

Backpropagation on the loss function

Kietikul Jearanaitanakij

Department of Computer Engineering, KMITL

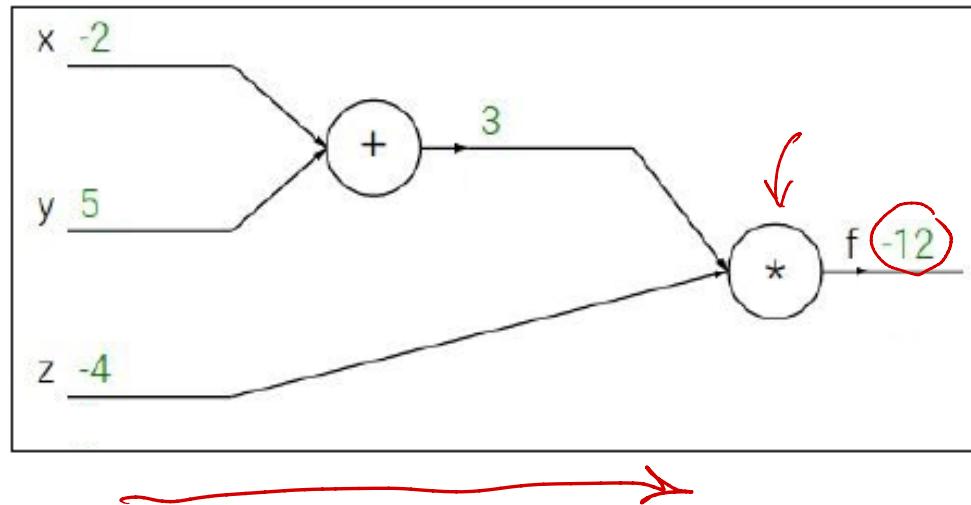
(Slides are adapted from cs231n @Stanford University)

Backpropagation(Revisited)

Backpropagation: a simple example

$$f(x, y, z) = (\underline{x} + \underline{y})z$$

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



Forward Computation

Backpropagation: a simple example

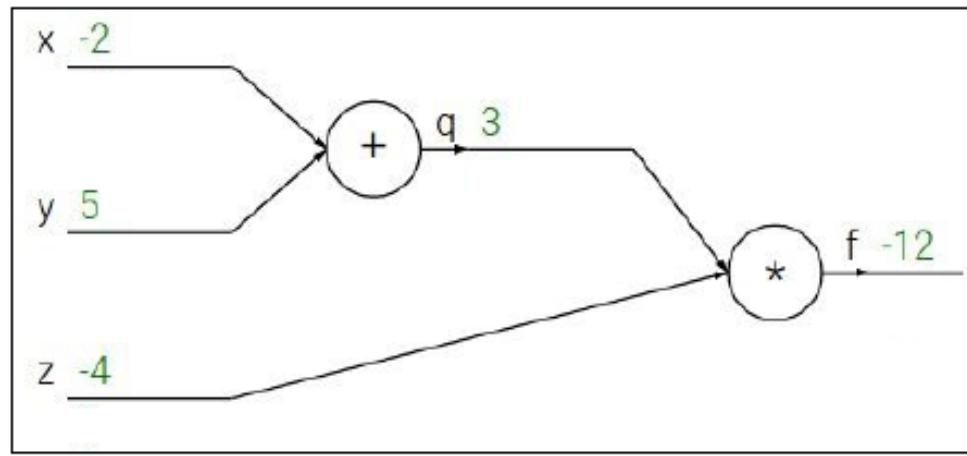
$$f(x, y, z) = (x + y)z$$

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

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

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

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



Backpropagation: a simple example

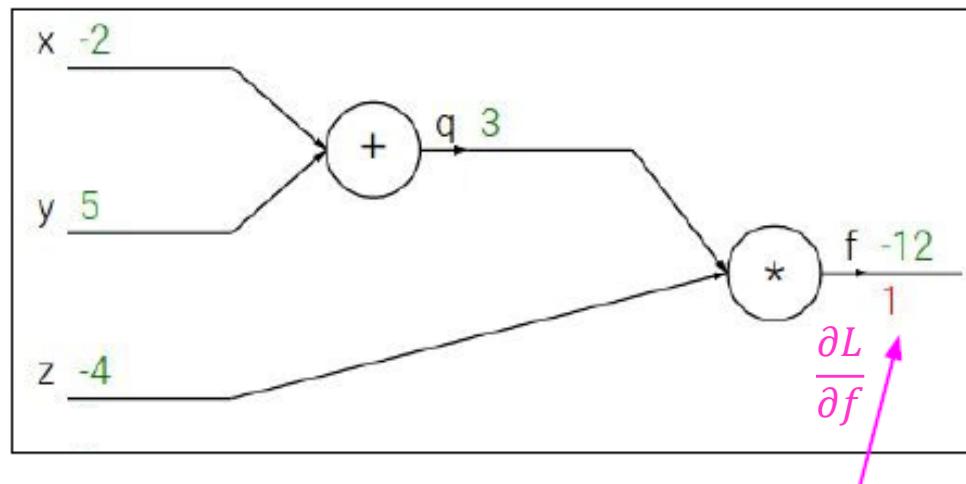
$$f(x, y, z) = (x + y)z$$

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

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

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

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



Suppose
gradient of
Loss wrt f
equals to 1

Backpropagation: a simple example

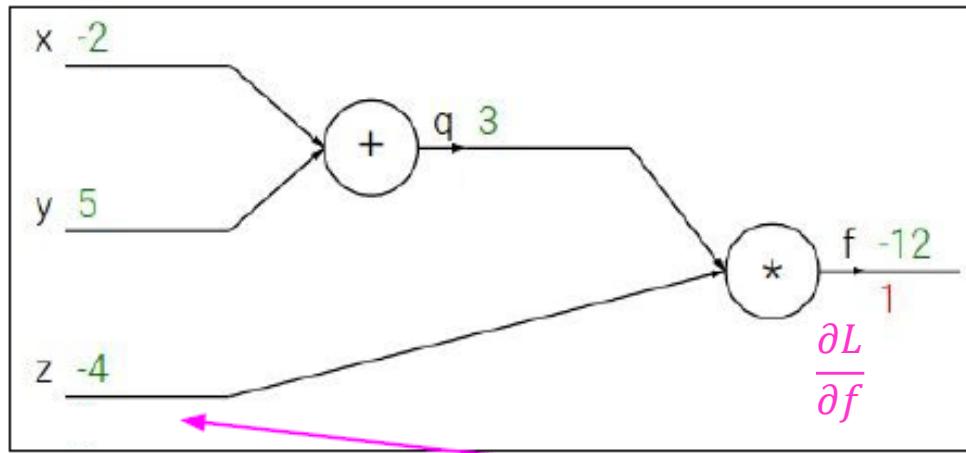
$$f(x, y, z) = (x + y)z$$

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

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

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

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



$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial z}$$

Chain rule:

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

Backpropagation: a simple example

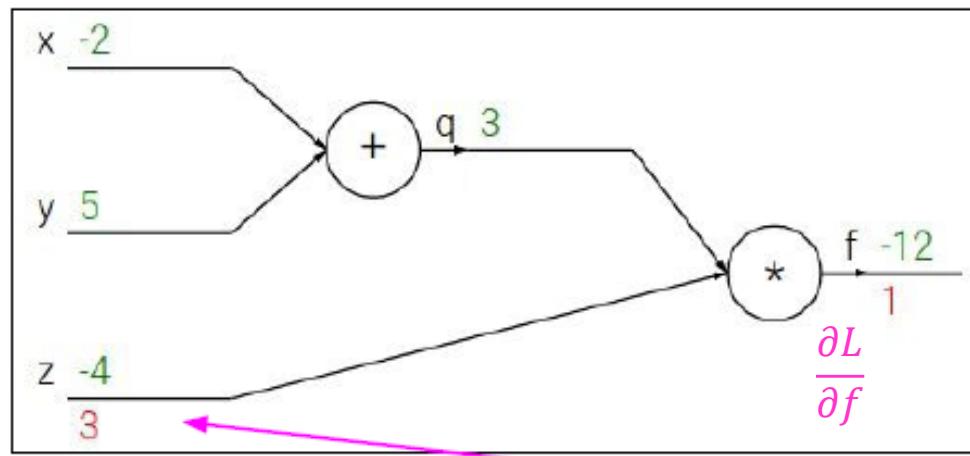
$$f(x, y, z) = (x + y)z$$

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

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

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

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



$$\frac{\partial L}{\partial z} = 1 \cdot \frac{\partial f}{\partial z}$$

Backpropagation: a simple example

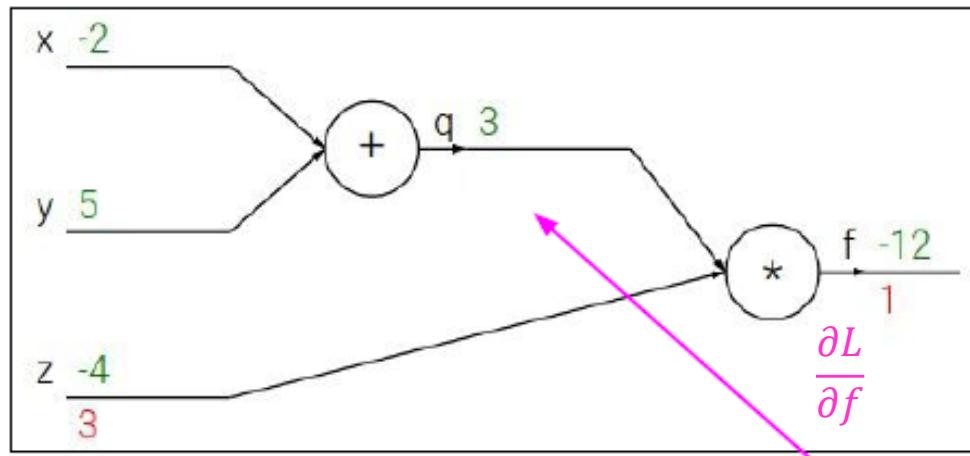
$$f(x, y, z) = (x + y)z$$

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

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

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

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



$$\frac{\partial L}{\partial q} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial q}$$

Backpropagation: a simple example

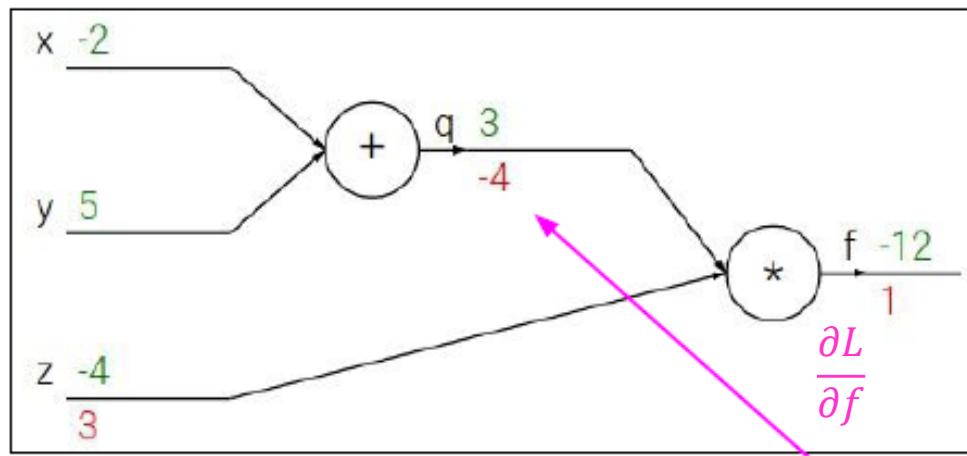
$$f(x, y, z) = (x + y)z$$

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

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

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

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



$$\frac{\partial L}{\partial q} = 1 \cdot \frac{\partial f}{\partial q}$$

Backpropagation: a simple example

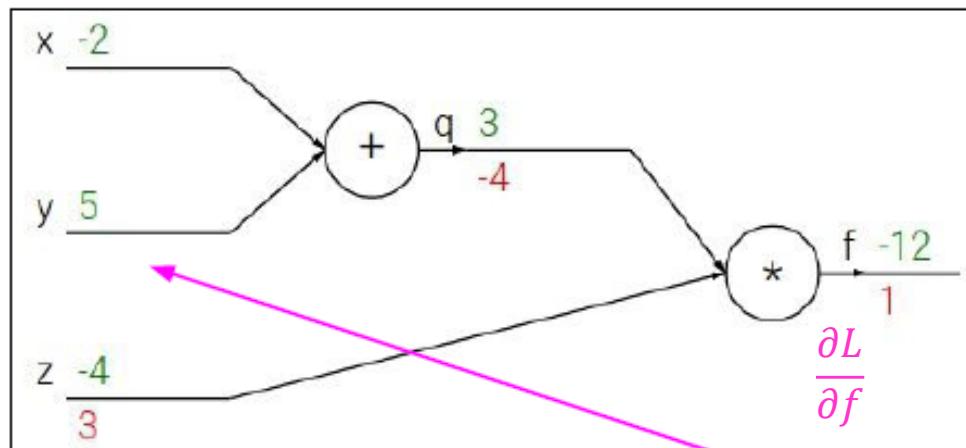
$$f(x, y, z) = (x + y)z$$

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

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

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

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



$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial q} \cdot \frac{\partial q}{\partial y}$$

