# Machine Learning

## Neural Network

Dr. Shuang LIANG

# Today's Topics

- Neural Network Introduction

- Neural Network Structure

- How Neural Network Works

- Backpropagation

# Today's Topics

- ***Neural Network Introduction***

- Neural Network Structure

- How Neural Network Works

- Backpropagation

# Neural networks are a hot topic

# A bit of history

The **Mark I Perceptron** machine was the first implementation of the perceptron algorithm.
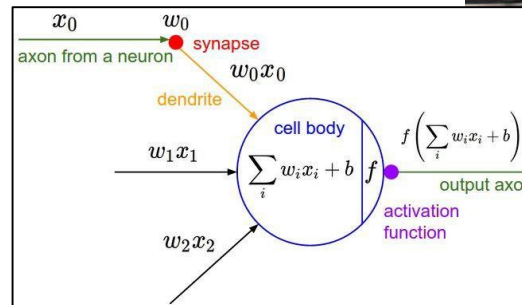
The machine was connected to a camera that used 20×20 cadmium sulfide photocells to produce a 400-pixel image.

recognized
letters of the alphabet

$$f(x) = \begin{cases} 1 & \text{if } w \cdot x + b > 0 \\ 0 & \text{otherwise} \end{cases}$$

update rule:

$$w_i(t+1) = w_i(t) + \alpha(d_j - y_j(t))x_{j,i},$$



$x_0$   $w_0$   synapse
axon from a neuron
$w_0 x_0$
dendrite
cell body
$w_1 x_1$
$\sum_i w_i x_i + b$   $f$
$f\left(\sum_i w_i x_i + b\right)$
output axon
activation function
$w_2 x_2$

Frank Rosenblatt, ~1957: Perceptron

# A bit of history



Widrow and Hoff, ~1960: Adaline/Madaline

These figures are reproduced from Widrow 1960, Stanford Electronics Laboratories Technical Report with permission from Stanford University Special Collections.

# A bit of history

$$\frac{\partial E_p}{\partial w_{ji}} = \frac{\partial E_p}{\partial o_{pj}} \frac{\partial o_{pj}}{\partial w_{ji}}$$



input pattern

$w_{ji}$

$o_{pj}$

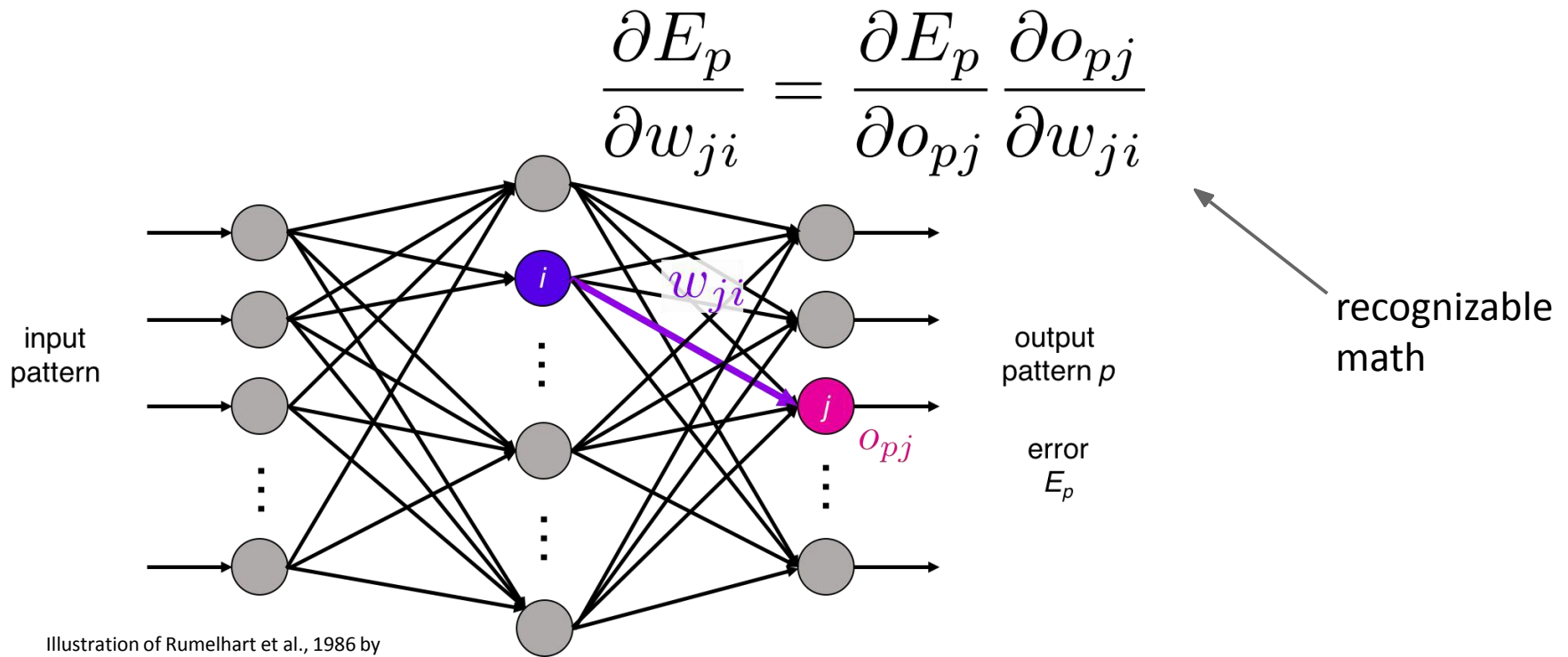output pattern $p$

error $E_p$

recognizable math

Illustration of Rumelhart et al., 1986 by
Lane McIntosh, copyright CS231n 2017

Rumelhart et al., 1986: First time back-propagation became popular

# A bit of history

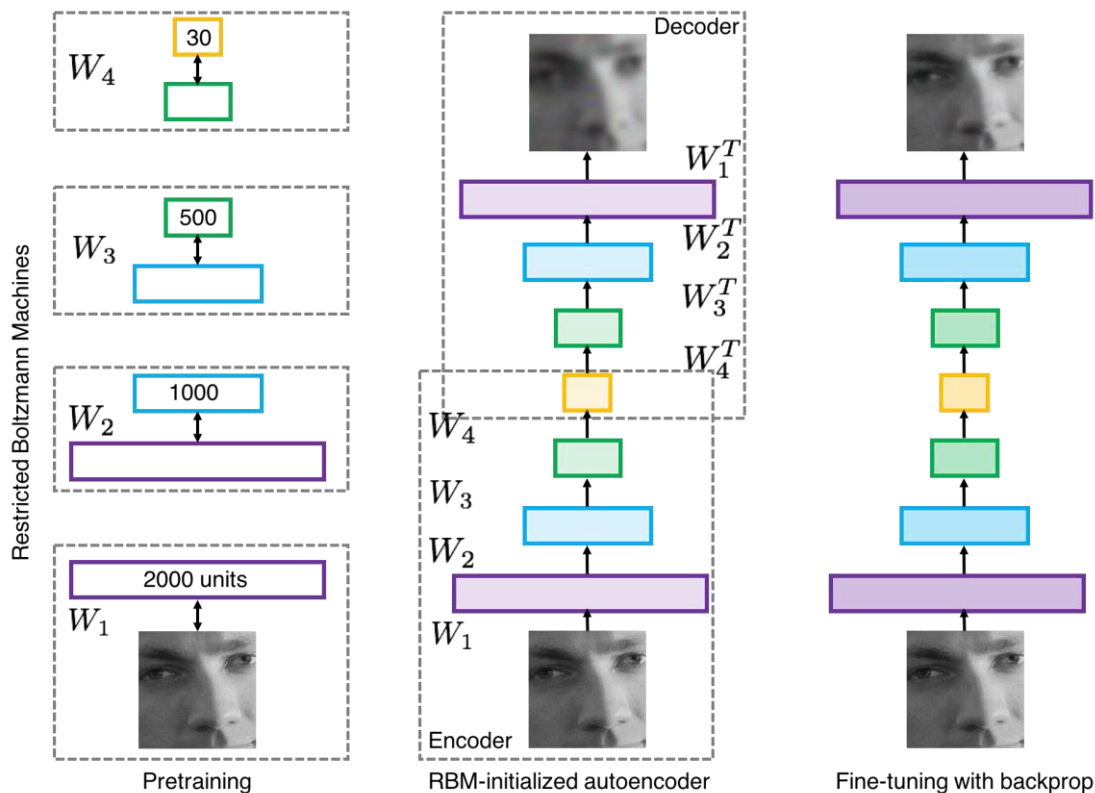[Hinton and Salakhutdinov 2006]

Reinvigorated research in Deep Learning



Illustration of Hinton and Salakhutdinov 2006 by Lane McIntosh, copyright CS231n 2017

# A bit of history

**First strong results**

***Acoustic Modeling using Deep Belief Networks***
*Abdel-rahman Mohamed, George Dahl, Geoffrey Hinton, 2010*

***Context-Dependent Pre-trained Deep Neural Networks for Large Vocabulary Speech Recognition***
George Dahl, Dong Yu, Li Deng, Alex Acero, 2012



Illustration of Dahl et al. 2012 by Lane McIntosh, copyright CS231n 2017

***Imagenet classification with deep convolutional neural networks***
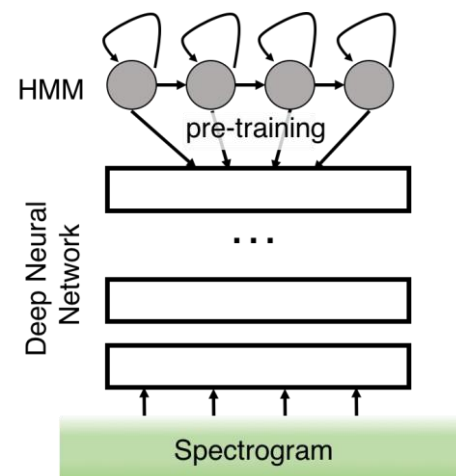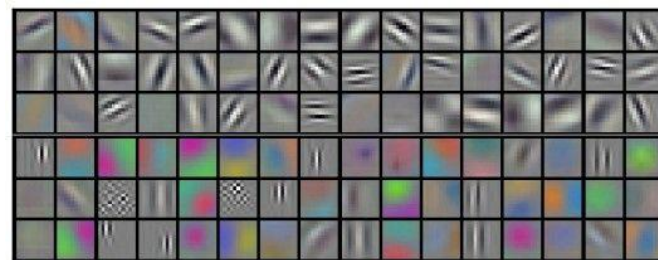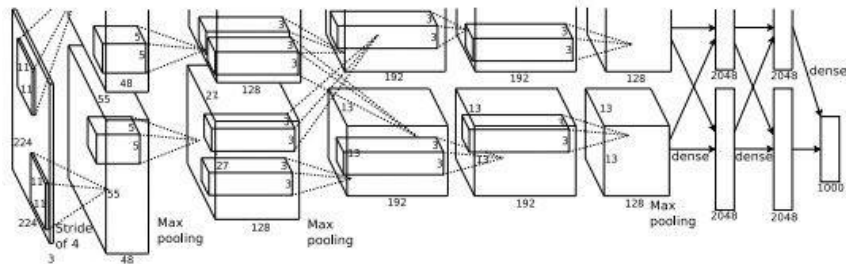Alex Krizhevsky, Ilya Sutskever, Geoffrey E Hinton, 2012



Figures copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

# Ups and downs of Neural Networks

- 1958: Perceptron (linear model)
- 1969: Perceptron has limitation
- 1980s: Multi-layer perceptron
- 1986: Backpropagation
- 1989: 1 hidden layer is "good enough", why deep?
- 2006: RBM initialization

