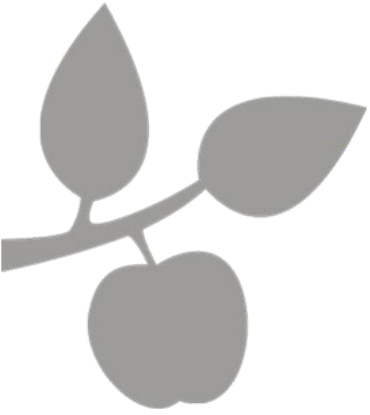


SDU Summer School

Deep Learning

Summer 2022

Welcome to the Summer School



Backpropagation

- **Function Principle**
- **Generalization to Vectors**

Chain Rule of Calculus

- If g is differentiable at x and f is differentiable at $g(x)$, then the composite function $F = f \circ g$ defined by $F(x) = f(g(x))$ is differentiable at x and F' is given by

$$F' = f'(g(x)) \cdot g'(x)$$

- In Leibnitz notation, if $y = f(u)$ and $u = g(x)$, then

$$\frac{dy}{dx} = \frac{dy}{du} \frac{du}{dx}$$

Forward vs. Backward Mode

■ Forward Accumulation:

- One first fixes the **independent variable** with respect to which differentiation is performed and computes the derivative of each sub-expression recursively

$$\frac{dy}{dx} = \frac{dy}{dw_{n-1}} \frac{dw_{i-1}}{dx} = \frac{dy}{dw_{n-1}} \left(\frac{dw_{n-1}}{dw_{n-2}} \frac{dw_{i-2}}{dx} \right) = \dots$$

■ Backward Accumulation:

- One first fixes the **dependent variable** to be differentiated and computes the derivative with respect to each sub-expression recursively

$$\frac{dy}{dx} = \frac{dy}{dw_1} \frac{dw_1}{dx} = \left(\frac{dy}{dw_2} \frac{dw_2}{dw_1} \right) \frac{dw_1}{dx} = \dots$$

Forward Accumulation Example

$$\begin{aligned}
 z &= f(x_1, x_2) \\
 &= x_1 x_2 + \sin x_1 \\
 &= w_1 w_2 + \sin w_1 \\
 &= w_3 + w_4 \\
 &= w_5
 \end{aligned}$$

The choice of the independent variable defines the used seed. For example, we want to differentiate with respect to x_1 :

$$\dot{w}_1 = \frac{dx_1}{dx_1} = 1 \quad \text{and} \quad \dot{w}_2 = \frac{dx_2}{dx_1} = 0$$

$$\frac{df}{dx_1} :$$

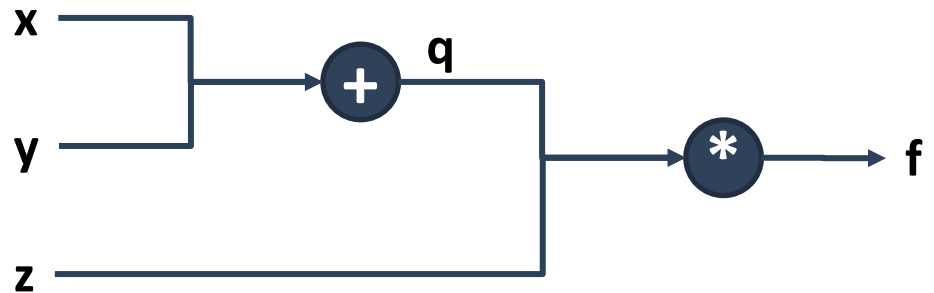
Compute Value	Compute derivative
$w_1 = x_1$	$\dot{w}_1 = 1$
$w_2 = x_2$	$\dot{w}_2 = 0$
$w_3 = w_1 \cdot w_2$	$\dot{w}_3 = w_2 \cdot \dot{w}_1 + w_1 \cdot \dot{w}_2$
$w_4 = \sin w_1$	$\dot{w}_4 = \cos w_1 \cdot \dot{w}_1$
$w_5 = w_3 + w_4$	$\dot{w}_5 = \dot{w}_3 + \dot{w}_4$

Observations

- In order to also get a derivative $\frac{df}{dx_2}$ another run would be required to receive the gradient
- Forward accumulation is good for functions
$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$$
with $m \gg n$
- In deep learning: Normally millions of weights (n) to optimize with only one output ($m = 1$), the costs
 - => Therefore, Backward accumulation, or Backpropagation

Backpropagation: A Simple Example

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

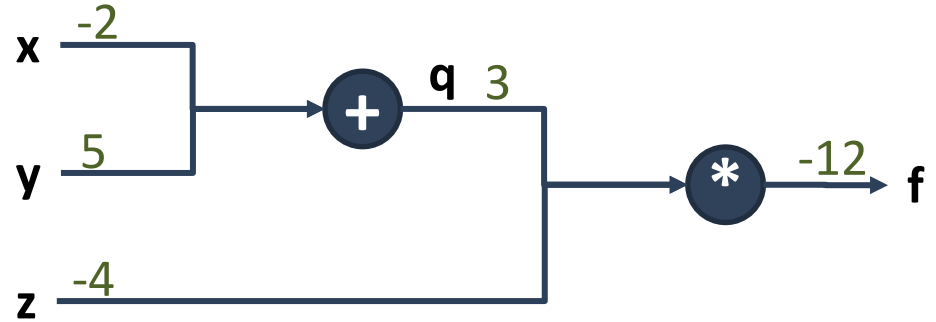


Backpropagation: A Simple Example

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

■ With

- $x = -2$
- $y = 5$
- $z = -4$

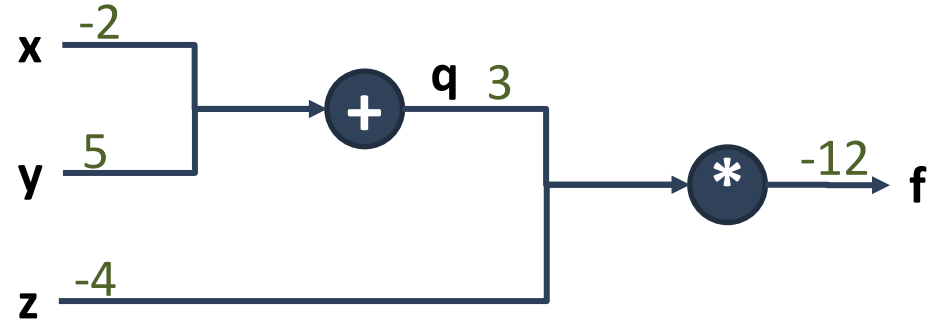


Backpropagation: A Simple Example

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

■ With

- $x = -2$
- $y = 5$
- $z = -4$



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

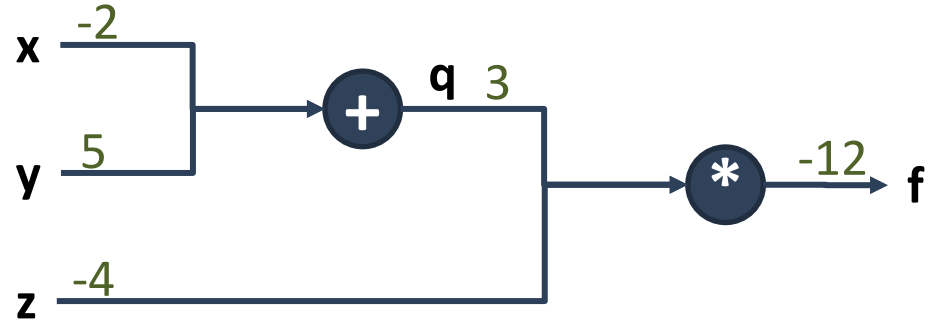
$$f = qz; \quad \frac{\partial f}{\partial q} = z; \quad \frac{\partial f}{\partial z} = q;$$

Backpropagation: A Simple Example

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

■ With

- $x = -2$
- $y = 5$
- $z = -4$



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

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

Looking for:

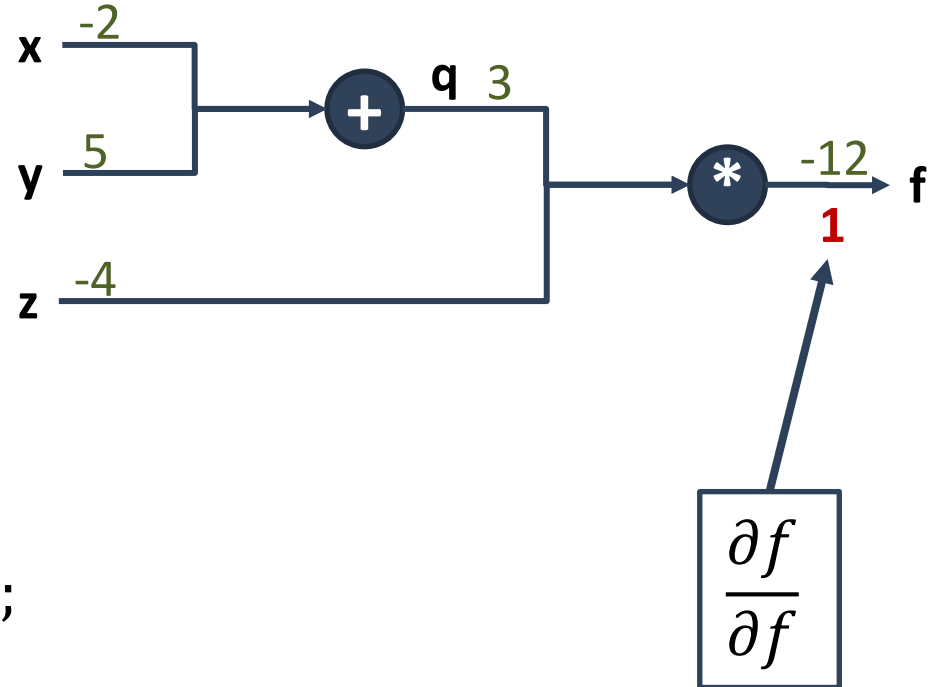
$$\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \quad \frac{\partial f}{\partial z}$$

Backpropagation: A Simple Example

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

■ With

- $x = -2$
- $y = 5$
- $z = -4$



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

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

Looking for:

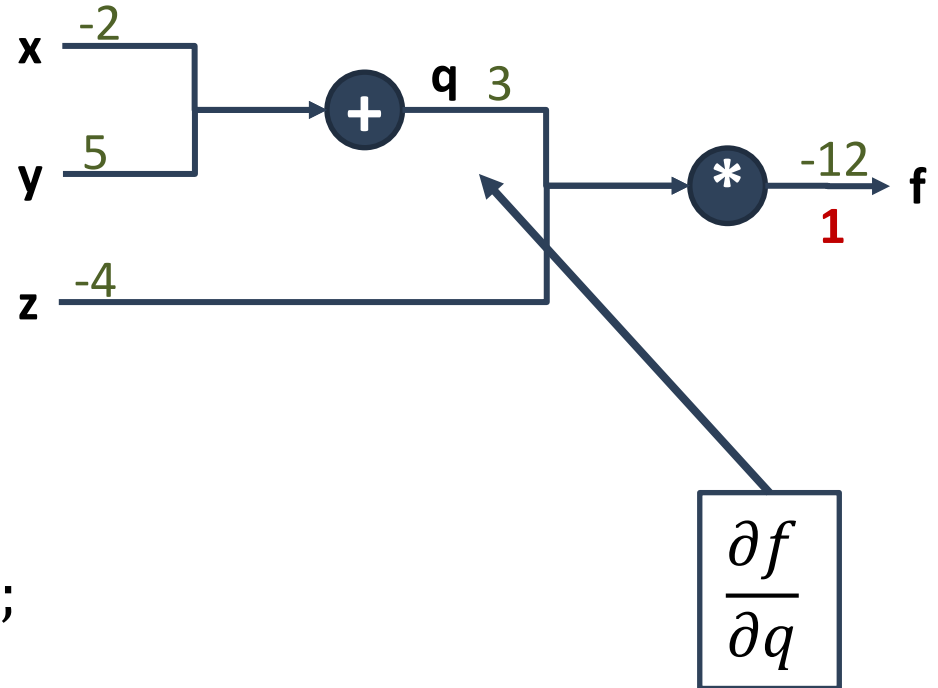
$$\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \quad \frac{\partial f}{\partial z}$$

