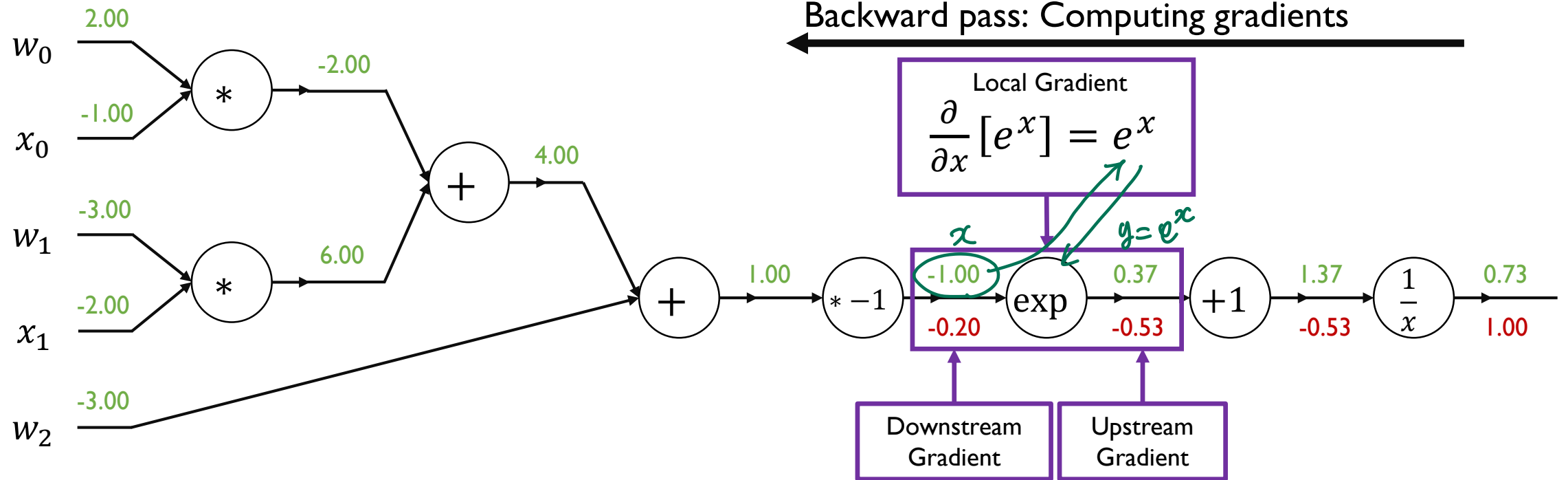


**Backward pass: Computing gradients**

Local Gradient

$$\frac{\partial}{\partial x}[x+y] = 1; \frac{\partial}{\partial y}[x+y] = 1$$

Downstream
Gradient

# Backpropagation of Logistic Regression



다른 예시: $f(x, w) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$

Backward pass: Computing gradients

Local Gradient

$$\frac{\partial}{\partial x}[x + y] = 1 ; \frac{\partial}{\partial y}[x + y] = 1$$
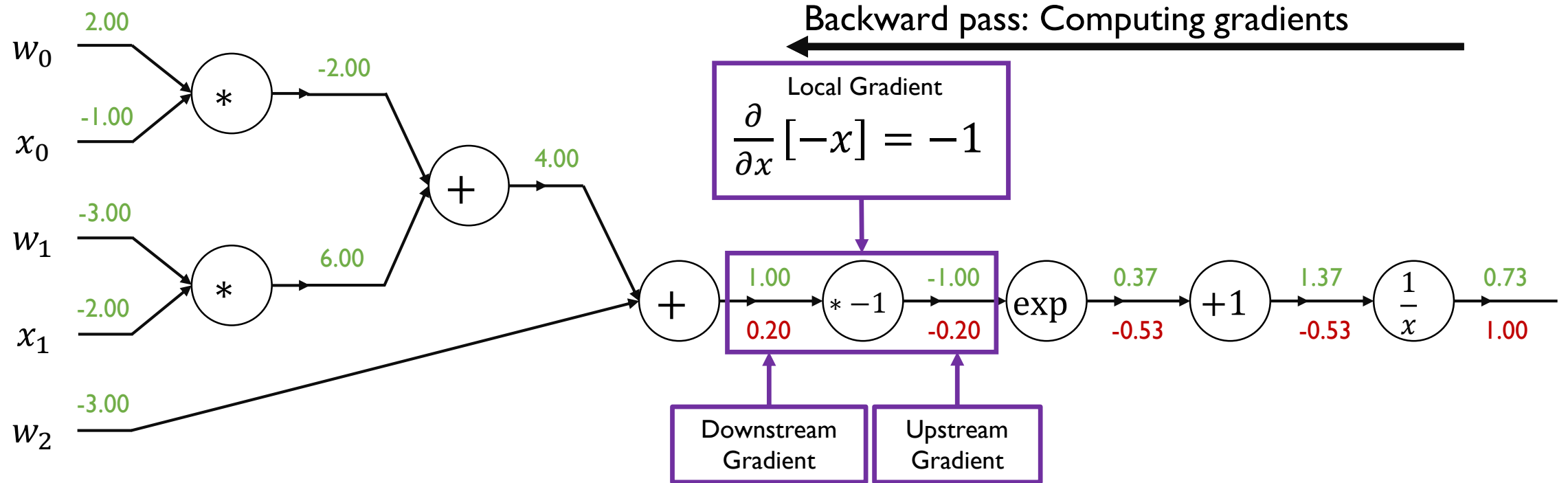
Downstream Gradient