# Ups and downs of Neural Networks

- 2009: GPU
- 2011: Start to be popular in speech recognition
- 2012: win ILSVRC image competition
- 2015.2: Image recognition surpassing human-level performance
- 2016.3: Alpha GO beats Lee Sedol
- 2016.10: Speech recognition system as good as humans
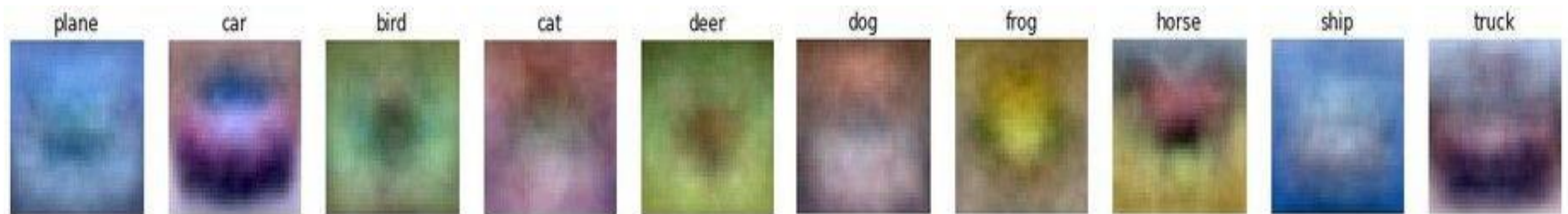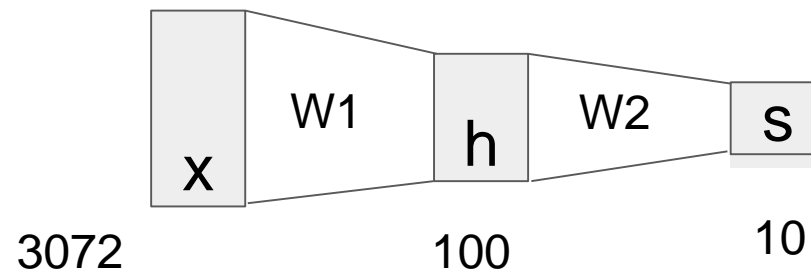- Now: Transformer, BERT,  Autopilot…
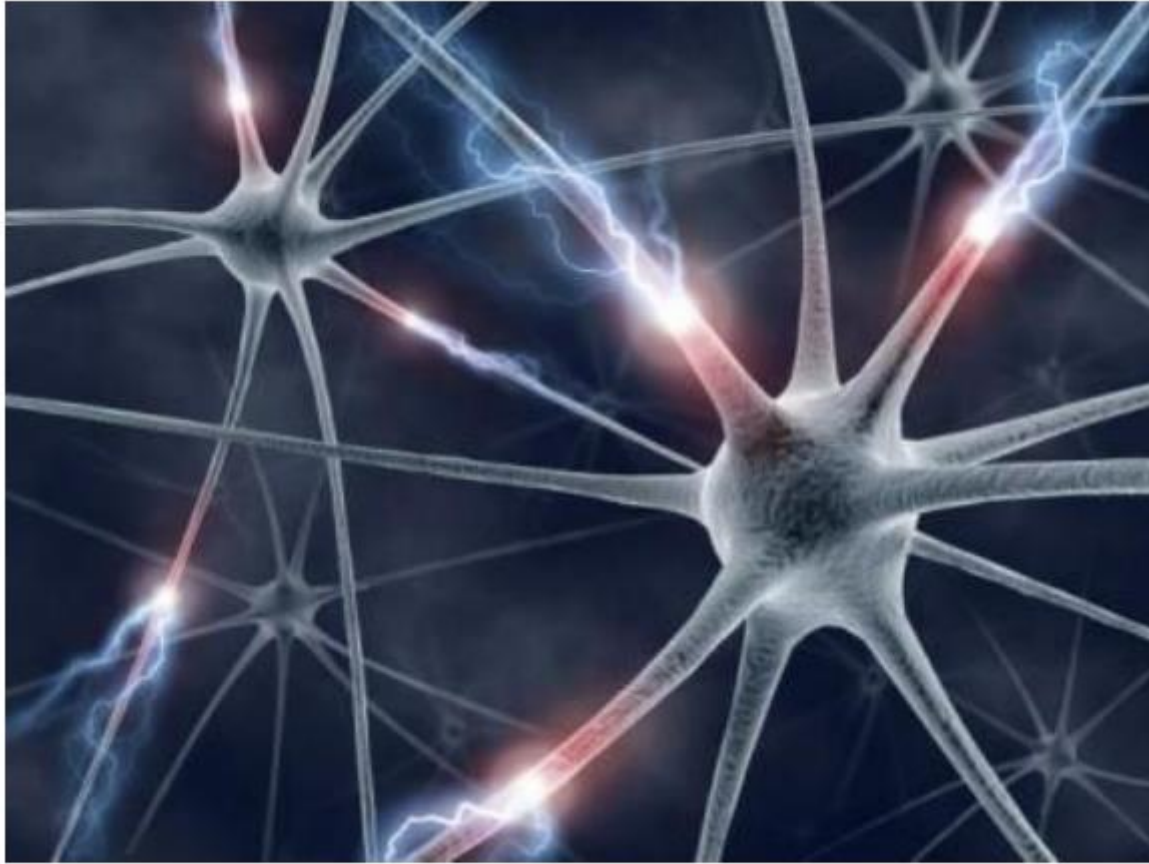
# Neural Networks

(**Before**) Linear score function

(**Now**) 2-layer Neural Network

$$f = Wx$$

$$f = W_2 \max(0, W_1 x)$$
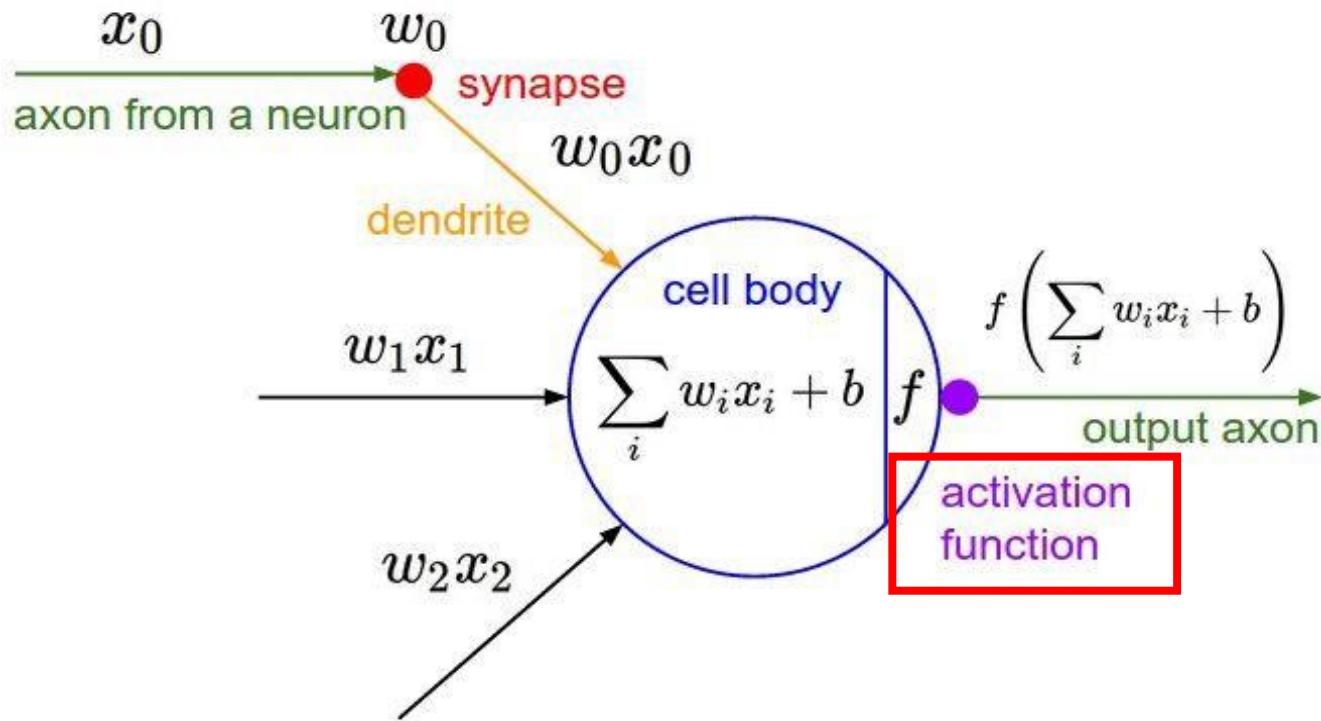
# Biological neuron



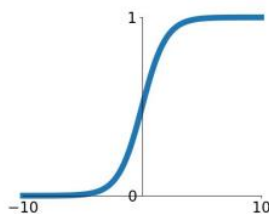This image by Fotis Bobolas is
licensed under CC-BY 2.0

# Artificial neuron

# Activation Function

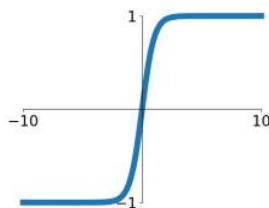**Sigmoid**

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



**tanh**

$\tanh(x)$
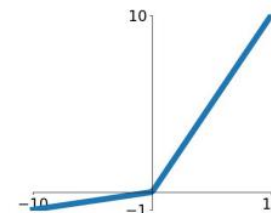


**ReLU**

$\max(0, x)$



**commonly used**

**Leaky ReLU**

$\max(0.1x, x)$



**Maxout**

$\max(w_1^T x + b_1, w_2^T x + b_2)$

**ELU**

$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$

# Biological v.s. Artificial

- From biological neuron to artificial neuron (perceptron)

# Biological v.s. Artificial

- From biological neuron network to artificial neuron networks

# Be very careful with your brain analogies!

- **Biological Neurons**
  - Many different types
  - Dendrites can perform complex non-linear computations
  - Synapses are not a single weight but a complex non-linear dynamical system
  - Rate code may not be adequate
- We can ignore whether the neural network actually simulates a biological neural network, and just think of a neural network as **a mathematical model with many parameters**

# Today's Topics

- Neural Network Introduction

- *Neural Network Structure*

- How Neural Network Works

- Backpropagation

# Perceptron

- It consists of two layers of neurons, the **input layer** and the **output layer**.



- Can realize logical AND, OR, NOT operation
- Only the neurons in the output layer perform activation function processing, and the learning ability is very limited
- Can't solve problems that are not linear separable, like XOR.

# Multi-layer Network



"2-layer Neural Net", or
"1-hidden-layer Neural Net"

"3-layer Neural Net", or
"2-hidden-layer Neural Net"

**"Fully-connected" layers**

Hidden layer and output layer neurons are
*functional neurons* with activation functions

# Multi-layer Network

**Weights**          **Weights**



input layer

hidden layer 1    hidden layer 2

output layer

The learning process of the neural network is to adjust the ***"connection weight"*** between neurons and the ***threshold*** of each functional neuron according to the training data

# Deep Neural Network

Deep = Many hidden layers

Now the commonly used **ResNet** has reached **152** layers

8 layers

19 layers

22 layers

7.3%

6.7%

16.4%

http://cs231n.stanford.edu/slides/winter1516_lecture8.pdf

AlexNet (2012)                VGG (2014)                GoogleNet (2014)

# Today's Topics

- Neural Network Introduction

- Neural Network Structure

- *How Neural Network Works*

- Backpropagation

# Matrix Operation



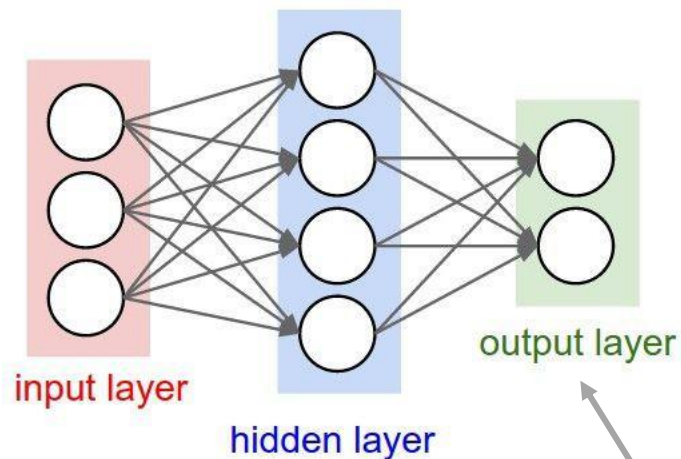$$\sigma\left(\begin{bmatrix} 1 & -2 \\ -1 & 1 \end{bmatrix}\begin{bmatrix} 1 \\ -1 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix}\right) = \begin{bmatrix} 0.98 \\ 0.12 \end{bmatrix}$$

$$\begin{bmatrix} 4 \\ -2 \end{bmatrix}$$

# Function Nesting

# Function Nesting



$$y = f(\ x\ )$$

Using parallel computing techniques to speed up matrix operation

$$= \sigma(\ W^L\ \cdots\ \sigma(\ W^2\ \sigma(\ W^1\ x\ +\ b^1\ )\ +\ b^2\ )\cdots +\ b^L\ )$$

# As a multi-class classifier



Feature extractor replacing feature engineering

get the classification probability

$x_1$

$x_2$

$x_K$

Softmax

$x$

$y_1$

$y_2$

$y_M$

**Input Layer**

**Hidden Layers**

**Output Layer**

= Multi-class Classifier

# Case study
# Handwriting Digit Recognition

# Case study
# Handwriting Digit Recognition

**Input**

**Output**



16 x 16 = 256

Ink → 1

No ink → 0

$x_1$

$x_2$

$x_{256}$

0.1    is 1

0.7    is 2

0.2    is 0

The image is "2"

Each dimension represents the confidence of a digit.