Backpropagation: A Simple Example

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

■ With

- $x = -2$
- $y = 5$
- $z = -4$



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

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

Looking for:

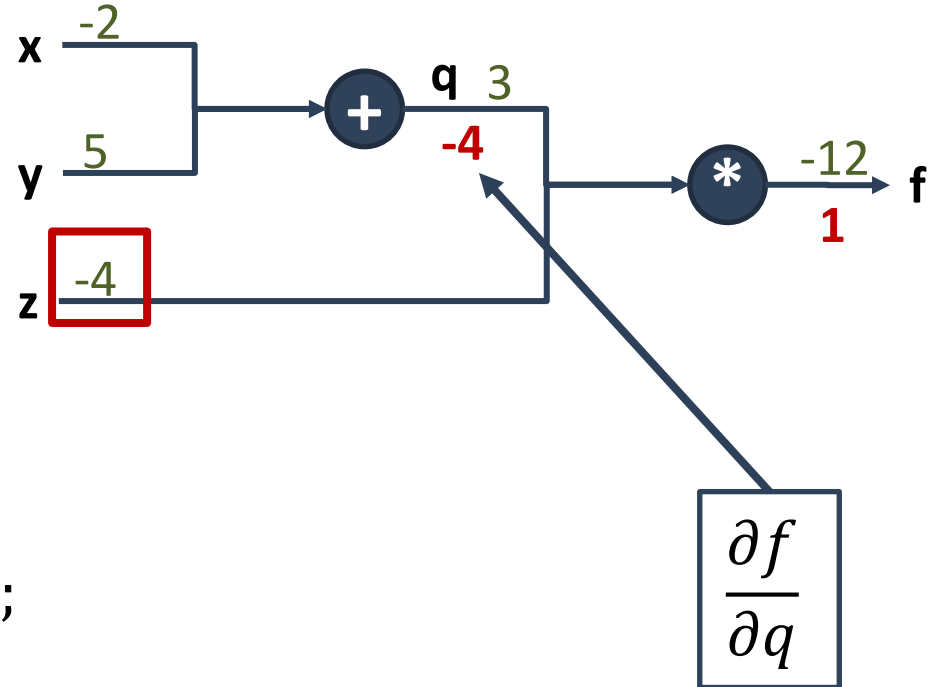
$$\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \quad \frac{\partial f}{\partial z}$$

Backpropagation: A Simple Example

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

■ With

- $x = -2$
- $y = 5$
- $z = -4$



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

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

Looking for:

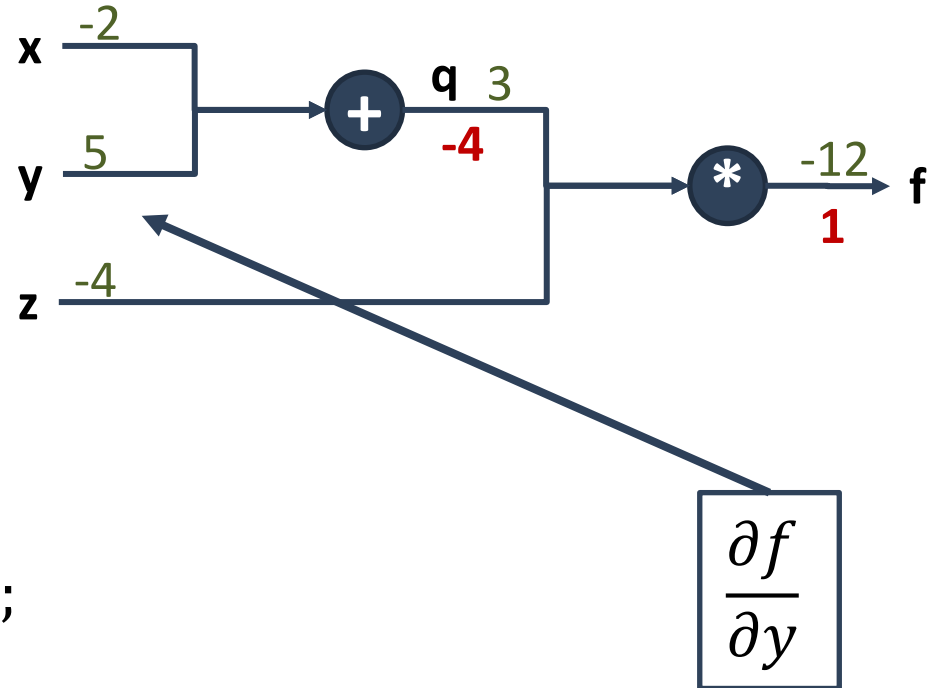
$$\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \quad \frac{\partial f}{\partial z}$$

Backpropagation: A Simple Example

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

■ With

- $x = -2$
- $y = 5$
- $z = -4$



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

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

Looking for:

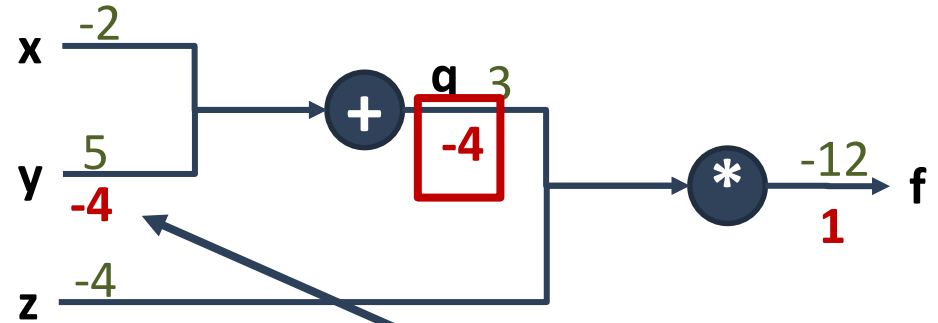
$$\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \quad \frac{\partial f}{\partial z}$$

Backpropagation: A Simple Example

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

■ With

- $x = -2$
- $y = 5$
- $z = -4$



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

Chain Rule:

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

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

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

Looking for:

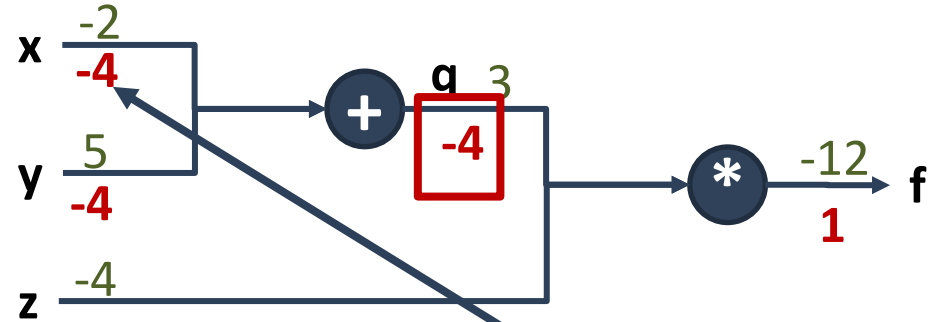
$$\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \quad \frac{\partial f}{\partial z}$$

Backpropagation: A Simple Example

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

■ With

- $x = -2$
- $y = 5$
- $z = -4$



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

Chain Rule:

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

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

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

Looking for:

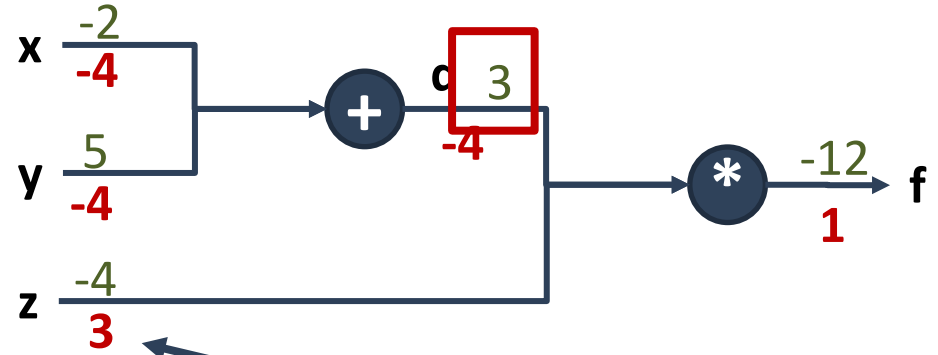
$$\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \quad \frac{\partial f}{\partial z}$$

Backpropagation: A Simple Example

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

■ With

- $x = -2$
- $y = 5$
- $z = -4$



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

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

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

Looking for:

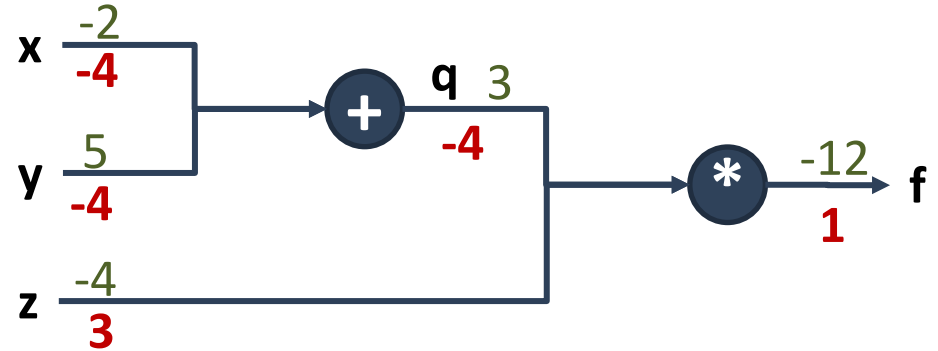
$$\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \quad \frac{\partial f}{\partial z}$$

