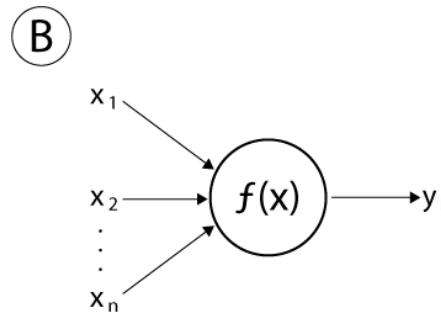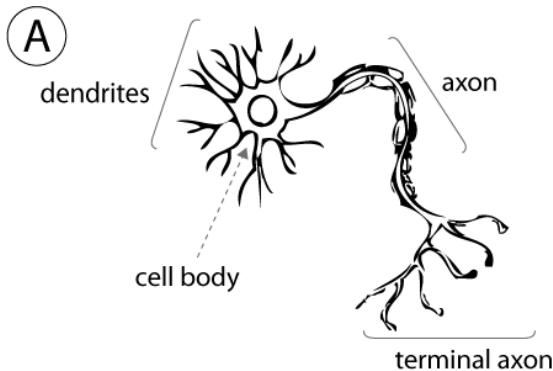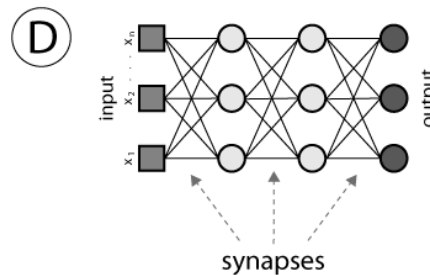# Deep learning

Hung-Hsuan Chen

# Biological neural network vs artificial neural network (ANN)



A **_synapse_** is a structure that permits a **_neuron_** to pass an electrical or chemical signal to another neuron.

A: human neuron
B: artificial neuron
C: biological synapse
D: ANN synapse

# Simple ANN vs deep learning



**Simple Neural Network**

**Deep Learning Neural Network**

🔴 Input Layer　🟠 Hidden Layer　🔵 Output Layer

- Simple neural network is also called single layer neural network
  - Single layer ➔ single hidden layer

# Why deep learning?

# Deep learning and big data

- Small model: tend to have high bias
- Deep learning (large model) and small data: tend to have high variance
- Deep learning (large model) and big data: small bias and small variance
  - Computation cost is an issue



Very small model          Very large model

# Learning representations

- Deep learning is part of a broader family of machine learning methods based on ***learning data representations***

- E.g., in CV, traditionally handcraft features (e.g., SIFT, HOG, etc.) are used as the features, but in deep learning, "raw pixels" used as the features

Input → Feature Extractor → Features → Traditional ML Algorithm → Output

Traditional Machine Learning Flow

Input → Deep Learning Algorithm → Output

Deep Learning Flow

# An example of learning data representation



Deep neural networks learn hierarchical feature representations

# Computation of a neuron



$O = \sigma(I)$

$$I = \sum_{i=1}^{d} w_i x_i^{(m)} = \boldsymbol{w}^T \boldsymbol{x}^{(m)}$$

$w_1 \quad w_2 \quad w_d$

$x_1^{(m)} \quad x_2^{(m)} \quad \dots \quad x_d^{(m)}$

m instance with d features

- Function: $O = f\big(\boldsymbol{x}^{(m)}\big) = \sigma\big(\boldsymbol{w}^T \boldsymbol{x}^{(m)}\big)$
  - We ignore the bias term
  - $\sigma$: a non-linear activation function
    - E.g., $\sigma(z) = \frac{1}{1+\exp(-z)}$
- How to find $\boldsymbol{w}$ such that the loss is minimized?
  - (Stochastic) gradient descent
- Logistic regression

# Review: gradient of the logistic function

- Lemma

  If $f(z) = \frac{1}{1+\exp(-z)}$, then $f'(z) = f(z)\big(1 - f(z)\big)$

- Proof

$$f'(z) = -\left(\frac{1}{1 + exp(-z)}\right)^2 (-\exp(-z))$$

$$= \frac{1 + \exp(-z) - 1}{(1 + \exp(-z))^2} = \frac{1}{1 + \exp(-z)}\left(1 - \frac{1}{(1 + \exp(-z))}\right)$$

$$= f(z)\big(1 - f(z)\big)$$

# Review: logistic regression model training (by SGD)

- Predicting function (ignore the bias term)
$$\hat{y}^{(m)} = f\big(\boldsymbol{x}^{(m)}\big) = \sigma\big(\boldsymbol{w}^T \boldsymbol{x}^{(m)}\big)$$

- Loss for the $m^{\text{th}}$ instance:
$$loss := -y^{(m)} \log \hat{y}^{(m)} - \big(1 - y^{(m)}\big) \log\big(1 - \hat{y}^{(m)}\big)$$

- Weight update rule
$$w_i \leftarrow w_i - \eta \frac{\partial loss}{\partial w_i}$$

# Review: gradient of logistic regression

$$\frac{\partial loss}{\partial w_i}$$

$$= \frac{\partial - y^{(m)} \log\left(\sigma\left(\boldsymbol{w}^T \boldsymbol{x}^{(m)}\right)\right)}{\partial \log\left(\sigma(\boldsymbol{w}^T \boldsymbol{x}^{(m)})\right)} \cdot \frac{\partial \log\left(\sigma\left(\boldsymbol{w}^T \boldsymbol{x}^{(m)}\right)\right)}{\partial \sigma(\boldsymbol{w}^T \boldsymbol{x}^{(m)})} \cdot \frac{\partial \sigma\left(\boldsymbol{w}^T \boldsymbol{x}^{(m)}\right)}{\partial \boldsymbol{w}^T \boldsymbol{x}^{(m)}} \cdot \frac{\partial \boldsymbol{w}^T \boldsymbol{x}^{(m)}}{\cancel{\partial} \; \partial w_i}$$

$$+ \frac{\partial(1 - y^{(m)}) \log\left(1 - \sigma\left(\boldsymbol{w}^T \boldsymbol{x}^{(m)}\right)\right)}{\partial \log\left(1 - \sigma(\boldsymbol{w}^T \boldsymbol{x}^{(m)})\right)} \cdot \frac{\partial \log\left(1 - \sigma\left(\boldsymbol{w}^T \boldsymbol{x}^{(m)}\right)\right)}{\partial \left(1 - \sigma(\boldsymbol{w}^T \boldsymbol{x}^{(m)})\right)} \cdot \frac{\partial(1 - \sigma(\boldsymbol{w}^T \boldsymbol{x}^{(m)}))}{\partial \sigma(\boldsymbol{w}^T \boldsymbol{x}^{(m)})}$$

$$\cdot \frac{\partial \sigma\left(\boldsymbol{w}^T \boldsymbol{x}^{(m)}\right)}{\partial \boldsymbol{w}^T \boldsymbol{x}^{(m)}} \cdot \frac{\partial \boldsymbol{w}^T \boldsymbol{x}^{(m)}}{\partial(w_i)}$$

$$= y^{(m)} \cdot \frac{1}{\sigma(\boldsymbol{w}^T \boldsymbol{x}^{(m)})} \cdot \sigma\left(\boldsymbol{w}^T \boldsymbol{x}^{(m)}\right) \cdot \left(1 - \sigma\left(\boldsymbol{w}^T \boldsymbol{x}^{(m)}\right)\right) \cdot x_i^{(m)}$$

$$+ \left(1 - y^{(m)}\right) \cdot \frac{1}{1 - \sigma(\boldsymbol{w}^T \boldsymbol{x}^{(m)})} \cdot (-1) \cdot \sigma\left(\boldsymbol{w}^T \boldsymbol{x}^{(m)}\right) \cdot \left(1 - \sigma\left(\boldsymbol{w}^T \boldsymbol{x}^{(m)}\right)\right) \cdot x_i^{(m)}$$

$$= x_i^{(m)}\left(y^{(m)} - \hat{y}^{(m)}\right)$$

# Feedforward neural network

Back-prop: adjust
weights to
minimize error

Forward: predict
$f(\boldsymbol{x}^{(m)})$ based on
current weights

$$\boldsymbol{O}^{(\ell)} = \begin{bmatrix} O_1^{(\ell)} & \dots & O_I^{(\ell)} \end{bmatrix}^T$$

$$\boldsymbol{W}^{(\ell)} = \begin{bmatrix} w_{11}^{(\ell)} & \cdots & w_{1J}^{(\ell)} \\ \vdots & \ddots & \vdots \\ w_{I1}^{(\ell)} & \cdots & w_{IJ}^{(\ell)} \end{bmatrix}$$



$$\boldsymbol{O}^{(2)} = \begin{bmatrix} O_1^{(2)} \end{bmatrix}^T$$

$$\boldsymbol{W}^{(2)} = \begin{bmatrix} w_{11}^{(2)} & w_{12}^{(2)} & w_{13}^{(2)} \end{bmatrix}$$

$$\boldsymbol{O}^{(1)} = \begin{bmatrix} O_1^{(1)} & O_2^{(1)} & O_3^{(1)} \end{bmatrix}^T$$

$$\boldsymbol{W}^{(1)} = \begin{bmatrix} w_{11}^{(1)} & \cdots & w_{14}^{(1)} \\ \vdots & \ddots & \vdots \\ w_{31}^{(1)} & \cdots & w_{34}^{(1)} \end{bmatrix}$$

The initial weights ($\boldsymbol{W}^{(1)}$ and $\boldsymbol{W}^{(2)}$) are randomly assigned

# Forward (predict)

- Forward (predict)
  - $I^{(1)} = W^{(1)} x$
  - $O^{(1)} = \sigma(I^{(1)})$
  - $I^{(2)} = W^{(2)} O^{(1)}$
  - $O^{(2)} = \sigma(I^{(2)})$
- Initial weights $w_{ij}^{(\ell)}$: randomly assign
  - E.g., $w_{ij}^{(\ell)} \leftarrow N(0,1)$

Forward: predict $f(x^{(m)})$ based on current weights

# Forward example

- Input feature and target
  - $x = [1 \quad -1 \quad 3 \quad -2]^T, y = 1$
- Initial weights (random assign)
  - $W^{(2)} = [-1 \quad 1 \quad -1]$
  - $W^{(1)} = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & -1 & 0.5 & -1 \\ -1 & -1 & 0 & 0.5 \end{bmatrix}$
- Forward
  - $I^{(1)} = W^{(1)}x = [4 \quad 0.5 \quad 1]^T$
  - $O^{(1)} = \sigma(I^{(1)}) = [0.98 \quad 0.62 \quad 0.73]^T$
  - $I^{(2)} = W^{(2)}O^{(1)} = [-1.09]$
  - $O^{(2)} = \sigma(I^{(2)}) = [0.25]$

# Back-propagation

Slides are taken from Prof. Fei-Fei Li

# Computational graph

- $f(x, y, z) = (x + y)z$
  - E.g., $x = -2, y = 5, z = -4$

$x: -2$

$\partial q / \partial x = 1$

$y: 5$

$\partial q / \partial y = 1$

$\partial f / \partial z$
$= q = 3$

$z: -4$

$q = x + y = 3$

$\partial f / \partial q = z$
$= -4$

$f = q * z = -12$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \times \frac{\partial q}{\partial x} = (-4) \times 1 = -4$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \times \frac{\partial q}{\partial y} = (-4) \times 1 = -4$$

$$\frac{\partial f}{\partial z} = 3$$

# Local gradient and chain rule

- Given $f(x, y) = z$ and the global loss is $L$

  1. Compute local gradients $\frac{\partial z}{\partial x}$ and $\frac{\partial z}{\partial y}$

  2. The value of $\frac{\partial L}{\partial z}$ is computed by other components in the computational graph

  3. Based on chain rule, we can get $\frac{\partial L}{\partial x}$ and $\frac{\partial L}{\partial y}$ by $\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \times \frac{\partial z}{\partial x}$, and $\frac{\partial L}{\partial y} = \frac{\partial L}{\partial z} \times \frac{\partial z}{\partial y}$

# Compute local gradient for every operator and apply chain rule

# Example: logistic function

$$f(\boldsymbol{w}, \boldsymbol{x}) = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2)}}$$

Observing $(x_1, x_2, y) = (-1, -2, 1)$

$w_0: -3$

$w_1: 2$

$x_1: -1$

$w_2: -3$

$x_2: -2$

$g_1$ $-2$

$g_2$ $6$

$g_3$ $+$ $-5$

$g_4$ $+$ $1$

$g_5$ $*(-1)$ $-1$

$g_6$ $exp$ $0.37$

$g_7$ $+1$ $1.37$

$g_8$ $\frac{1}{x}$ $0.73$

$L$ $0.315$

By cross entropy loss (here $\hat{y} = g_8$):
$L(y, g_8) = -y \log g_8 - (1 - y) \log(1 - g_8) = 0.315$

# Example: logistic function

$$f(\boldsymbol{w}, \boldsymbol{x}) = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2)}}$$  Observing $(x_1, x_2, y) = (-1, -2, 1)$

$w_0: -3$

$g_3$

$w_1: 2$    $g_1$ $-2$    $+$    $-5$ $g_4$    $g_5$    $g_6$    $g_7$    $g_8$

$x_1: -1$    $*$

$g_4$ $+$    $1$    $g_5$ $*(-1)$    $-1$    $g_6$ $exp$    $0.37$    $g_7$ $+1$    $1.37$    $g_8$ $\frac{1}{x}$    $0.73$    $L$    $0.315$