# Case study
# Handwriting Digit Recognition

$x_1$

$x_2$

$\vdots$

$x_{256}$

**Neural Network**

What is needed is a function ……

$y_1$ — is 1

$y_2$ — is 2

$\vdots$

$y_{10}$ — is 0

Input: 256-dim vector

output: 10-dim vector

# Case study
# Handwriting Digit Recognition



Input  Layer 1  Layer 2  Layer L  Output

$x_1$

$x_2$

$x_N$

**Input Layer**

**Hidden Layers**

**Output Layer**

A function set containing the candidates for
Handwriting Digit Recognition

$y_1$  is 1

$y_2$  is 2

$y_{10}$  is 0

**You need to decide the network structure to let a good function in your function set.**

# Case study
# Handwriting Digit Recognition

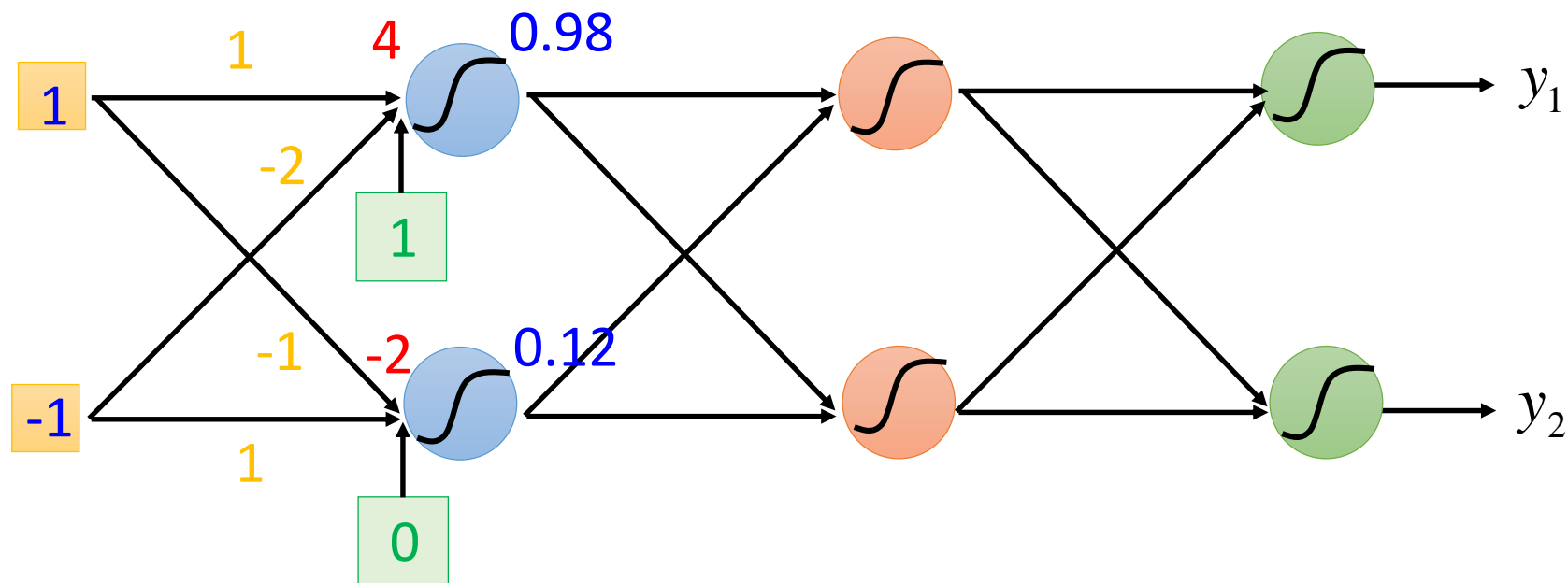- How many layers? How many neurons for each layer?

| Trial and Error | + | Intuition |
|---|---|---|

# Today's Topics

- Neural Network Introduction

- Neural Network Structure
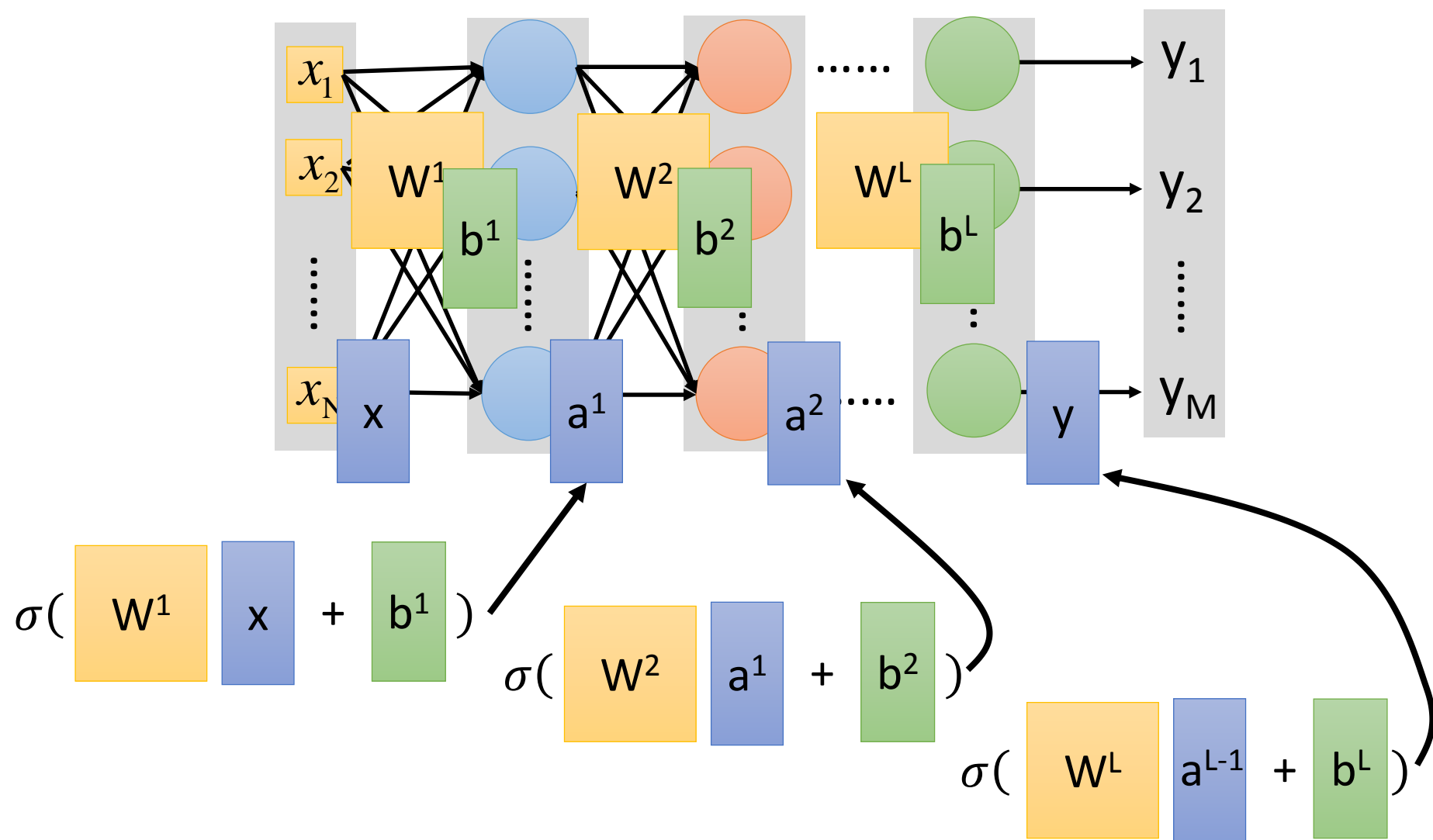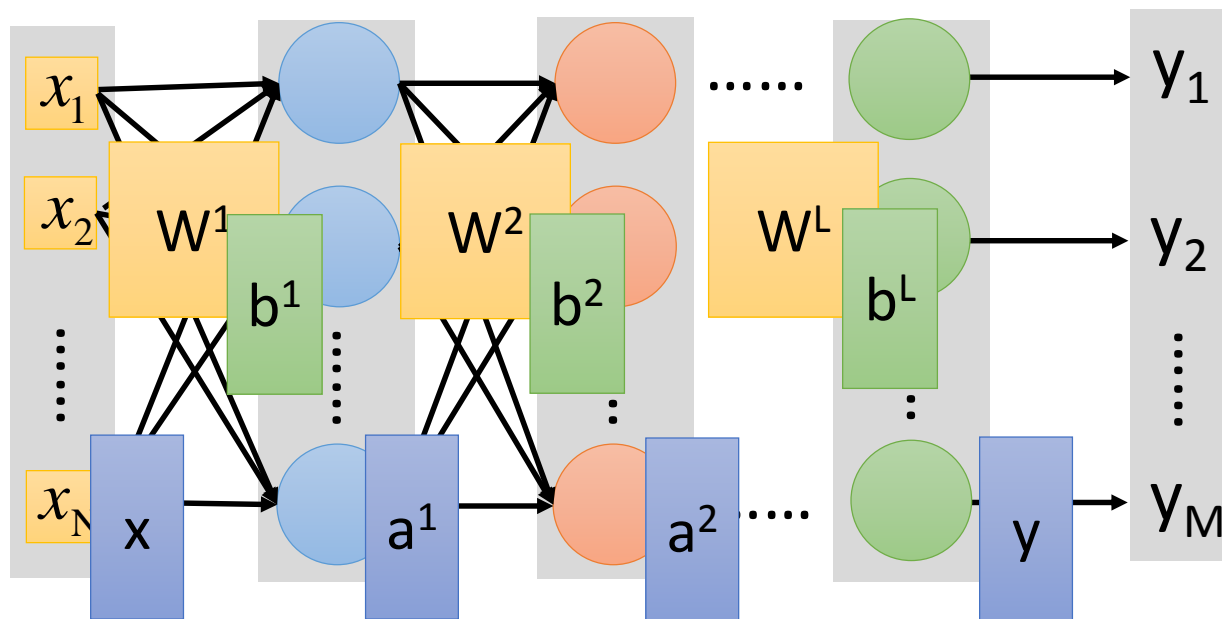
- How Neural Network Works

- *Backpropagation*

# Neural Network Optimization

- We have already learned to optimize the learner using the gradient descent method

- Can neural networks also be optimized using gradient descent?

# Case: CNN (AlexNet)

input image

weights

loss



Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

# Case: Neural Turing Machine

input image

loss



Figure reproduced with permission from a Twitter post by Andrej Karpathy.

# Why we need BP

- If we use gradient descent directly

  Network parameters $\theta = \{w_1, w_2, \cdots, b_1, b_2, \cdots\}$

Starting Parameters

$$\theta^0 \longrightarrow \theta^1 \longrightarrow \theta^2 \longrightarrow \ \cdots\cdots$$

$$\nabla L(\theta) = \begin{bmatrix} \partial L(\theta)/\partial w_1 \\ \partial L(\theta)/\partial w_2 \\ \vdots \\ \partial L(\theta)/\partial b_1 \\ \partial L(\theta)/\partial b_2 \\ \vdots \end{bmatrix}$$

$Compute \ \nabla L(\theta^0) \qquad \theta^1 = \theta^0 - \eta \nabla L(\theta^0)$

$Compute \ \nabla L(\theta^1) \qquad \theta^2 = \theta^1 - \eta \nabla L(\theta^1)$

Millions of parameters ……

To compute the gradients efficiently, we use ***backpropagation***.