Backpropagation: A Simple Example

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

■ With

- $x = -2$
- $y = 5$
- $z = -4$



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

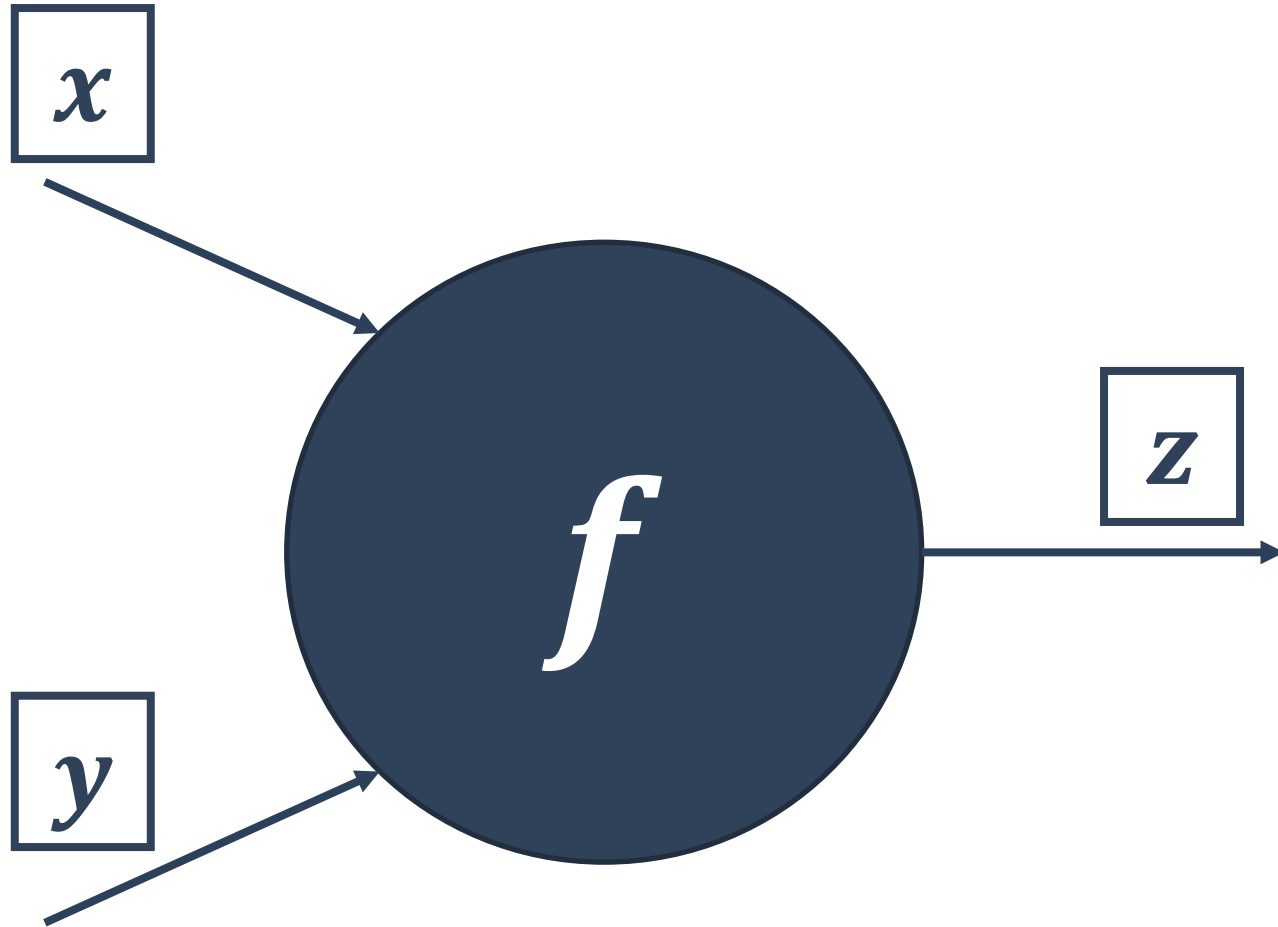
$$f = qz; \quad \frac{\partial f}{\partial q} = z; \quad \frac{\partial f}{\partial z} = q;$$

$$\nabla_{x,y,z} f = \begin{pmatrix} -4 \\ -4 \\ 3 \end{pmatrix}$$

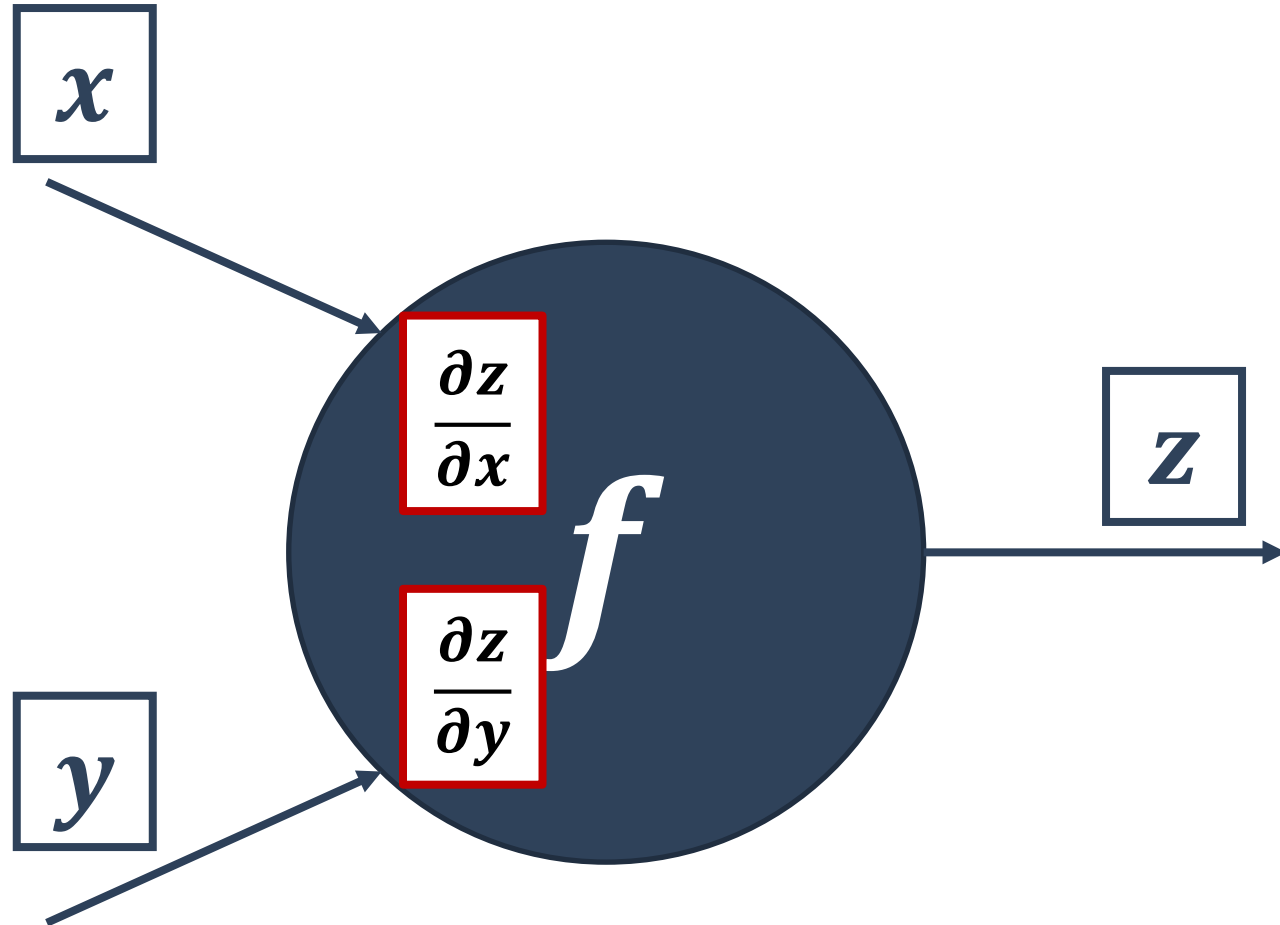
Looking for:

$$\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \quad \frac{\partial f}{\partial z}$$