Backpropagation: a simple example

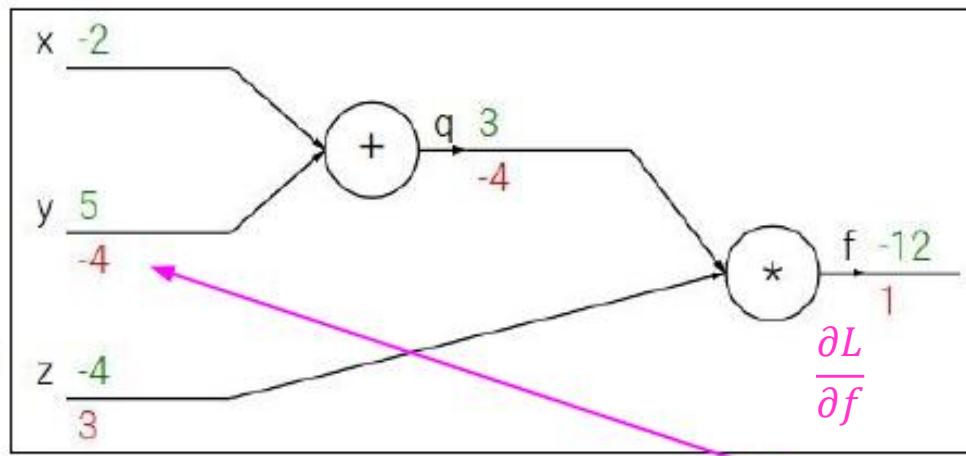
$$f(x, y, z) = (x + y)z$$

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

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

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

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



$$\frac{\partial L}{\partial y} = -4 \cdot \frac{\partial q}{\partial y}$$

Backpropagation: a simple example

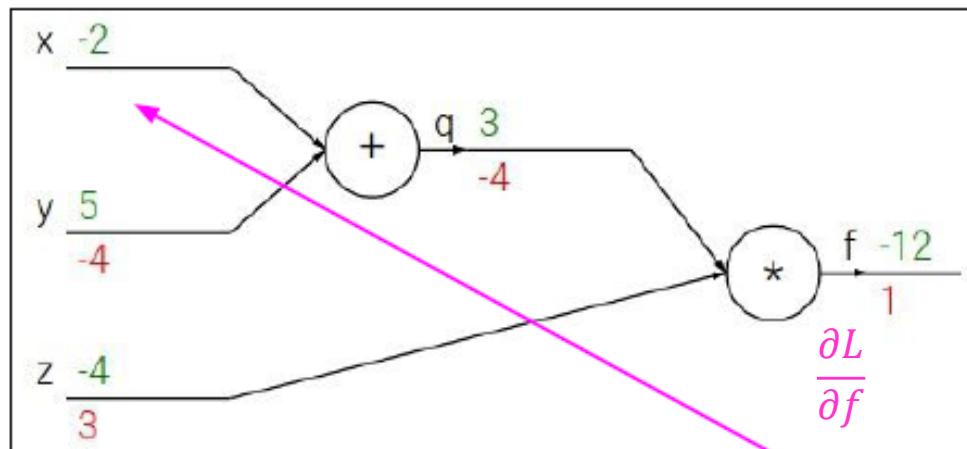
$$f(x, y, z) = (x + y)z$$

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

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

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

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



$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial q} \cdot \frac{\partial q}{\partial x}$$

Backpropagation: a simple example

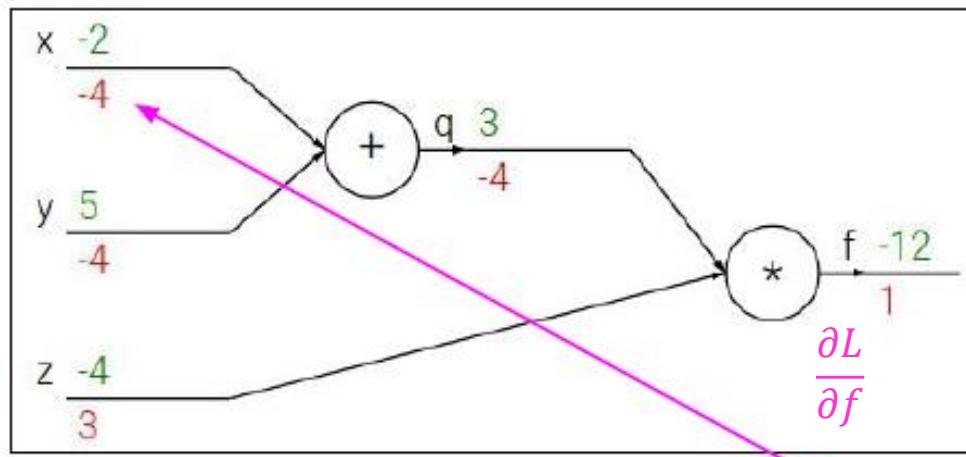
$$f(x, y, z) = (x + y)z$$

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

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

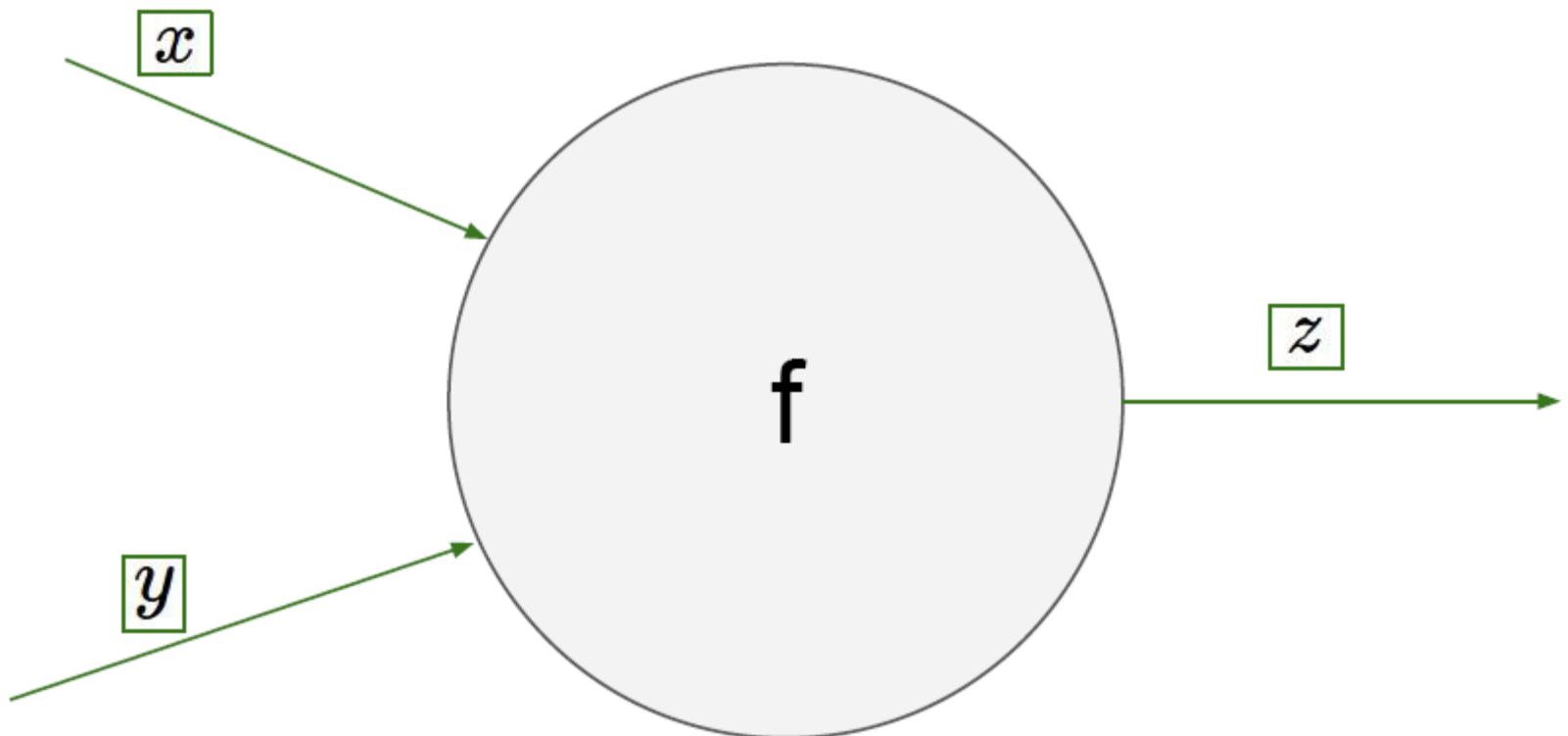
$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