# BP: A simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

# BP: A simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$

# BP: A simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$

# BP: A simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial f}$$

# BP: A simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



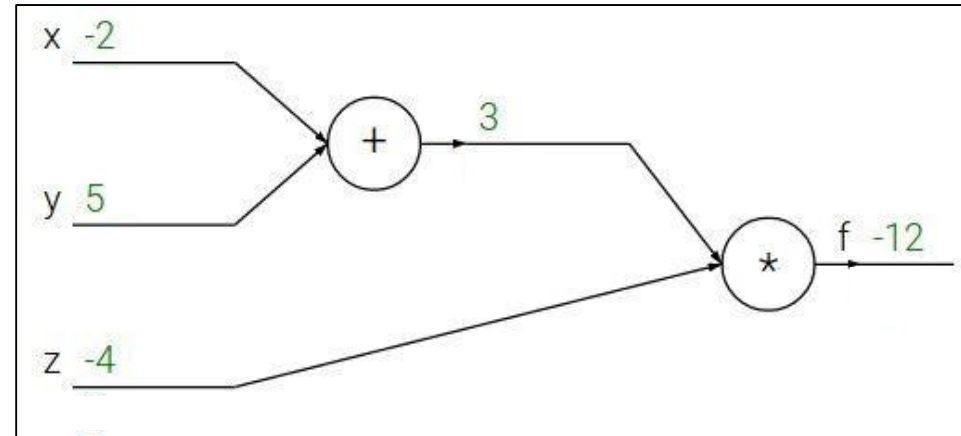$$\frac{\partial f}{\partial z}$$

# BP: A simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$



x  -2

y  5

q  3

z  -4

+

*

f  -12

1

3

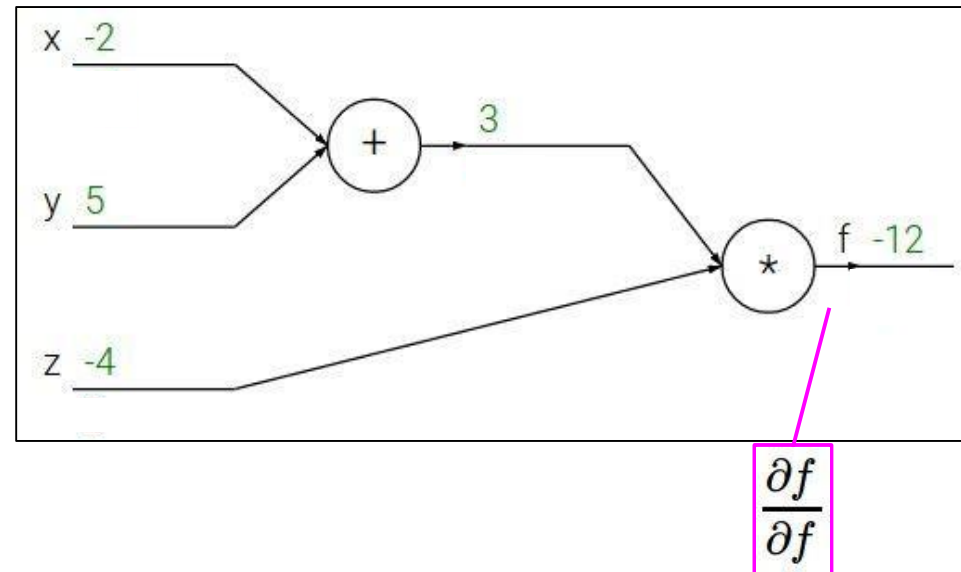$\dfrac{\partial f}{\partial q}$

# BP: A simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
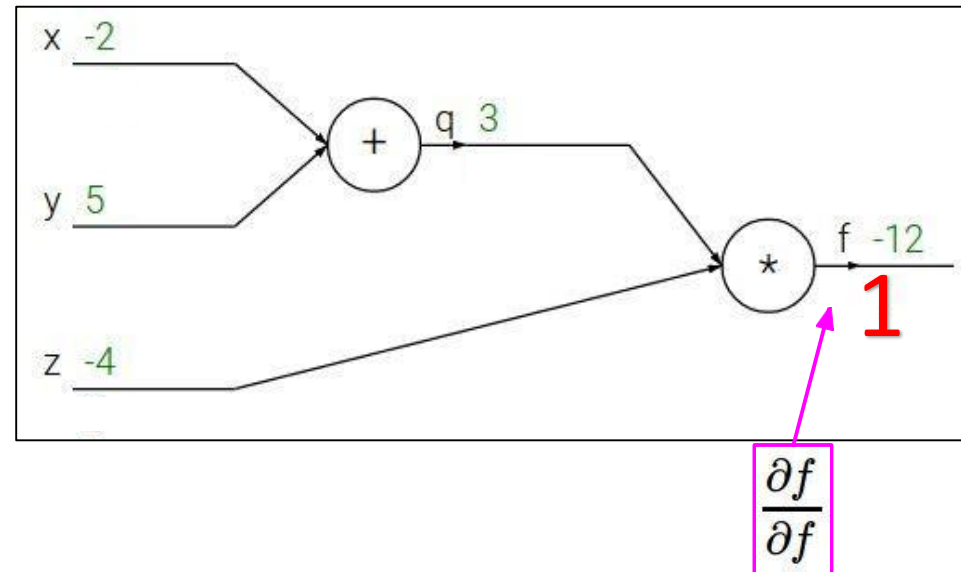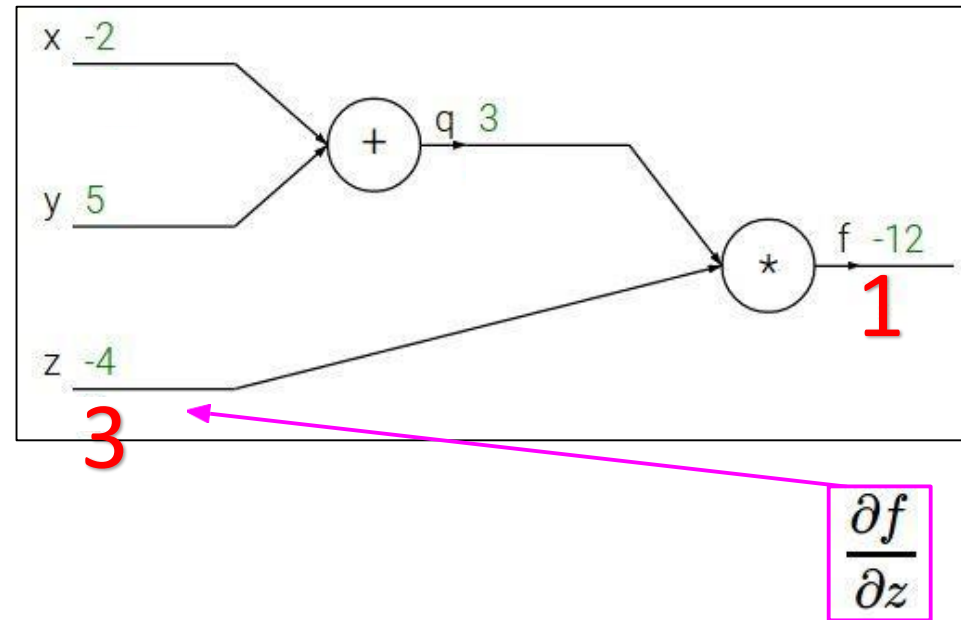
# BP: A simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$



x   -2

y   5

z   -4

q   3

-4

f   -12

1

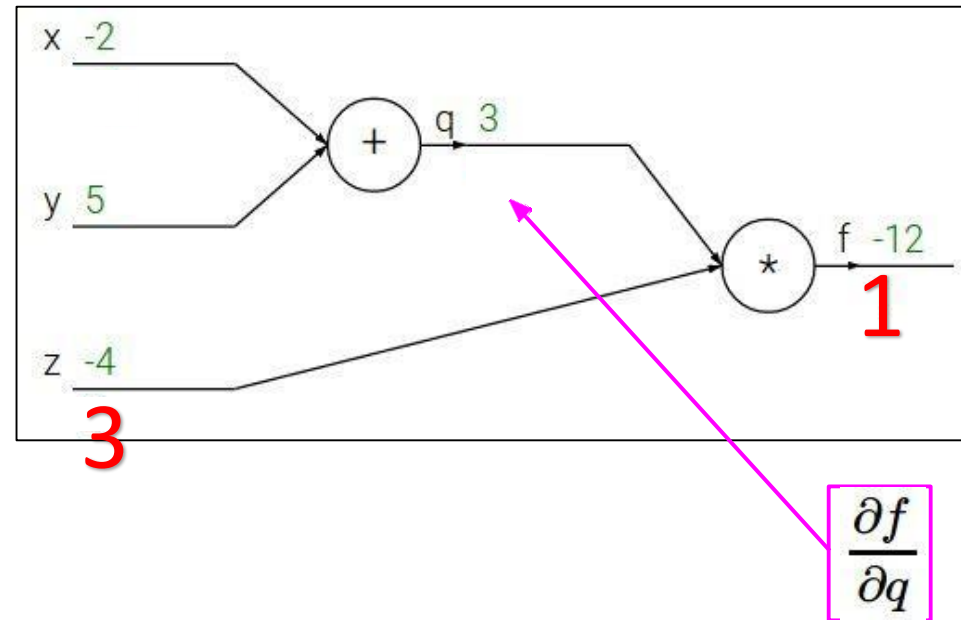3

$\dfrac{\partial f}{\partial y}$

# BP: A simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$

$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$



Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

Chain rule:

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$
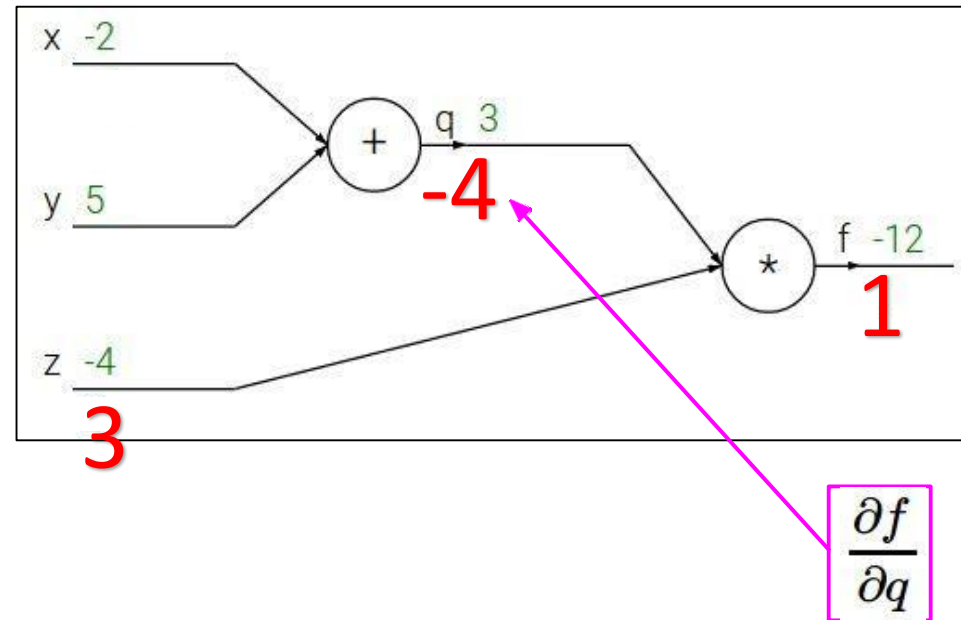
# BP: A simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
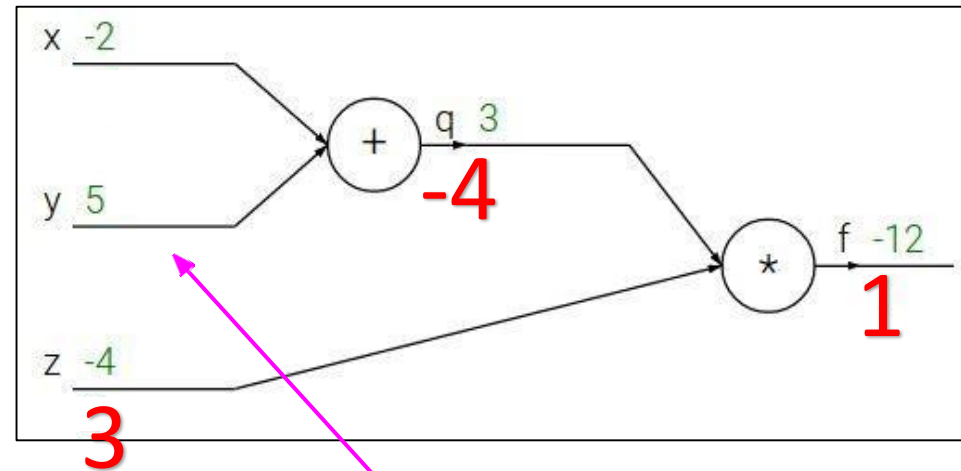
# BP: A simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\quad \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$
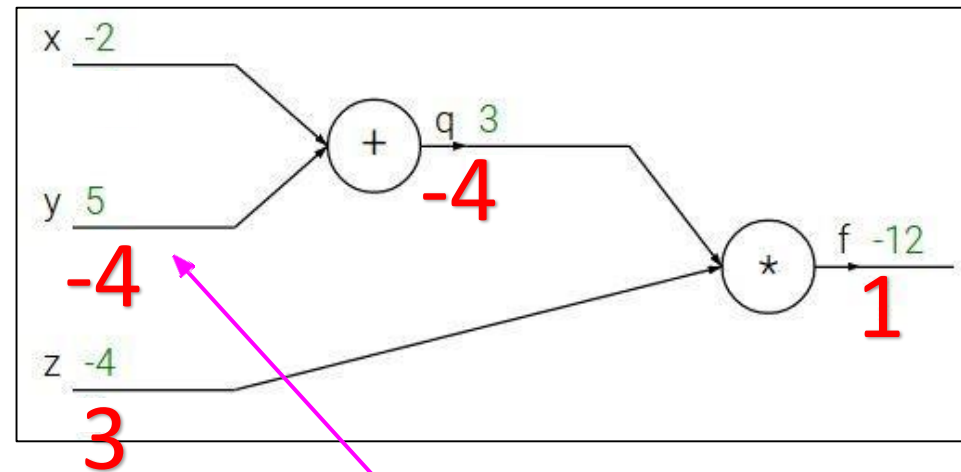
Chain rule:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

# Backpropagation



$x$

$y$

f

$z$

# Backpropagation

# Backpropagation



$x$

"local gradient"

$\dfrac{\partial z}{\partial x}$

$\dfrac{\partial z}{\partial y}$

$y$

f

$z$

$\dfrac{\partial L}{\partial z}$

gradients

# Backpropagation



$x$

$\dfrac{\partial L}{\partial x} = \dfrac{\partial L}{\partial z}\dfrac{\partial z}{\partial x}$

"local gradient"

$\dfrac{\partial z}{\partial x}$

$\dfrac{\partial z}{\partial y}$

f

$z$

$\dfrac{\partial L}{\partial z}$

$y$

gradients

# Backpropagation



$x$

$$\boxed{\frac{\partial L}{\partial x}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x}$$

"local gradient"

$\boxed{\dfrac{\partial z}{\partial x}}$

f

$z$

$\boxed{\dfrac{\partial z}{\partial y}}$

$\boxed{\dfrac{\partial L}{\partial z}}$

$y$

$$\boxed{\frac{\partial L}{\partial y}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial y}$$

gradients

# Backpropagation



"local gradient"