$-1.37$

$w_2: -3$    $g_2$ $6$

$x_2: -2$    $*$

$$L(y, g_8) = -y \log g_8 - (1 - y) \log(1 - g_8)$$
$$\Rightarrow \frac{\partial L}{\partial g_8} = -\frac{y}{g_8} + \frac{1 - y}{1 - g_8}$$
$$\frac{\partial L}{\partial g_8}\bigg|_{g_8 = 0.73} = -\frac{1}{0.73} = -1.37$$

12/8/20    20

# Example: logistic function

$$f(\boldsymbol{w}, \boldsymbol{x}) = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2)}}$$

Observing $(x_1, x_2, y) = (-1, -2, 1)$



$w_0: -3$

$w_1: 2$

$x_1: -1$

$g_1 \quad -2$

$g_3$

$-5 \quad g_4$

$g_5 \quad -1$

$g_6 \quad 0.37$

$g_7 \quad 1.37$

$g_8 \quad 0.73$

$0.315$

$1 \quad *(-1) \quad exp \quad +1 \quad \frac{1}{x} \quad L$

$0.73 \quad -1.37$

$w_2: -3$

$g_2 \quad 6$

$x_2: -2$

$$g_8(g_7) = \frac{1}{g_7}$$

$$\Rightarrow \frac{\partial g_8}{\partial g_7} = -\frac{1}{g_7^2}$$

$$\left.\frac{\partial g_8}{\partial g_7}\right|_{g_7 = 1.37} = -\frac{1}{1.37^2} = -0.53$$

$$\frac{\partial L}{\partial g_7} = \frac{\partial L}{\partial g_8} \times \frac{\partial g_8}{\partial g_7} = -1.37 \times -0.53 = 0.73$$

# Example: logistic function

$$f(\boldsymbol{w}, \boldsymbol{x}) = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2)}}$$

Observing $(x_1, x_2, y) = (-1, -2, 1)$



$g_7(g_6) = g_6 + 1$

$$\Rightarrow \frac{\partial g_7}{\partial g_6} = 1$$

$$\frac{\partial g_7}{\partial g_6}\bigg|_{g_6 = 0.37} = 1$$

$$\frac{\partial L}{\partial g_6} = \frac{\partial L}{\partial g_7} \times \frac{\partial g_7}{\partial g_6} = 0.73 \times 1 = 0.73$$
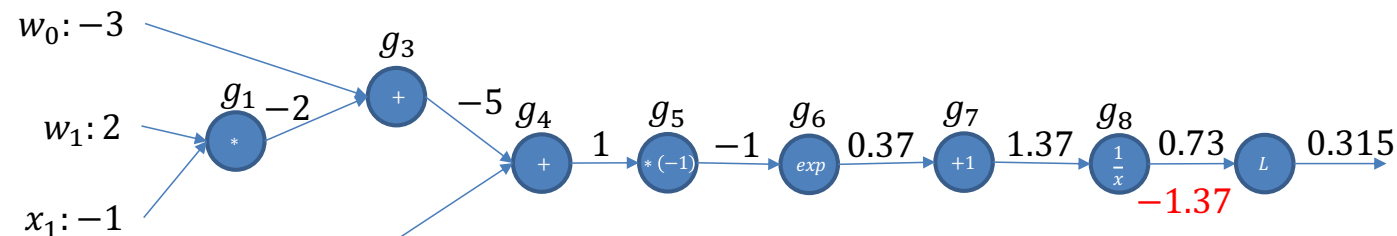
# Example: logistic function

$$f(\boldsymbol{w}, \boldsymbol{x}) = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2)}}$$

Observing $(x_1, x_2, y) = (-1, -2, 1)$



$w_0: -3$

$w_1: 2$

$x_1: -1$

$w_2: -3$

$x_2: -2$

$g_3$

$g_1$  $-2$

$-5$  $g_4$

$g_5$  $-1$

$g_6$  $0.37$

$g_7$  $1.37$

$g_8$  $0.73$

$0.315$

$1$

$0.27$

$0.73$

$0.73$

$-1.37$

$g_2$  $6$

$g_6(g_5) = \exp(g_5)$

$$\Rightarrow \frac{\partial g_6}{\partial g_5} = \exp(g_5)$$

$$\left. \frac{\partial g_6}{\partial g_5} \right|_{g_5 = -1} = \exp(-1) = 0.37$$

$$\frac{\partial L}{\partial g_5} = \frac{\partial L}{\partial g_6} \times \frac{\partial g_6}{\partial g_5} = 0.73 \times 0.37 = 0.27$$

# Example: logistic function

$$f(\boldsymbol{w}, \boldsymbol{x}) = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2)}}$$
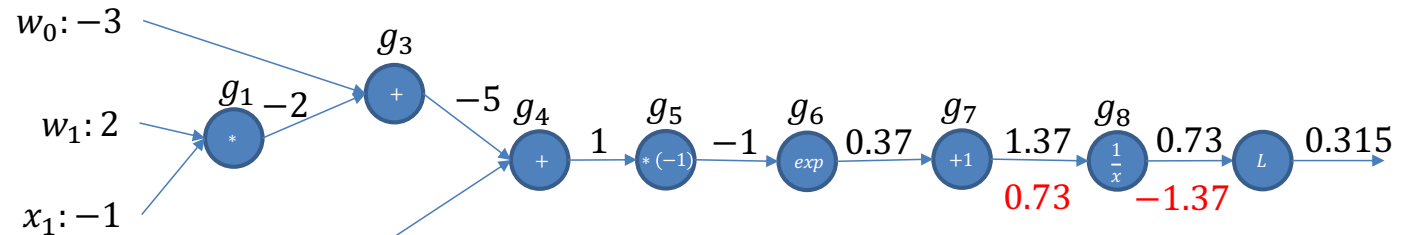
Observing $(x_1, x_2, y) = (-1, -2, 1)$



$w_0: -3$

$w_1: 2$

$x_1: -1$

$g_1$ $-2$

$g_3$

$-5$ $g_4$

$1$ $g_5$

$-1$ $g_6$

$0.37$ $g_7$

$1.37$ $g_8$

$0.73$

$0.315$

$-0.27$   $0.27$   $0.73$   $0.73$   $-1.37$

$w_2: -3$

$g_2$ $6$

$x_2: -2$

$$g_5(g_4) = -g_4$$
$$\Rightarrow \frac{\partial g_5}{\partial g_4} = -1$$
$$\left. \frac{\partial g_5}{\partial g_4} \right|_{g_4 = 1} = -1$$
$$\frac{\partial L}{\partial g_4} = \frac{\partial L}{\partial g_5} \times \frac{\partial g_5}{\partial g_4} = 0.27 \times (-1) = -0.27$$

# Example: logistic function

$$f(\boldsymbol{w}, \boldsymbol{x}) = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2)}}$$

Observing $(x_1, x_2, y) = (-1, -2, 1)$

$w_0: -3$

$w_1: 2$

$g_1 \; -2$

$*$

$g_3$

$+$

$-5 \; g_4$

$1 \; g_5$

$*(-1)$

$-1 \; g_6$

$exp$

$0.37 \; g_7$

$+1$

$1.37 \; g_8$

$\frac{1}{x}$

$0.73$

$L$

$0.315$

$-0.27$

$+$

$-0.27$    $0.27$    $0.73$    $0.73$    $-1.37$

$x_1: -1$

$6$

$-0.27$

$w_2: -3$

$g_2$

$*$

$x_2: -2$

有 input

$$g_4(g_2, g_3) = g_2 + g_3$$
$$\Rightarrow \frac{\partial g_4}{\partial g_2} = 1, \frac{\partial g_4}{\partial g_3} = 1$$
$$\left.\frac{\partial g_4}{\partial g_2}\right|_{(g_2, g_3)=(6,-5)} = 1, \left.\frac{\partial g_4}{\partial g_3}\right|_{(g_2, g_3)=(6,-5)} = 1$$

$$\frac{\partial L}{\partial g_2} = \frac{\partial L}{\partial g_4} \times \frac{\partial g_4}{\partial g_2} = -0.27 \times 1 = -0.27$$
$$\frac{\partial L}{\partial g_3} = \frac{\partial L}{\partial g_4} \times \frac{\partial g_4}{\partial g_3} = -0.27 \times 1 = -0.27$$
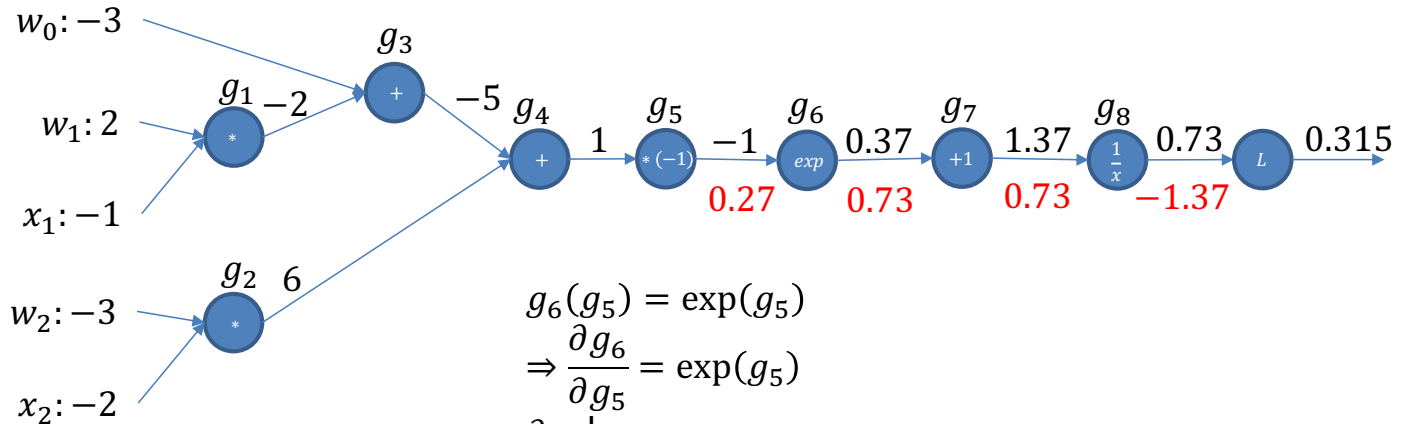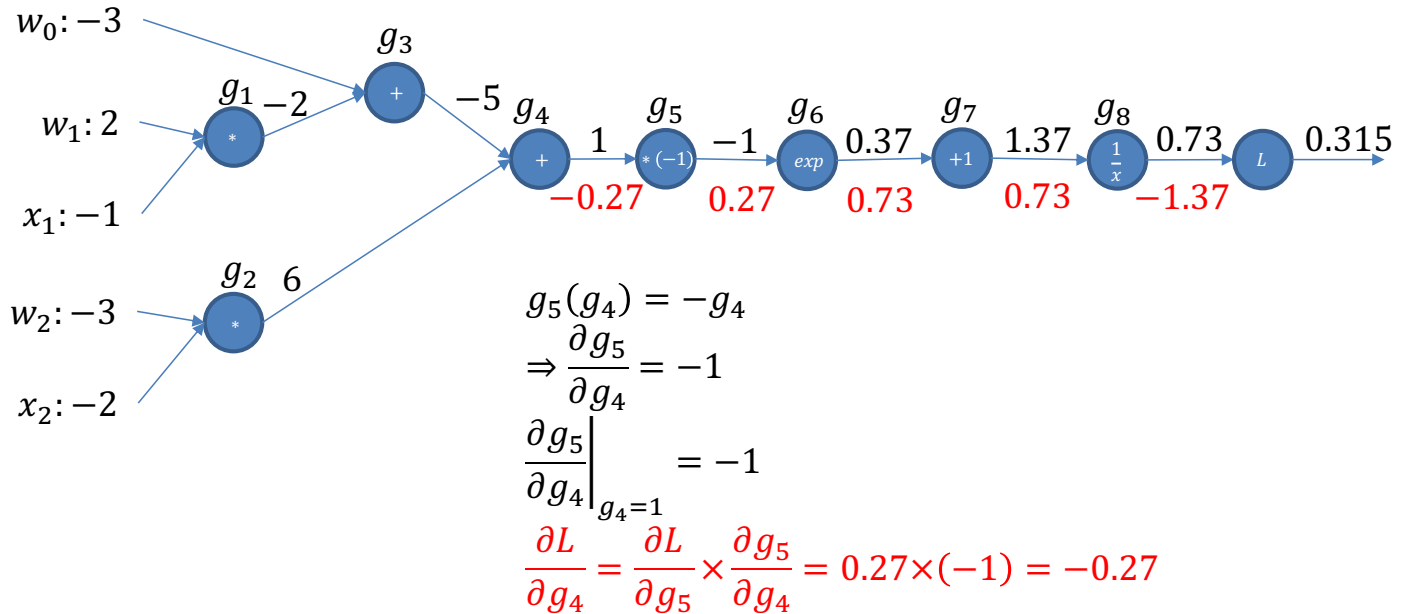
12/8/20

25

# Example: logistic function

$$f(\boldsymbol{w}, \boldsymbol{x}) = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2)}}$$