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

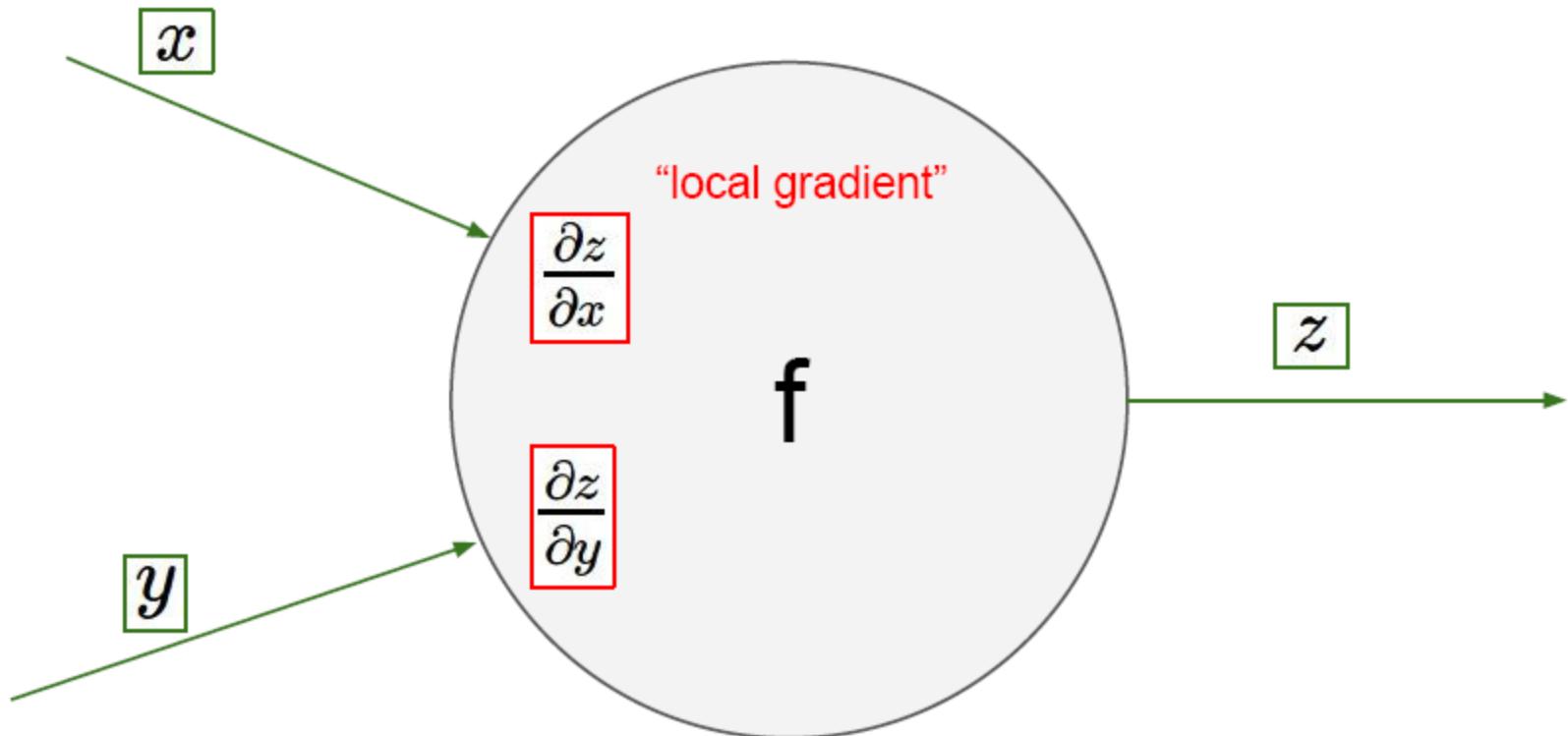


$$\frac{\partial L}{\partial x} = -4 \cdot \frac{\partial q}{\partial x}$$

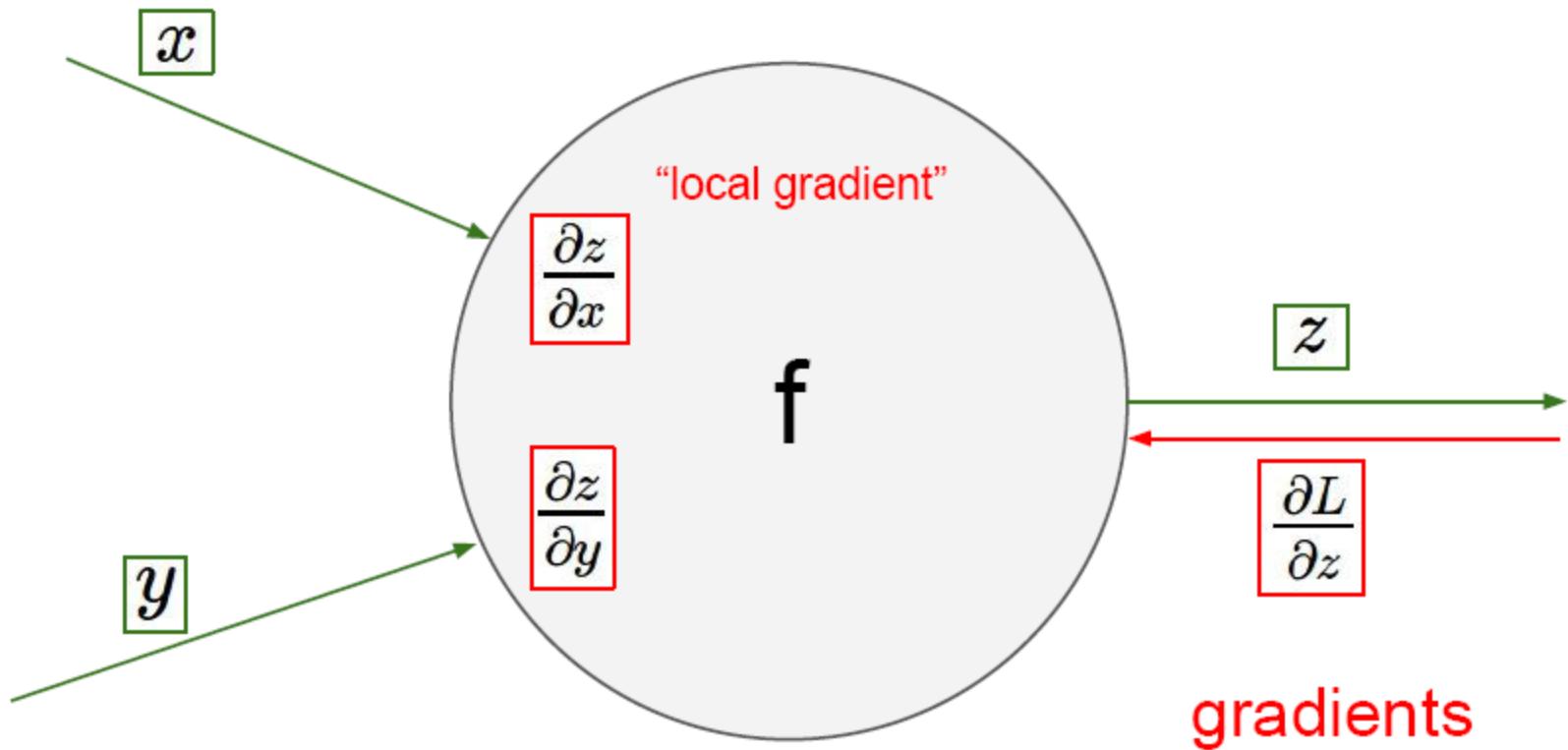
Backpropagation of gradient in general



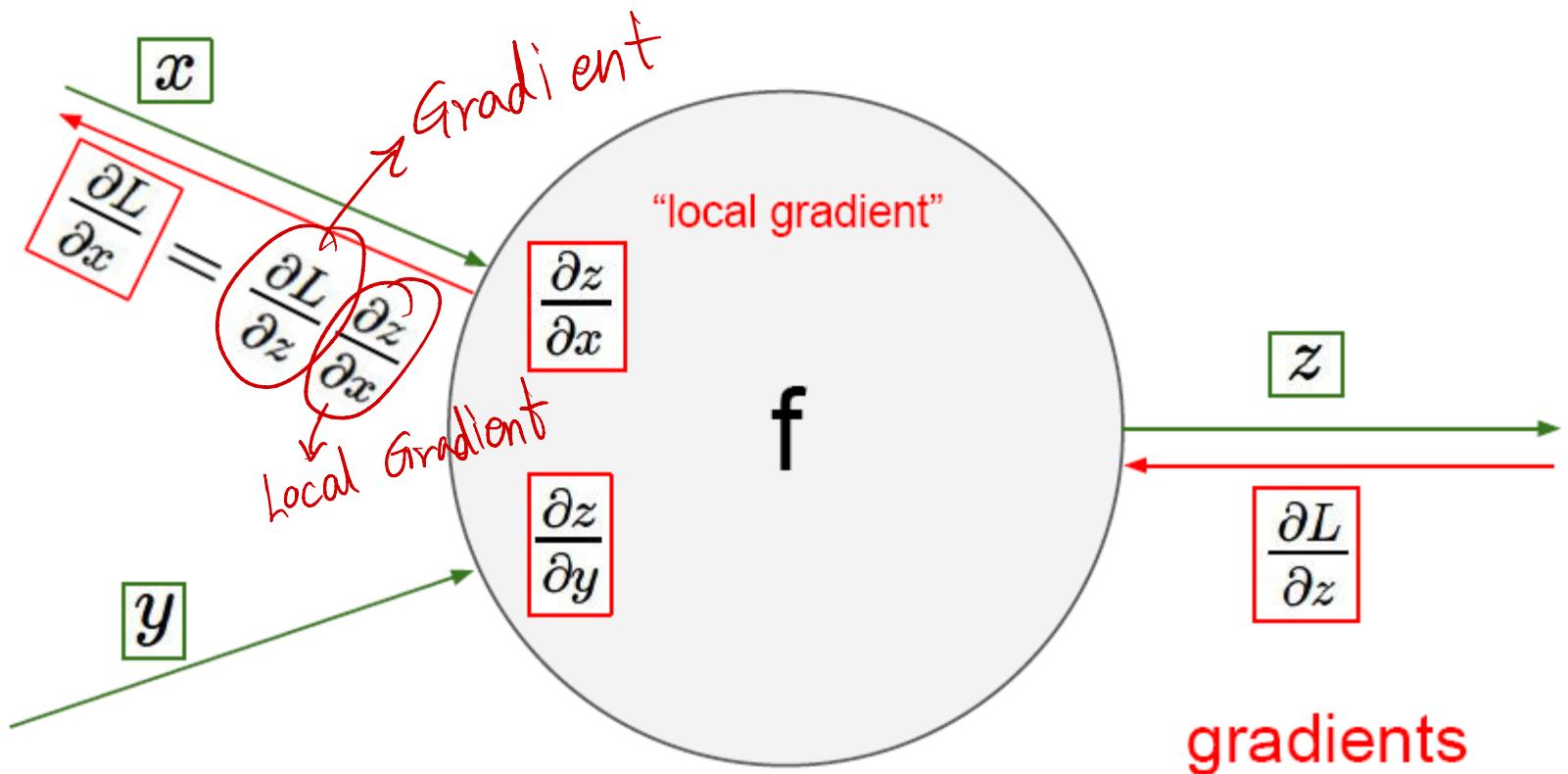
Backpropagation of gradient in general



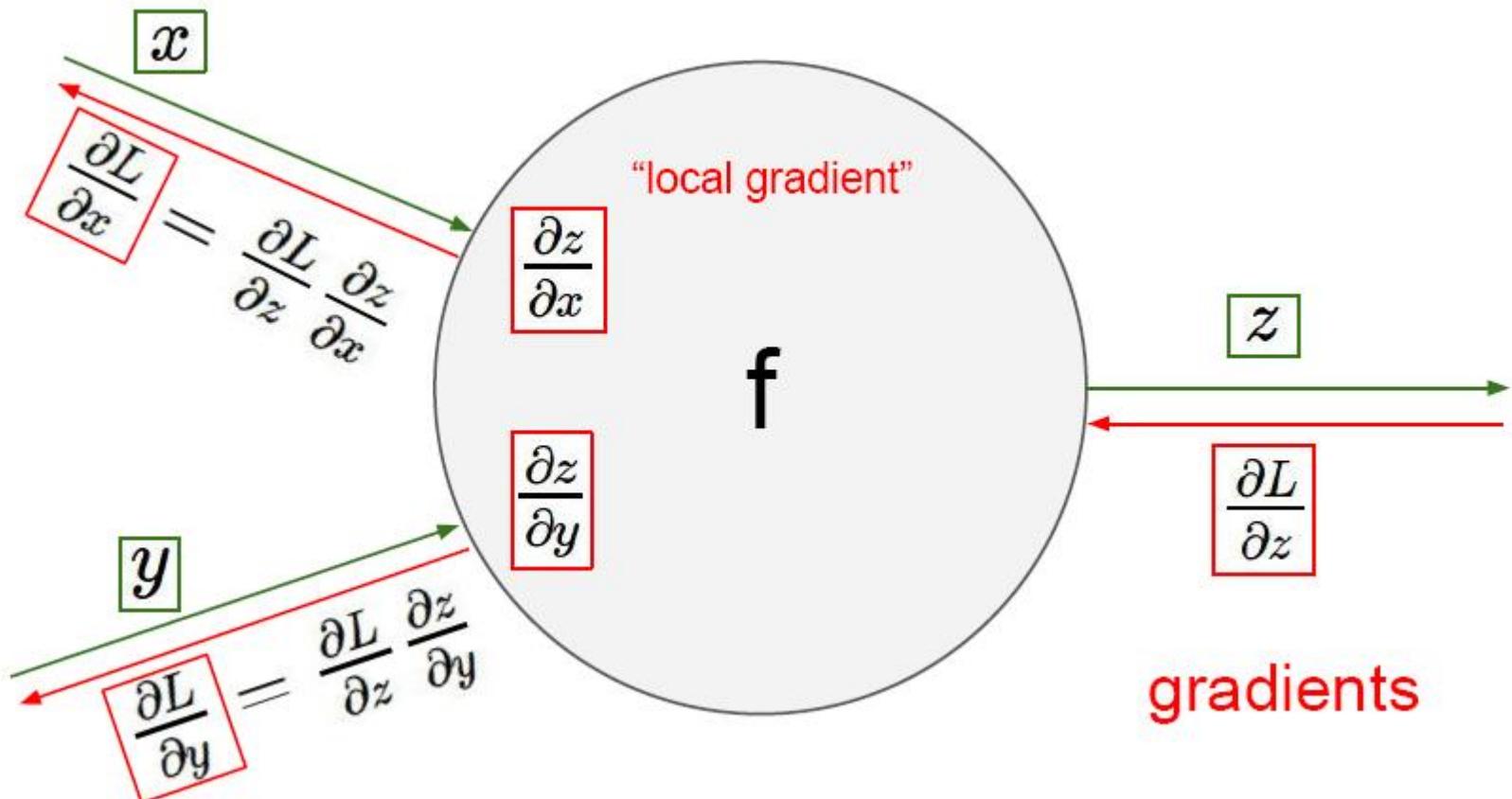
Backpropagation of gradient in general



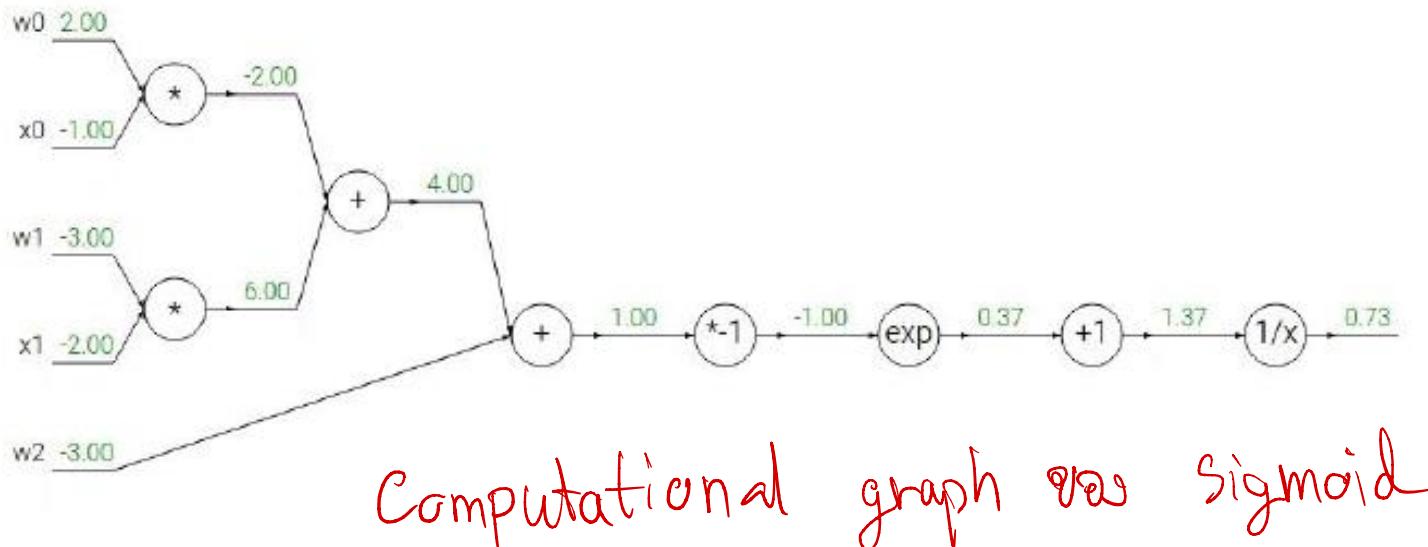
Backpropagation of gradient in general



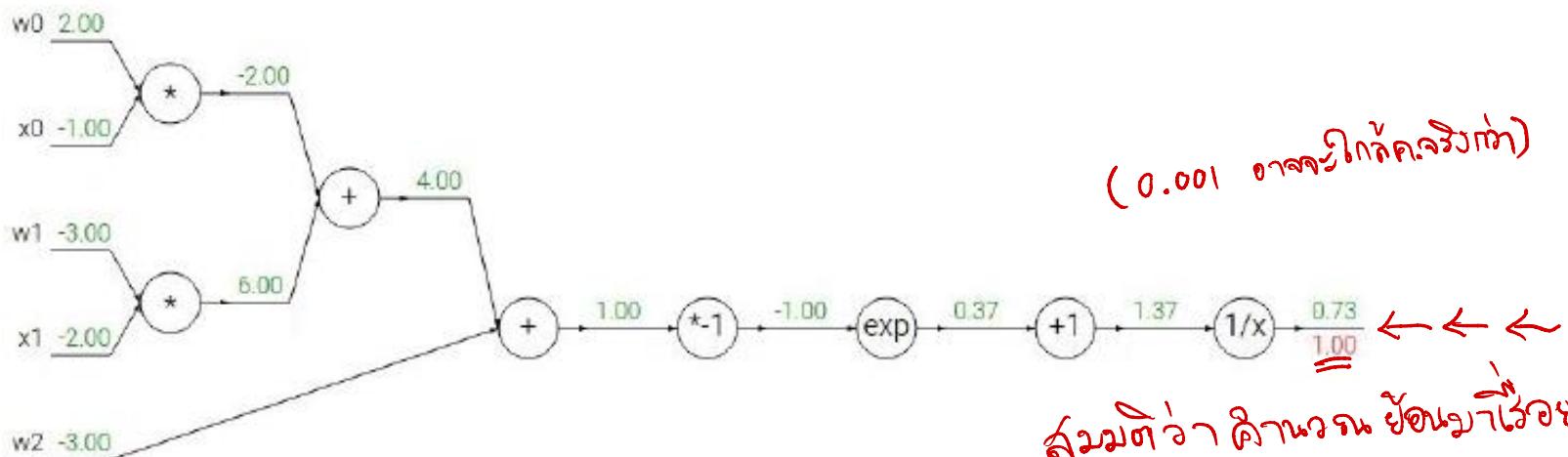
Backpropagation of gradient in general



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



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



$$f(x) = e^x$$

→

$$\frac{df}{dx} = e^x$$

$$f_a(x) = ax$$

→

$$\frac{df}{dx} = a$$

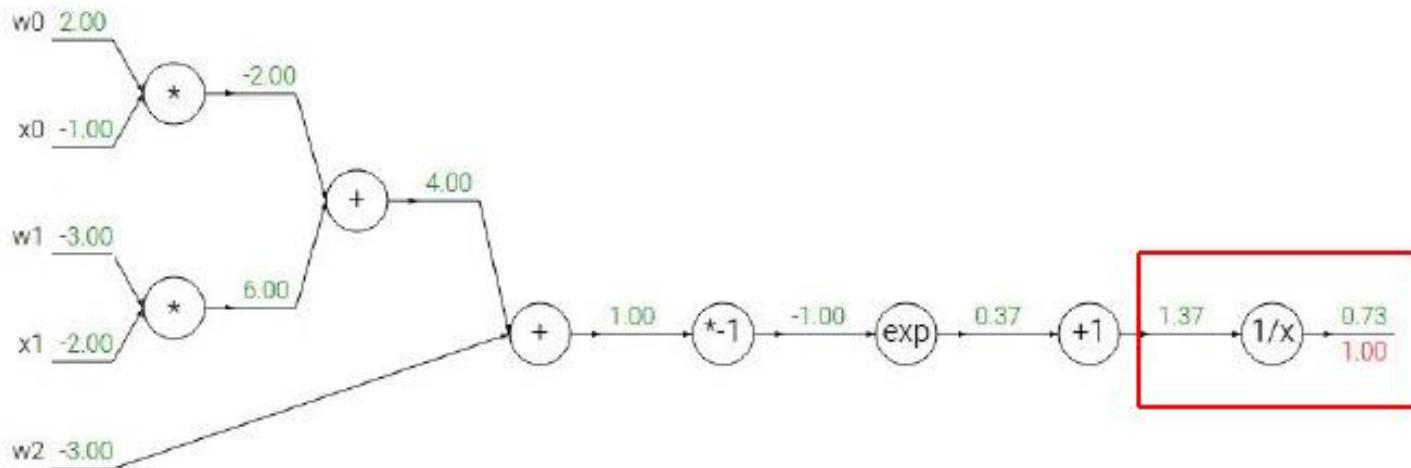
$$f(x) = \frac{1}{x}$$

$$f_c(x) = c + x$$

$$\frac{df}{dx} = -1/x^2$$

$$\frac{df}{dx} = 1$$

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



$$f(x) = e^x$$

\rightarrow

$$\frac{df}{dx} = e^x$$

$$f_a(x) = ax$$

\rightarrow

$$\frac{df}{dx} = a$$

$$f(x) = \frac{1}{x}$$

\rightarrow

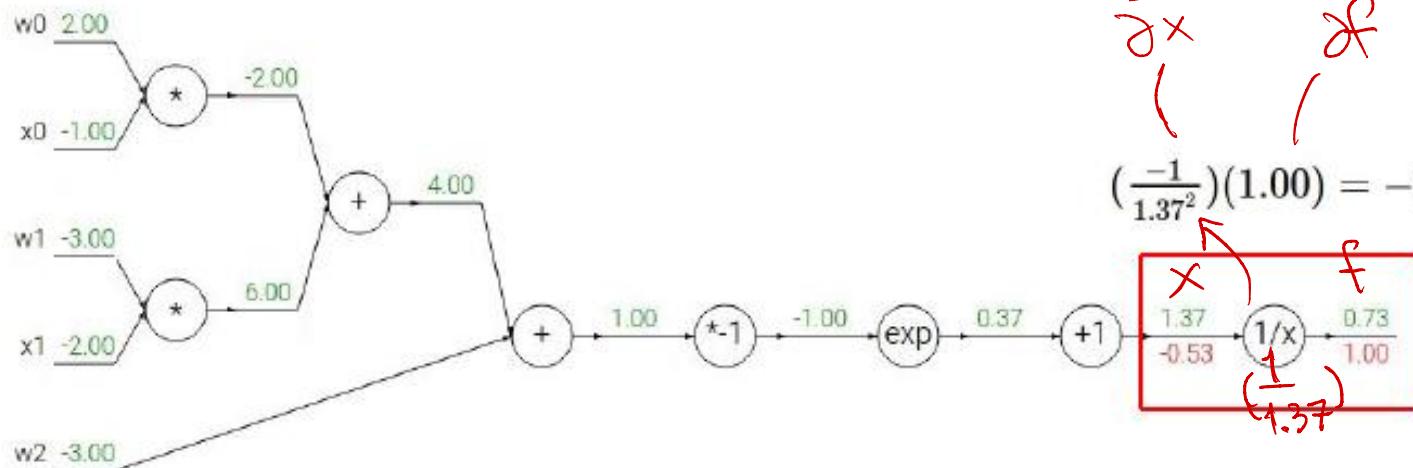
$$\frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x$$

\rightarrow

$$\frac{df}{dx} = 1$$

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