Upstream Gradient

# Backpropagation of Logistic Regression

다른 예시: $f(x, w) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$

# Backpropagation of Logistic Regression

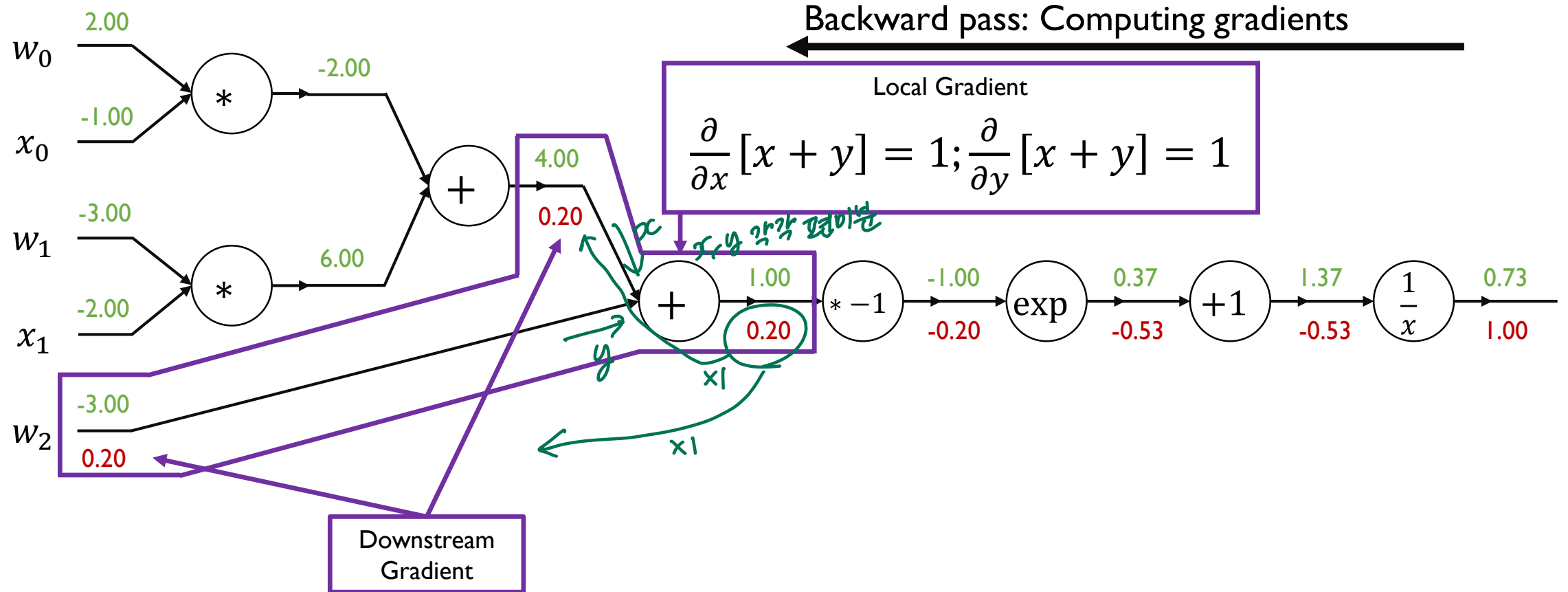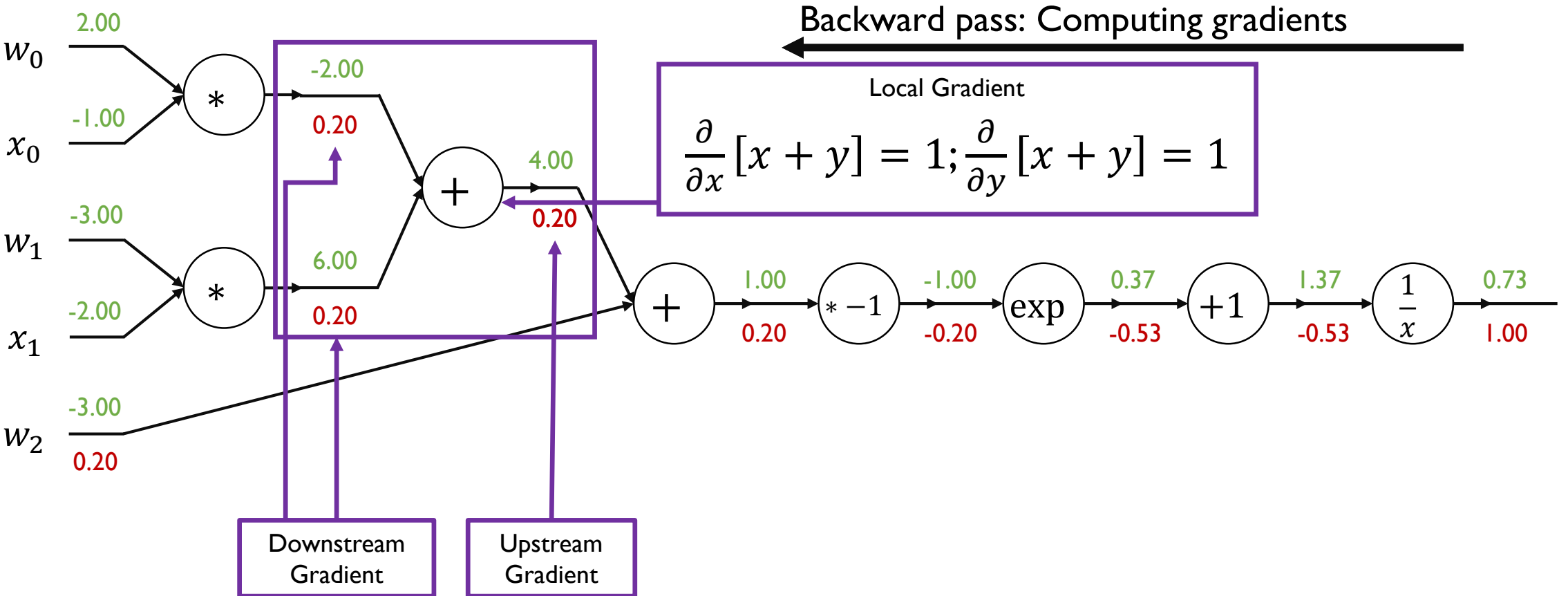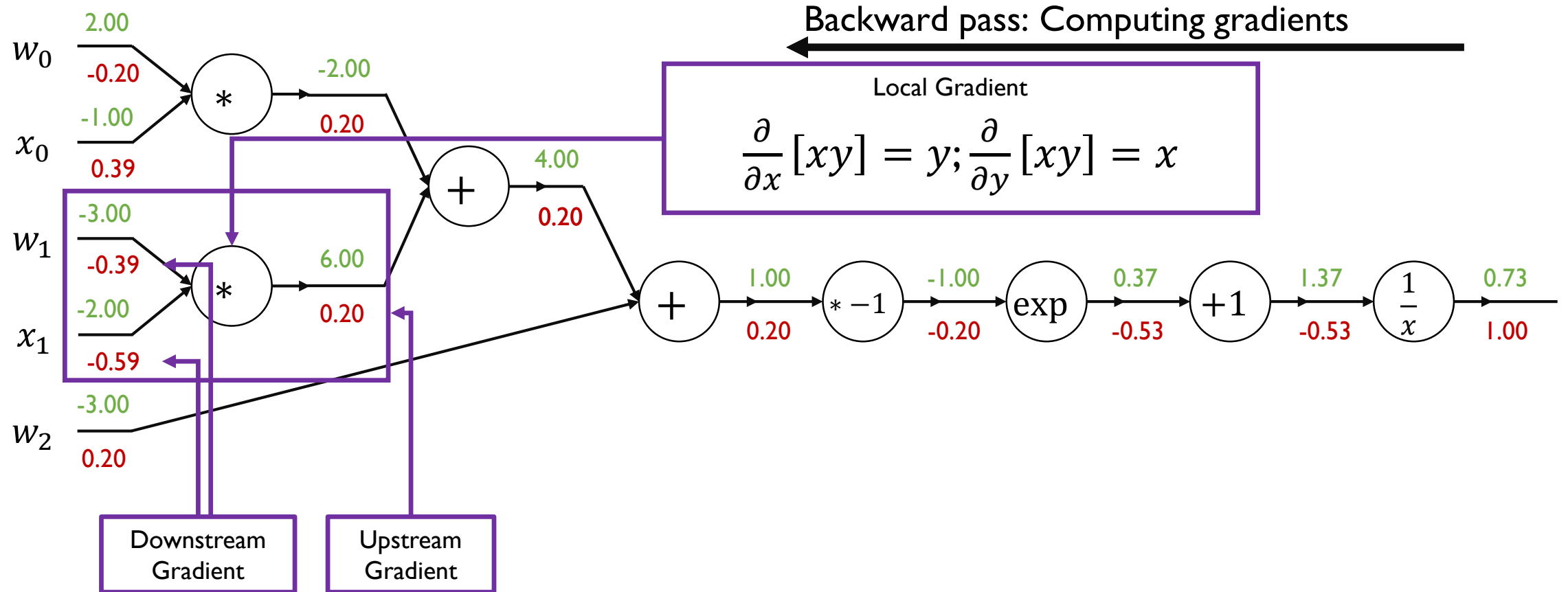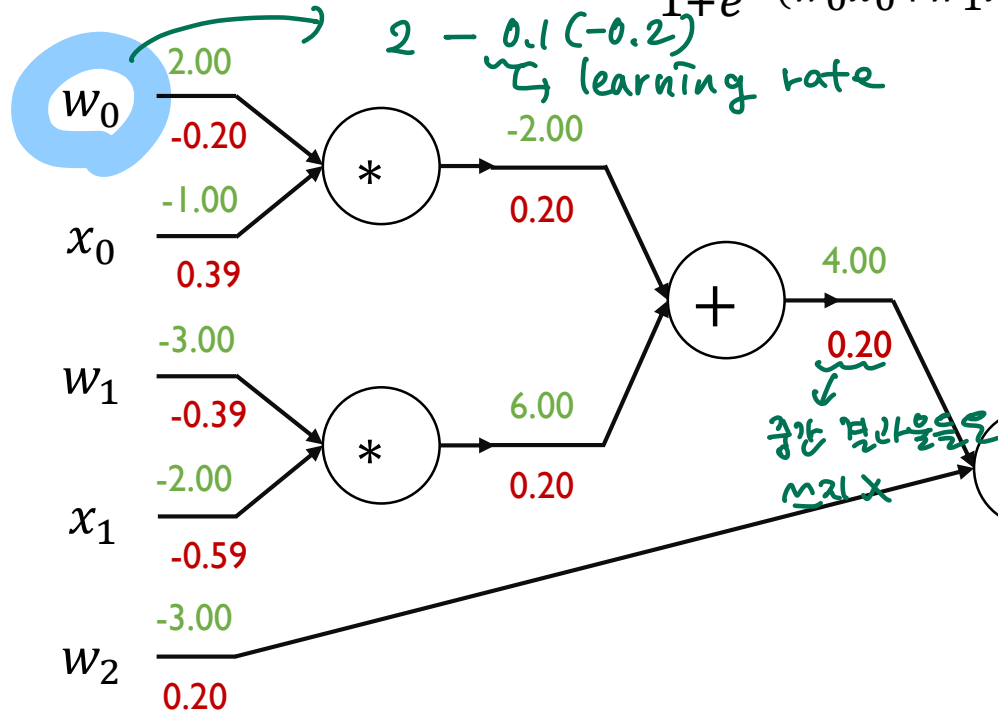