Observing $(x_1, x_2, y) = (-1, -2, 1)$

$w_0: -3$

$-0.27$

$g_3$

$g_1$ $-2$

$w_1: 2$

$*$ $-0.27$ $-0.27$

$-5$ $g_4$

$g_5$ $g_6$ $g_7$ $g_8$

$x_1: -1$

$+$ $1$ $*(-1)$ $-1$ $exp$ $0.37$ $+1$ $1.37$ $\frac{1}{x}$ $0.73$ $L$ $0.315$

$6$

$-0.27$ $-0.27$ $0.27$ $0.73$ $0.73$ $-1.37$

$w_2: -3$

$g_2$

$*$

$x_2: -2$

$$g_3(w_0, g_1) = w_0 + g_1$$
$$\Rightarrow \frac{\partial g_3}{\partial w_0} = 1, \frac{\partial g_3}{\partial g_1} = 1$$
$$\left.\frac{\partial g_3}{\partial w_0}\right|_{(w_0, g_1) = (-3, -2)} = 1, \left.\frac{\partial g_3}{\partial g_1}\right|_{(w_0, g_1) = (-3, -2)} = 1$$
$$\frac{\partial L}{\partial w_0} = \frac{\partial L}{\partial g_3} \times \frac{\partial g_3}{\partial w_0} = -0.27 \times 1 = -0.27$$
$$\frac{\partial L}{\partial g_1} = \frac{\partial L}{\partial g_3} \times \frac{\partial g_3}{\partial g_1} = -0.27 \times 1 = -0.27$$
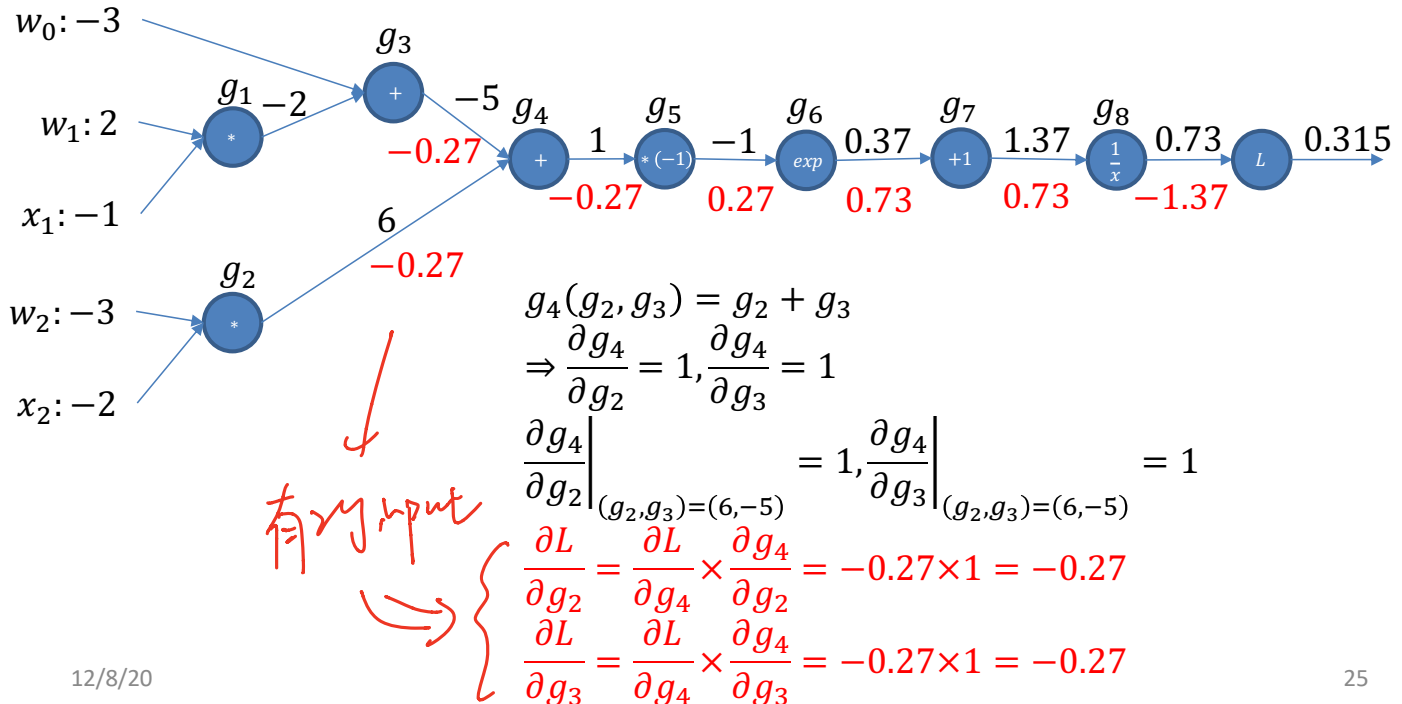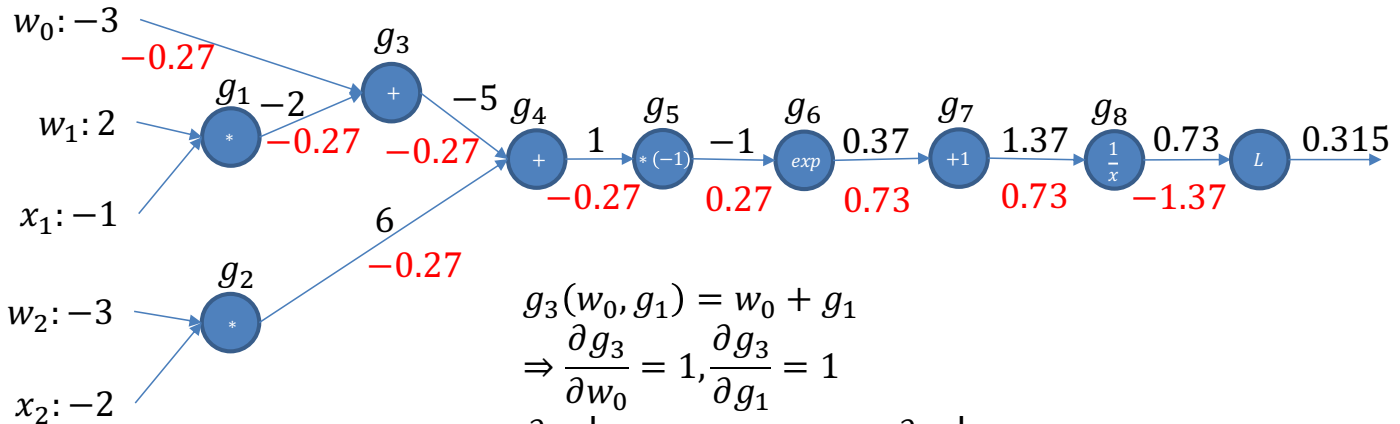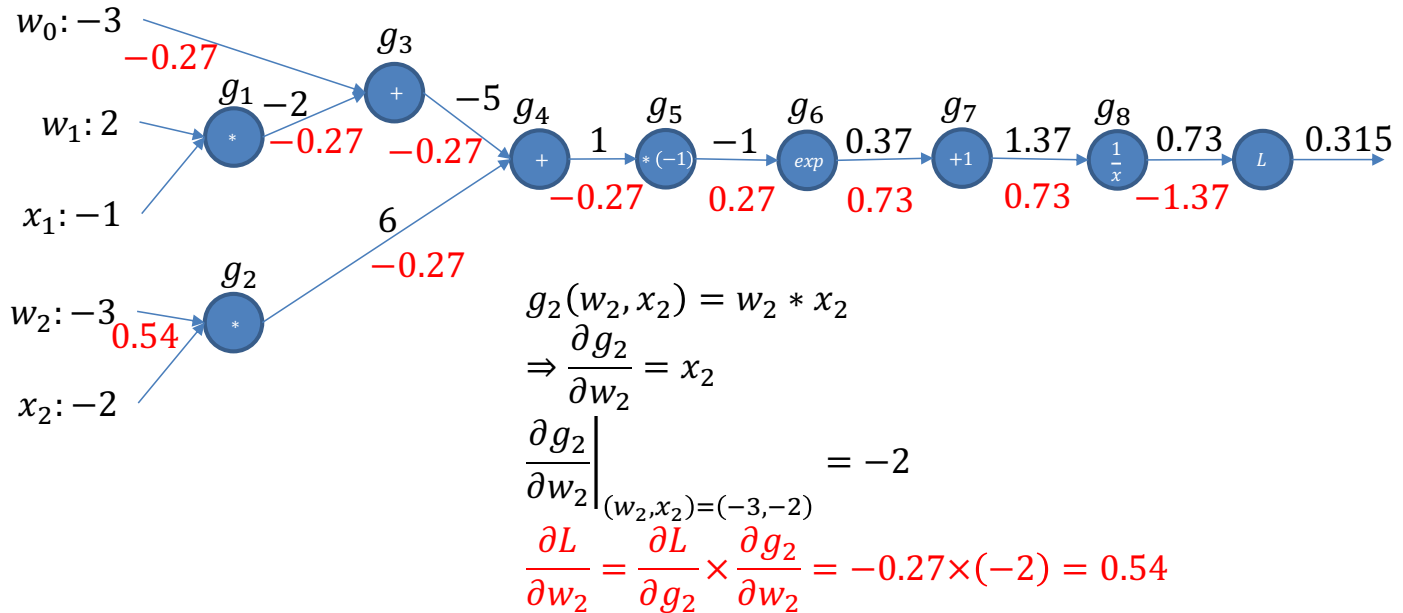
# Example: logistic function

$$f(\boldsymbol{w}, \boldsymbol{x}) = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2)}}$$

Observing $(x_1, x_2, y) = (-1, -2, 1)$



$w_0: -3$
$-0.27$

$g_3$

$g_1$  $-2$
$w_1: 2$  $-0.27$  $-0.27$  $-5$  $g_4$  $g_5$  $-1$  $g_6$  $0.37$  $g_7$  $1.37$  $g_8$  $0.73$  $0.315$

$x_1: -1$  $+$  $1$  $*(-1)$  $exp$  $+1$  $\frac{1}{x}$  $L$
$-0.27$  $0.27$  $0.73$  $0.73$  $-1.37$

$6$
$-0.27$

$g_2$
$w_2: -3$
$0.54$  $*$

$x_2: -2$

$g_2(w_2, x_2) = w_2 * x_2$

$$\Rightarrow \frac{\partial g_2}{\partial w_2} = x_2$$

$$\left.\frac{\partial g_2}{\partial w_2}\right|_{(w_2, x_2) = (-3, -2)} = -2$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial g_2} \times \frac{\partial g_2}{\partial w_2} = -0.27 \times (-2) = 0.54$$

# Example: logistic function

$$f(\boldsymbol{w}, \boldsymbol{x}) = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2)}}$$

Observing $(x_1, x_2, y) = (-1, -2, 1)$



$w_0: -3$
$-0.27$

$g_3$

$w_1: 2$
$0.27$
$g_1$ $-2$
$-0.27$ $-0.27$

$x_1: -1$

$g_3$ $+$ $-5$ $g_4$
$g_4$ $+$ $1$ $g_5$ $*(-1)$ $-1$ $g_6$ $exp$ $0.37$ $g_7$ $+1$ $1.37$ $g_8$ $\frac{1}{x}$ $0.73$ $L$ $0.315$
$-0.27$ $0.27$ $0.73$ $0.73$ $-1.37$