$$f(x) = e^x$$

→

$$\frac{df}{dx} = e^x$$

$$f_a(x) = ax$$

→

$$\frac{df}{dx} = a$$

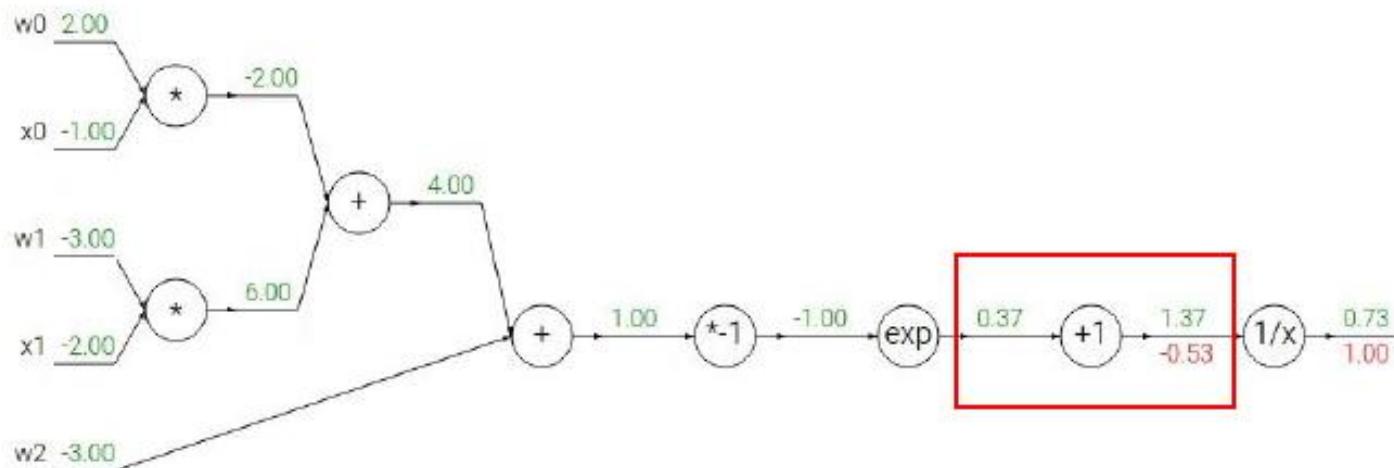
$$f(x) = \frac{1}{x}$$

$$f_c(x) = c + x$$

$$\frac{df}{dx} = -1/x^2$$

$$\frac{df}{dx} = 1$$

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



$$f(x) = e^x$$

\rightarrow

$$\frac{df}{dx} = e^x$$

$$f_a(x) = ax$$

\rightarrow

$$\frac{df}{dx} = a$$

$$f(x) = \frac{1}{x}$$

\rightarrow

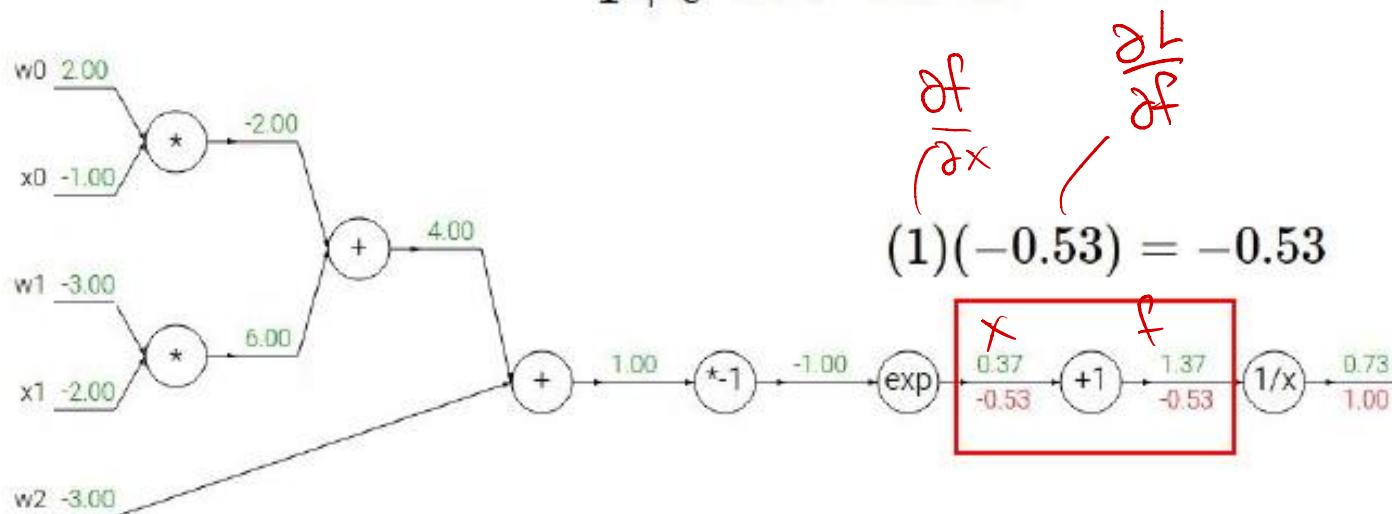
$$\frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x$$

\rightarrow

$$\frac{df}{dx} = 1$$

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



$$f(x) = e^x$$

→

$$\frac{df}{dx} = e^x$$

$$f_a(x) = ax$$

→

$$\frac{df}{dx} = a$$

$$f(x) = \frac{1}{x}$$

→

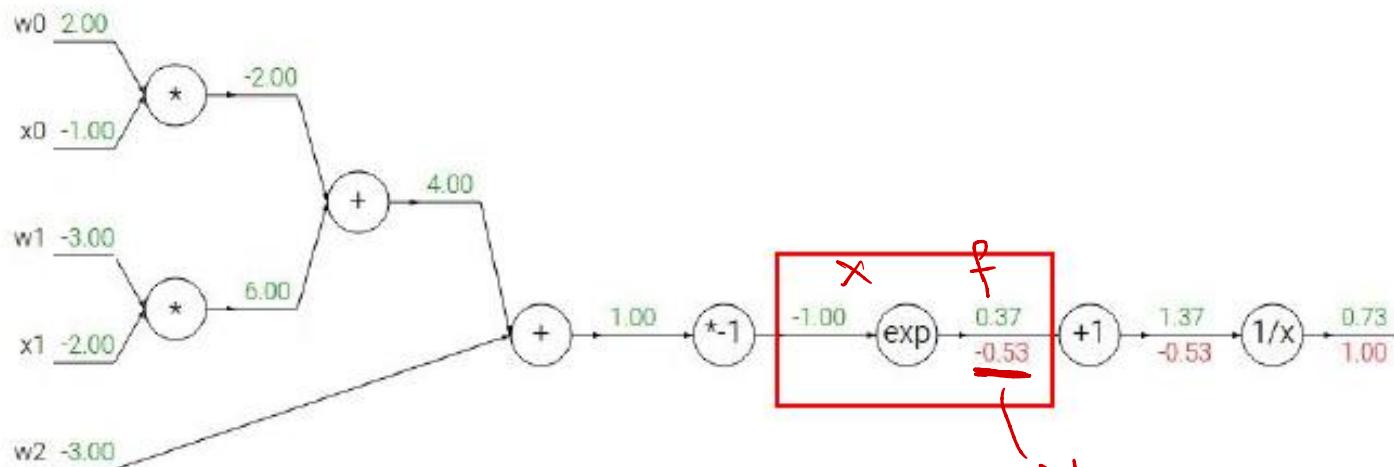
$$\frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x$$

→

$$\frac{df}{dx} = 1$$

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



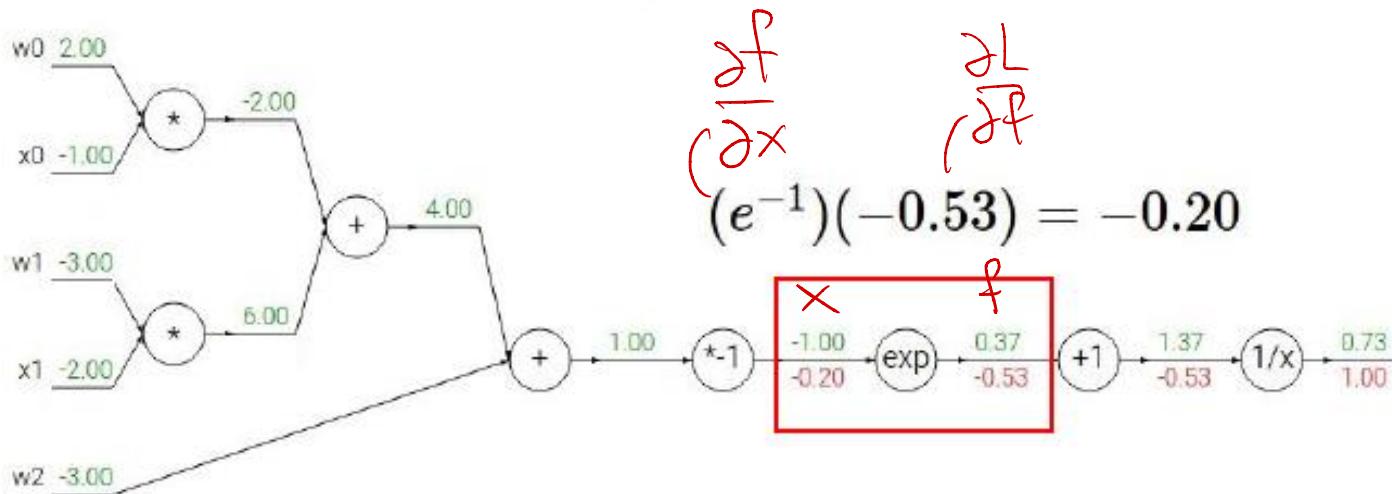
$$f(x) = e^x \rightarrow \frac{df}{dx} = e^x$$

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

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

$$\frac{df}{dx} = 1 \rightarrow \frac{df}{dx} = 1$$

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



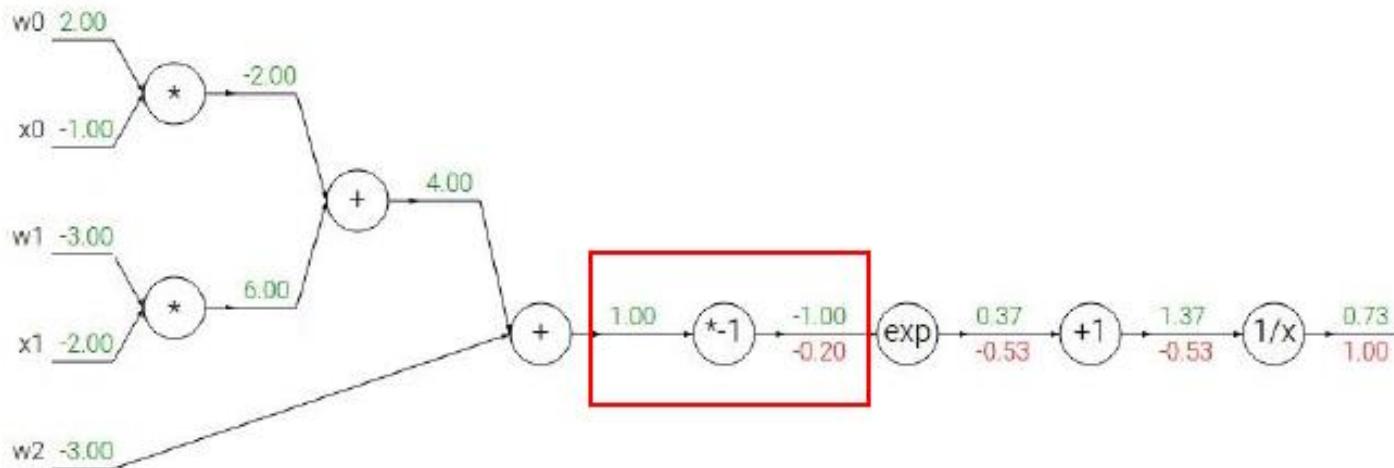
$$f(x) = e^x \rightarrow \frac{df}{dx} = e^x$$