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial x}$$

$$\frac{\partial z}{\partial x}$$

$$\frac{\partial z}{\partial y}$$
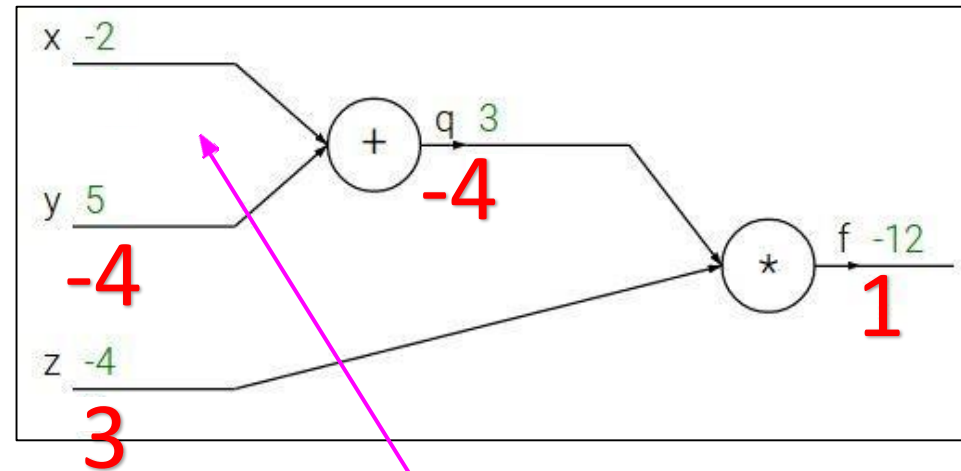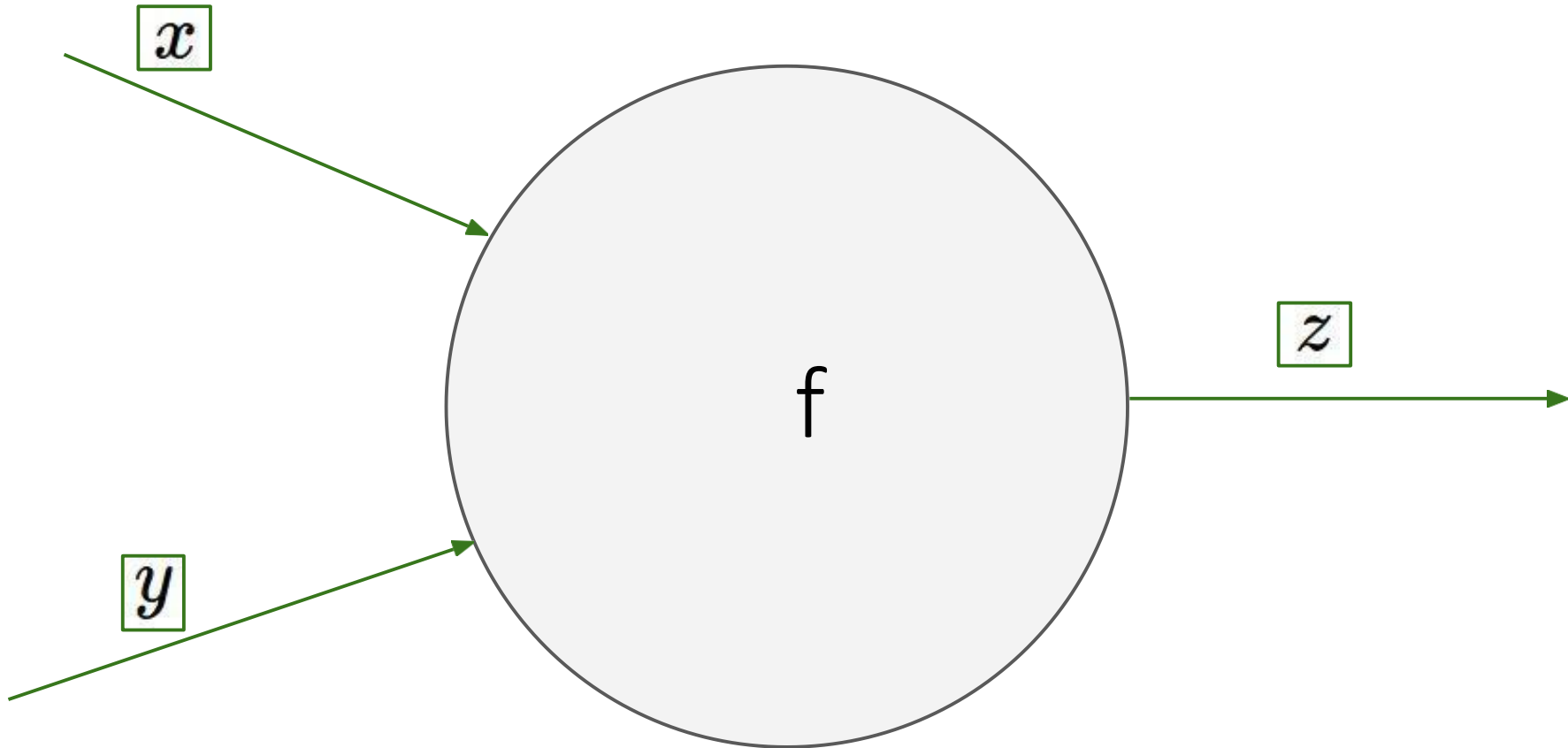
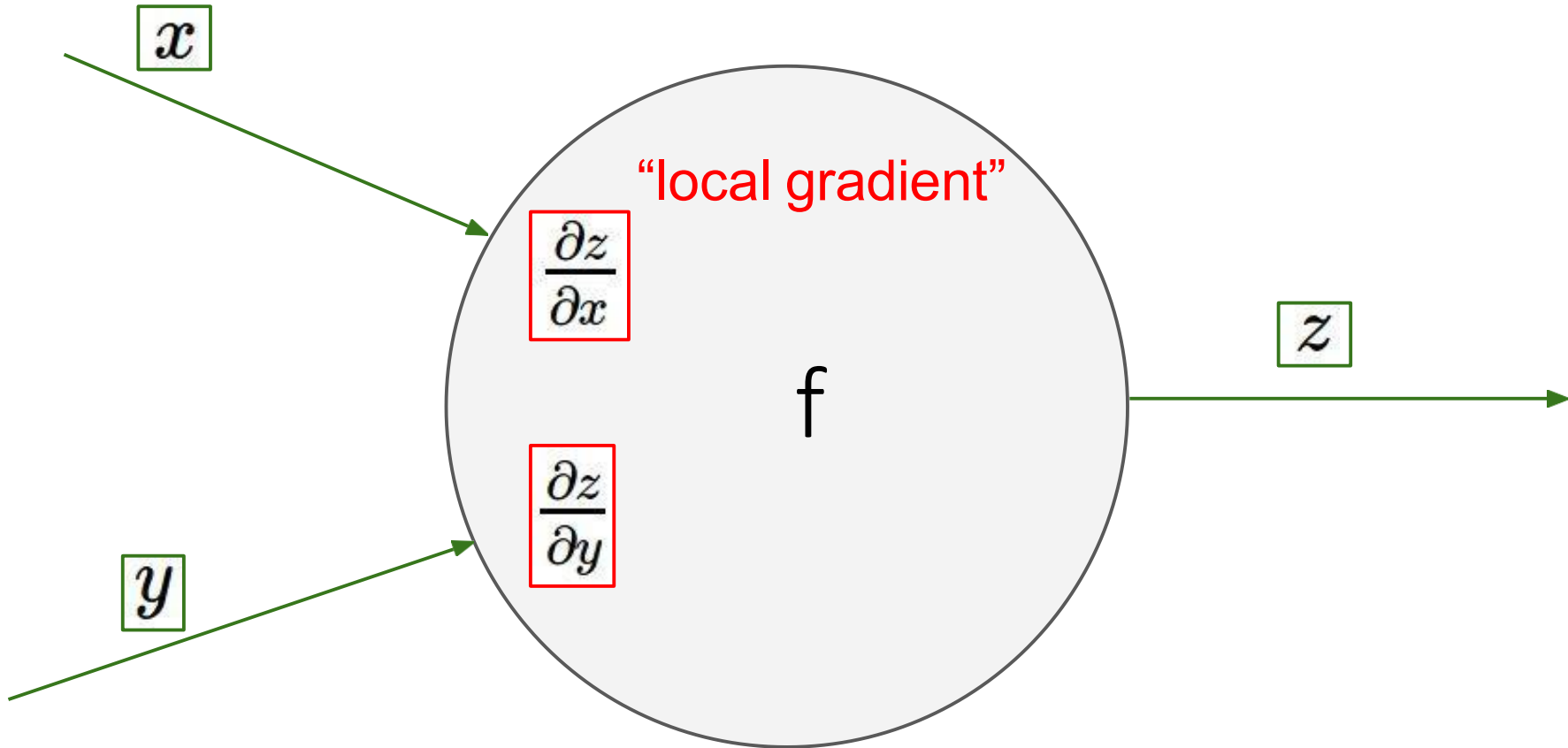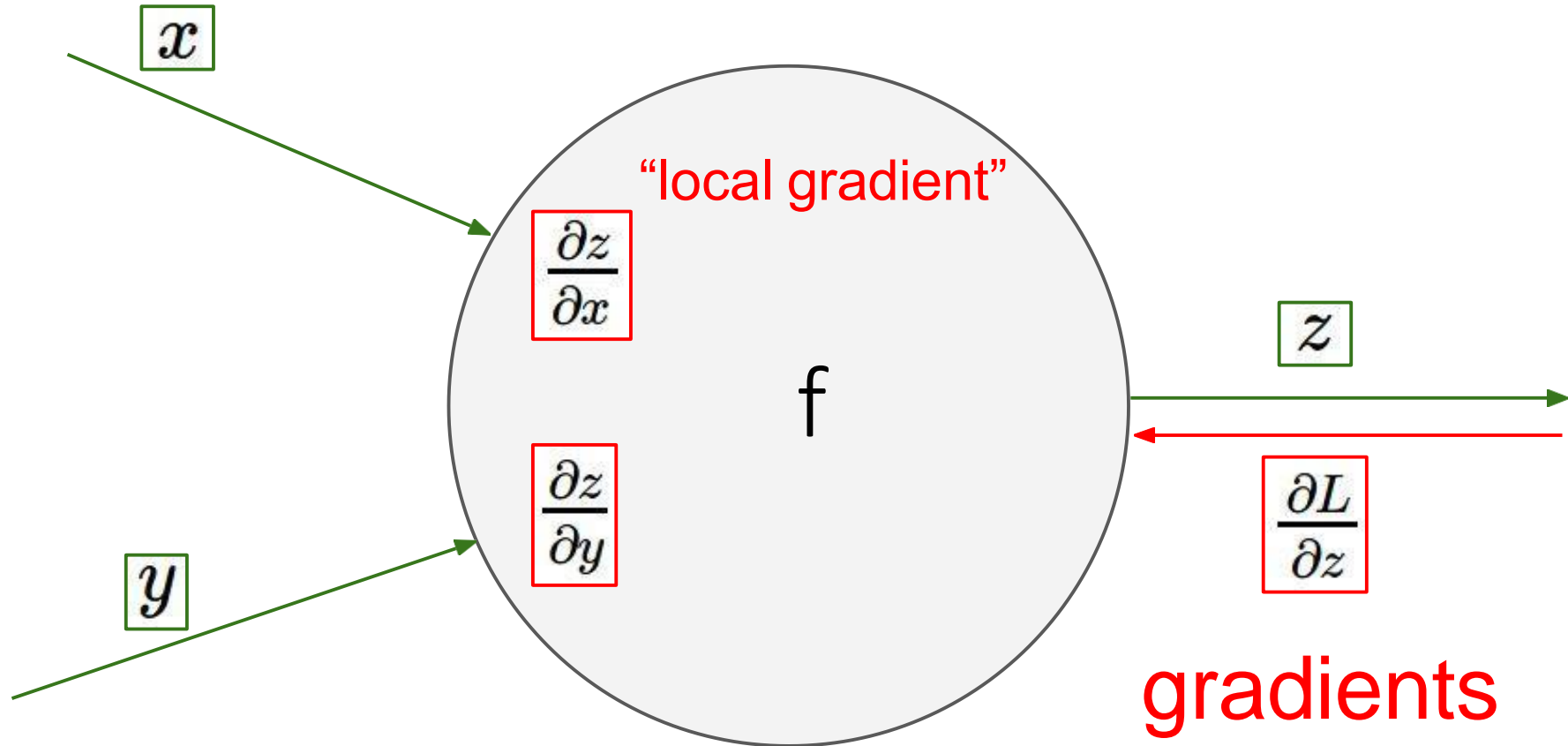$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial y}$$

$x$

$y$

$z$

$$\frac{\partial L}{\partial z}$$

gradients

f

# BP: Another example

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



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \qquad \Bigg| \qquad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \Bigg| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# BP: Another example