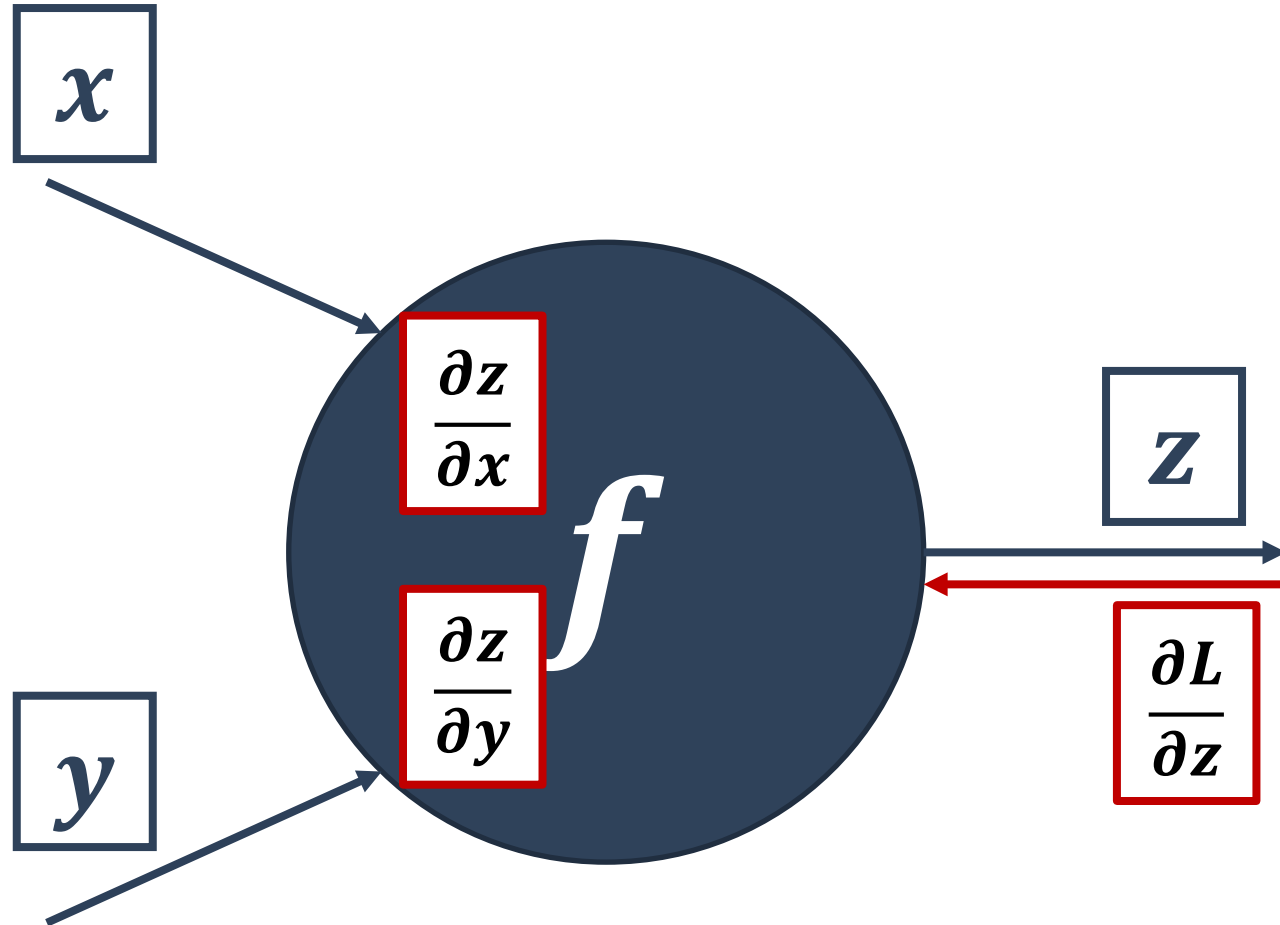
General Principle



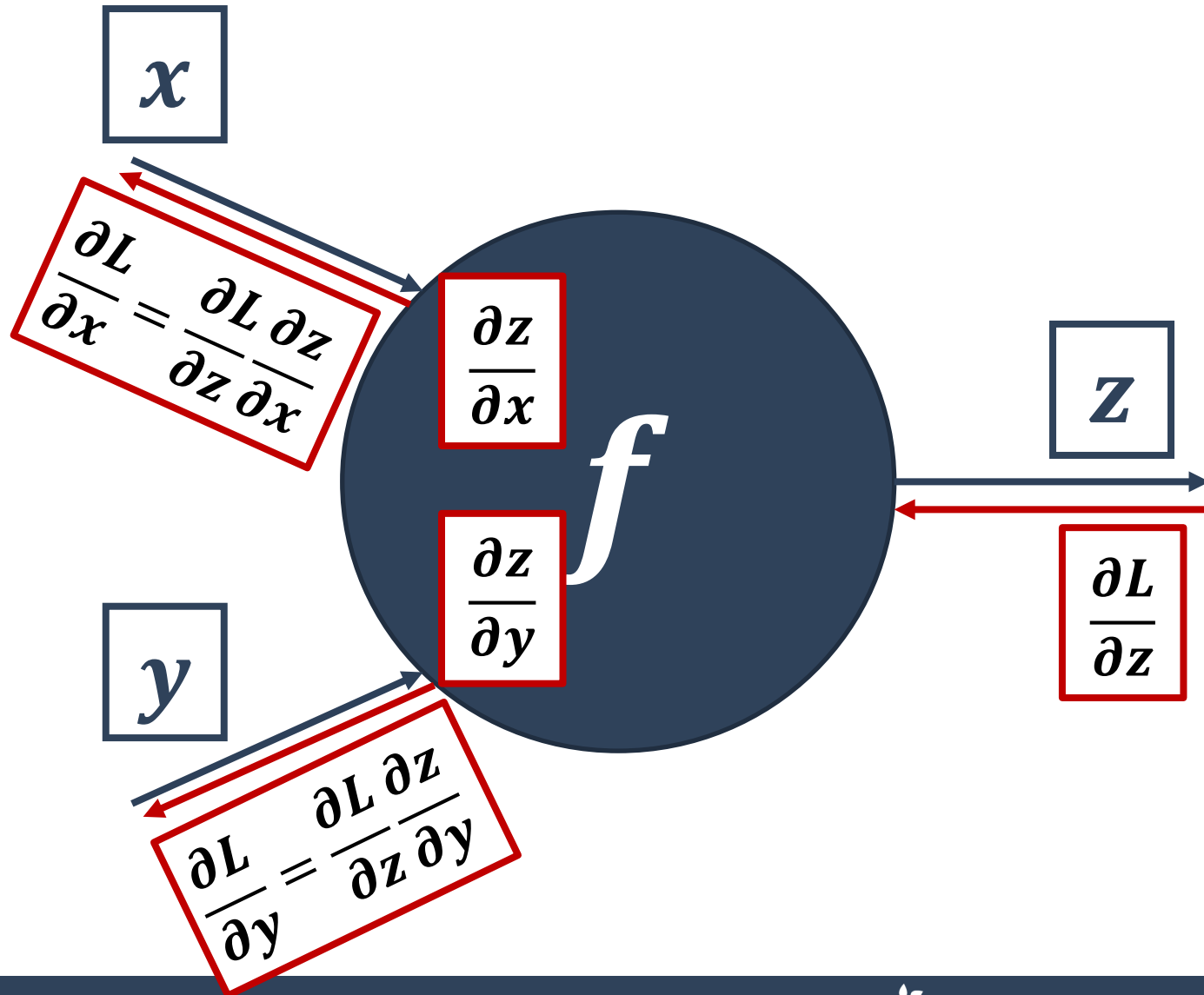
General Principle



General Principle

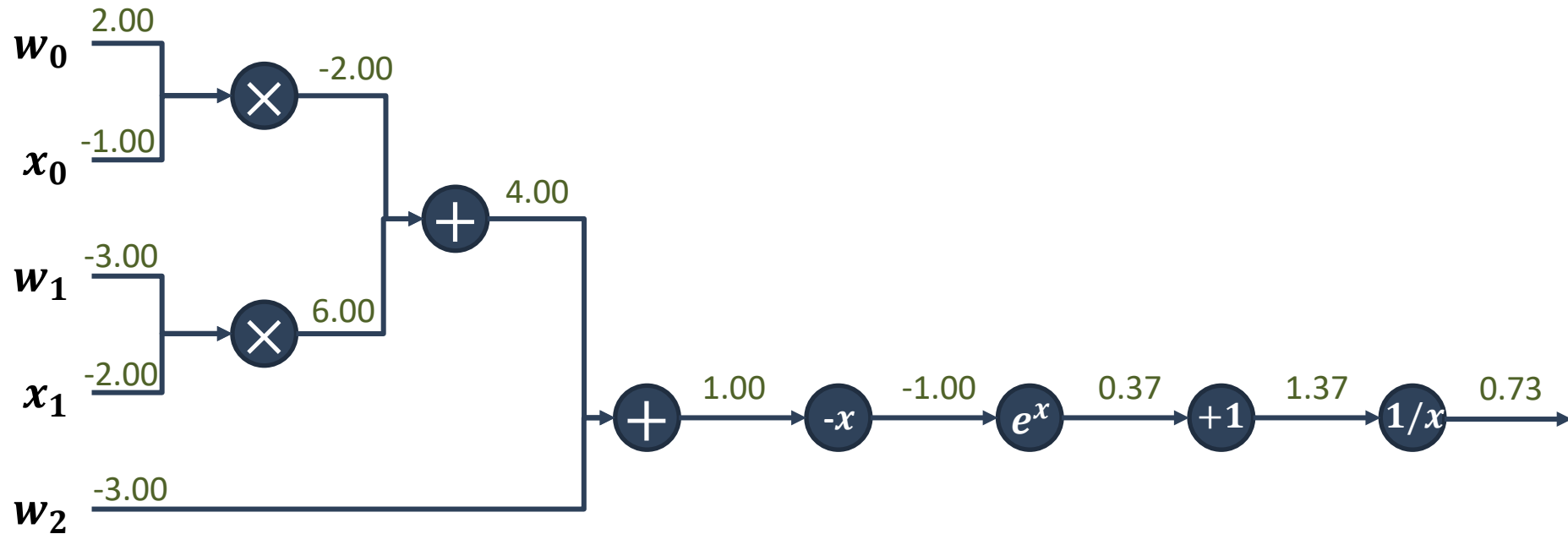


General Principle



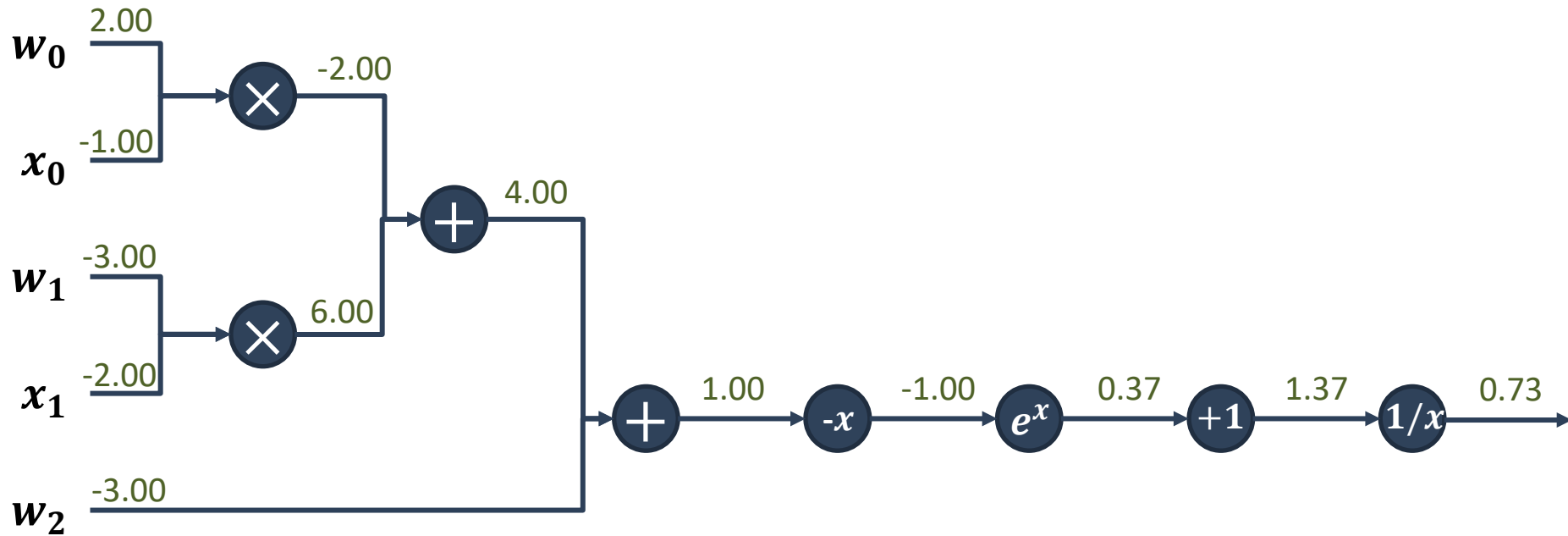
Different Example

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



Different Example

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



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

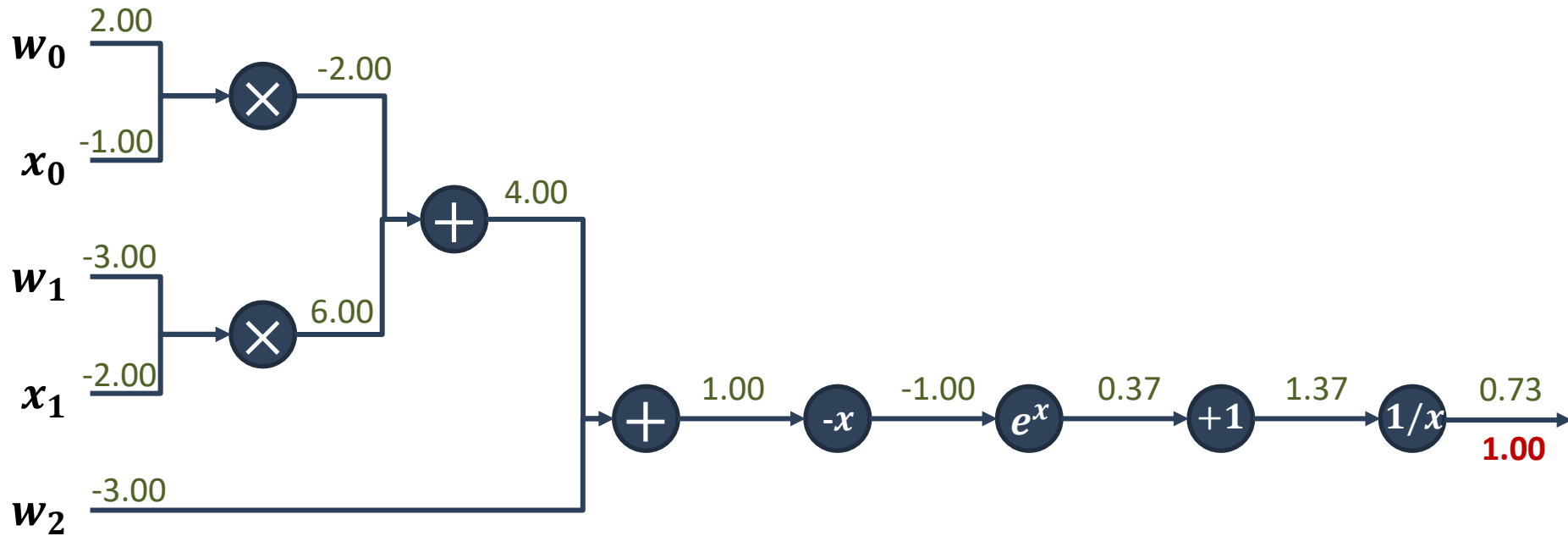
$$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$$