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

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

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

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



$$f(x) = e^x$$

\rightarrow

$$\frac{df}{dx} = e^x$$

$$f_a(x) = ax$$

\rightarrow

$$\frac{df}{dx} = a$$

$$f(x) = \frac{1}{x}$$

\rightarrow

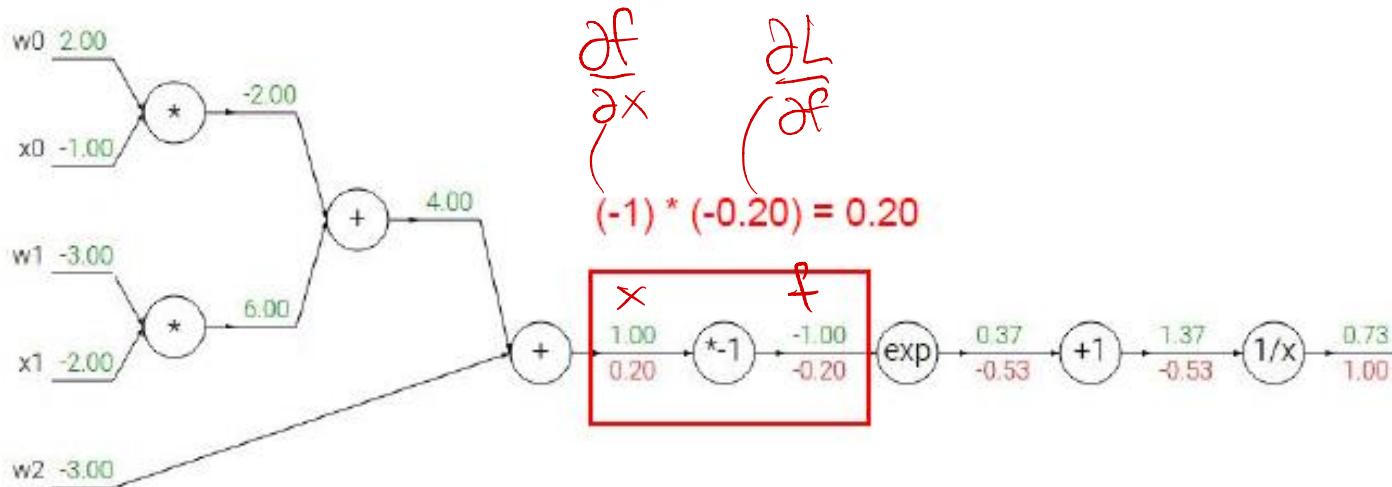
$$\frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x$$

\rightarrow

$$\frac{df}{dx} = 1$$

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



$$f(x) = e^x$$

→

$$\frac{df}{dx} = e^x$$

$$f_a(x) = ax$$

→

$$\frac{df}{dx} = a$$

$$f(x) = \frac{1}{x}$$

→

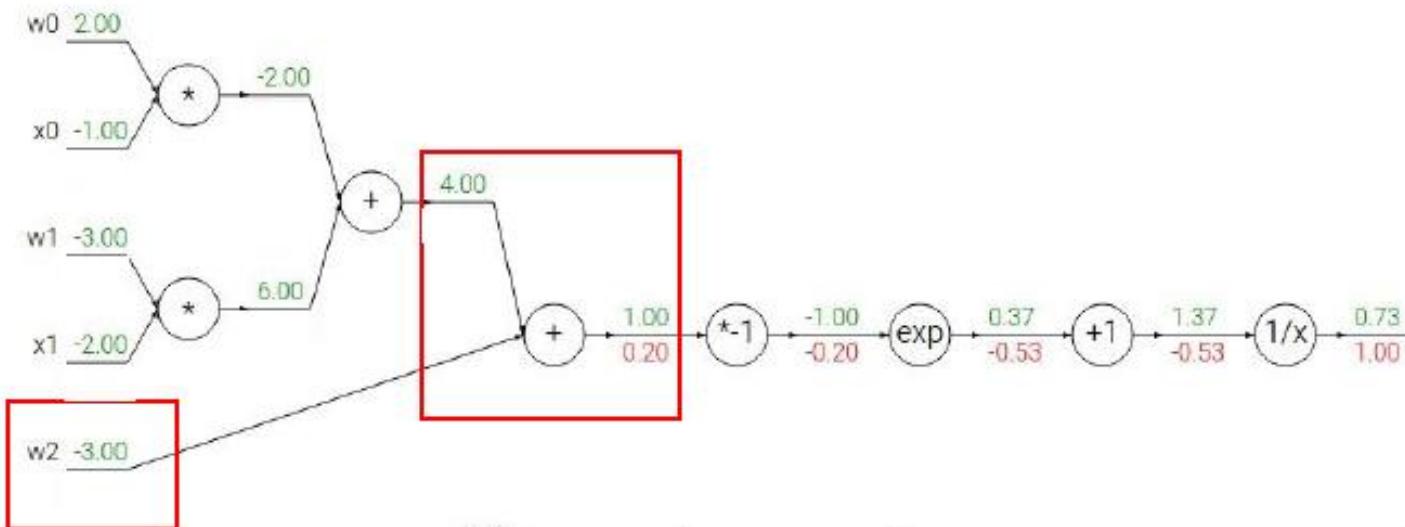
$$\frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x$$

→

$$\frac{df}{dx} = 1$$

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



$$f(x) = e^x$$

\rightarrow

$$\frac{df}{dx} = e^x$$

$$f_a(x) = ax$$

\rightarrow

$$\frac{df}{dx} = a$$

$$f(x) = \frac{1}{x}$$

\rightarrow

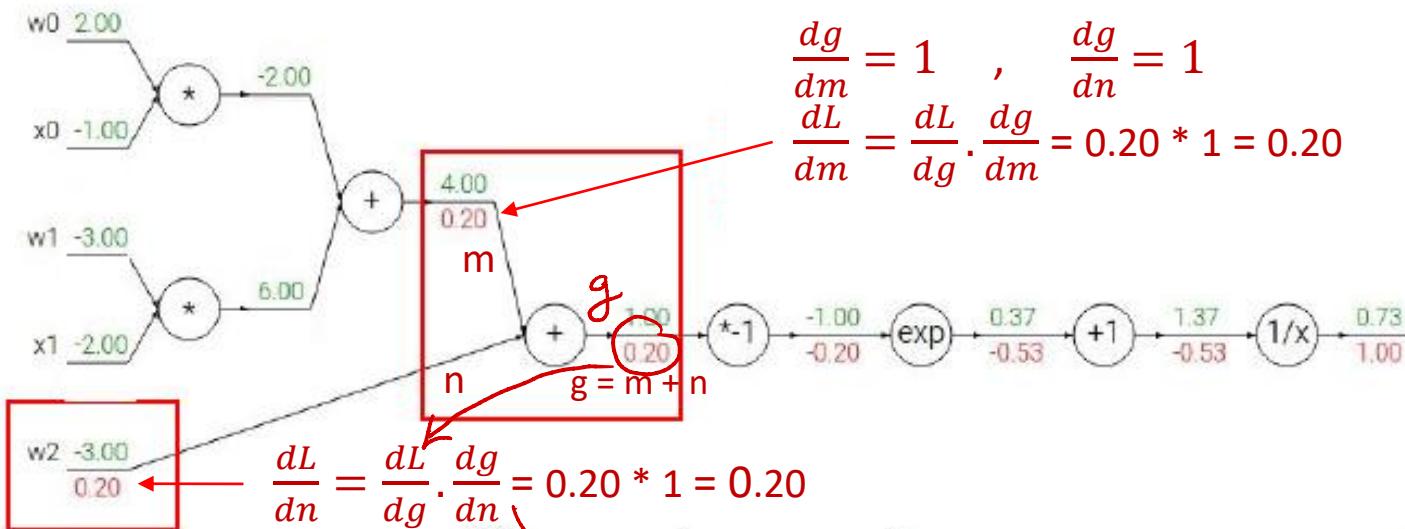
$$\frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x$$

\rightarrow

$$\frac{df}{dx} = 1$$

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



$$f(x) = e^x$$

$$f_a(x) = ax$$

\rightarrow

$$\frac{df}{dx} = e^x$$

$$\frac{df}{dx} = a$$

$$f(x) = \frac{1}{x}$$

$$f_c(x) = c + x$$

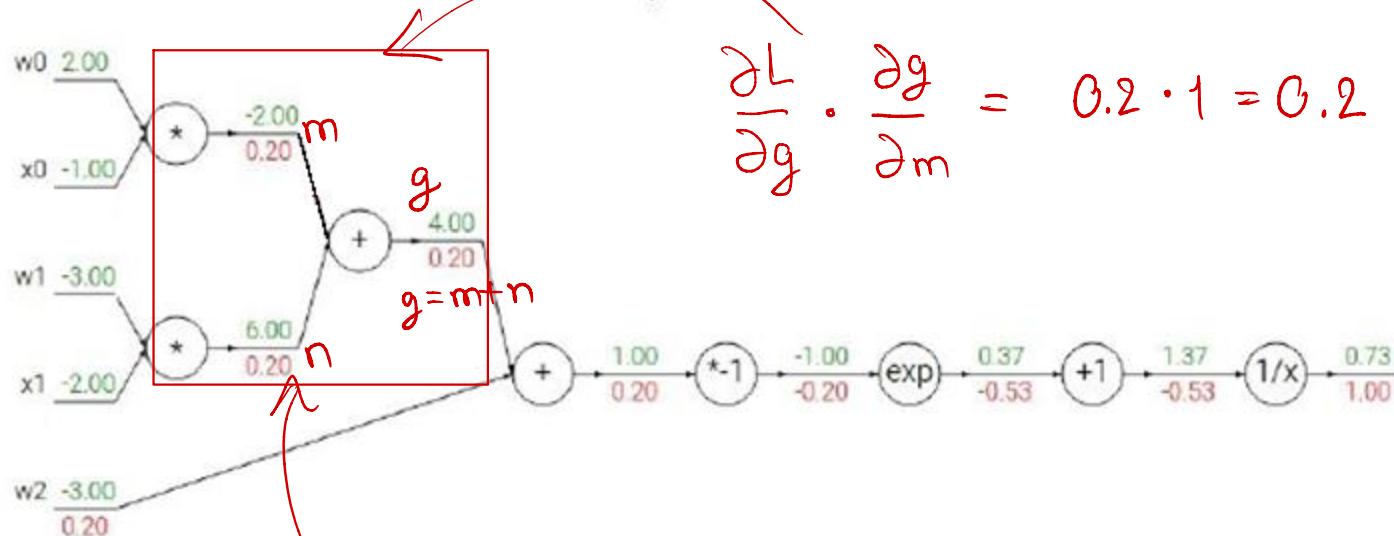
\rightarrow

$$\frac{df}{dx} = -1/x^2$$

$$\frac{df}{dx} = 1$$

Another example:

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



$$\frac{\partial L}{\partial g} \cdot \frac{\partial g}{\partial m} = 0.2 \cdot 1 = 0.2$$

$$f(x) = e^x$$

$$f_a(x) = ax$$

$$\rightarrow \frac{df}{dx} = e^x$$

$$\rightarrow \frac{\partial L}{\partial g} \cdot \frac{\partial g}{\partial n} \frac{df}{dx} = a$$