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



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

Dr. Shuang LIANG, Tongji

# BP: Another example

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



$$\left(\frac{-1}{1.37^2}\right)(1.00) = -0.53$$

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

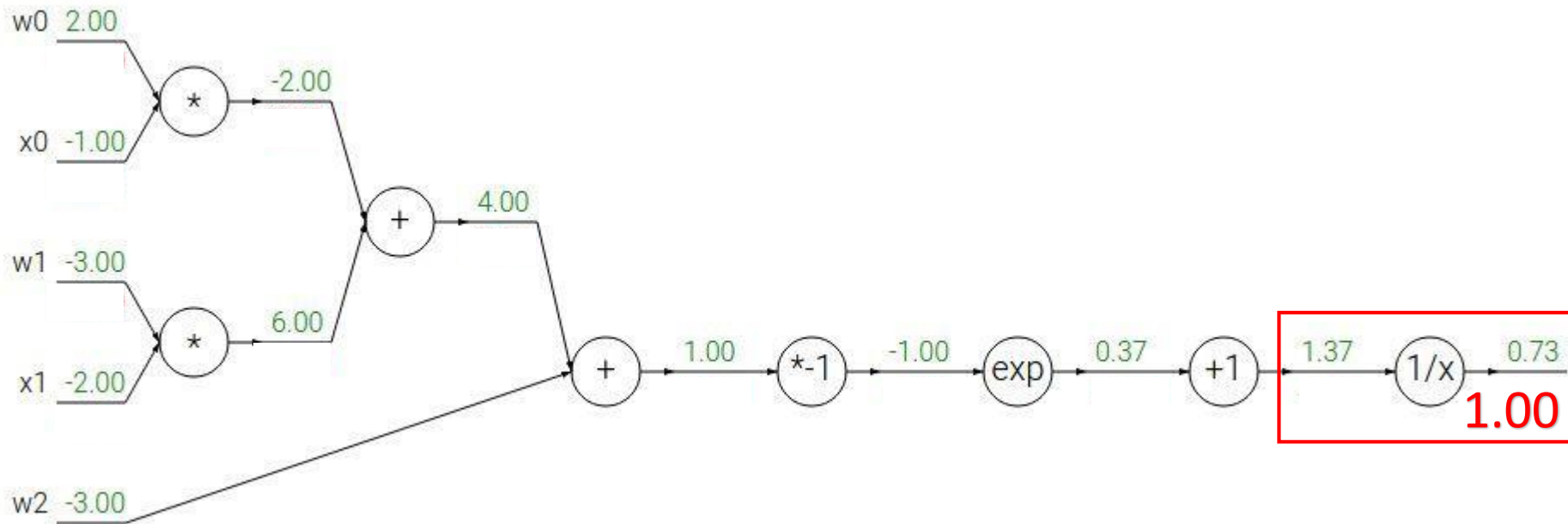$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# BP: Another example

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



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# BP: Another example

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



$$(1)(-0.53) = -0.53$$

$$f(x) = e^x \quad \rightarrow \quad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \quad \rightarrow \quad \frac{df}{dx} = a$$

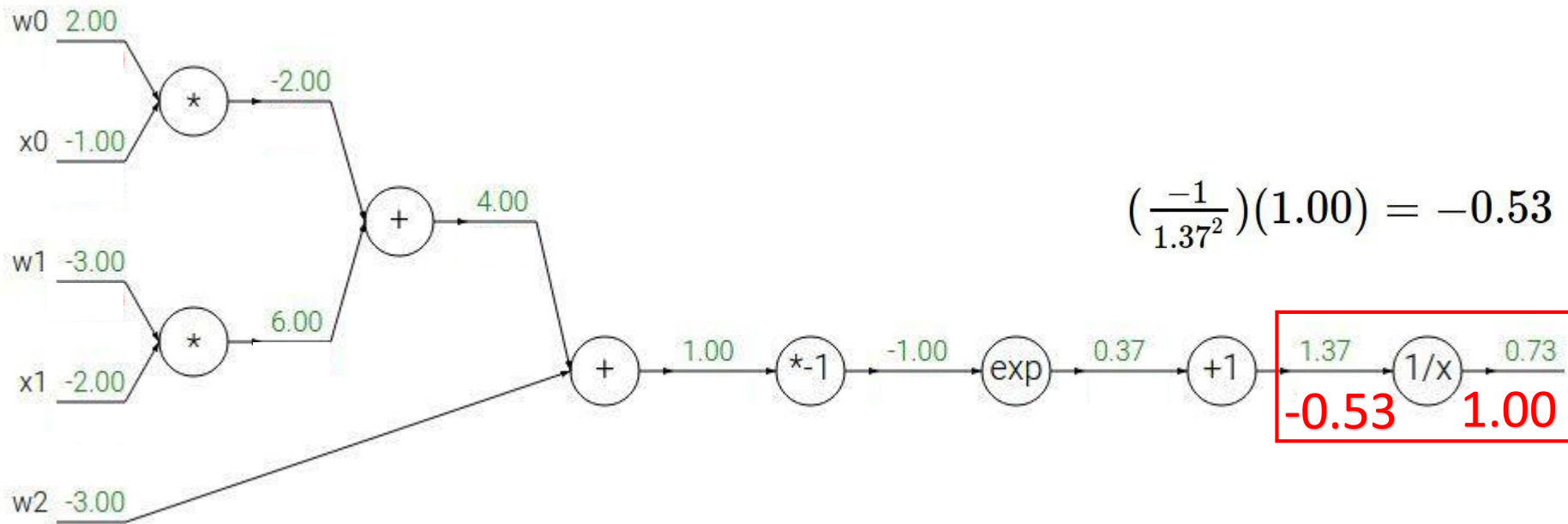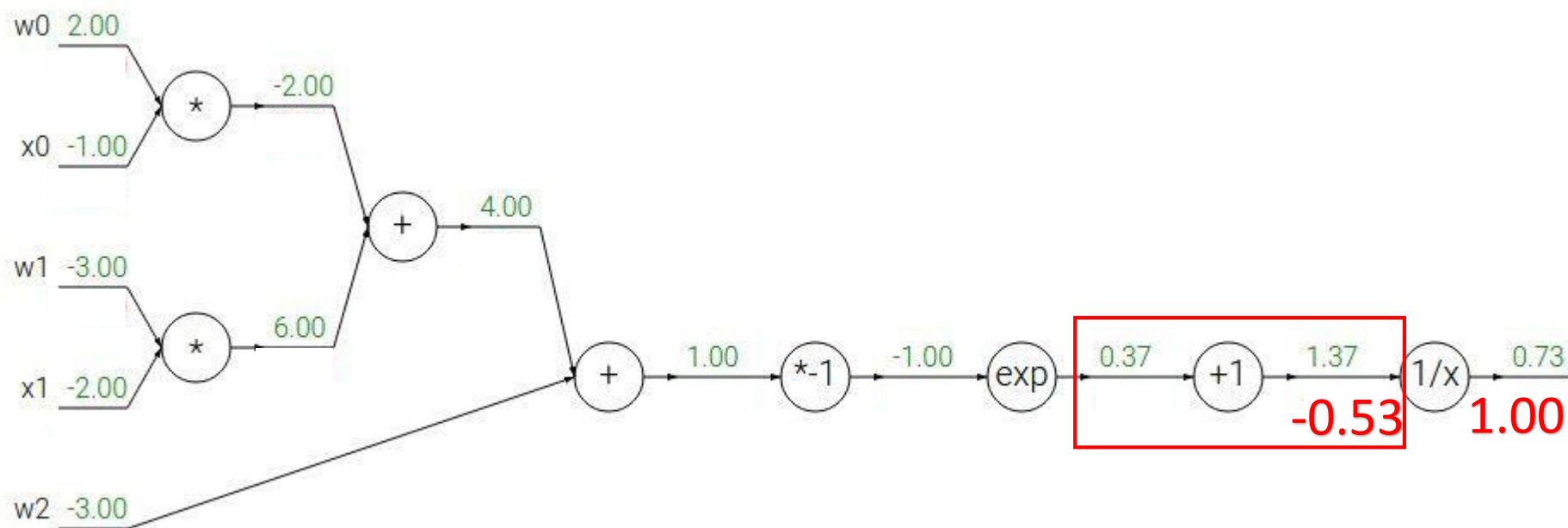$$f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \quad \rightarrow \quad \frac{df}{dx} = 1$$

# BP: Another example

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



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$
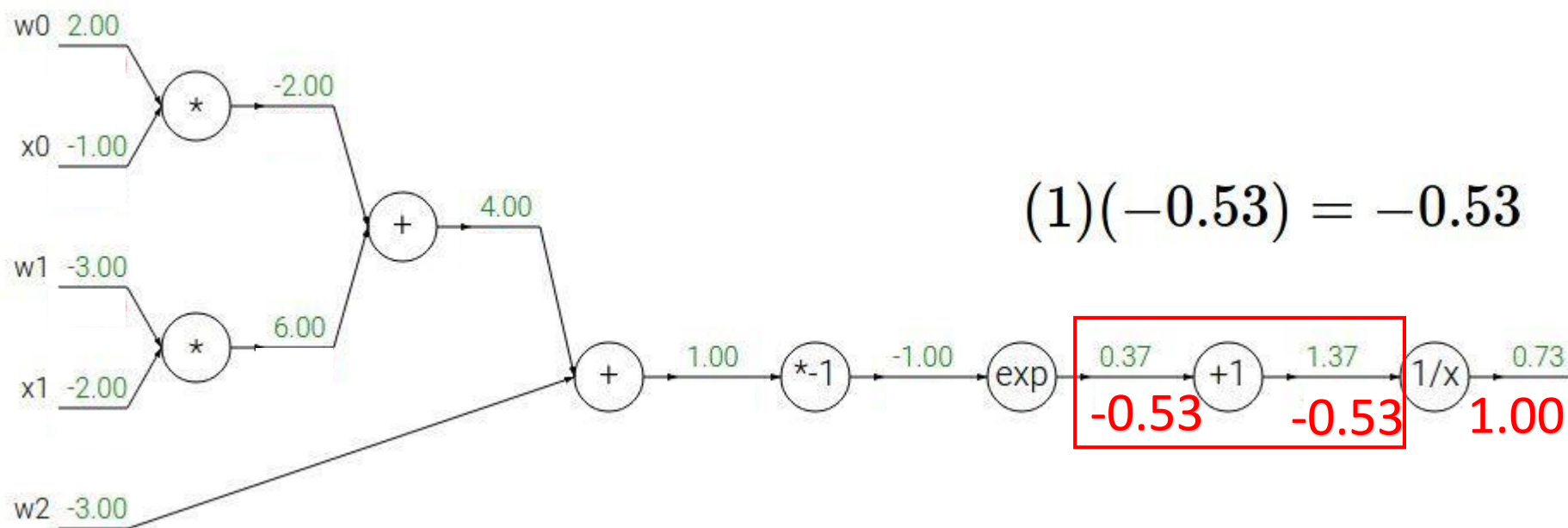
# BP: Another example

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



$$(e^{-1})(-0.53) = -0.20$$

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$
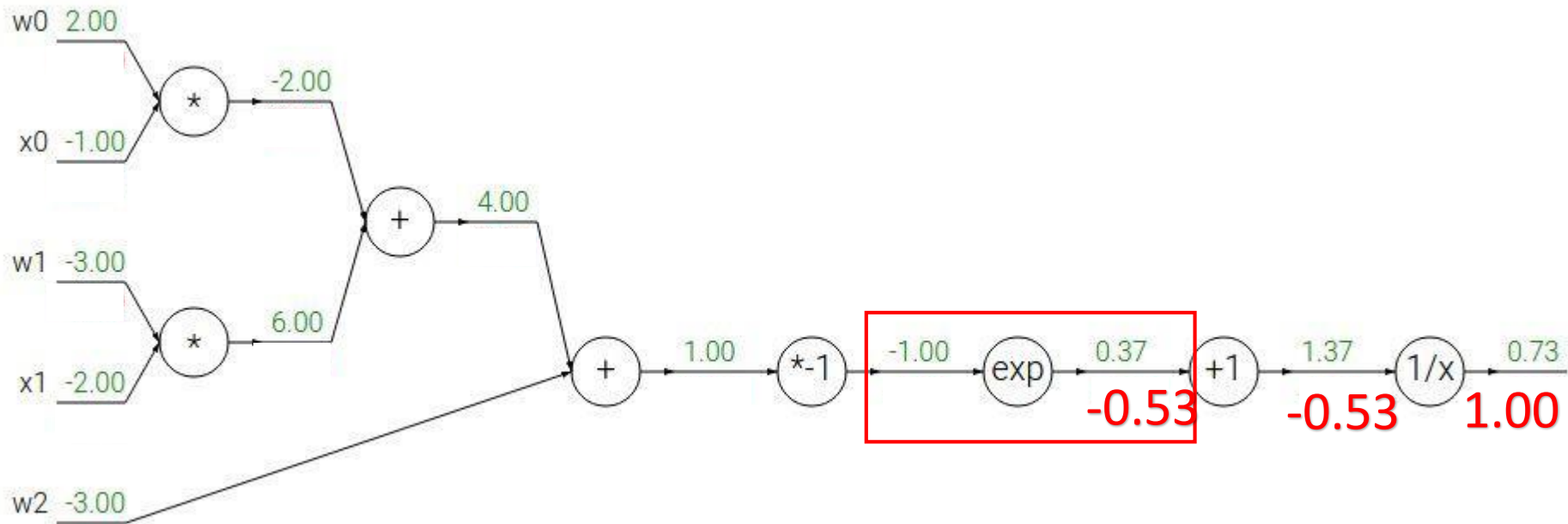
# BP: Another example