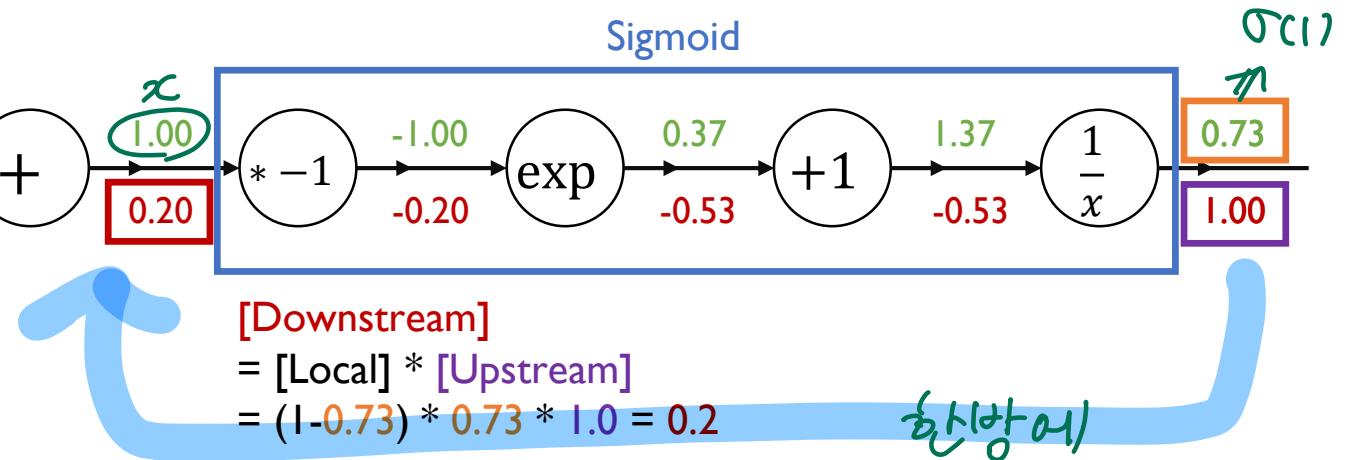다른 예시: $f(x, w) = \dfrac{1}{1+e^{-(w_0x_0 + w_1x_1 + w_2)}} = \boxed{\sigma(w_0x_0 + w_1x_1 + w_2)}$