$$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$$

Different Example

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



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

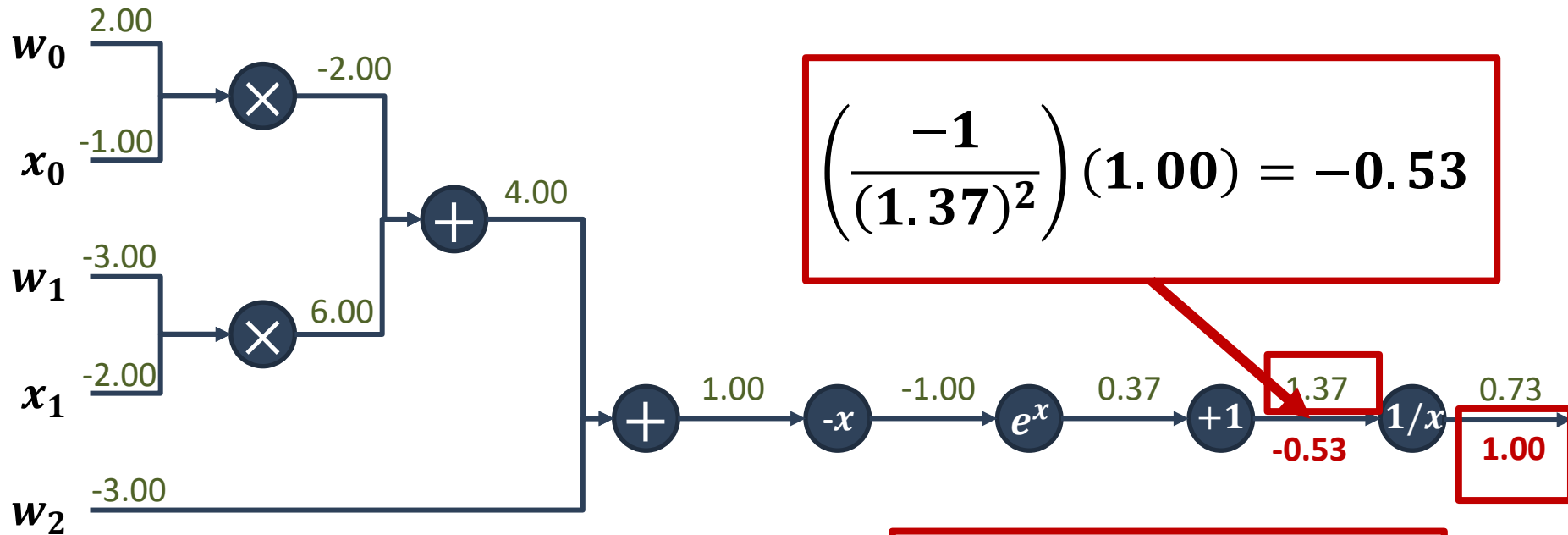
$$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$$

$$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$$

Different Example

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



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

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

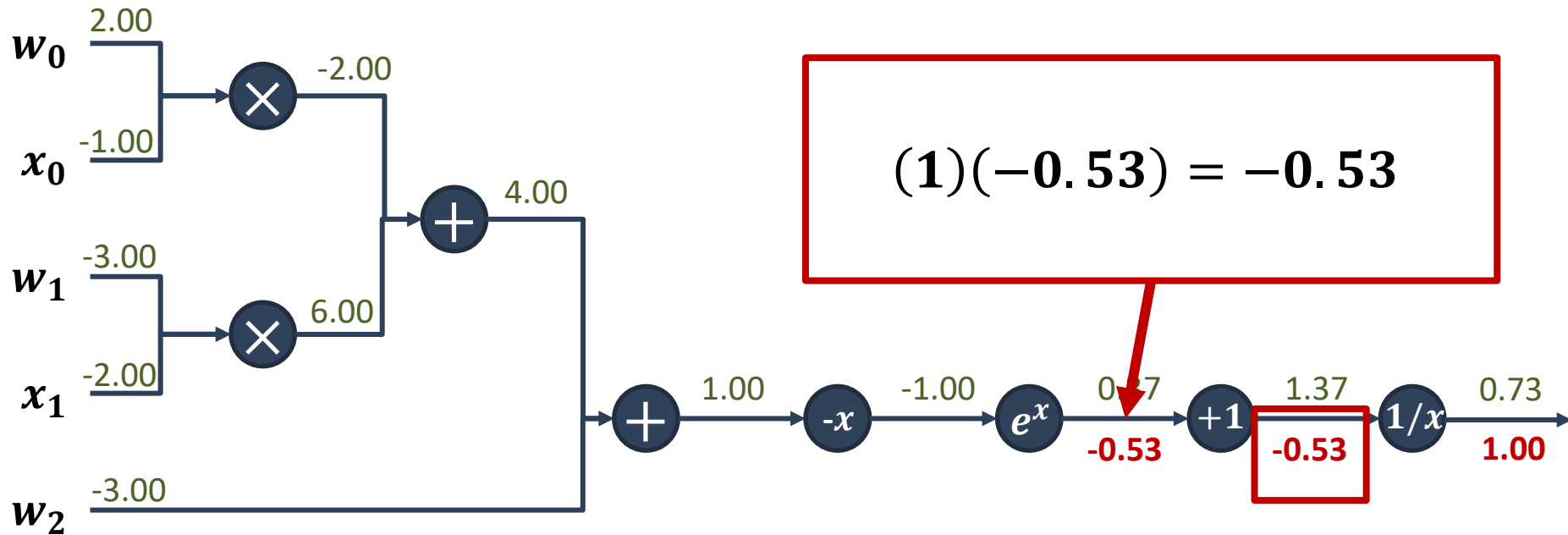
$$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$$

$$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$$

Different Example

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



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

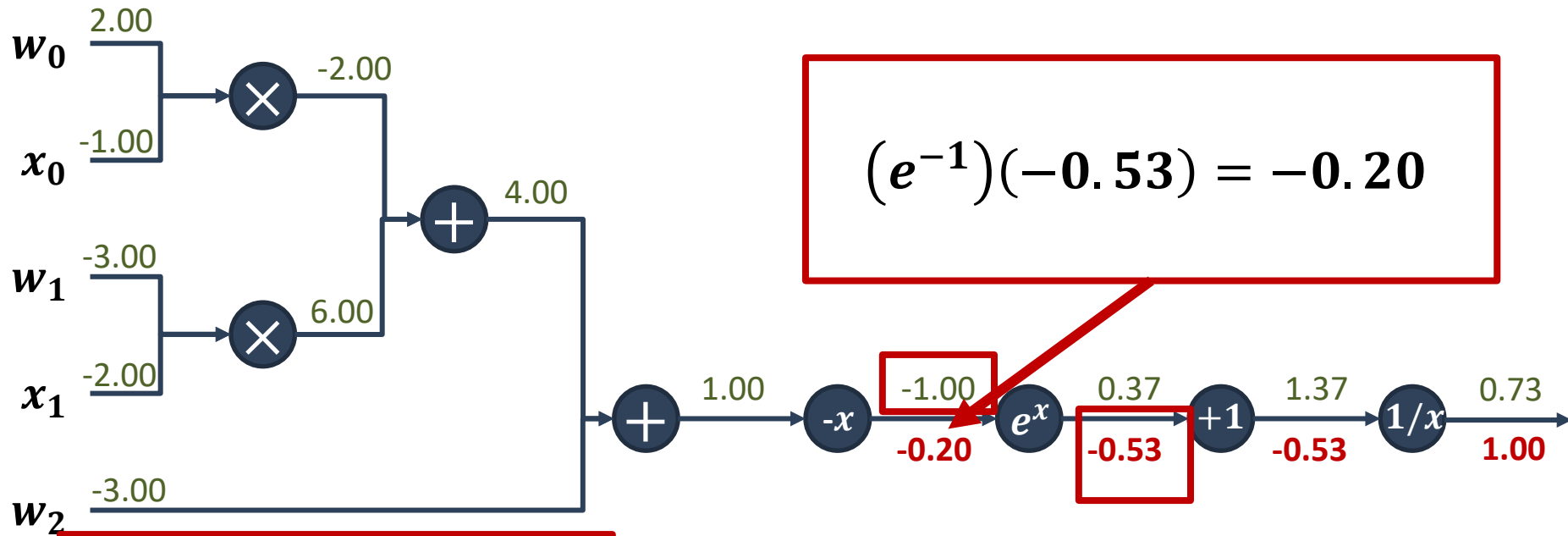
$$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$$

$$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$$

Different Example

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



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

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

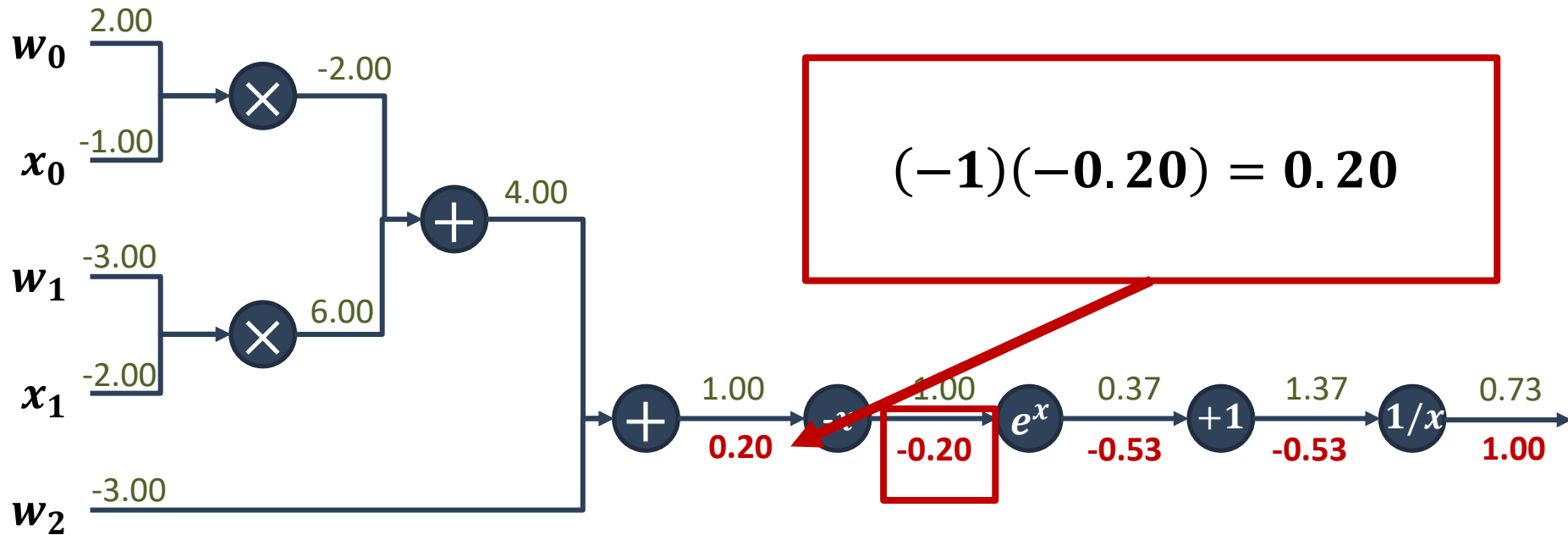
$$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$$