$6$
$-0.27$

$w_2: -3$
$0.54$
$g_2$ $*$

$x_2: -2$

$g_1(w_1, x_1) = w_1 * x_1$

$$\Rightarrow \frac{\partial g_1}{\partial w_1} = x_1$$

$$\left. \frac{\partial g_1}{\partial w_1} \right|_{(w_1, x_1) = (2, -1)} = -1$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial g_1} \times \frac{\partial g_1}{\partial w_1} = -0.27 \times (-1) = 0.27$$
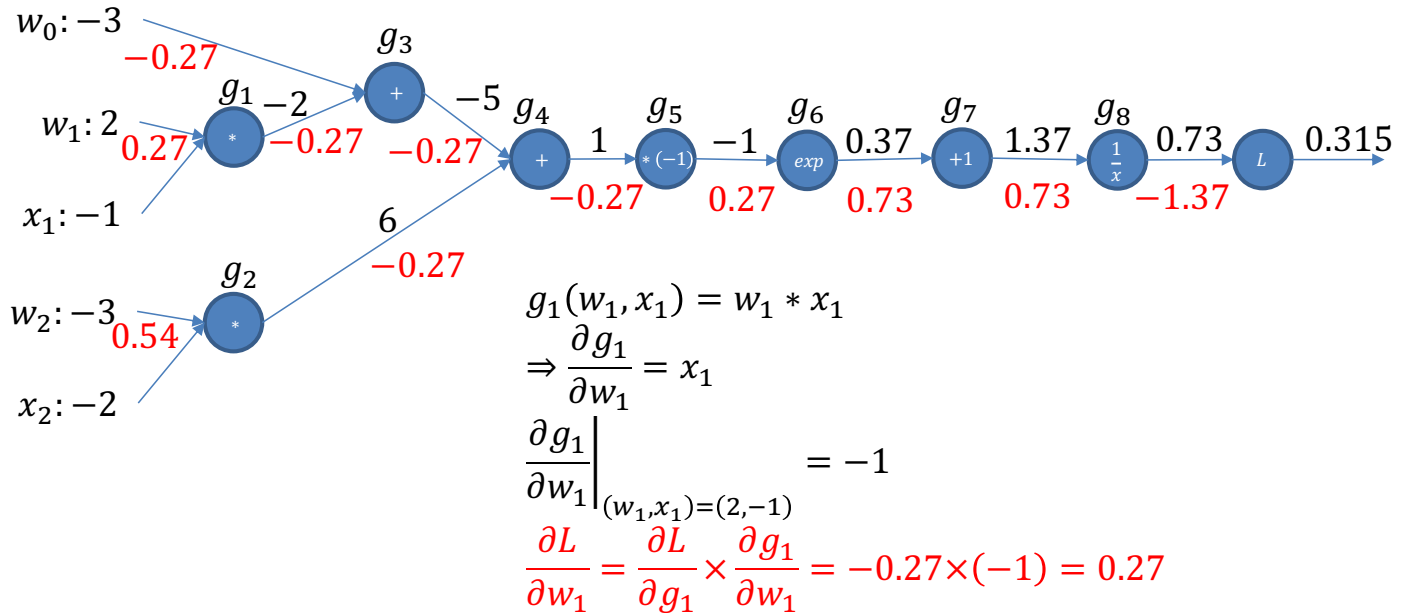
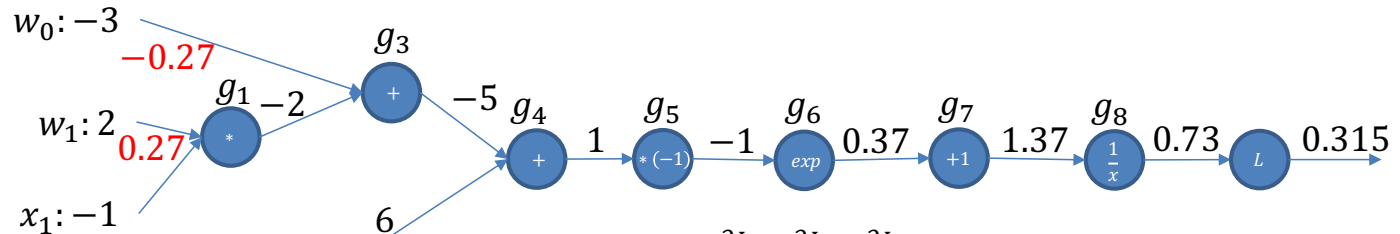# Example: logistic function

$$f(\boldsymbol{w}, \boldsymbol{x}) = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2)}}$$

Observing $(x_1, x_2, y) = (-1, -2, 1)$

$w_0: -3$
$-0.27$

$w_1: 2$
$0.27$

$g_1$ $-2$

$g_3$ $+$

$-5$ $g_4$ $+$ $1$ $g_5$ $*(-1)$ $-1$ $g_6$ $exp$ $0.37$ $g_7$ $+1$ $1.37$ $g_8$ $\frac{1}{x}$ $0.73$ $L$ $0.315$

$x_1: -1$

$6$

$g_2$ $*$

$w_2: -3$
$0.54$

$x_2: -2$

Now we've get $\frac{\partial L}{\partial w_0}, \frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}$

Gradient descent: $\theta^{(k+1)} = \theta^{(k)} - \eta \frac{\partial L}{\partial \theta^{(k)}}$

If we set $\eta = 0.001$:

$$\begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} -3 \\ 2 \\ -3 \end{bmatrix} - 0.001 \begin{bmatrix} -0.27 \\ 0.27 \\ 0.54 \end{bmatrix} = \begin{bmatrix} -2.9973 \\ 1.9973 \\ -3.00054 \end{bmatrix}$$

Compare the old $\boldsymbol{w}$ and new $\boldsymbol{w}$:
$f(w_0 = -3, w_1 = 2, w_2 = -3, \boldsymbol{x}) = 0.7311$
$f(w_0 = -2.9973, w_1 = 1.9973, w_2 = -3.00054, \boldsymbol{x}) =$
0.7323, indeed getting closer to 1

$\pi$ ... $\sigma$

$\rightarrow \Delta \times \sigma(1) \times [1 - \sigma(1)]$   $\frac{\partial C}{\partial} = \Delta$   使 compuwe easy。

$\because C = a + b$

$\therefore \frac{\partial C}{\partial a} = 1, \frac{\partial C}{\partial b} = 1$

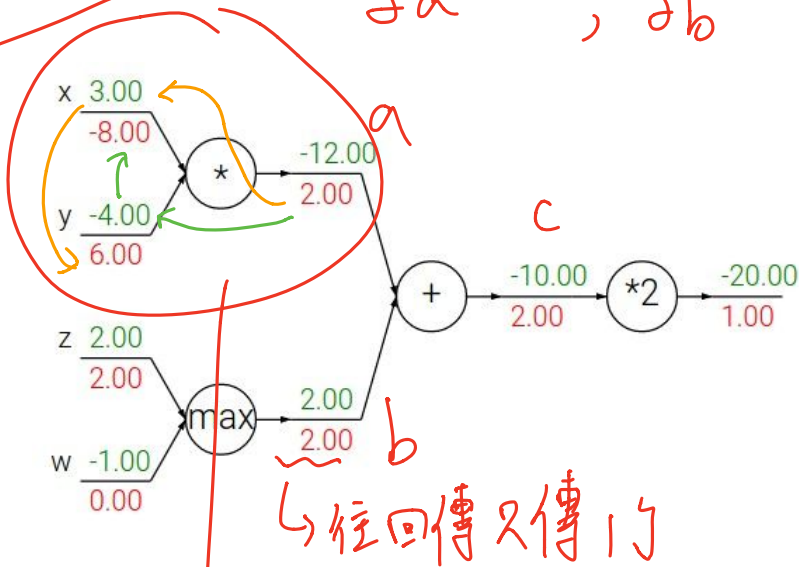# Patterns in backward flow

**add** gate: gradient distributor

**max** gate: gradient router

**mul** gate: gradient switcher

**sigmoid** gate: $\sigma'(z) = \sigma(z)(1 - \sigma(z))$

$b = \max(z, w)$

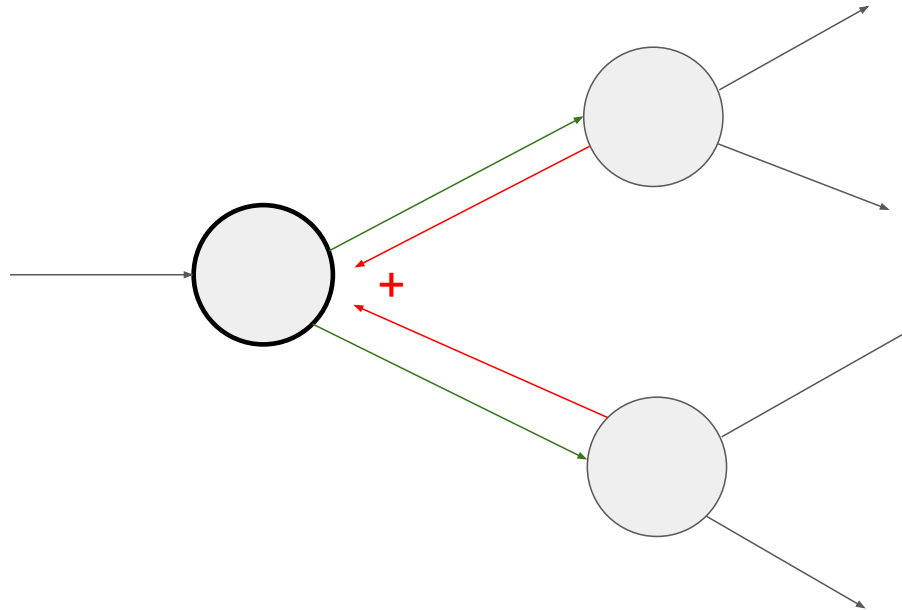$\rightarrow \frac{\partial b}{\partial z} = \begin{cases} 1 & \text{if } z > w \\ 0 & \text{if } z < w \end{cases}$



a

c

b

↳往回傳只傳1了

$\frac{\partial b}{\partial w} = \begin{cases} 0 & \text{if } z > w \\ 1 & \text{if } z < w \end{cases}$