$2 - 0.1\,(-0.2)$
$\hookrightarrow$ learning rate

Backward pass: Computing gradients $\longleftarrow$

$\sigma(x) = \dfrac{1}{1+e^{-x}}$  One can represent multiple computational steps as a single computational module

$w_0$   2.00
  -0.20

$x_0$   -1.00
  0.39

$*$   -2.00
  0.20

$+$   4.00
  0.20

$w_1$   -3.00
  -0.39

$x_1$   -2.00
  -0.59

$*$   6.00
  0.20

$w_2$   -3.00
  0.20

중간 변수들을 외킨 $x$

$+$   $x$   1.00
  0.20

Sigmoid

$*-1$   -1.00
  -0.20

exp   0.37
  -0.53

$+1$   1.37
  -0.53

$\frac{1}{x}$

$\sigma(1)$

0.73

1.00

[Downstream]
= [Local] * [Upstream]
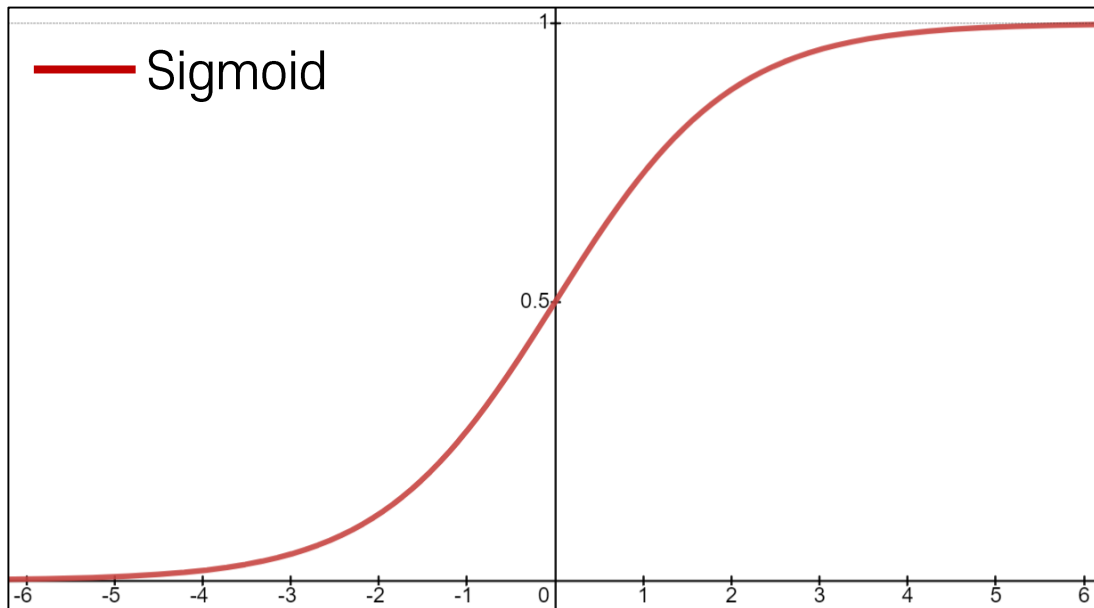= (1-0.73) * 0.73 * 1.0 = 0.2

한방에!

**Sigmoid local gradient**   $\dfrac{\partial}{\partial x}[\sigma(x)] = \dfrac{e^{-x}}{(1+e^{-x})^2} = \left(\dfrac{1+e^{-x}-1}{1+e^{-x}}\right)\left(\dfrac{1}{1+e^{-x}}\right) = (1-\sigma(x))\sigma(x)$