$$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$$

Different Example

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



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

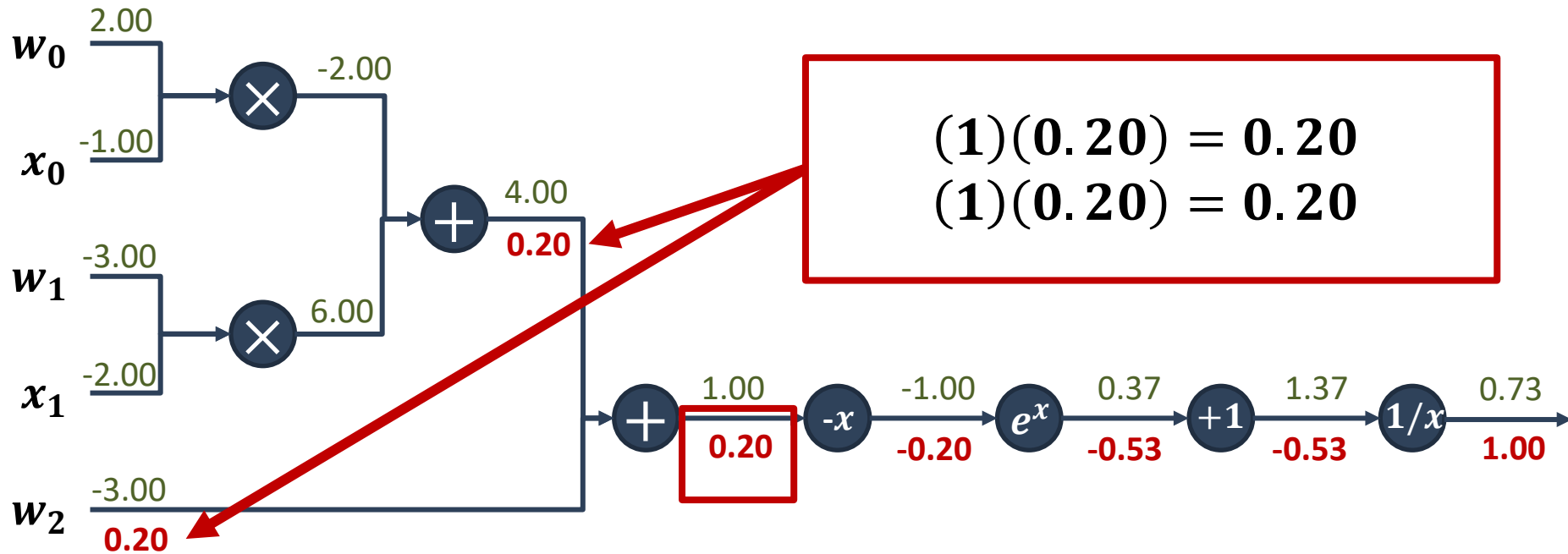
$$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$$

$$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$$

Different Example

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



$$(1)(0.20) = 0.20$$

$$(1)(0.20) = 0.20$$

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

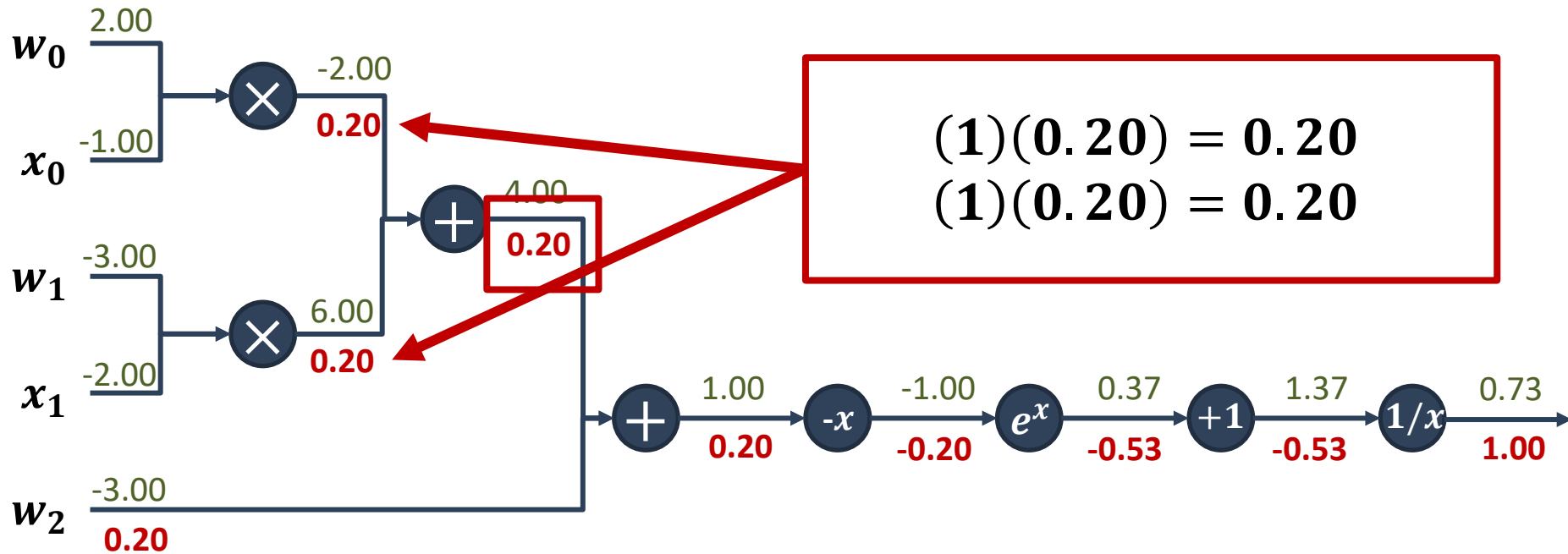
$$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$$

$$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$$

Different Example

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



$$(1)(0.20) = 0.20$$

$$(1)(0.20) = 0.20$$

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

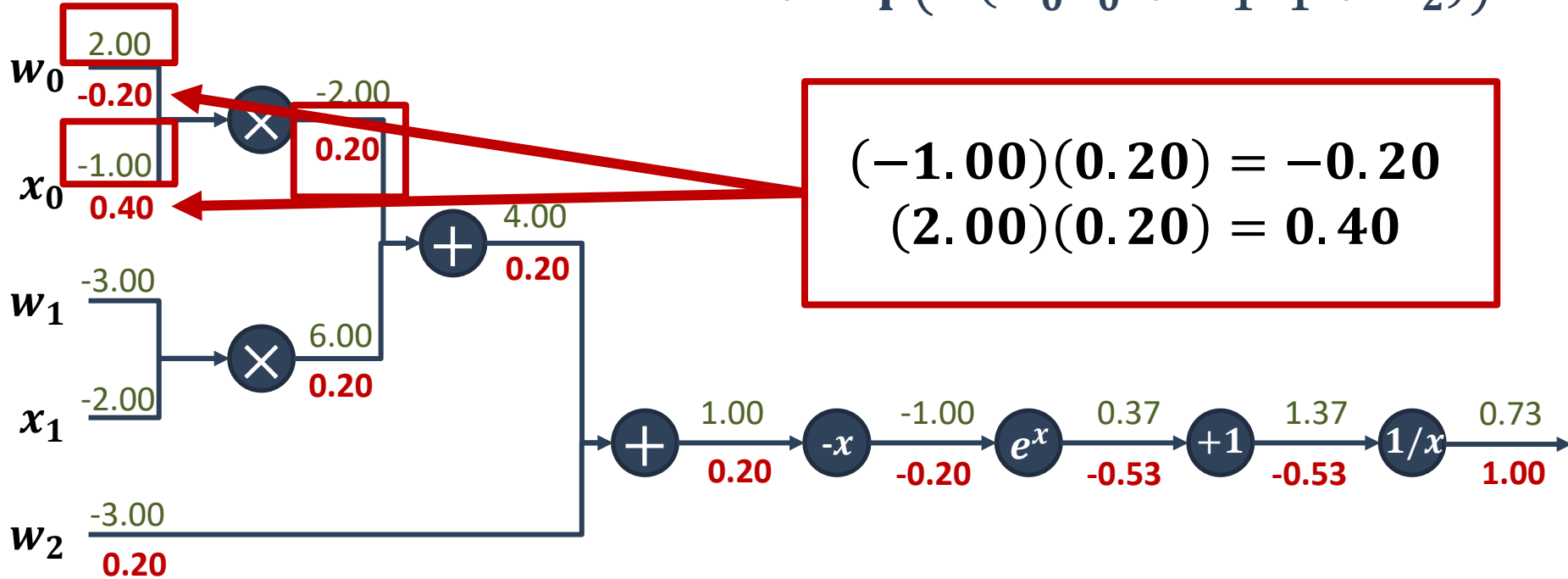
$$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$$

$$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$$

Different Example

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



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

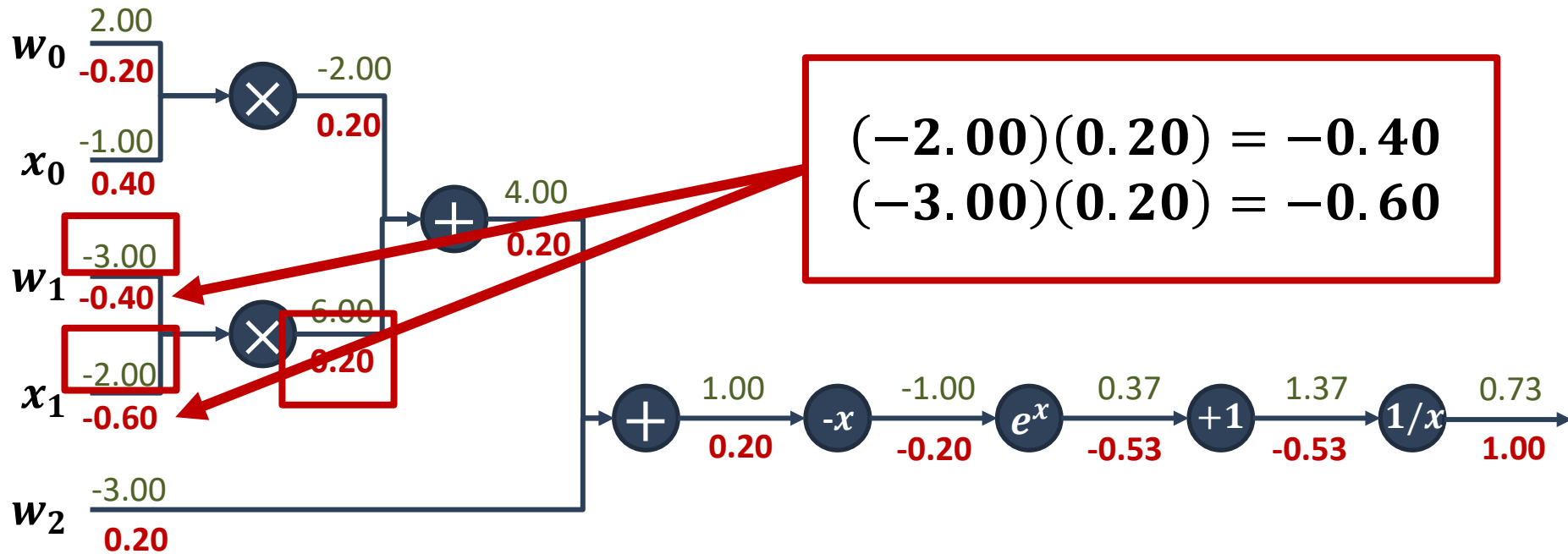
$$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$$