$\rightarrow a = x - y$

Gradients add at branches

# Forward backprop example



$x_1$   3

$g_1$

$w_{11}^{(1)}$   0.2   $\times$   0.6

$x_2$   2

$g_2$

$w_{12}^{(1)}$   0.1   $\times$   0.2

$g_5$   $+$   0.8   $g_7$   $\sigma$   0.69

$w_{11}^{(2)}$   $-1$

$g_9$   $\times$   $-0.69$

$x_1$   3

$g_3$

$w_{21}^{(1)}$   $-0.7$   $\times$   $-2.1$

$x_2$   2

$g_4$

$w_{22}^{(1)}$   1.2   $\times$   2.4

$g_6$   $+$   0.3   $g_8$   $\sigma$   0.57

$g_{10}$   $\times$   0.57

$w_{12}^{(2)}$   1

$g_{11}$   $+$   $-0.12$   $g_{12}$   $\sigma$   0.47   $L$   $L = 0.755$

$L$
$$= -y \log g_{12} - (1-y) \log(1 - g_{12})$$
$$= -\log 0.47 = 0.755$$

# Forward backprop example



$x_1$ 3

$g_1$ $\times$ 0.6

$w_{11}^{(1)}$ 0.2

$x_2$ 2

$g_2$ $\times$ 0.2

$w_{12}^{(1)}$ 0.1

$g_5$ + 0.8

$g_7$ $\sigma$ 0.69

$w_{11}^{(2)}$ $-1$

$g_9$ $\times$ $-0.69$

$x_1$ 3

$g_3$ $\times$ $-2.1$

$w_{21}^{(1)}$ $-0.7$

$x_2$ 2

$g_4$ $\times$ 2.4

$w_{22}^{(1)}$ 1.2

$g_6$ + 0.3

$g_8$ $\sigma$ 0.57

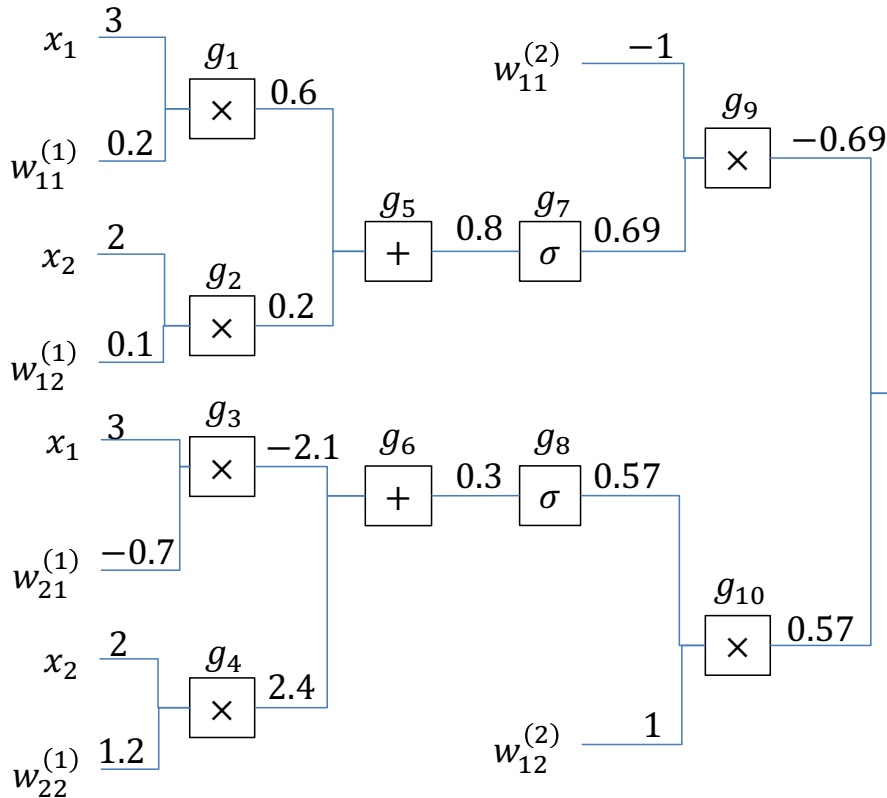$g_{10}$ $\times$ 0.57

$w_{12}^{(2)}$ 1
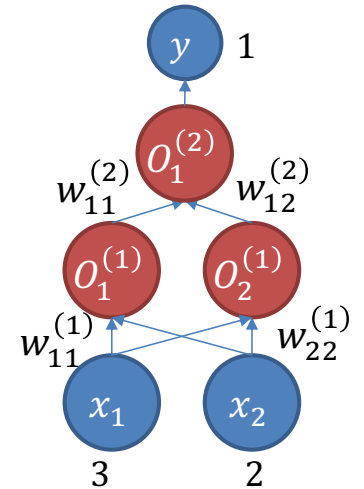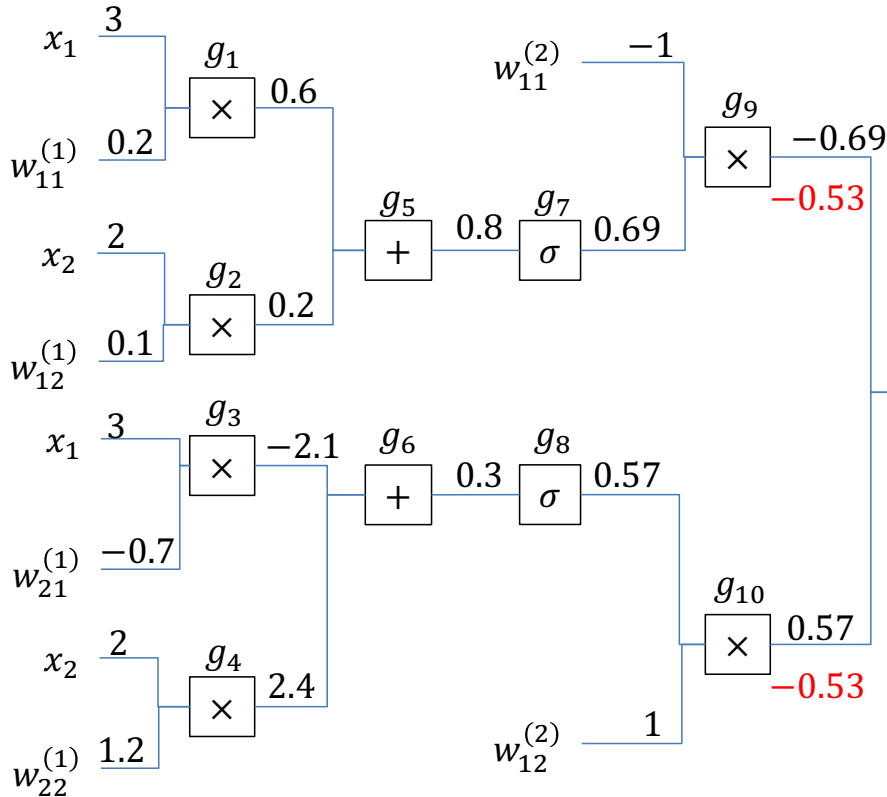
$g_{11}$ + $-0.12$

$g_{12}$ $\sigma$ 0.47

$L$    $L = 0.755$

$-2.128$

$$\frac{\partial L}{\partial g_{12}} = -\frac{y}{g_{12}} + \frac{1-y}{1-g_{12}} = -\frac{1}{0.47}$$
$$= -2.128$$

$y$  1

$O_1^{(2)}$
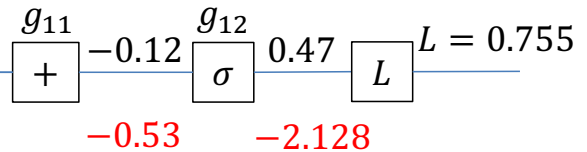
$w_{11}^{(2)}$    $w_{12}^{(2)}$

$O_1^{(1)}$    $O_2^{(1)}$

$w_{11}^{(1)}$    $w_{22}^{(1)}$

$x_1$    $x_2$

3    2

# Forward backprop example

$x_1$   3

$w_{11}^{(1)}$   0.2

$g_1$ $\times$   0.6

$x_2$   2

$w_{12}^{(1)}$   0.1

$g_2$ $\times$   0.2

$g_5$ $+$   0.8

$g_7$ $\sigma$   0.69

$w_{11}^{(2)}$   $-1$

$g_9$ $\times$   $-0.69$

$x_1$   3

$w_{21}^{(1)}$   $-0.7$

$g_3$ $\times$   $-2.1$

$x_2$   2

$w_{22}^{(1)}$   1.2

$g_4$ $\times$   2.4

$g_6$ $+$   0.3

$g_8$ $\sigma$   0.57

$g_{10}$ $\times$   0.57

$w_{12}^{(2)}$   1

$g_{11}$ $+$   $-0.12$

$g_{12}$ $\sigma$   0.47

$L$   $L = 0.755$

$-0.53$     $-2.128$

$g_{12} = \sigma(g_{11})$

$\dfrac{\partial g_{12}}{\partial g_{11}} = \sigma(g_{11})\big(1 - \sigma(g_{11})\big)$

$= (0.47)(0.53) = 0.2491$

$\dfrac{\partial L}{\partial g_{11}} = \dfrac{\partial L}{\partial g_{12}} \times \dfrac{\partial g_{12}}{\partial g_{11}}$

$= -2.128 \times 0.2491 = -0.53$

$y$   1

$O_1^{(2)}$

$w_{11}^{(2)}$    $w_{12}^{(2)}$

$O_1^{(1)}$    $O_2^{(1)}$

$w_{11}^{(1)}$    $w_{22}^{(1)}$

$x_1$    $x_2$

3    2

# Forward backprop example



$x_1$ 3

$g_1$ [×] 0.6

$w_{11}^{(1)}$ 0.2

$x_2$ 2

$g_2$ [×] 0.2

$w_{12}^{(1)}$ 0.1

$g_5$ [+] 0.8 $g_7$ [σ] 0.69

$w_{11}^{(2)}$ −1

$g_9$ [×] −0.69
−0.53

$x_1$ 3

$g_3$ [×] −2.1

$w_{21}^{(1)}$ −0.7

$x_2$ 2

$g_4$ [×] 2.4

$w_{22}^{(1)}$ 1.2

$g_6$ [+] 0.3 $g_8$ [σ] 0.57

$g_{10}$ [×] 0.57
−0.53

$w_{12}^{(2)}$ 1

$g_{11}$ [+] −0.12 $g_{12}$ [σ] 0.47 [L] $L = 0.755$
−0.53 −2.128

$g_{11} = g_9 + g_{10}$

$$\frac{\partial g_{11}}{\partial g_9} = \frac{\partial g_{11}}{\partial g_{10}} = 1$$

$$\frac{\partial L}{\partial g_9} = \frac{\partial L}{\partial g_{11}} \times \frac{\partial g_{11}}{\partial g_9} = -0.53$$

$$\frac{\partial L}{\partial g_{10}} = \frac{\partial L}{\partial g_{11}} \times \frac{\partial g_{11}}{\partial g_{10}} = -0.53$$

# Forward backprop example



$x_1$   3

$g_1$

$\times$   0.6

$w_{11}^{(1)}$   0.2

$x_2$   2

$g_2$

$\times$   0.2

$w_{12}^{(1)}$   0.1

$g_5$

$+$   0.8

$g_7$

$\sigma$   0.69

$w_{11}^{(2)}$   $-1$

$-0.366$

$g_9$

$\times$   $-0.69$

$-0.53$

$0.53$

$x_1$   3

$g_3$

$\times$   $-2.1$

$w_{21}^{(1)}$   $-0.7$

$x_2$   2

$g_4$

$\times$   2.4

$w_{22}^{(1)}$   1.2

$g_6$

$+$   0.3

$g_8$

$\sigma$   0.57

$g_{10}$

$\times$   0.57

$-0.53$

$w_{12}^{(2)}$   1

$g_{11}$

$+$   $-0.12$

$-0.53$

$g_{12}$

$\sigma$   0.47
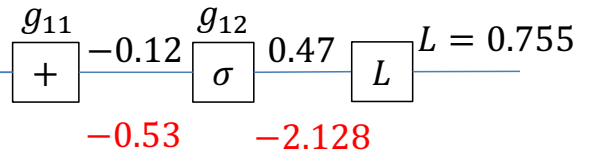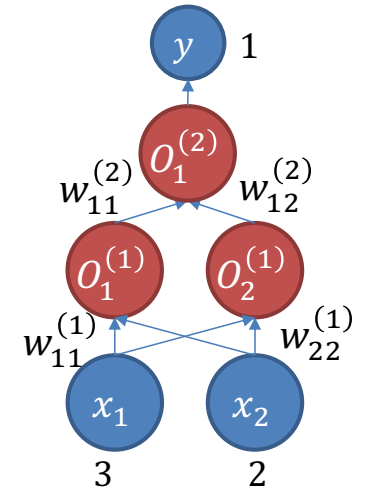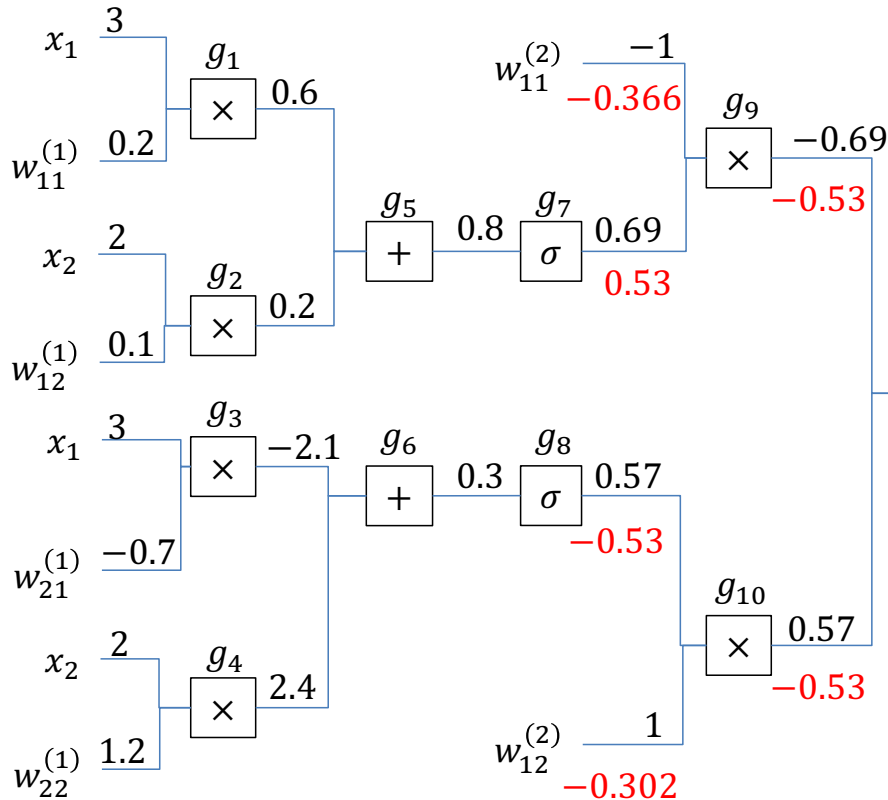
$-2.128$

$L$   $L = 0.755$

$$g_9 = w_{11}^{(2)} \times g_7$$
$$\frac{\partial g_9}{\partial w_{11}^{(2)}} = g_7, \frac{\partial g_9}{\partial g_7} = w_{11}^{(2)}$$
$$\frac{\partial L}{\partial w_{11}^{(2)}} = \frac{\partial L}{\partial g_9} \times \frac{\partial g_9}{\partial w_{11}^{(2)}} = -0.366$$
$$\frac{\partial L}{\partial g_7} = \frac{\partial L}{\partial g_9} \times \frac{\partial g_9}{\partial g_7} = 0.53$$

12/8/20

36

# Forward backprop example



$x_1$ 3

$g_1$ ×  0.6

$w_{11}^{(1)}$ 0.2

$x_2$ 2

$g_2$ ×  0.2

$w_{12}^{(1)}$ 0.1

$g_5$ +  0.8

$g_7$ σ  0.69
0.53

$w_{11}^{(2)}$ −1
−0.366

$g_9$ ×  −0.69
−0.53

$x_1$ 3

$g_3$ ×  −2.1

$w_{21}^{(1)}$ −0.7

$x_2$ 2

$g_4$ ×  2.4

$w_{22}^{(1)}$ 1.2

$g_6$ +  0.3

$g_8$ σ  0.57
−0.53

$g_{10}$ ×  0.57
−0.53

$w_{12}^{(2)}$ 1
−0.302

$g_{11}$ +  −0.12
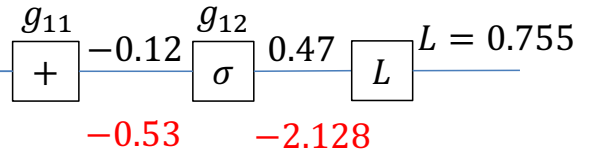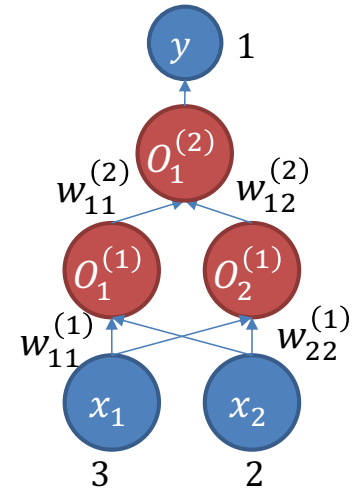−0.53

$g_{12}$ σ  0.47
−2.128

$L$   $L = 0.755$

$$\frac{\partial L}{\partial w_{12}^{(2)}} = \frac{\partial L}{\partial g_{10}} \times \frac{\partial g_{10}}{\partial w_{12}^{(2)}}$$
$$= -0.53 \times 0.57 = -0.302$$
$$\frac{\partial L}{\partial g_7} = \frac{\partial L}{\partial g_{10}} \times \frac{\partial g_{10}}{\partial g_8} = -0.53 \times 1$$
$$= -0.53$$