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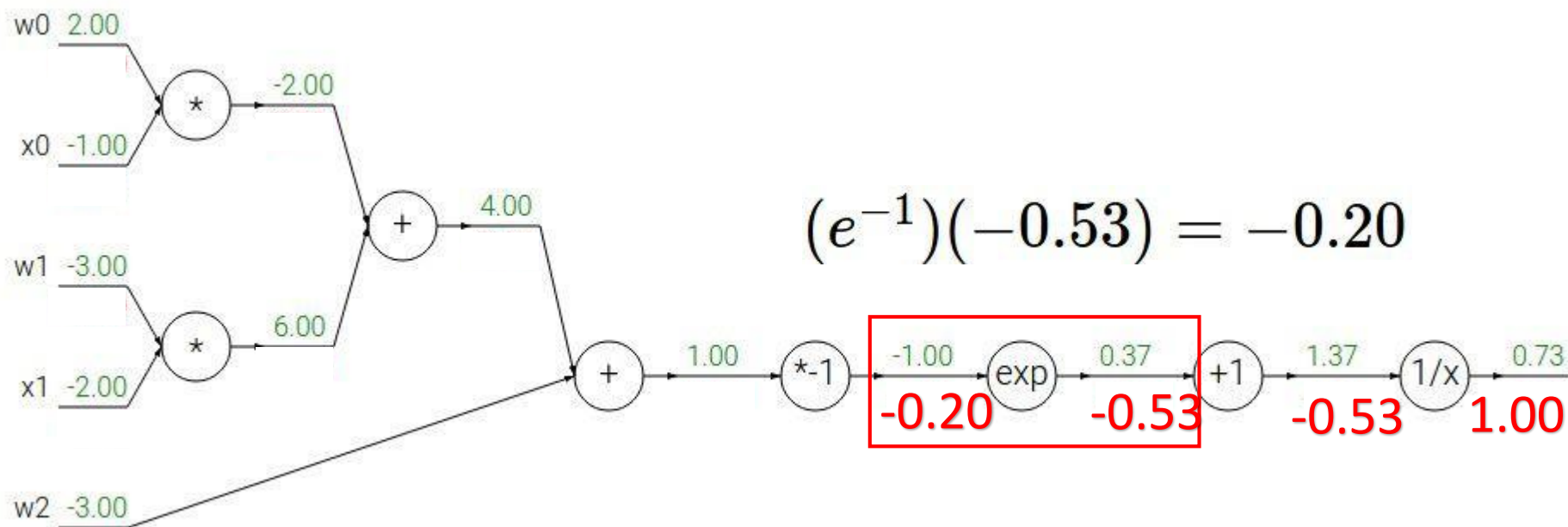


$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

Dr. Shuang LIANG, Tongji

# BP: Another example

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



(-1) * (-0.20) = 0.20

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \qquad \bigg| \qquad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

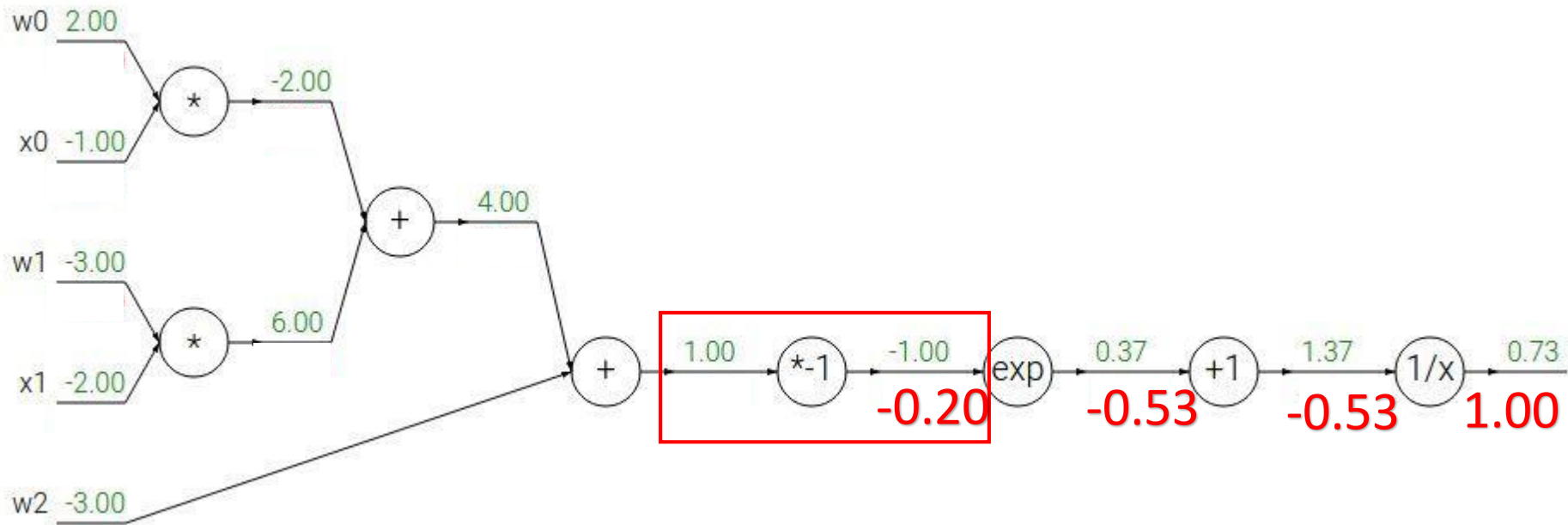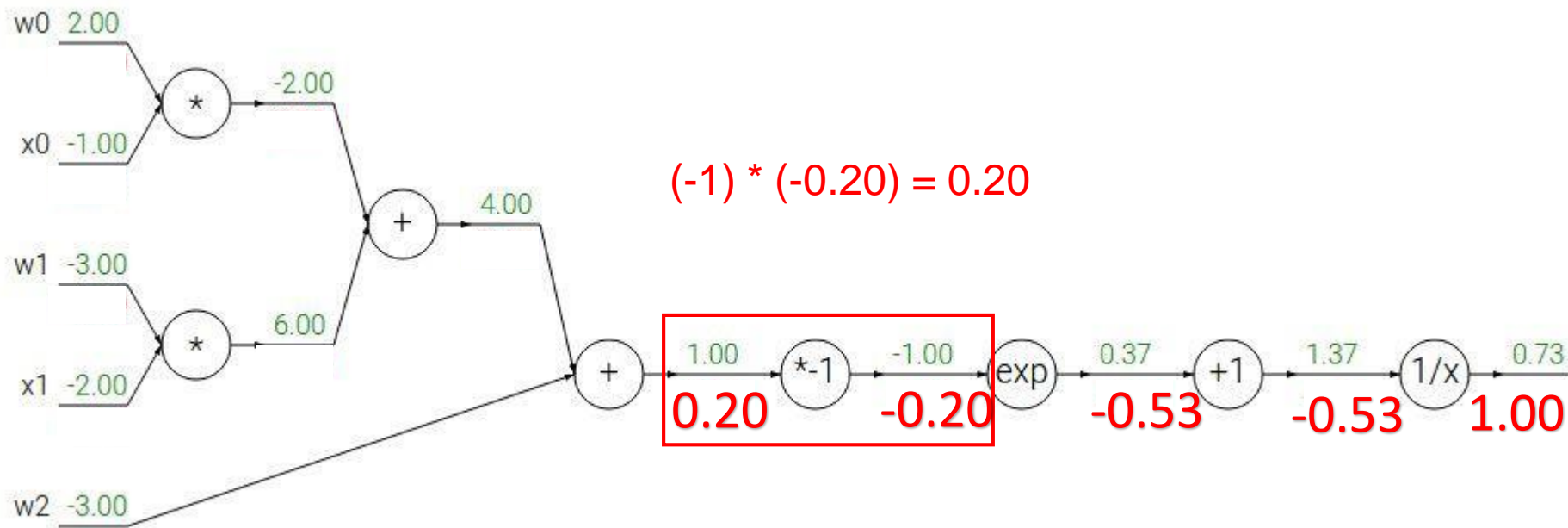$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \bigg| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# BP: Another example

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



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# BP: Another example

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

[local gradient] x [upstream gradient]

[1] x [0.2] = 0.2

[1] x [0.2] = 0.2  (both inputs!)



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$
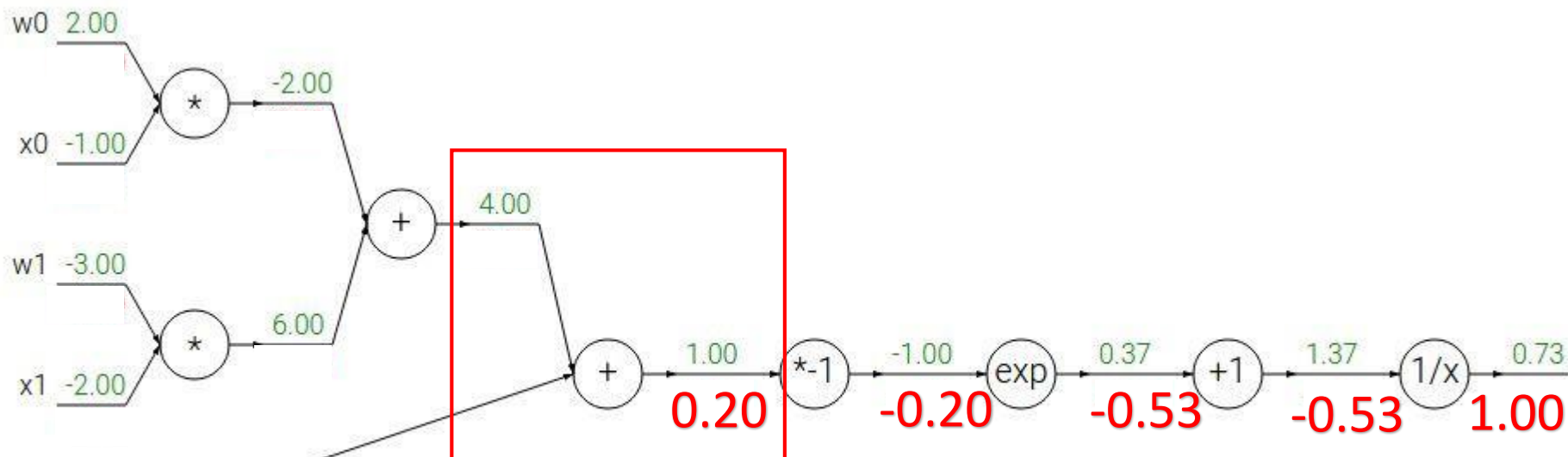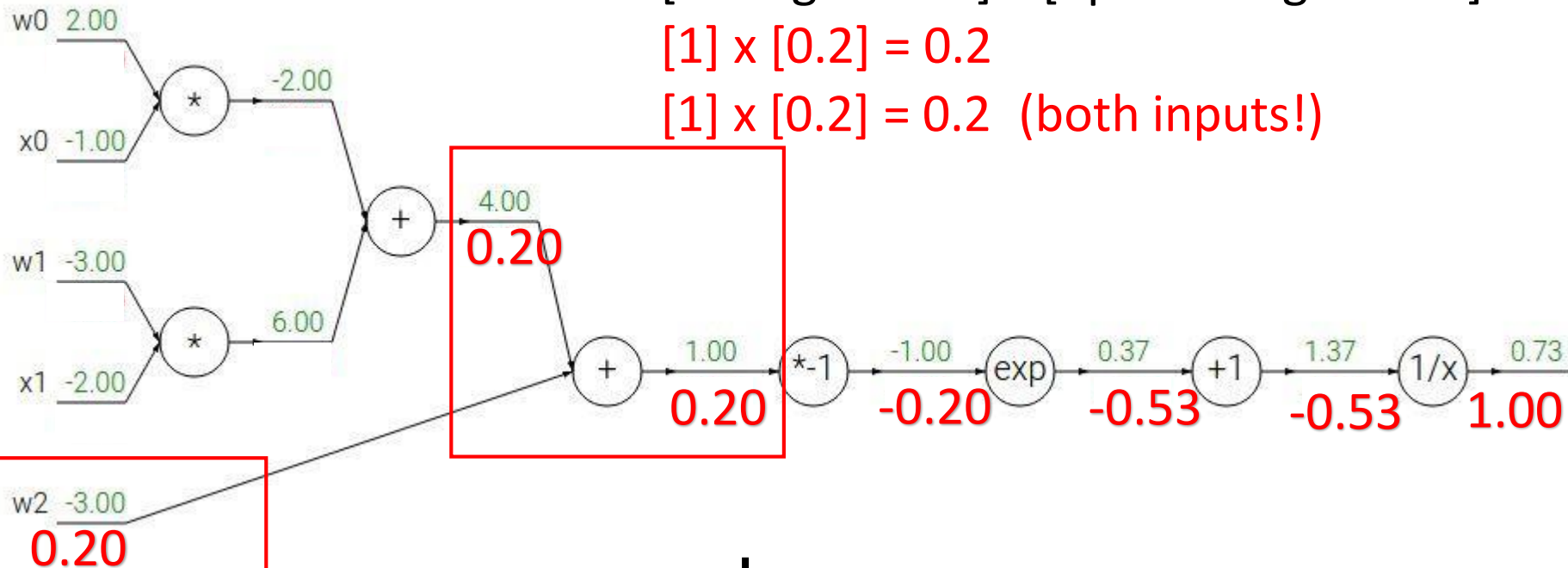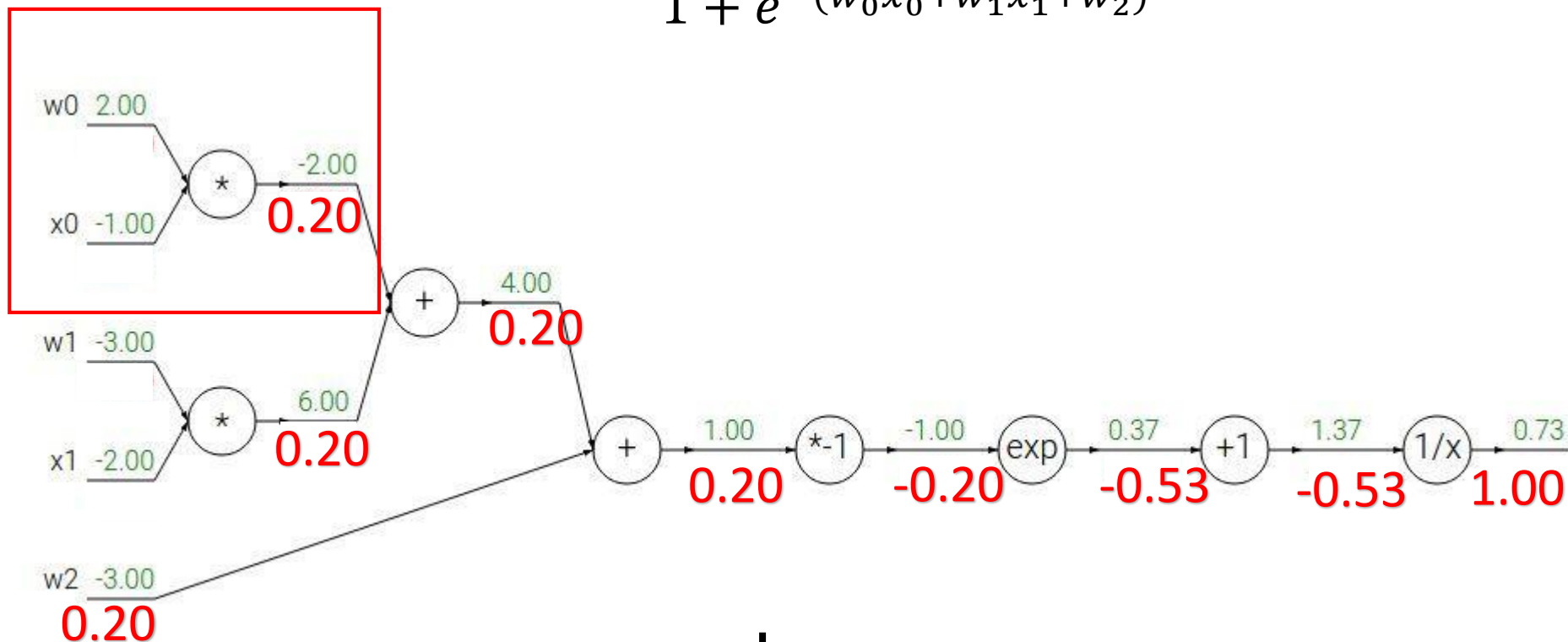
# BP: Another example