$$= 0.2 \cdot 1 = 0.2$$

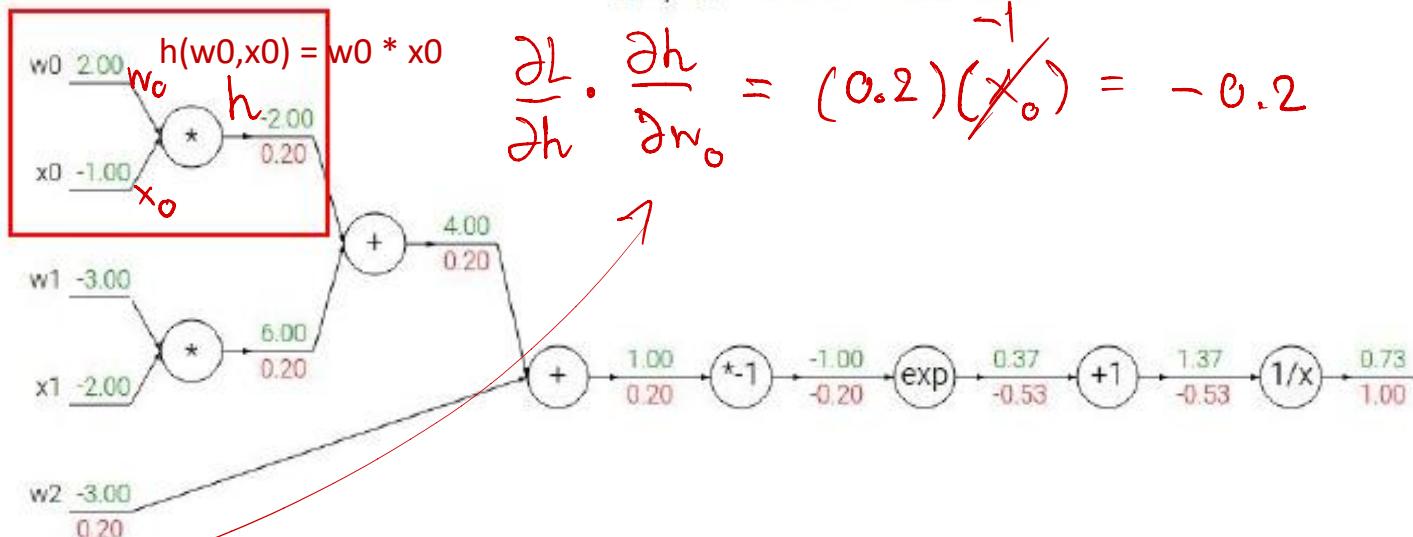
$$f(x) = \frac{1}{x}$$

$$f_c(x) = c + x$$

$$\frac{df}{dx} = -1/x^2$$

$$\frac{df}{dx} = 1$$

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



$$f(x) = e^x$$

$$\underline{f_a(x) = ax}$$

$$\frac{df}{dx} = e^x$$

$$\frac{df}{dx} = a$$

$$f(x) = \frac{1}{x}$$

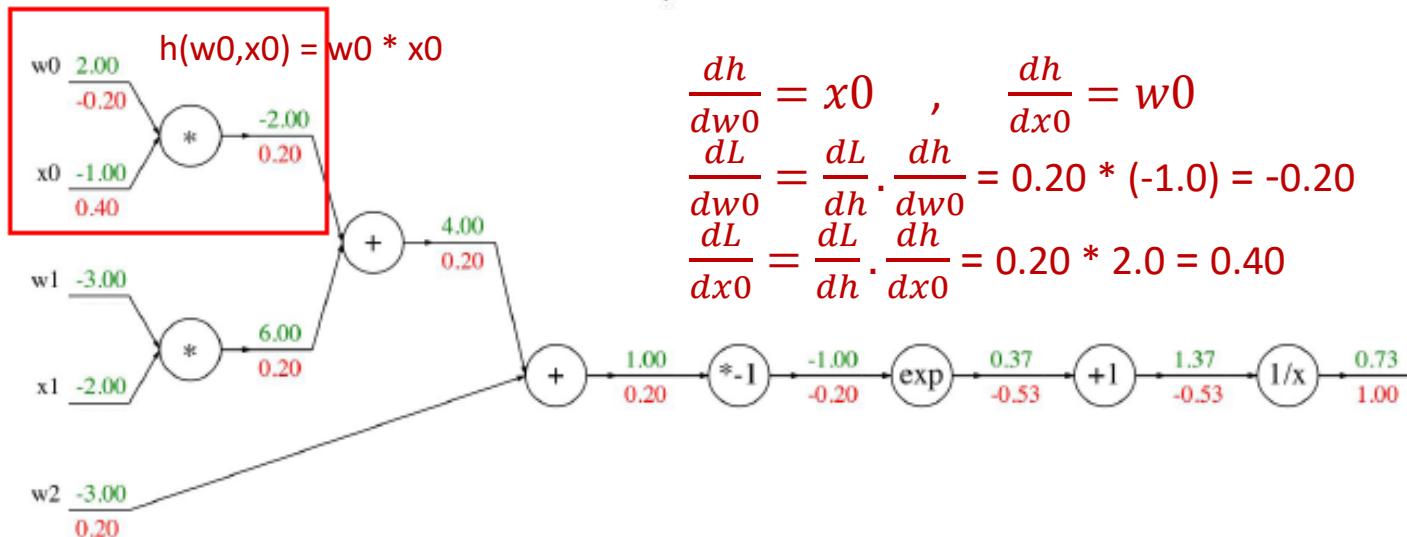
$$f_c(x) = c + x$$

$$\frac{df}{dx} = -1/x^2$$

$$\frac{df}{dx} = 1$$

$$\frac{\partial L}{\partial h} \cdot \frac{\partial h}{\partial x_0} = (0.2)(\cancel{w_0}) = 0.4$$

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



$$f(x) = e^x$$

→

$$\frac{df}{dx} = e^x$$

$$f_a(x) = ax$$

→

$$\frac{df}{dx} = a$$

$$f(x) = \frac{1}{x}$$

→

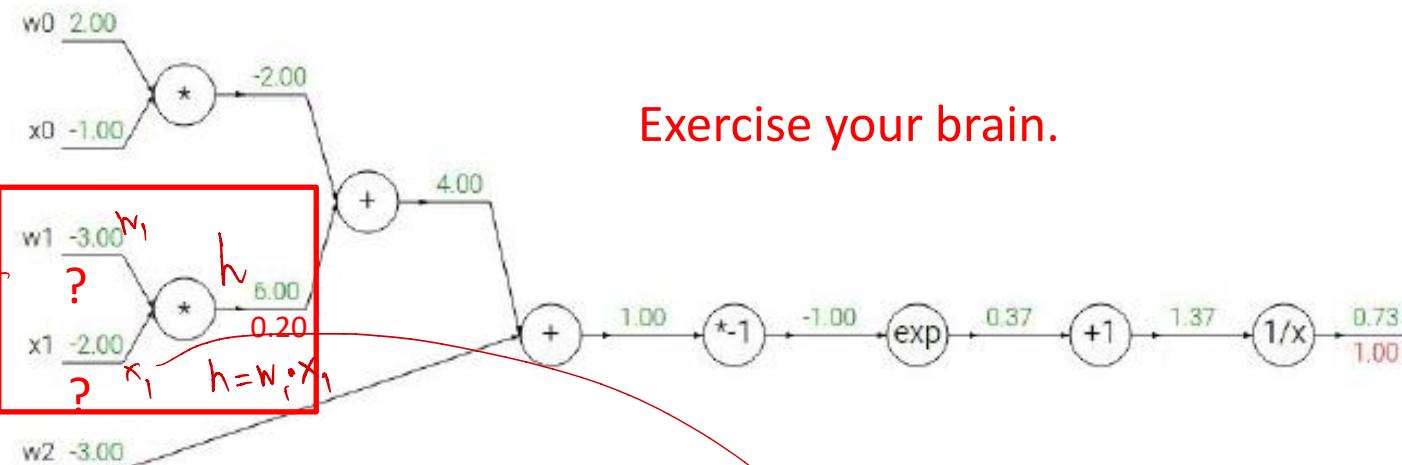
$$\frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x$$

→

$$\frac{df}{dx} = 1$$

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



$$f(x) = e^x$$

$$f_a(x) = ax$$

$$\frac{df}{dx} = e^x$$

$$\frac{df}{dx} = a$$

$$f(x) = \frac{1}{x}$$

$$f_c(x) = c + x$$

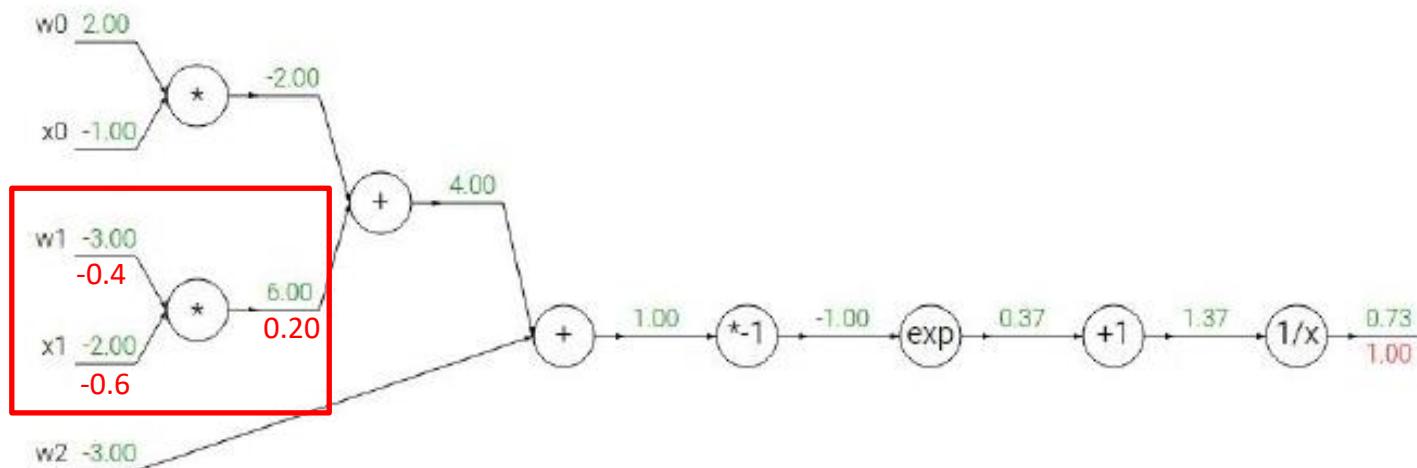
$$\frac{df}{dx} = -1/x^2$$

$$\frac{df}{dx} = 1$$

$$\frac{\partial L}{\partial h} \cdot \frac{\partial h}{\partial w_1} = (0.2)(x_1) = -0.4$$

$$\frac{\partial L}{\partial h} \cdot \frac{\partial h}{\partial x_1} = (0.2)(w_1) = -0.6$$

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



$$f(x) = e^x$$

→

$$\frac{df}{dx} = e^x$$

$$f_a(x) = ax$$

→

$$\frac{df}{dx} = a$$

$$f(x) = \frac{1}{x}$$

→

$$\frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x$$

→

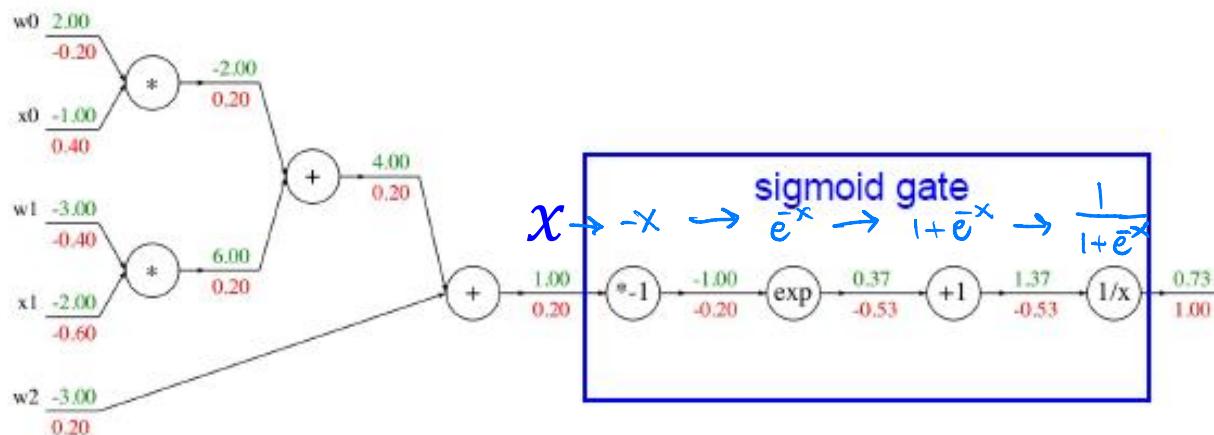
$$\frac{df}{dx} = 1$$

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

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

sigmoid function

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

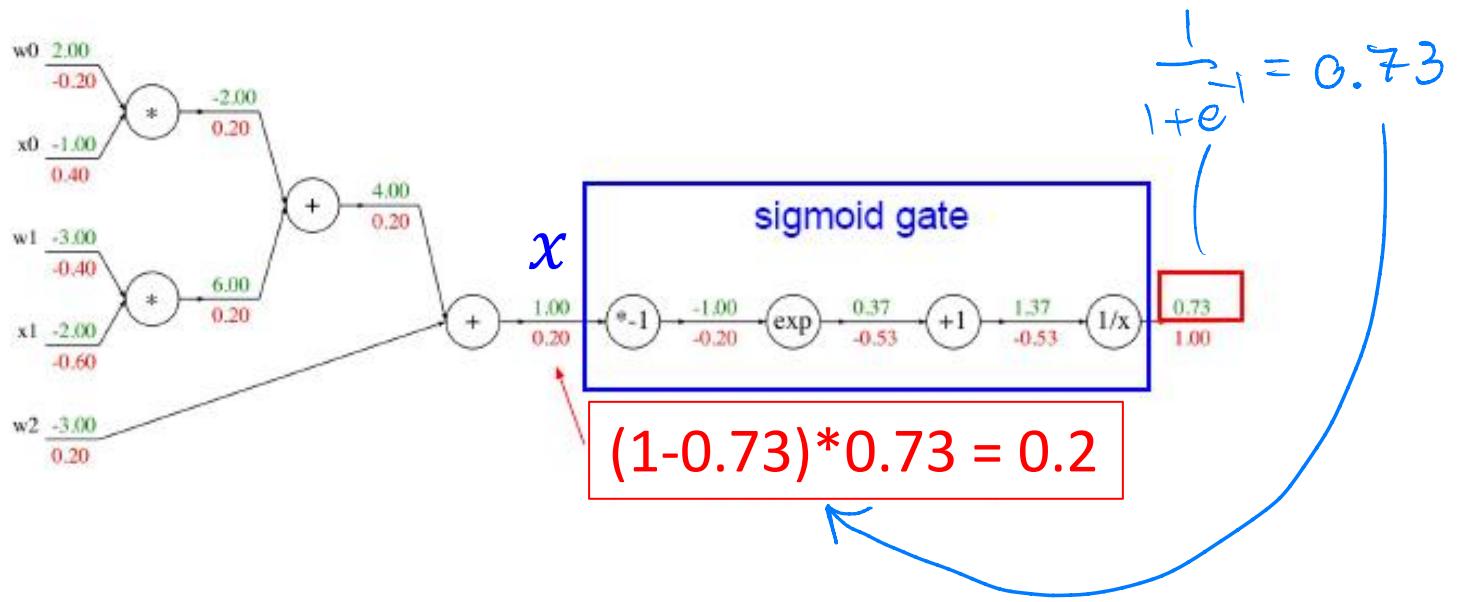


$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

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

sigmoid function

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



Exercise

In the context of artificial neural networks, the **rectifier** is an activation function defined as the positive part of its argument:

$$f(x) = \max(0, x) \quad ; \text{this is also known as ReLU function}$$

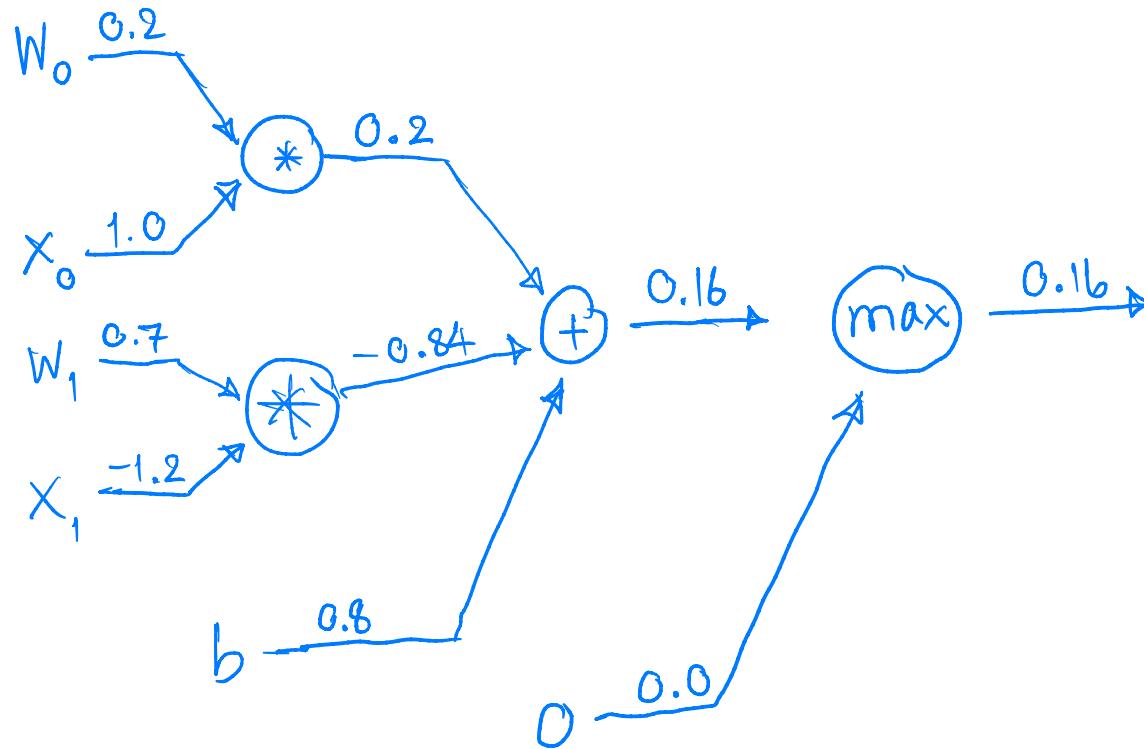
Let $x = w_0x_0 + w_1x_1 + b$

$$w_0 = 0.2, w_1 = 0.7, x_0 = 1.0, x_1 = -1.2, b = 0.8$$

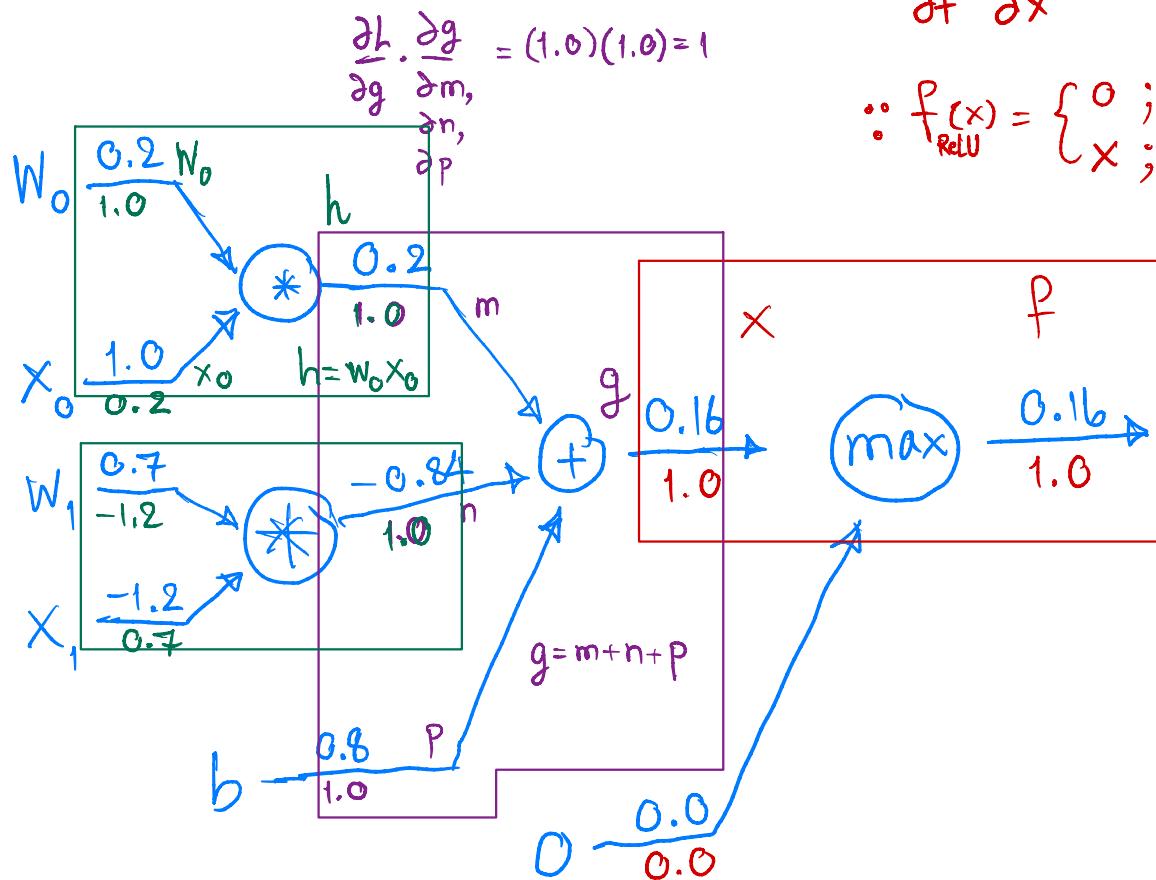
1. Draw the complete computational graph for $f(x)$.
2. Assume that $\frac{\partial L}{\partial f} = 1.0$
3. Calculate backpropagate gradient of L at every branch of the graph.

/

Computational graph



Gradients



$$\frac{\partial L}{\partial p} \cdot \frac{\partial p}{\partial x} = (1)(1) = 1$$

$$\therefore f_{\text{ReLU}}(x) = \begin{cases} 0 & ; x < 0 \\ x & ; x > 0 \end{cases}$$

$$\frac{\partial L}{\partial h} \cdot \frac{\partial h}{\partial w_0} = (1.0)(1.0) = 1.0$$

$$\frac{\partial L}{\partial h} \cdot \frac{\partial h}{\partial x_0} = (1.0)(0.2) = 0.2$$

$$\frac{\partial L}{\partial h} \cdot \frac{\partial h}{\partial w_1} = (1.0)(-0.84) = -0.84$$

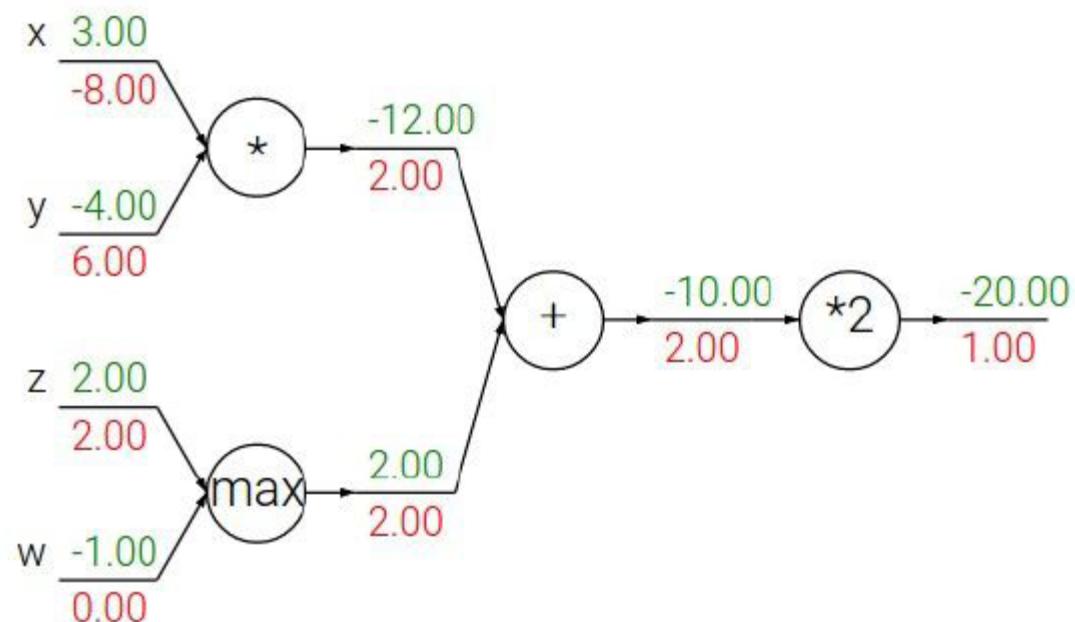
$$\frac{\partial L}{\partial h} \cdot \frac{\partial h}{\partial x_1} = (1.0)(1.0) = 1.0$$