$y$  1

$O_1^{(2)}$

$w_{11}^{(2)}$     $w_{12}^{(2)}$

$O_1^{(1)}$     $O_2^{(1)}$

$w_{11}^{(1)}$     $w_{22}^{(1)}$

$x_1$     $x_2$

3     2

# Forward backprop example



$x_1$ 3

$g_1$

$w_{11}^{(1)}$ 0.2

$\times$ 0.6

$x_2$ 2

$g_2$

$w_{12}^{(1)}$ 0.1

$\times$ 0.2

$g_5$

$+$ 0.8

0.113

$g_7$

$\sigma$ 0.69

0.53

$w_{11}^{(2)}$ $-1$

$-0.366$

$g_9$

$\times$ $-0.69$

$-0.53$

$x_1$ 3

$g_3$

$w_{21}^{(1)}$ $-0.7$

$\times$ $-2.1$

$g_6$

$+$ 0.3

$-0.13$

$g_8$

$\sigma$ 0.57

$-0.53$

$x_2$ 2

$g_4$

$w_{22}^{(1)}$ 1.2

$\times$ 2.4

$g_{10}$

$\times$ 0.57

$-0.53$

$w_{12}^{(2)}$ 1

$-0.302$

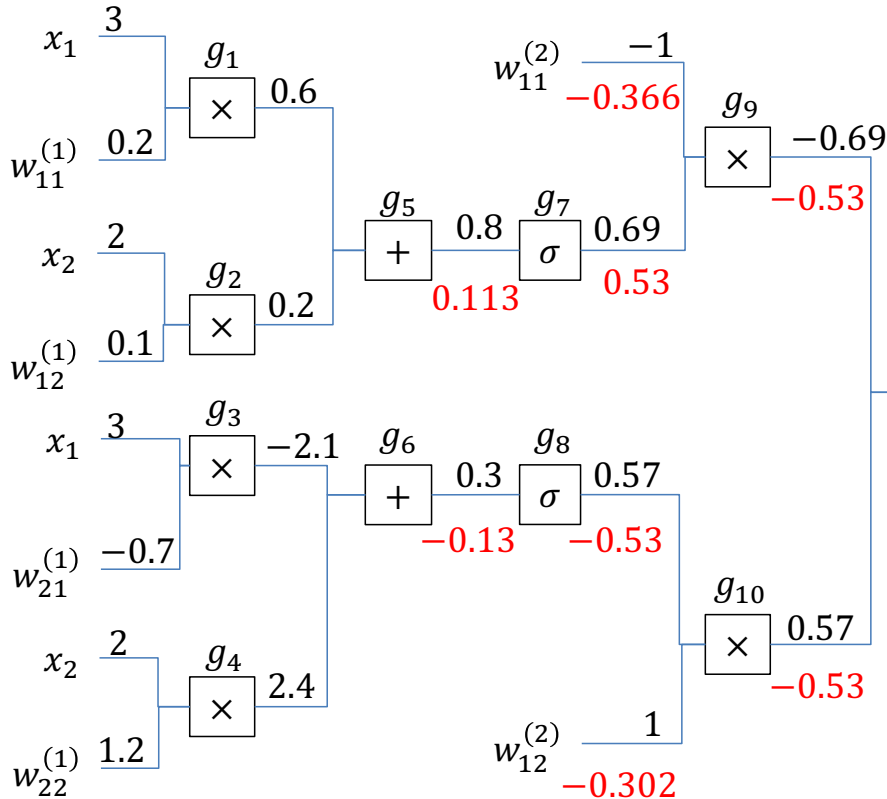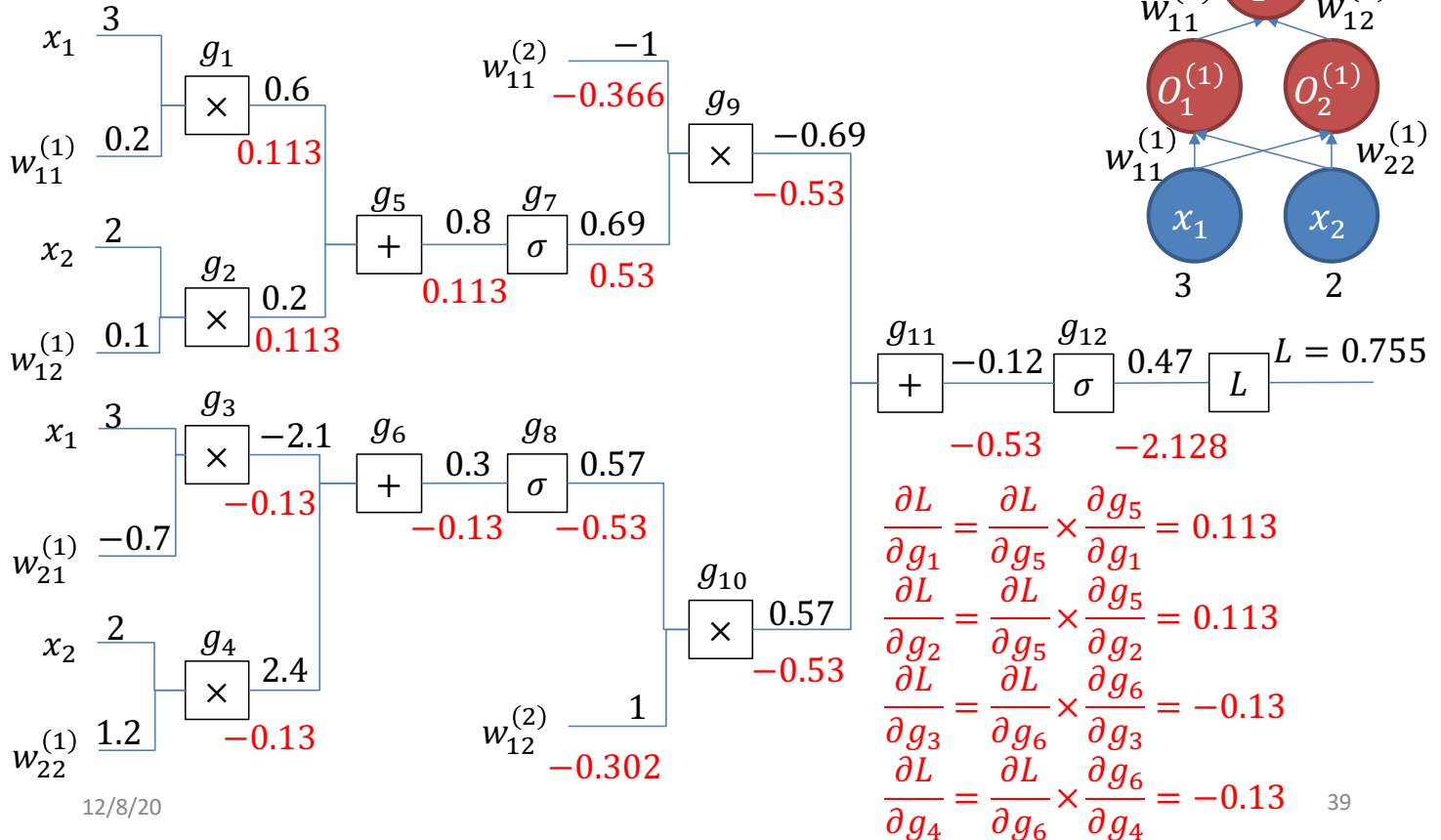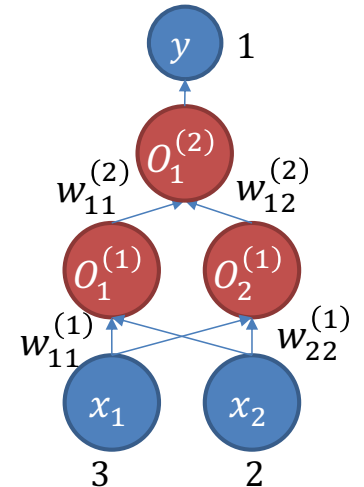$g_{11}$

$+$ $-0.12$
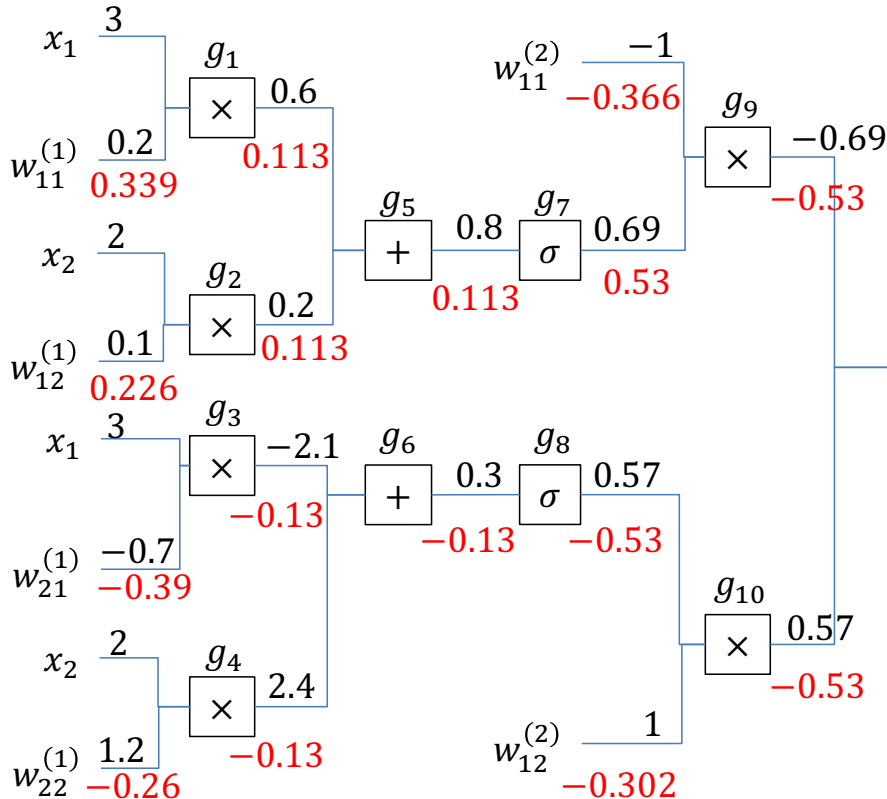
$-0.53$

$g_{12}$

$\sigma$ 0.47

$-2.128$

$L$

$L = 0.755$

$$\frac{\partial L}{\partial g_5} = \frac{\partial L}{\partial g_7} \times \frac{\partial g_7}{\partial g_5}$$
$$= 0.53 \times \big(0.69(1 - 0.69)\big) = 0.113$$
$$\frac{\partial L}{\partial g_6} = \frac{\partial L}{\partial g_8} \times \frac{\partial g_8}{\partial g_6}$$
$$= -0.53 \times \big(0.57(1 - 0.57)\big)$$
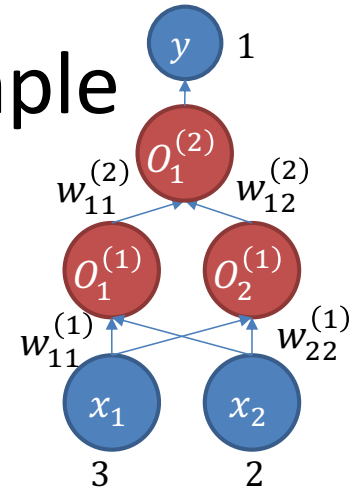$$= -0.13$$

# Forward backprop example



$x_1$ 3

$g_1$ 0.6

$w_{11}^{(1)}$ 0.2

0.113

$x_2$ 2

$g_2$ 0.2

$w_{12}^{(1)}$ 0.1

0.113

$w_{11}^{(2)}$ −1

−0.366

$g_5$ + 0.8

$g_7$ σ 0.69

0.113

0.53

$g_9$ × −0.69

−0.53

$x_1$ 3

$g_3$ × −2.1

−0.13

$w_{21}^{(1)}$ −0.7

$g_6$ + 0.3

−0.13

$g_8$ σ 0.57

−0.53

$x_2$ 2

$g_4$ × 2.4

$w_{22}^{(1)}$ 1.2

−0.13

$g_{10}$ × 0.57

−0.53

$w_{12}^{(2)}$ 1

−0.302

$g_{11}$ + −0.12

−0.53

$g_{12}$ σ 0.47

−2.128

$L$   $L = 0.755$

$y$   1

$O_1^{(2)}$

$w_{11}^{(2)}$   $w_{12}^{(2)}$

$O_1^{(1)}$   $O_2^{(1)}$

$w_{11}^{(1)}$   $w_{22}^{(1)}$

$x_1$   $x_2$

3   2

$$\frac{\partial L}{\partial g_1} = \frac{\partial L}{\partial g_5} \times \frac{\partial g_5}{\partial g_1} = 0.113$$

$$\frac{\partial L}{\partial g_2} = \frac{\partial L}{\partial g_5} \times \frac{\partial g_5}{\partial g_2} = 0.113$$

$$\frac{\partial L}{\partial g_3} = \frac{\partial L}{\partial g_6} \times \frac{\partial g_6}{\partial g_3} = -0.13$$

$$\frac{\partial L}{\partial g_4} = \frac{\partial L}{\partial g_6} \times \frac{\partial g_6}{\partial g_4} = -0.13$$