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



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \qquad \bigg| \qquad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

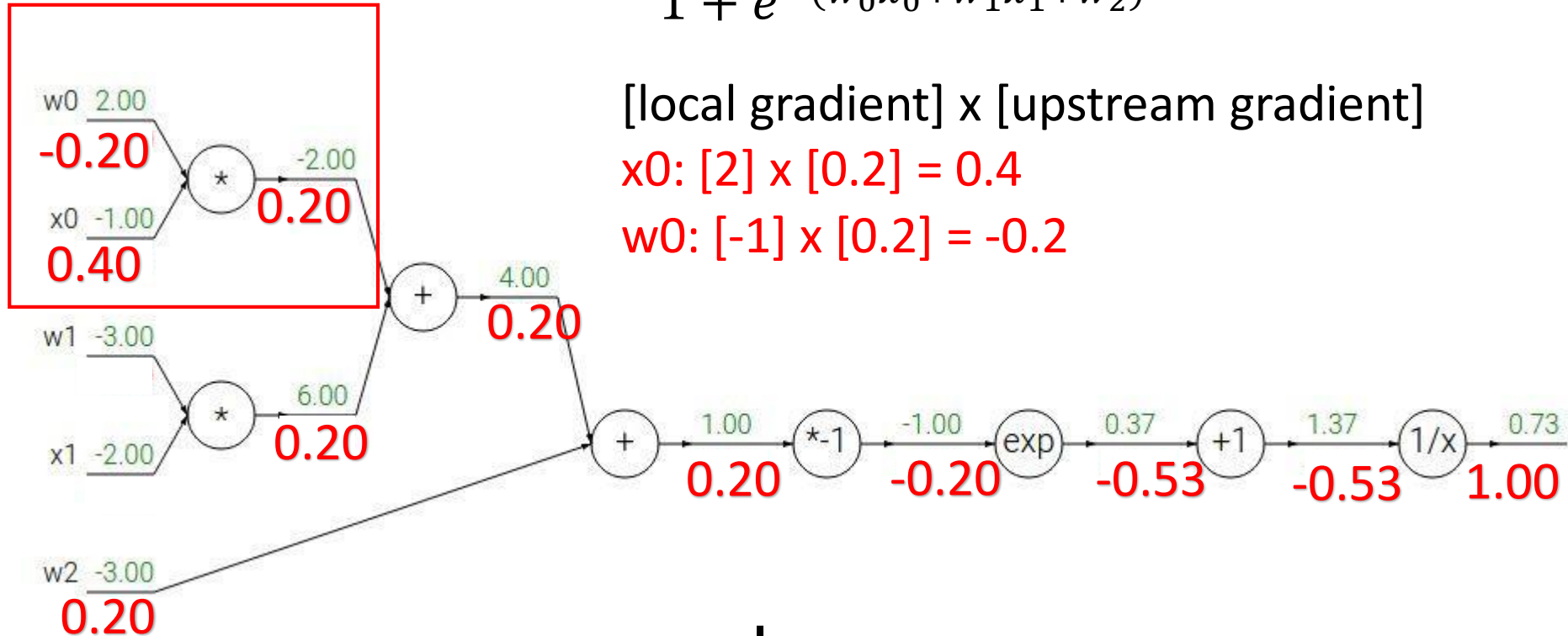$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \bigg| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# BP: Another example

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

[local gradient] x [upstream gradient]

x0: [2] x [0.2] = 0.4

w0: [-1] x [0.2] = -0.2

w0 2.00
-0.20

x0 -1.00
0.40

$*$   -2.00
0.20

w1 -3.00

x1 -2.00

$*$   6.00
0.20

$+$   4.00
0.20

$+$   1.00
0.20

$*-1$   -1.00
-0.20

exp   0.37
-0.53

$+1$   1.37
-0.53

$1/x$   0.73
1.00

w2 -3.00
0.20

$f(x) = e^x$  $\rightarrow$  $\dfrac{df}{dx} = e^x$  |  $f(x) = \dfrac{1}{x}$  $\rightarrow$  $\dfrac{df}{dx} = -1/x^2$

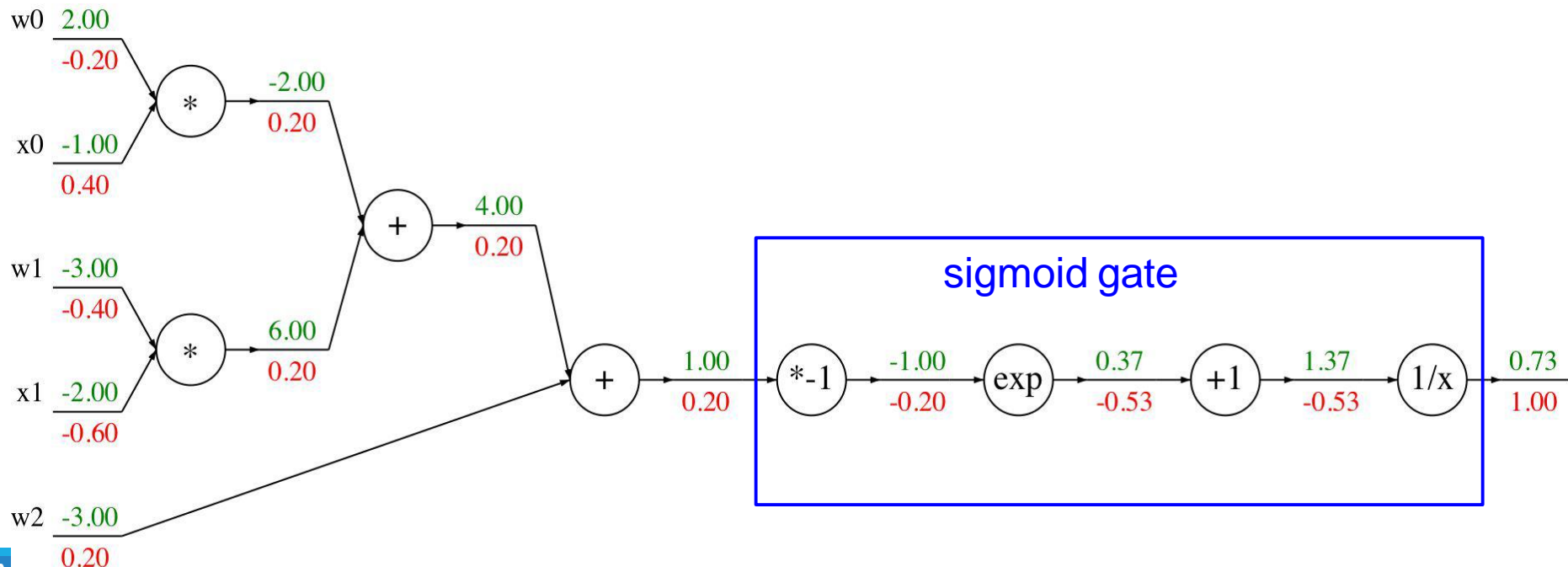$f_a(x) = ax$  $\rightarrow$  $\dfrac{df}{dx} = a$  |  $f_c(x) = c + x$  $\rightarrow$  $\dfrac{df}{dx} = 1$

# BP: Another example

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

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

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left( \frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \left( \frac{1}{1 + e^{-x}} \right) = (1 - \sigma(x)) \sigma(x)$$

# BP: Another example

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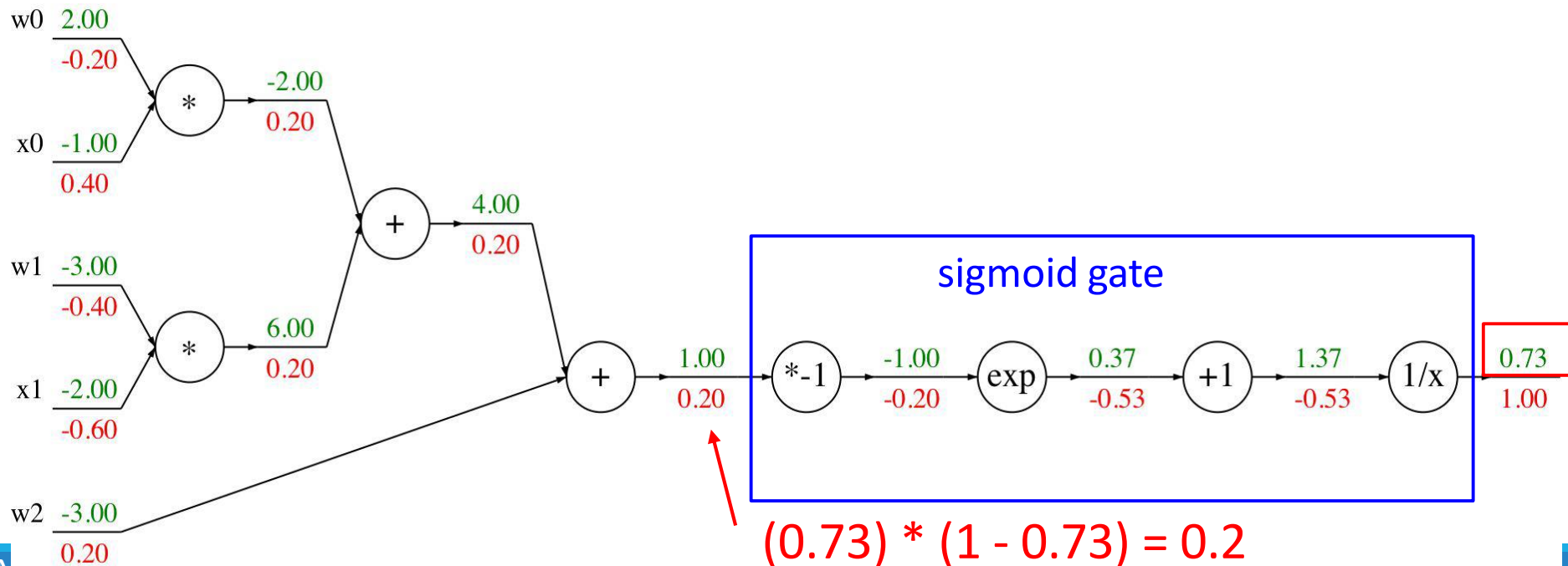
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$ sigmoid function

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left( \frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \left( \frac{1}{1 + e^{-x}} \right) = (1 - \sigma(x))\sigma(x)$$



sigmoid gate

(0.73) * (1 - 0.73) = 0.2

# Summary

- **Neuron**
  - Input and output

- **Neural Networks**
  - Perceptron
  - Multi-layer network
  - Deep neural networks

- **How Neural Network Works**
  - Calculation process
  - Work as a multi-class classifier

- **Backpropagation**

# Thinking

- Are more layers in a neural network better?

- Why shouldn't we use linear function as activation function of neural network?