Patterns in backward flow

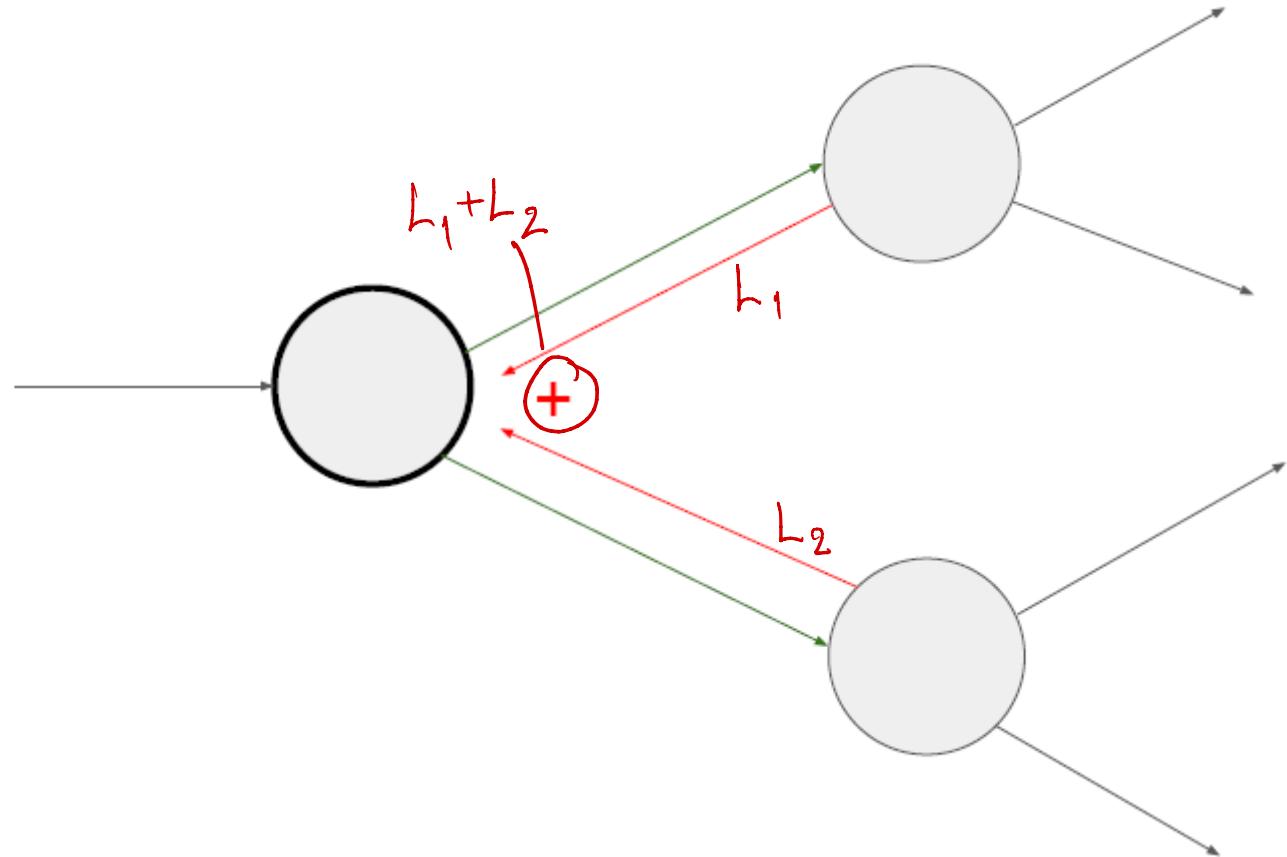
add gate: gradient distributor

max gate: gradient router

mul gate: gradient switcher



Gradients add at branches

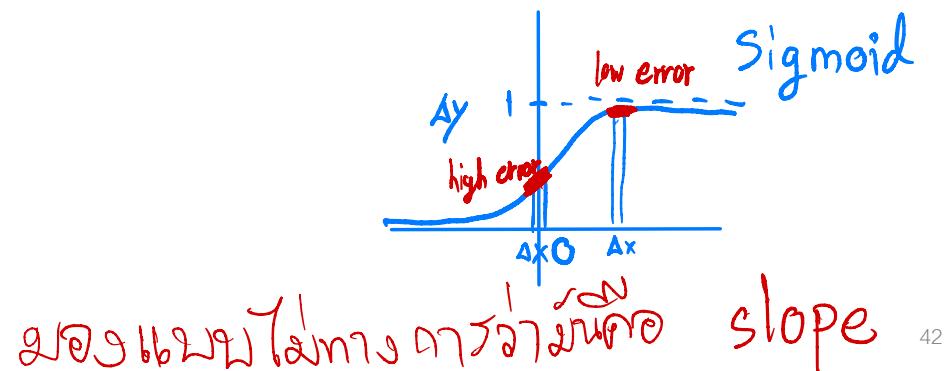


Revisit: Backpropagation for multilayer NN

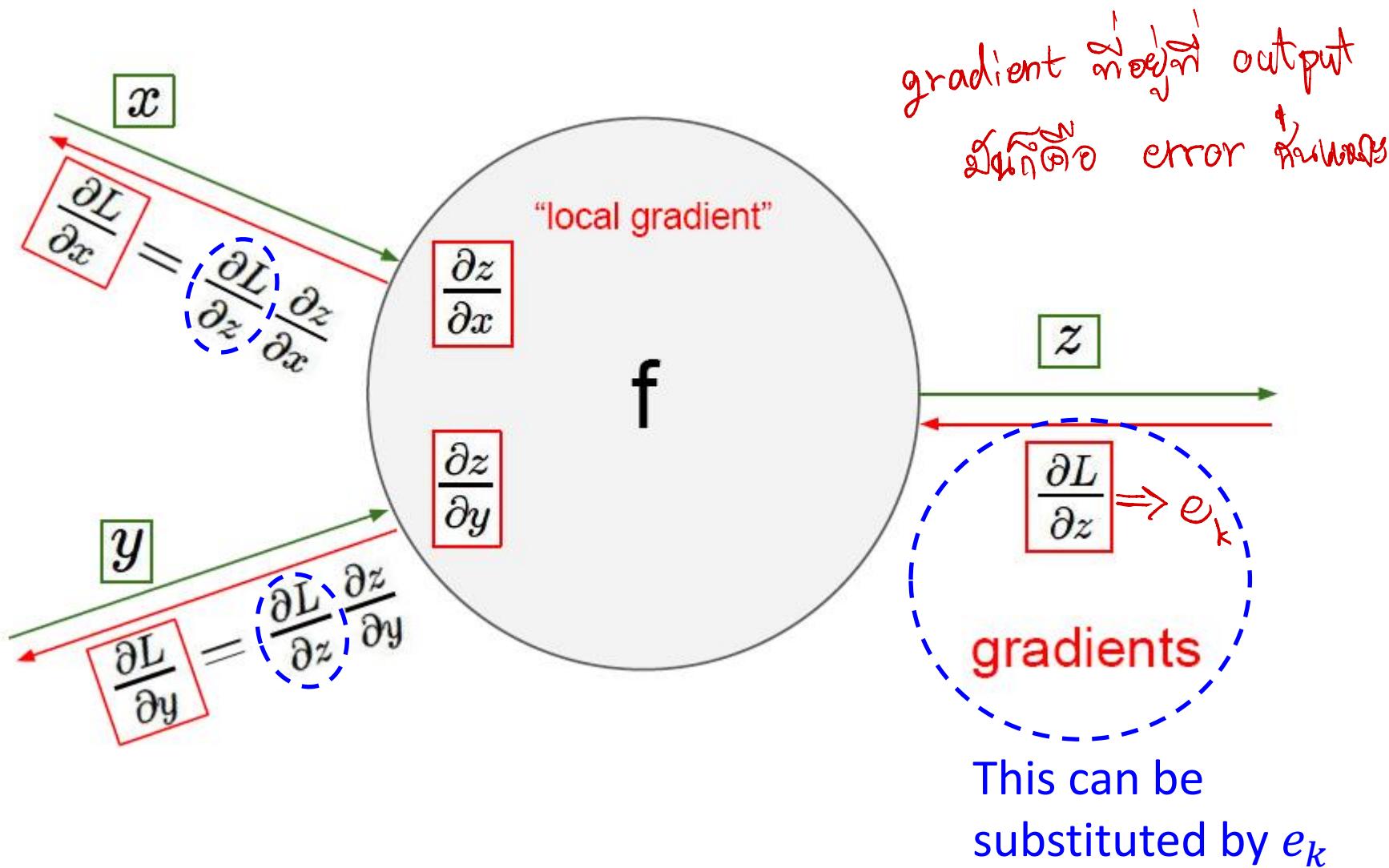
- To propagate error signals, we start at the output layer and work backward to the hidden layer.
- The error signal at the output of neuron k at iteration p is defined by:

$$e_k(p) = \text{Target} - \text{Actual}$$
$$e_k(p) = yd_k(p) - y_k(p)$$

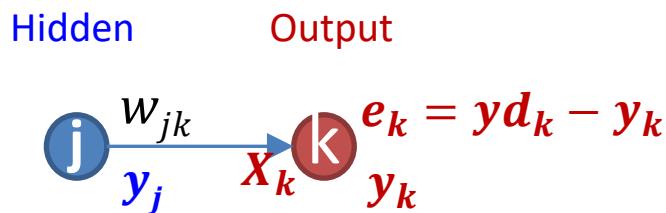
Where $yd_k(p)$ is the desired output of neuron k at iteration p.



Recall gradients backpropagation from page 18.



- To update weights, the weight correction (Δw) is adjusted to the previous weights.



$$w_{jk}(p + 1) = w_{jk}(p) - \Delta w_{jk}(p),$$

$$\Delta w_{jk}(p) = \alpha \cdot y_j(p) \cdot \delta_k(p),$$

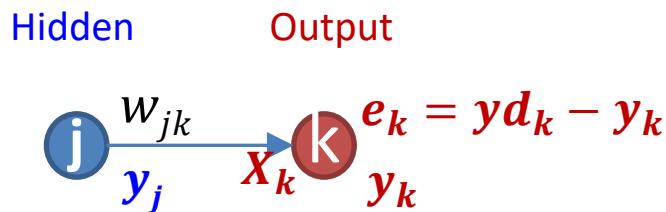
$$\delta_k(p) = \frac{\partial y_k(p)}{\partial x_k(p)} \cdot e_k(p),$$

How do we get these terms?

ව්‍යුහ තැබූ ඇති

- To update weights, the weight correction (Δw) is adjusted to the previous weights.

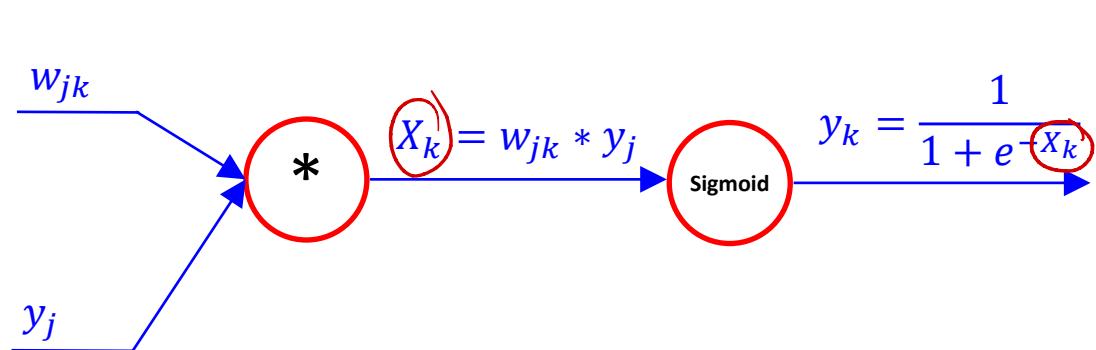
$$w_{jk}(p + 1) = w_{jk}(p) - \Delta w_{jk}(p),$$



$$\Delta w_{jk}(p) = \alpha \cdot y_j(p) \cdot \delta_k(p),$$

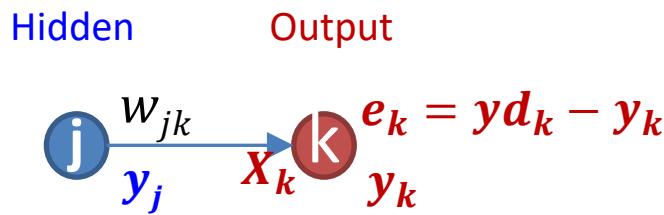
$$\delta_k(p) = \frac{\partial y_k(p)}{\partial x_k(p)} \cdot e_k(p),$$

How do we get these terms?



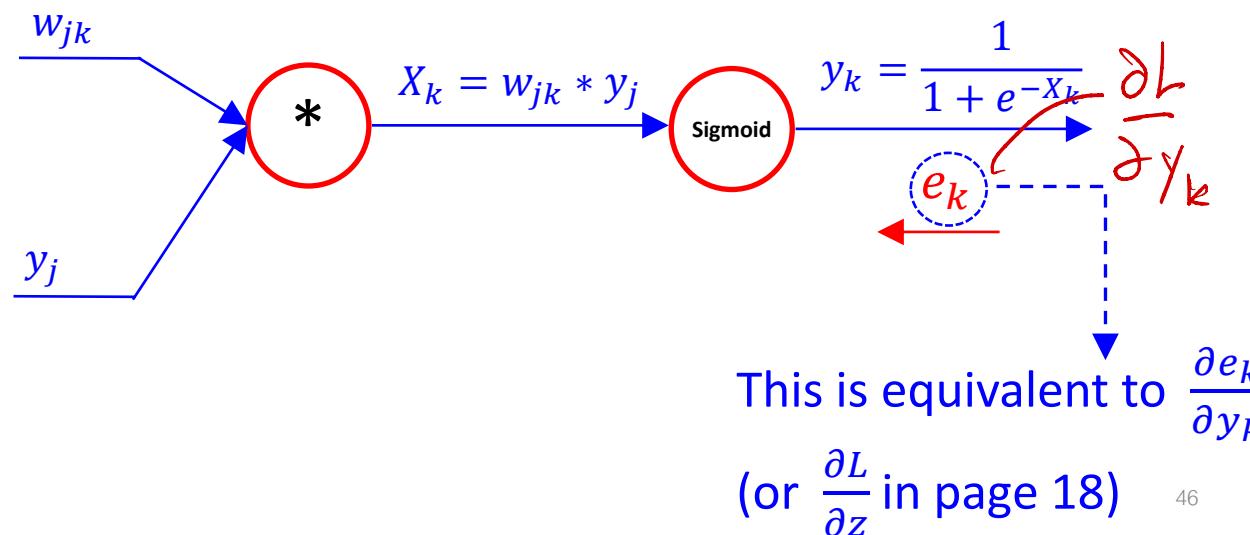
- To update weights, the weight correction (Δw) is adjusted to the previous weights.

$$w_{jk}(p + 1) = w_{jk}(p) - \Delta w_{jk}(p),$$



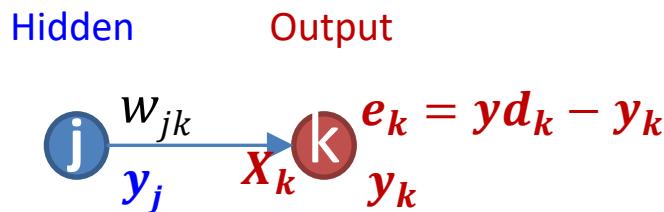
$$\Delta w_{jk}(p) = \alpha \cdot y_j(p) \cdot \delta_k(p),$$

$$\delta_k(p) = \frac{\partial y_k(p)}{\partial x_k(p)} \cdot e_k(p),$$



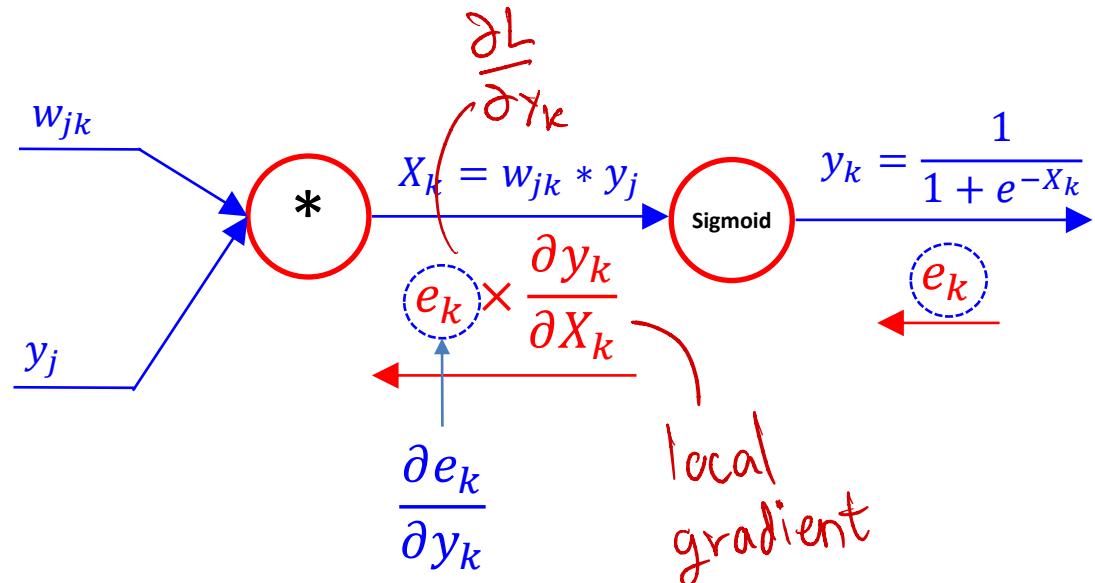
- To update weights, the weight correction (Δw) is adjusted to the previous weights.

$$w_{jk}(p + 1) = w_{jk}(p) - \Delta w_{jk}(p),$$



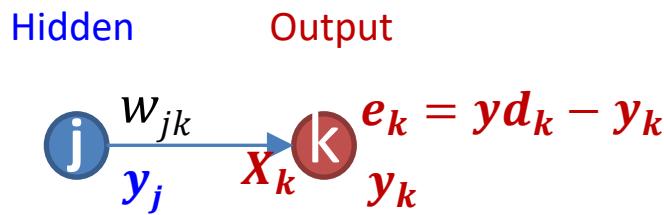
$$\Delta w_{jk}(p) = \alpha \cdot y_j(p) \cdot \delta_k(p),$$

$$\delta_k(p) = \frac{\partial y_k(p)}{\partial x_k(p)} \cdot e_k(p),$$



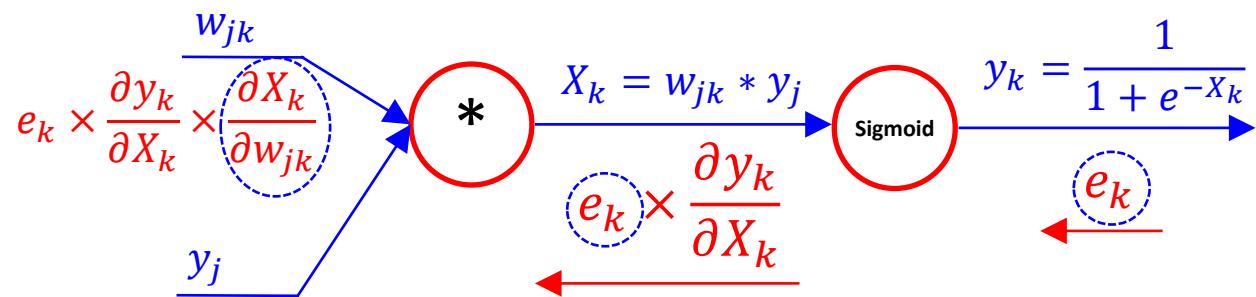
- To update weights, the weight correction (Δw) is adjusted to the previous weights.

$$w_{jk}(p + 1) = w_{jk}(p) - \Delta w_{jk}(p),$$



$$\Delta w_{jk}(p) = \alpha \cdot y_j(p) \cdot \delta_k(p),$$

$$\delta_k(p) = \frac{\partial y_k(p)}{\partial X_k(p)} \cdot e_k(p),$$



- To update weights, the weight correction (Δw) is adjusted to the previous weights.

$$w_{jk}(p + 1) = w_{jk}(p) - \Delta w_{jk}(p),$$