12/8/20

39

# Forward backprop example



$x_1$ 3

$g_1$ 0.6
× 0.113

$w_{11}^{(1)}$ 0.2
0.339

$x_2$ 2

$g_2$ 0.2
× 0.113

$w_{12}^{(1)}$ 0.1
0.226

$w_{11}^{(2)}$ −1
−0.366

$g_5$ 0.8
+ 0.113

$g_7$ 0.69
σ 0.53

$g_9$ −0.69
× −0.53

$g_{11}$ −0.12
+ −0.53

$g_{12}$ 0.47
σ −2.128

$L$ = 0.755

$x_1$ 3

$g_3$ −2.1
× −0.13

$w_{21}^{(1)}$ −0.7
−0.39

$g_6$ 0.3
+ −0.13

$g_8$ 0.57
σ −0.53

$x_2$ 2

$g_4$ 2.4
× −0.13

$w_{22}^{(1)}$ 1.2
−0.26

$g_{10}$ 0.57
× −0.53

$w_{12}^{(2)}$ 1
−0.302

$y$ 1

$O_1^{(2)}$

$w_{11}^{(2)}$  $w_{12}^{(2)}$

$O_1^{(1)}$  $O_2^{(1)}$

$w_{11}^{(1)}$  $w_{22}^{(1)}$

$x_1$  $x_2$

3  2

$$\frac{\partial L}{\partial w_{11}^{(1)}} = \frac{\partial L}{\partial g_1} \times \frac{\partial g_1}{\partial w_{11}^{(1)}} = 0.339$$

$$\frac{\partial L}{\partial w_{12}^{(1)}} = \frac{\partial L}{\partial g_2} \times \frac{\partial g_2}{\partial w_{12}^{(1)}} = 0.226$$

$$\frac{\partial L}{\partial w_{21}^{(1)}} = \frac{\partial L}{\partial g_3} \times \frac{\partial g_3}{\partial w_{21}^{(1)}} = -0.39$$

$$\frac{\partial L}{\partial w_{22}^{(1)}} = \frac{\partial L}{\partial g_4} \times \frac{\partial g_4}{\partial w_{22}^{(1)}} = -0.26$$

# Forward backprop example



- Use original $\boldsymbol{w}^{(1)}, \boldsymbol{w}^{(2)}$:
  - $\hat{y} = 0.4711$
- Use the new $\boldsymbol{w}^{(1)}, \boldsymbol{w}^{(2)}$ computed by gradient descent ($\eta = 0.001$)
  - $\boldsymbol{w}^{(\ell)} = \boldsymbol{w}^{(\ell)} - \eta \nabla_{\boldsymbol{w}^{(\ell)}} L$
  - $\hat{y} = 0.4714$
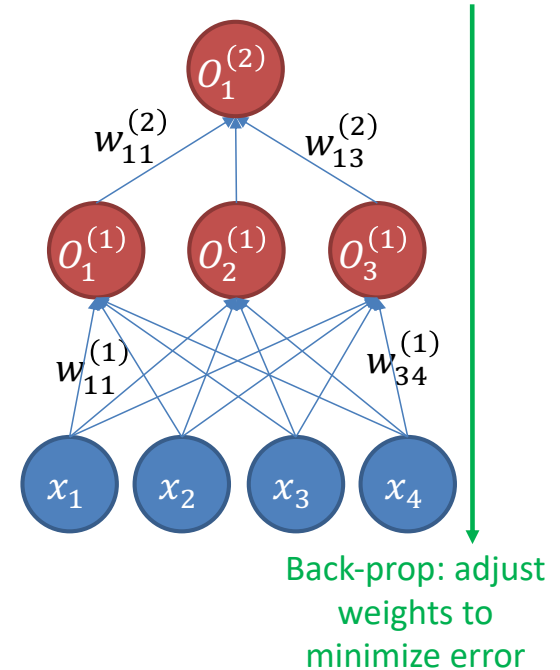- The new $\hat{y}$ indeed closer to 1

# Back-prop update



- Back-prop

  - $$\frac{\partial loss}{\partial I^{(2)}} = \frac{\partial loss}{\partial O^{(2)}} \frac{\partial O^{(2)}}{\partial I^{(2)}}$$

  - $$\frac{\partial loss}{\partial W^{(2)}} = \frac{\partial loss}{\partial I^{(2)}} \frac{\partial I^{(2)}}{\partial W^{(2)}}$$

  - $$\frac{\partial loss}{\partial I^{(1)}} = \frac{\partial loss}{\partial I^{(2)}} \frac{\partial I^{(2)}}{\partial O^{(1)}} \frac{\partial O^{(1)}}{\partial I^{(1)}}$$

  - $$\frac{\partial loss}{\partial W^{(1)}} = \frac{\partial loss}{\partial I^{(1)}} \frac{\partial I^{(1)}}{\partial W^{(1)}}$$

Back-prop: adjust weights to minimize error

Red: computed in higher layers
Blue: depends on the activation function
Green: depends on the loss function

- $I^{(1)} = W^{(1)} x$
- $O^{(1)} = \sigma(I^{(1)})$
- $I^{(2)} = W^{(2)} O^{(1)}$
- $O^{(2)} = \sigma(I^{(2)})$

# Weight updates

- Update $\boldsymbol{W}^{(\ell)}$ by (stochastic) gradient descent

  - Compute $\frac{\partial loss}{\partial \boldsymbol{W}^{(\ell)}}$ for all $\boldsymbol{W}^{(\ell)}$

  - $\boldsymbol{W}^{(\ell)} \leftarrow \boldsymbol{W}^{(\ell)} - \eta \frac{\partial loss}{\partial \boldsymbol{W}^{(\ell)}}$

# Loss functions

- K-nary multi-class classification: cross entropy loss
  - $Loss = -\sum_{k=1}^{K} I(y_k = 1) \log \hat{y}$
    - Example: 3 classes
      - $p(\hat{y}_i = k) = [1/4 \quad 1/4 \quad 1/2]$
      - $p(y_i = k) = [0 \quad 1 \quad 0]$
      - $-\sum_k p(y_i = k) \log[p(\hat{y}_i = k)] = -\left[0 \log\frac{1}{4} + 1 \log\frac{1}{4} + 0 \log\frac{1}{2}\right] = 2$
  - When $K = 2 \Rightarrow$ back to binary classification ($y \in \{0,1\}$)
    - $Loss = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$
- Regression
  - $Loss = \frac{1}{2}(y - \hat{y})^2$
- Multi-label classification (assuming $L$ possible labels $(\ell_1, \ldots, \ell_L)$ for each instance)
  - $Loss = \sum_{i=1}^{L}\left(-\ell_i \log \hat{\ell}_i - (1 - \ell_i) \log(1 - \hat{\ell}_i)\right)$

# Common activation functions

- Logistic (sigmoid) function
  - $f(z) = \frac{1}{1+\exp(-z)}$
- Hyperbolic (tanh) function
  - $f(z) = \frac{2}{1+\exp(-2z)} - 1$
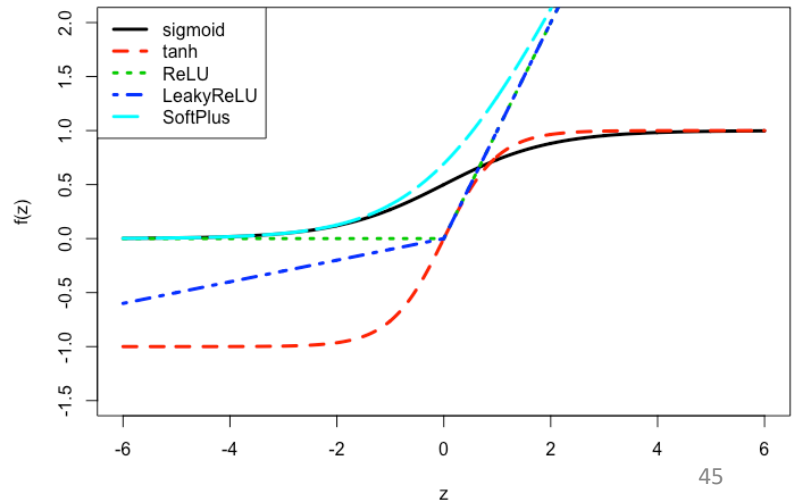- Rectified Linear Unit (ReLU)
  - $f(z) = \max(0, z)$
- Soft-plus
  - $f(z) = \log(1 + \exp(z))$

- Leaky ReLU
  - $f(z) = \begin{cases} z \text{ if } z \geq 0 \\ \alpha z \text{ if } z < 0 \end{cases}$
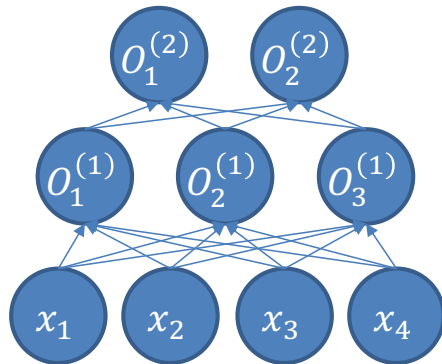    - $\alpha$ is small and positive (e.g., 0.01)

# Activation function for the output layer

- For the final layer, select the activation function based on the type of your task
  - Regression: return one value based on linear activation function
    - $f(z) = z$
  - $K$-ary multi-class classification: return a vector based on softmax function
    - $f(\mathbf{z}) = \left[\frac{e^{z_1}}{\sum_j e^{z_j}}, \frac{e^{z_2}}{\sum_j e^{z_j}}, \dots, \frac{e^{z_K}}{\sum_j e^{z_j}}\right]$    *or softmax*
      - $\text{Sum}(f(\mathbf{z})) = 1$
  - Multi-label classification: return a vector based on sigmoid function
    - $f(\mathbf{z}) = \left[\frac{1}{1+\exp(-z_1)}, \frac{1}{1+\exp(-z_2)}, \dots, \frac{1}{1+\exp(-z_K)}\right]$
      - $\text{Sum}(f(\mathbf{z})) \neq 1$

# Why non-linear activation function?

- If no (non-linear) activation function, multi-layer reduces to single layer

$O_1^{(2)}$ $O_2^{(2)}$

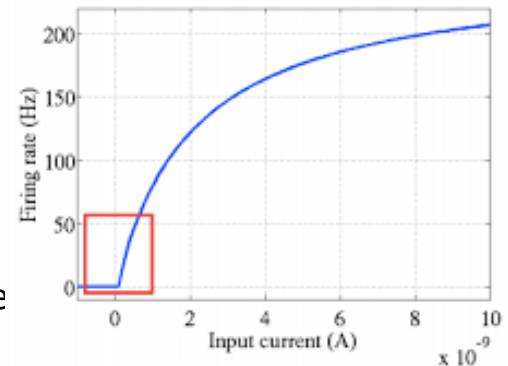$O_1^{(1)}$ $O_2^{(1)}$ $O_3^{(1)}$

$x_1$ $x_2$ $x_3$ $x_4$

$W^{(2)} \in R^{2\times3}$

$W^{(1)} \in R^{3\times4}$

$$\boldsymbol{O}^{(2)} = \boldsymbol{W}^{(2)}\boldsymbol{O}^{(1)}$$
$$= \boldsymbol{W}^{(2)}\left(\boldsymbol{W}^{(1)}\boldsymbol{x}\right)$$
$$= \left(\boldsymbol{W}^{(2)}\boldsymbol{W}^{(1)}\right)\boldsymbol{x}$$