$\hookrightarrow (1-0.73)\,0.73$

# Sigmoid Activation

## Sigmoid

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



- Maps real numbers in $(-\infty, \infty)$ into a range of $[0, 1]$

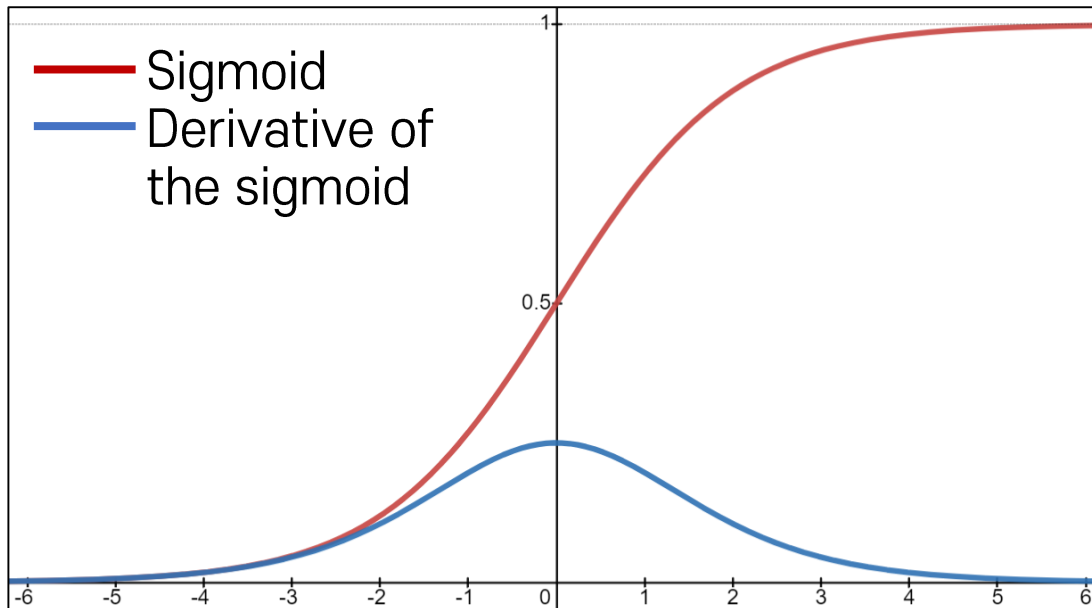- gives a probabilistic interpretation

확률값으로 해석

> Historically, sigmoid activation function gives nice interpretation of saturating firing rate of a neuron

# Problems of Sigmoid Activation

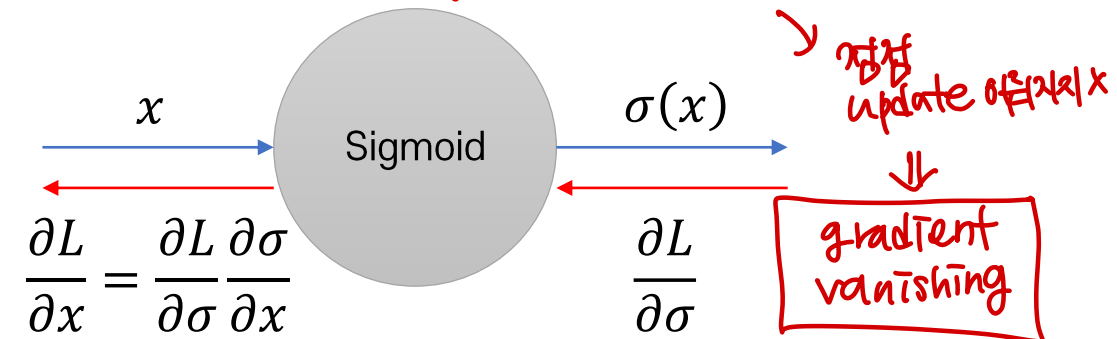$$0 < \sigma(x) \cdot (1 - \sigma(x)) \leq \frac{1}{4}$$

$$\boxed{0 \sim 1}$$

**Sigmoid**

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



Sigmoid
Derivative of the sigmoid

- Saturated neurons kills the gradients

- The gradient value $\sigma(x) \leq \frac{1}{4}$, which decreases the gradient during backpropagation, i.e., causing a gradient vanishing problem

back propagation 계속 곱해서
gradient 점차 작아짐.

정정
update 안되거나 x

gradient vanishing

$$x \rightarrow \boxed{\text{Sigmoid}} \rightarrow \sigma(x)$$

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial \sigma} \frac{\partial \sigma}{\partial x} \qquad \frac{\partial L}{\partial \sigma}$$

# Tanh Activation



-∞, ∞ ⇒ -1, 1

=0

## Tanh

- $\tanh(x) = 2 \times \mathrm{sigmoid}(x) - 1$

- Squashes numbers to range **[-1, 1]**

## Strength

학습이 더 빠르게 일어난다.

- Zero-centered (average is 0)

## Weakness

- Still kills gradients when saturated, i.e., still causing a gradient vanishing problem

$0 \sim \frac{1}{2}$ 임. → gradient vanishing 여전히 존재.