$$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$$

Different Example

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



$$(-2.00)(0.20) = -0.40$$

$$(-3.00)(0.20) = -0.60$$

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

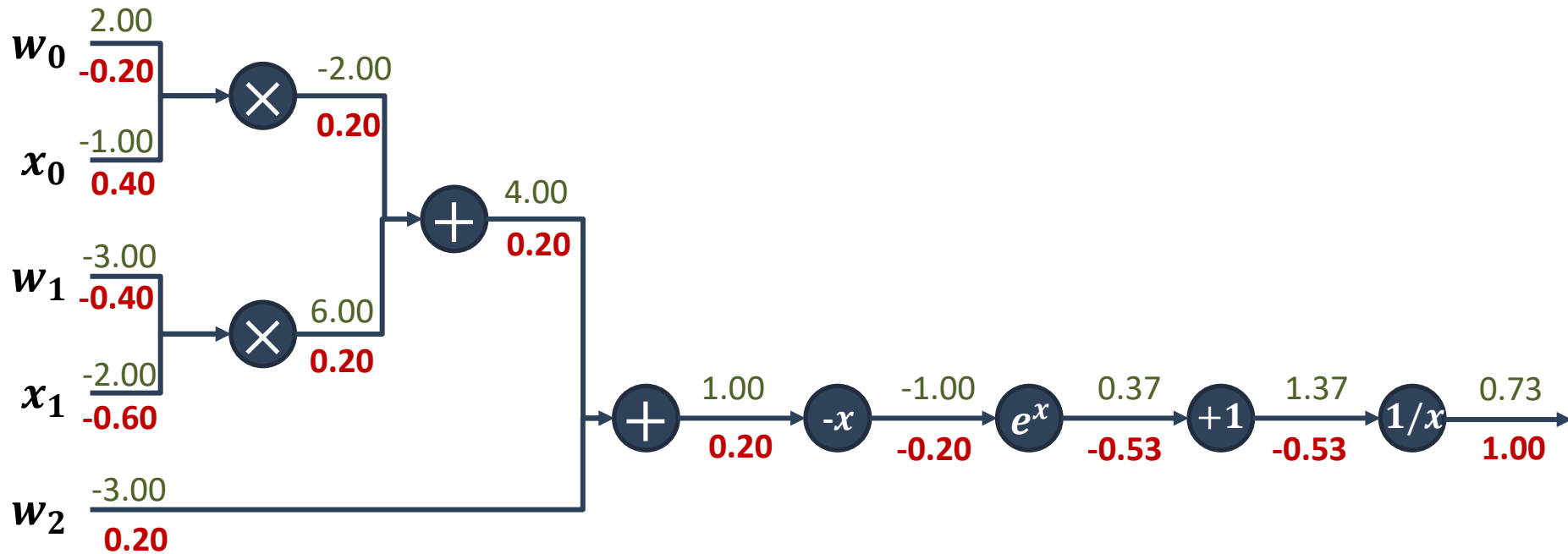
$$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$$

$$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$$

Different Example

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



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

$$f(x) = \frac{1}{x} \rightarrow \frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

$$f_a(x) = ax \rightarrow \frac{\partial f}{\partial x} = a$$

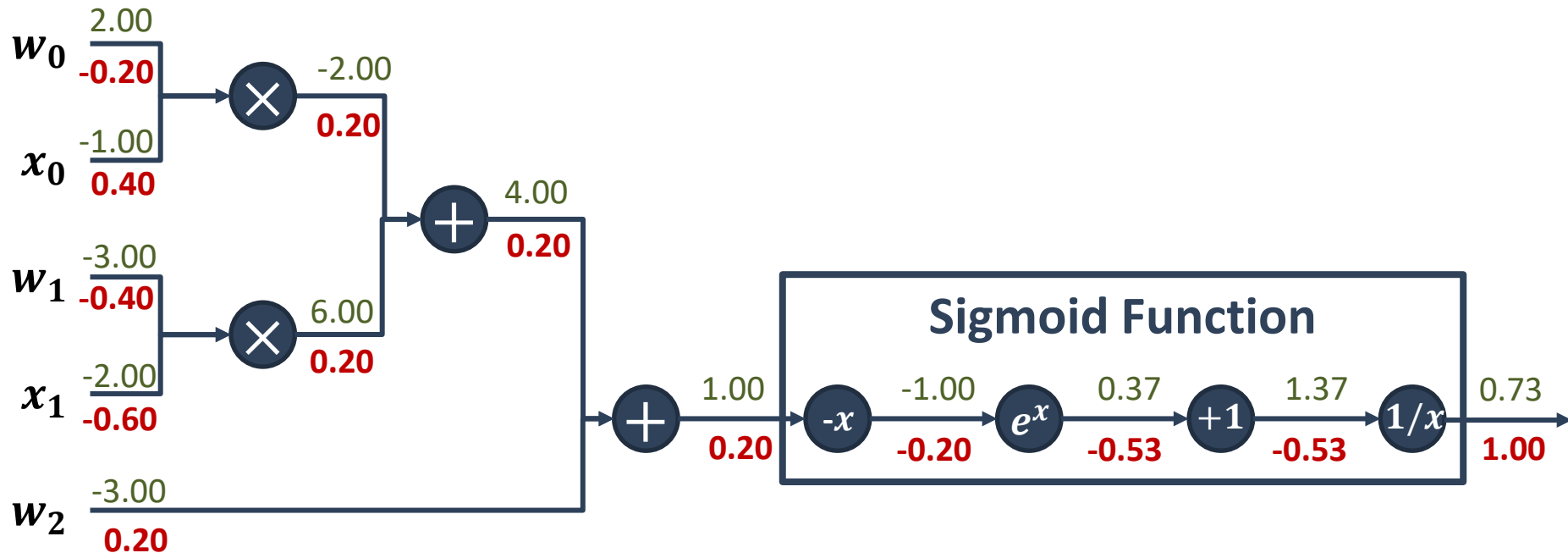
$$f_c(x) = c + x \rightarrow \frac{\partial f}{\partial x} = 1$$

Different Example: Sigmoid Function

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

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

$$\frac{\partial \sigma}{\partial x} = (1 - \sigma(x))\sigma(x)$$

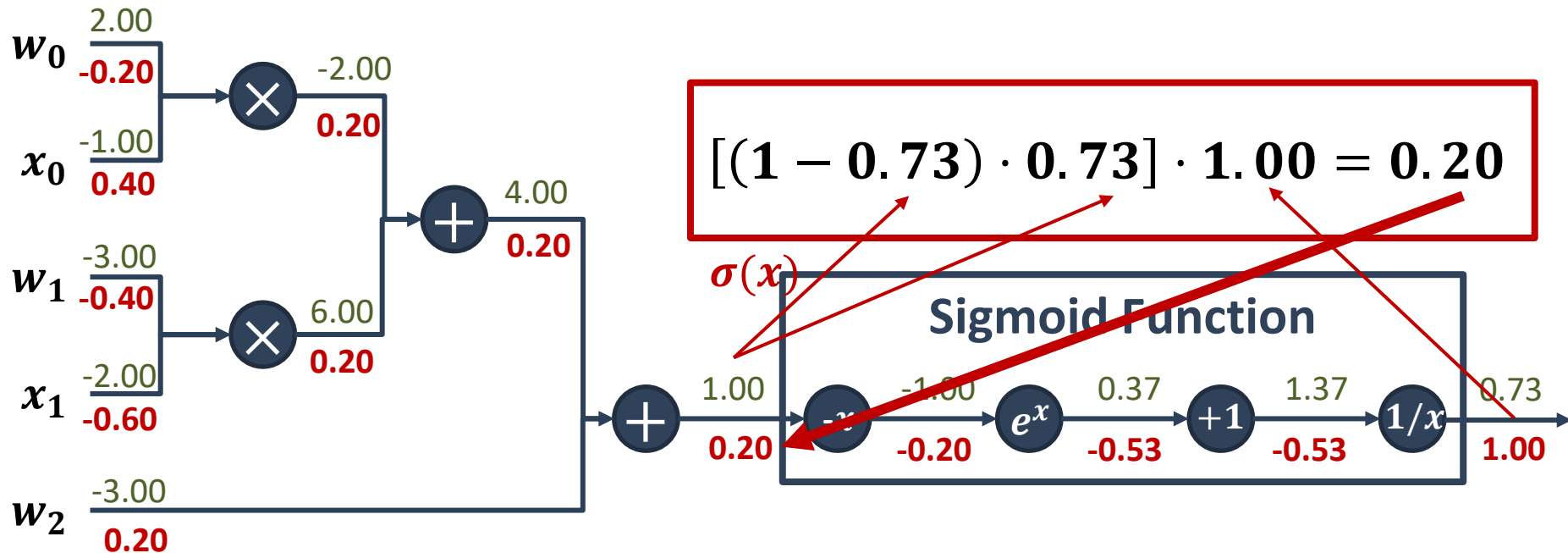


Different Example: Sigmoid Function

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

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

$$\frac{\partial \sigma}{\partial x} = (1 - \sigma(x))\sigma(x)$$



Patterns in Backflow of the Gradient

- **add**

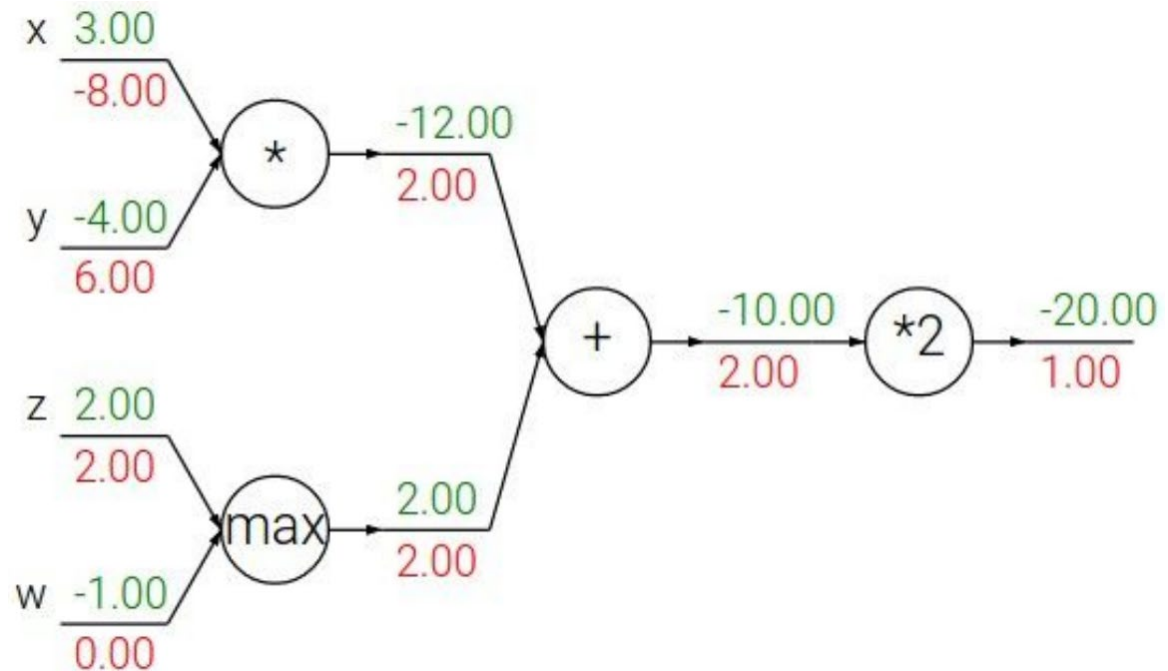
- Gradient distributor

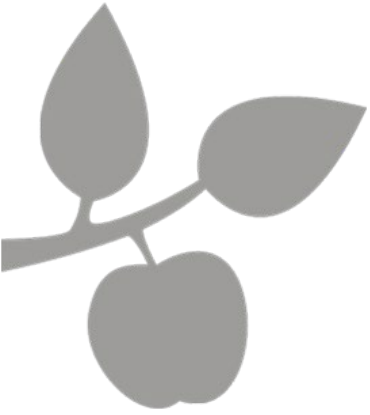
- **max**

- Gradient router

- **mul**

- Gradient switcher





Backpropagation

- Function Principle
- **Generalization to Vectors**

Generalization to Vectors

- Suppose $\mathbf{x} \in \mathbb{R}^m, \mathbf{y} \in \mathbb{R}^n$
 - g maps from \mathbb{R}^m to \mathbb{R}^n and
 - f maps from \mathbb{R}^n to \mathbb{R}