$$\Rightarrow \left(\boldsymbol{W}^{(2)}\boldsymbol{W}^{(1)}\right) \in R^{2\times4}$$
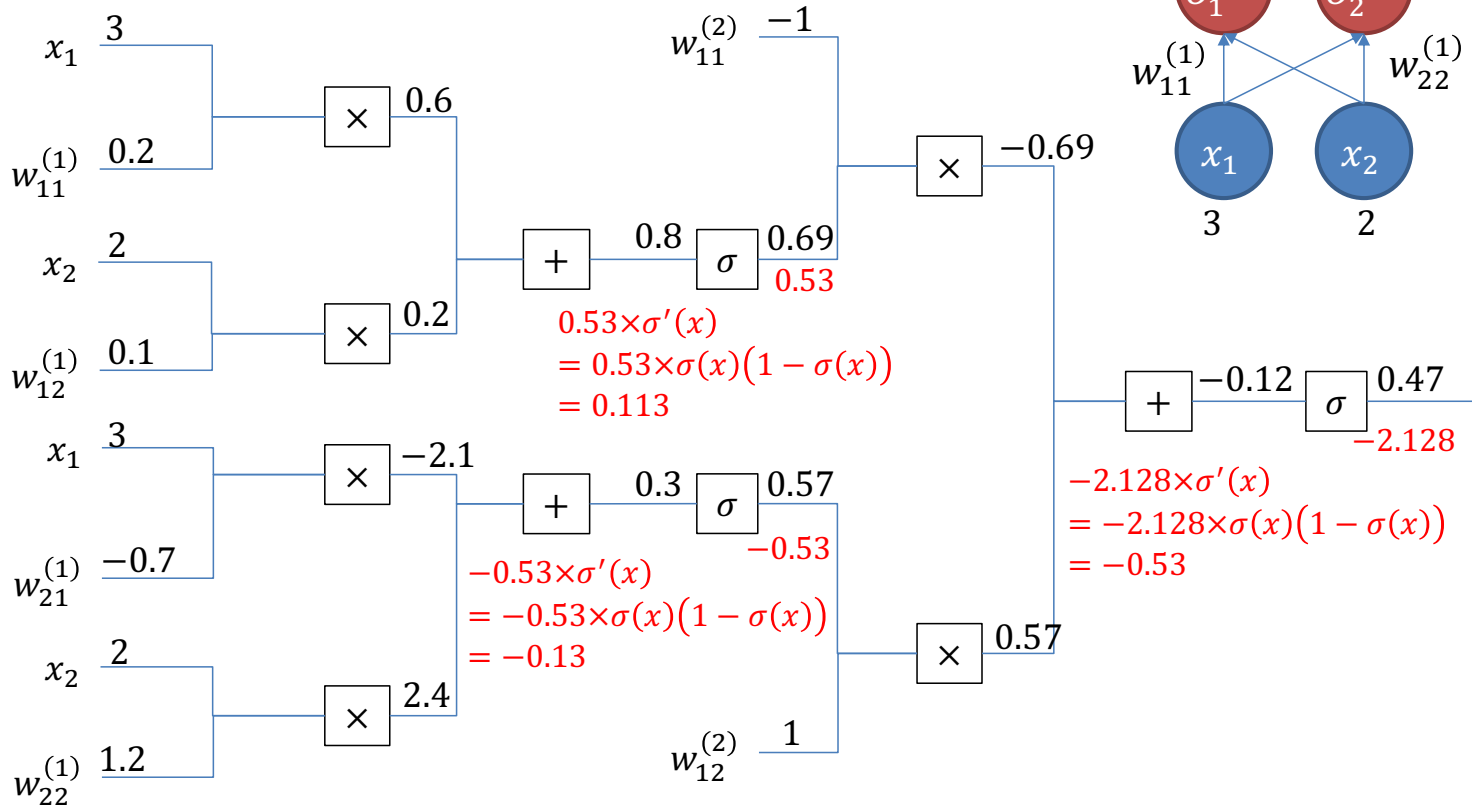
# ReLU: "default" activation function for the hidden layers

- Vanishing gradient problem in sigmoid and tanh
  - When $z$ is outside $[-4, 4]$, $f'_{sigmoid}(z) \approx 0$ and $f'_{tanh}(z) \approx 0$ ➔ new information cannot be back-propagated
- ReLU's sparse neuron outputs may prevent overfitting
  - When $z < 0$, $f_{ReLU}(z) = 0$ ➔ many neuron outputs are zero ➔ model becomes smaller
- ReLU's computation cost is small
  - No exponential computation
- ReLU's may resemble biological neurons
  - Firing rate > 0 only when input current is large enough (see the figure)
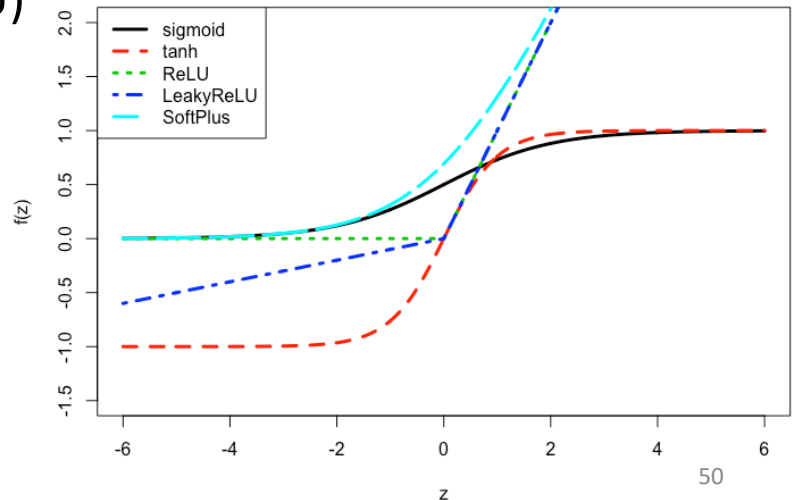
# Vanishing gradient of sigmoid function



$x_1$ 3

$\times$ 0.6

$w_{11}^{(1)}$ 0.2

$x_2$ 2

$w_{12}^{(1)}$ 0.1

$\times$ 0.2

$+$ 0.8 $\sigma$ 0.69

0.53

$0.53 \times \sigma'(x)$
$= 0.53 \times \sigma(x)(1 - \sigma(x))$
$= 0.113$

$w_{11}^{(2)}$ $-1$

$\times$ $-0.69$

$x_1$ 3

$\times$ $-2.1$

$w_{21}^{(1)}$ $-0.7$

$x_2$ 2

$w_{22}^{(1)}$ 1.2

$\times$ 2.4

$+$ 0.3 $\sigma$ 0.57

$-0.53$

$-0.53 \times \sigma'(x)$
$= -0.53 \times \sigma(x)(1 - \sigma(x))$
$= -0.13$

$w_{12}^{(2)}$ 1

$\times$ 0.57

$+$ $-0.12$ $\sigma$ 0.47

$-2.128$

$-2.128 \times \sigma'(x)$
$= -2.128 \times \sigma(x)(1 - \sigma(x))$
$= -0.53$

Network diagram:

$O_1^{(2)}$

$w_{11}^{(2)}$ $w_{12}^{(2)}$

$O_1^{(1)}$ $O_2^{(1)}$

$w_{11}^{(1)}$ $w_{22}^{(1)}$

$x_1$ $x_2$

3 2

# Derivative of activation functions

- Logistic (sigmoid) function
  - $f'(z) = f(z)\big(1 - f(z)\big)$
- Hyperbolic (tanh) function
  - $f'^{(z)} = 1 - \big(f(z)\big)^2$
- Rectified Linear Unit (ReLU)
  - $f(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{else} \end{cases}$
- Soft-plus
  - $f'(z) = \dfrac{1}{1+\exp(-z)}$

- Leaky ReLU
  - $f'(z) = \begin{cases} 1 \text{ if } z \geq 0 \\ \alpha \text{ if } z < 0 \end{cases}$
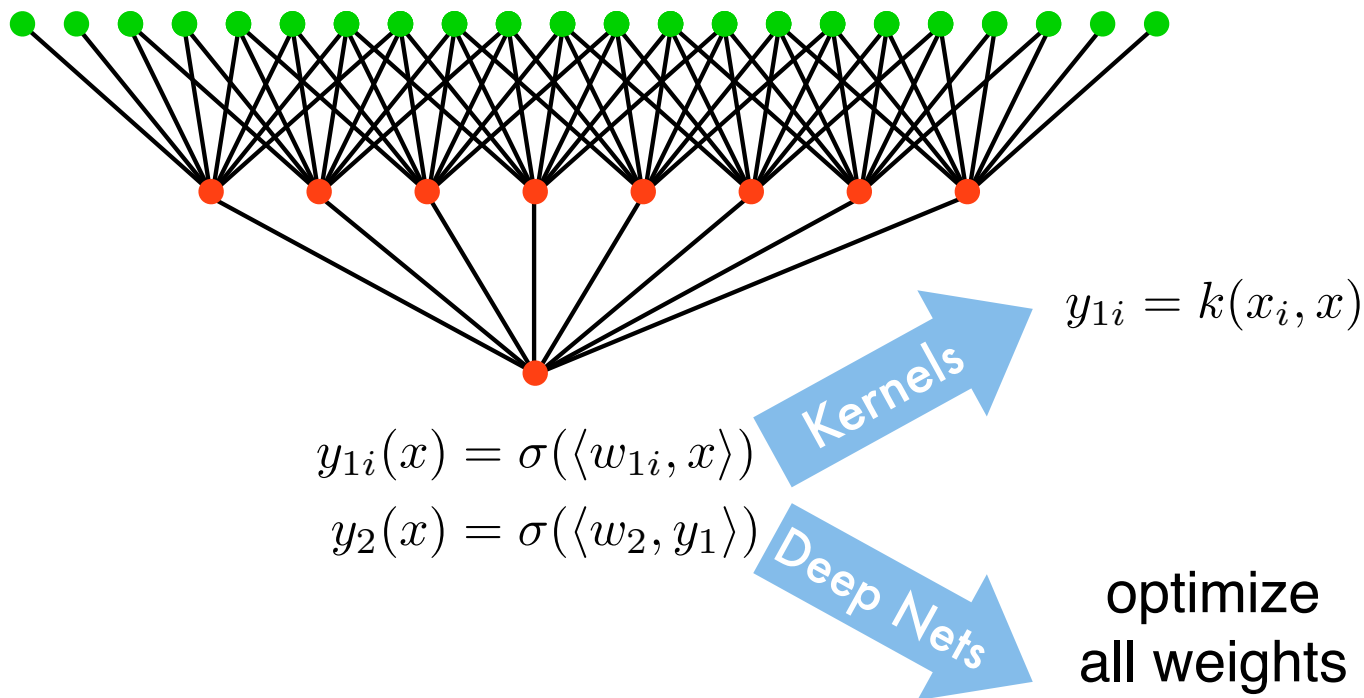    - $\alpha$ is small and positive (e.g., 0.01)

# Example: write your own activation operation

- Define autograd function (in PyTorch) by writing forward and backward for Tensors
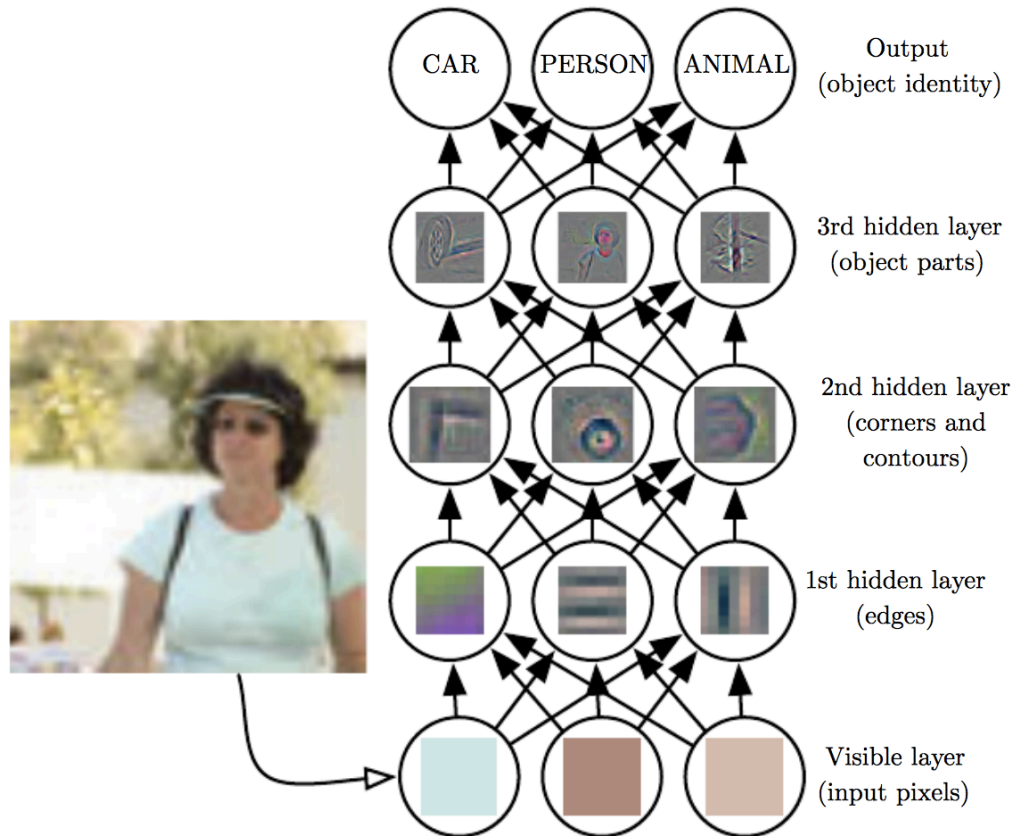
```python
class ReLU(torch.autograd.Function):
    def forward(self, x):
        self.save_for_backward(x)
        return x.clamp(min=0)

    def backward(self, grad_y):
        x, = self.saved_tensors
        grad_input = grad_y.clone()
        grad_input[x < 0] = 0
        return grad_input
```

- You may define your customized activation functions

- PyTorch performs chain rules for you if all operators on the graph are well-defined:
  - `loss.backward()`

# Nonlinearities via Layers



$$y_{1i} = k(x_i, x)$$

Kernels

$$y_{1i}(x) = \sigma(\langle w_{1i}, x \rangle)$$
$$y_2(x) = \sigma(\langle w_2, y_1 \rangle)$$

Deep Nets

optimize
all weights

Carnegie Mellon University

# Learning representations

# Popular deep learning framework

- Most popular
  - Tensorflow (Keras), PyTorch
- Others
  - MXNet, Caffe and Caffe2, Theano, Torch

# Other successful deep learning architecture

- Convolutional Neural Network (CNN)
- Recurrent Neural Network (RNN)

# Conclusion

- Supervised DNN is composed by many hidden layers
  - Learning representations
  - Interaction between features to form high-dimensional features
  - Non-linearity property
- Training a supervised DNN is based on SGD
  - View few (e.g., 32) instances in on batch
  - Chain rule and backpropagation are the tricks to perform derivative