*layer 쌓면 good?*

# ReLU Activation

$-\infty, \infty \Rightarrow 0, \infty$

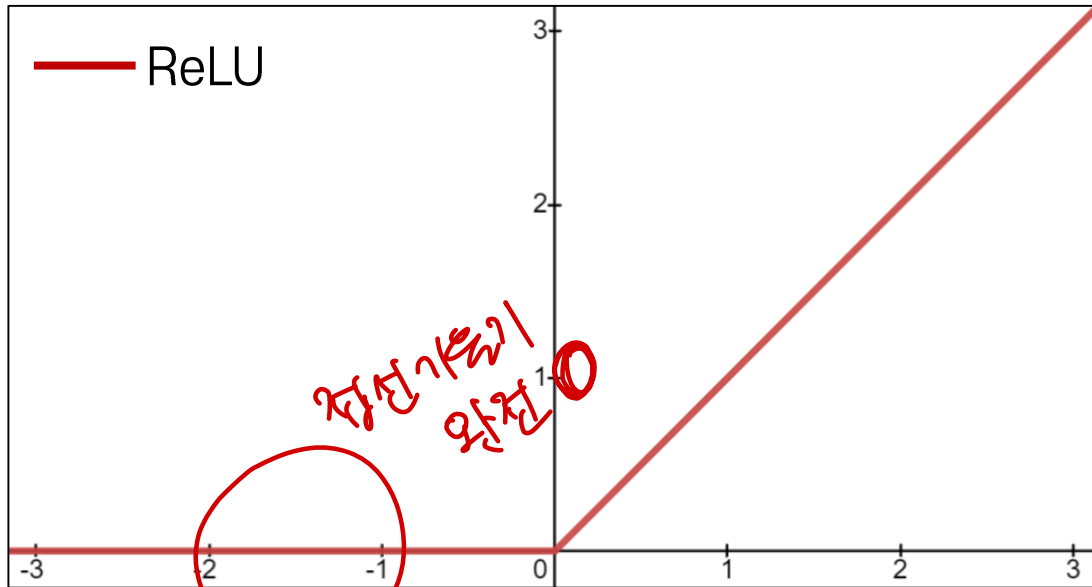**ReLU (Rectified Linear Unit)**

- $f(x) = \max(0, x)$

  *연산 간단*

**Strength**

- Does not saturate (in + region)
- Very computationally efficient
- Converge much faster than sigmoid/tanh

**Weakness**

- Not zero-centered output
- Gradient is completely zero for $x < 0$



ReLU

*정선 기울기 완전 0*

The slope of the function
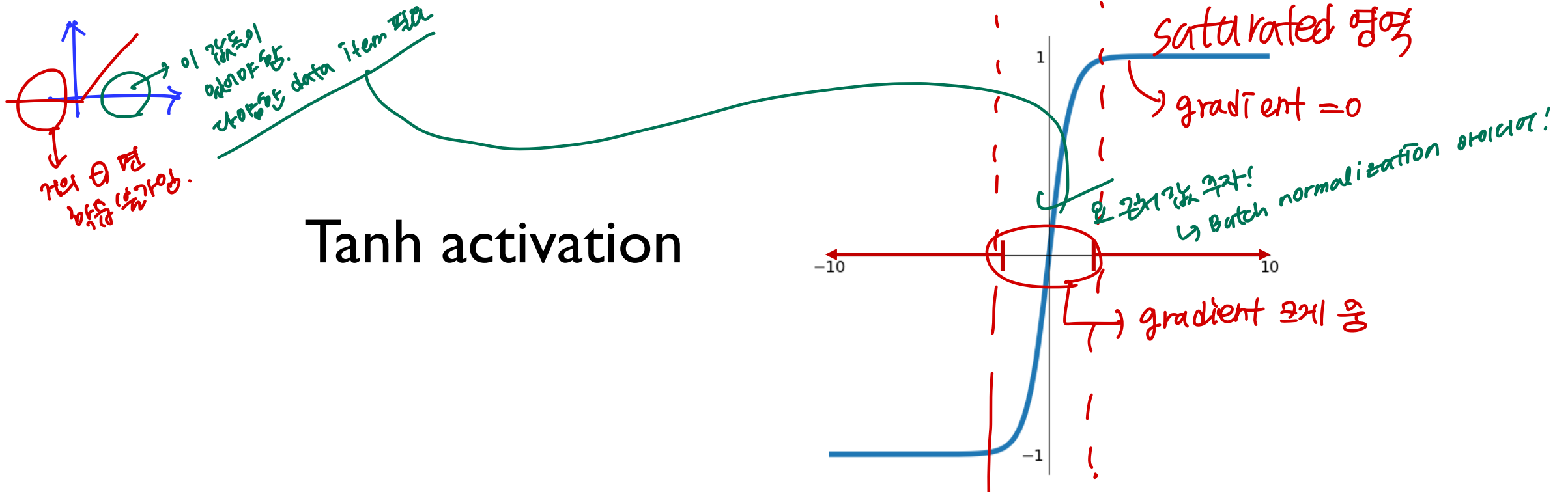$x \geq 0$: 1 (bypass)
$x < 0$: 0 (gating)

# Batch Normalization

- Saturated gradients when random initialization is done
- The parameters are not updated → Hard to optimize (in red region)

Tanh activation

# Definition of Batch Normalization

## "You want unit Gaussian activations? just make them so."

• We consider a batch of activations at some layer to make each dimension unit Gaussian

1. Compute the empirical mean $\mathbb{E}[x^{(k)}]$ and variance $\mathrm{Var}[x^{(k)}]$

    independently for each dimension $k$



2. Normalize

$$\hat{x}_i^{(k)} = \frac{x_i^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\mathrm{Var}[x^{(k)}]}}$$

This is a vanilla differentiable function

Ioffe, Sergey, and Christian, Szegedy. "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift.".
In Proceedings of the 32nd International Conference on Machine Learning (pp. 448–456). PMLR, 2015.

# Batch Normalization Process



FC

*fully-connected layer 나*
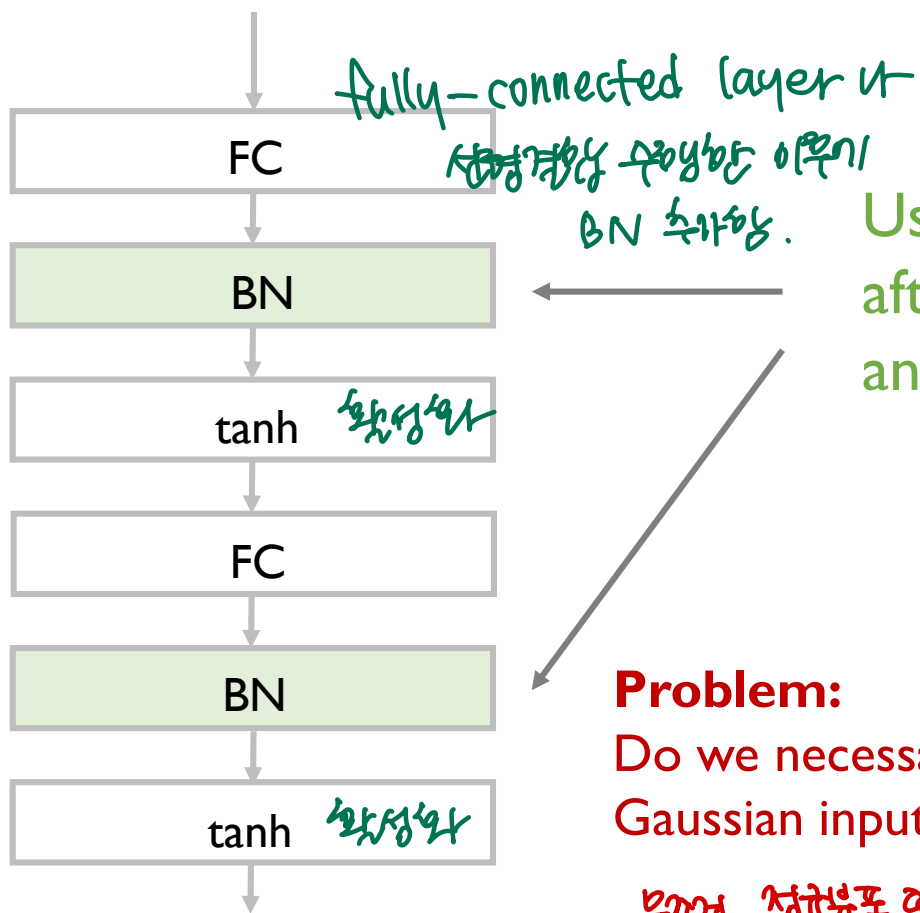
*컨볼루션 수행한 이후에*

*BN 추가함.*

BN

tanh *활성화*

FC

BN

tanh *활성화*

Usually inserted
after fully-connected or convolutional layers,
and before nonlinearity.

**Problem:**
Do we necessarily want a unit
Gaussian input to a tanh layer?

*무조건 정규분포 따르게 하는건*
*neural network가 잘 추출한 정보*
*무시하는 과정일 수 있음.*

$$\hat{x}_i^{(k)} = \frac{x_i^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

Ioffe, Sergey, and Christian, Szegedy. "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift.".
In Proceedings of the 32nd International Conference on Machine Learning (pp. 448–456). PMLR, 2015.

# Batch Normalization

↑
가중치 소실 문제 해결 할수 있는 장치.

본래 표준한 분산으로
복원 가능.
←

Normalize:

$$\hat{x}_i^{(k)} = \frac{x_i^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\mathrm{Var}[x^{(k)}]}}$$

And then allow the network to squash the range if it wants to:

$ax+b$

$$y_i^{(k)} = \gamma^{(k)} \hat{x}_i^{(k)} + \beta^{(k)}$$

→ 평균 b,
분산 $a^2$ 으로 바꿈.

원하는 정보로 복원

Note, the network can learn:

$$\gamma^{(k)} = \sqrt{\mathrm{Var}[x^{(k)}]}$$

$$\beta^{(k)} = \mathbb{E}[x^{(k)}]$$

to recover the identity mapping

← 평균 분산
스스로 결정하도록

## Determined while training neural network

Ioffe, Sergey, and Christian, Szegedy. "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift.".
In Proceedings of the 32nd International Conference on Machine Learning (pp. 448–456). PMLR, 2015.

# Batch Normalization

**Input:** Values of $x$ over a mini-batch: $B = \{x_{1...m}\}$;
          Parameters to be learned; $\gamma, \beta$

**Output:** $\{y_i = \text{BN}_{\gamma,\beta}(x_i)\}$

$\mu_B \leftarrow \frac{1}{m}\sum_{i=1}^{m} x_i$            // mini-batch mean

$\sigma_B^2 \leftarrow \frac{1}{m}\sum_{i=1}^{m}(x_i - \mu_B)^2$     // mini-batch variance

$\hat{x}_i \leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$             // normalize

$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma,\beta}(x_i)$     // scale and shift

- Improves gradient flow through the network

- Reduces the strong dependence on initialization

Ioffe, Sergey, and Christian, Szegedy. "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift.". In Proceedings of the 32nd International Conference on Machine Learning (pp. 448–456). PMLR, 2015.

THANK YOU