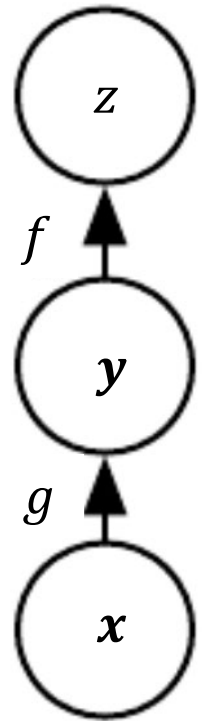
- If $\mathbf{y} = g(\mathbf{x})$ and $z = f(\mathbf{y})$, then

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \cdot \frac{\partial y_j}{\partial x_i}$$

- Or, in vector notation:

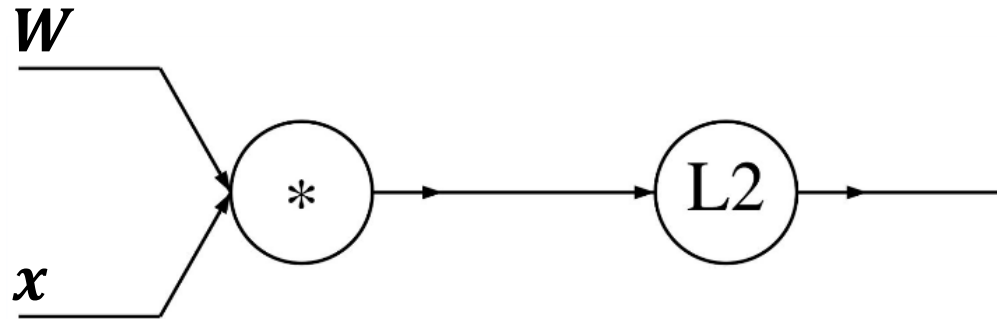
$$\nabla_{\mathbf{x}} z = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^T \nabla_{\mathbf{y}} z$$

- That is the product of the Jacobian matrix $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$ and the gradient vector $\nabla_{\mathbf{y}} z$.



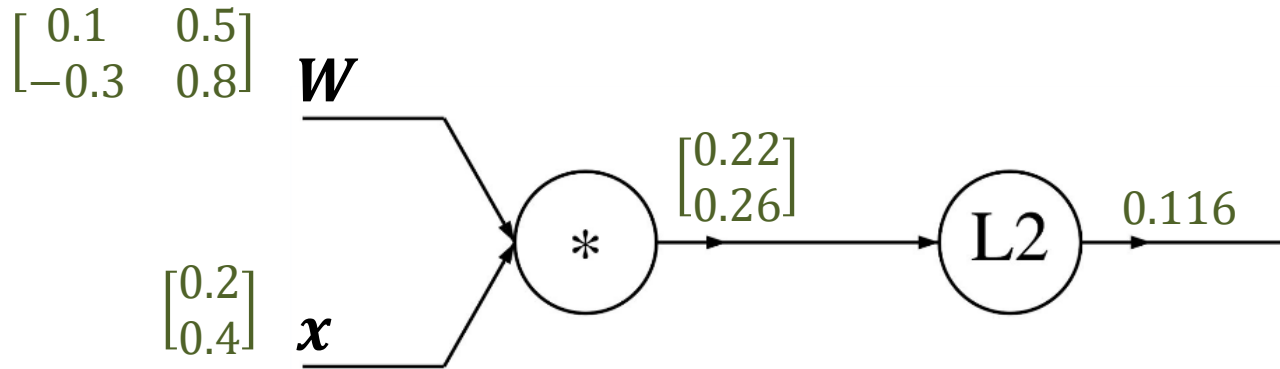
Vectorized Example

$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W}\mathbf{x}\|^2 = \sum_i^n (\mathbf{W}\mathbf{x})_i^2$$



Vectorized Example

$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W}\mathbf{x}\|^2 = \sum_i^n (\mathbf{W}\mathbf{x})_i^2$$

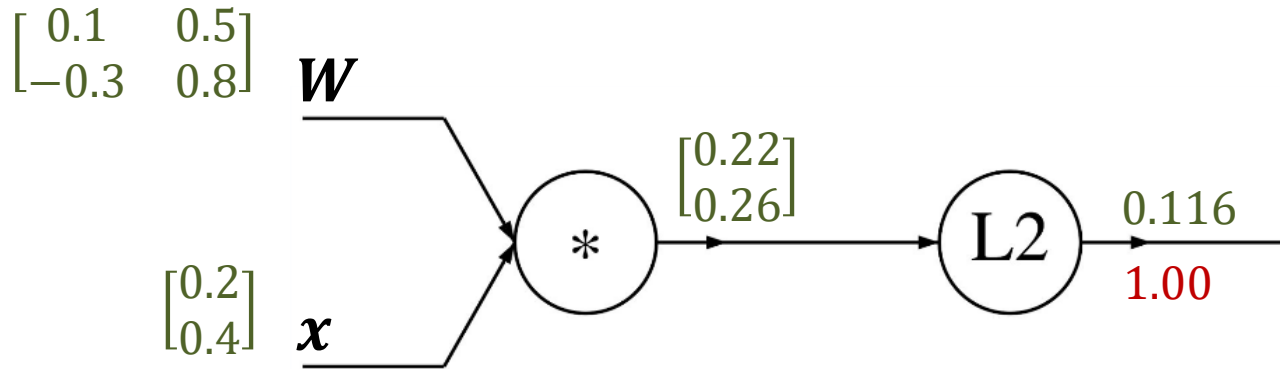


$$\mathbf{q} = \mathbf{W}\mathbf{x} = \begin{pmatrix} W_{1,1}x_1 & \cdots & W_{1,n}x_n \\ \vdots & \ddots & \vdots \\ W_{n,1}x_1 & \cdots & W_{n,n}x_n \end{pmatrix}$$

$$f(\mathbf{q}) = \|\mathbf{q}\|^2 = q_1^2 + \cdots + q_n^2$$

Vectorized Example

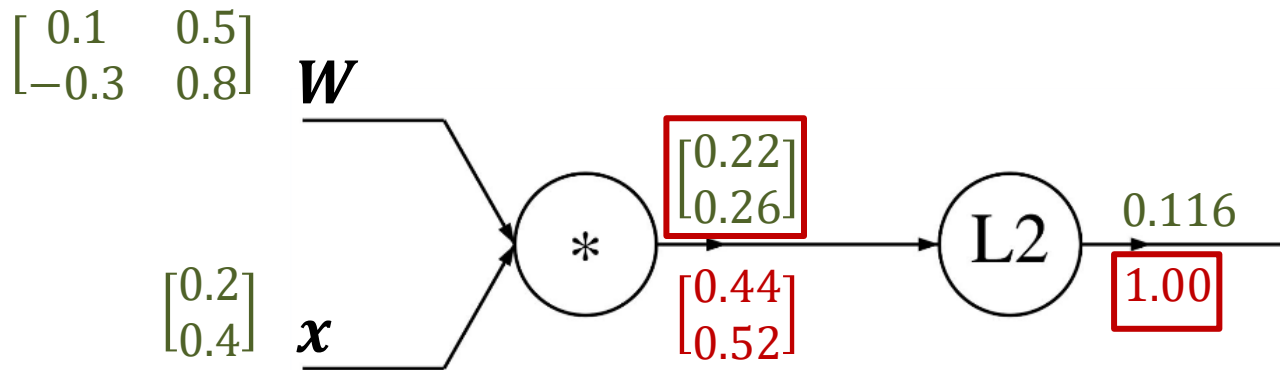
$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W}\mathbf{x}\|^2 = \sum_i^n (\mathbf{W}\mathbf{x})_i^2$$



$$\mathbf{q} = \mathbf{W}\mathbf{x} = \begin{pmatrix} W_{1,1}x_1 & \cdots & W_{1,n}x_n \\ \vdots & \ddots & \vdots \\ W_{n,1}x_1 & \cdots & W_{n,n}x_n \end{pmatrix}$$
$$f(\mathbf{q}) = \|\mathbf{q}\|^2 = q_1^2 + \cdots + q_n^2$$

Vectorized Example

$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W}\mathbf{x}\|^2 = \sum_i^n (\mathbf{W}\mathbf{x})_i^2$$



$$\mathbf{q} = \mathbf{W}\mathbf{x} = \begin{pmatrix} W_{1,1}x_1 & \cdots & W_{1,n}x_n \\ \vdots & \ddots & \vdots \\ W_{n,1}x_1 & \cdots & W_{n,n}x_n \end{pmatrix}$$

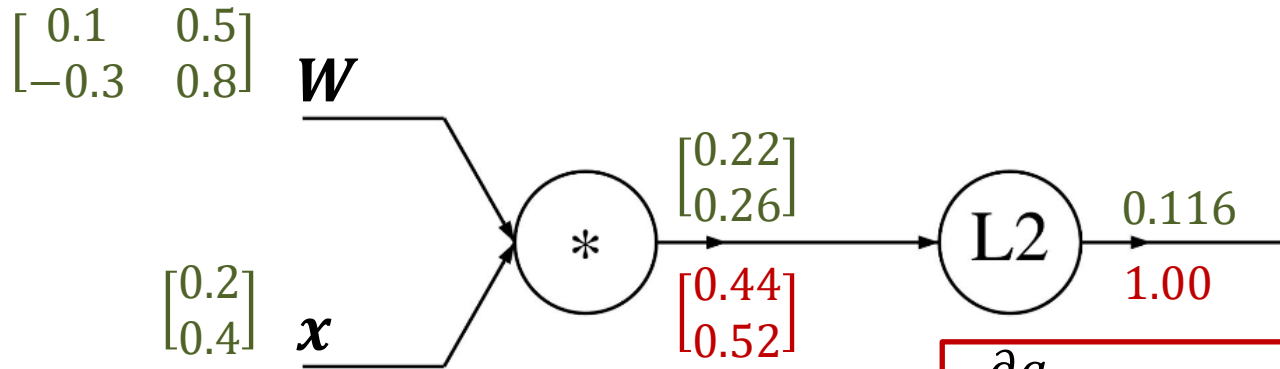
$$f(\mathbf{q}) = \|\mathbf{q}\|^2 = q_1^2 + \cdots + q_n^2$$

$$\frac{\partial f}{\partial q_i} = 2q_i$$

$$\nabla_{\mathbf{q}} f = 2\mathbf{q}$$

Vectorized Example

$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W}\mathbf{x}\|^2 = \sum_i^n (\mathbf{W}\mathbf{x})_i^2$$



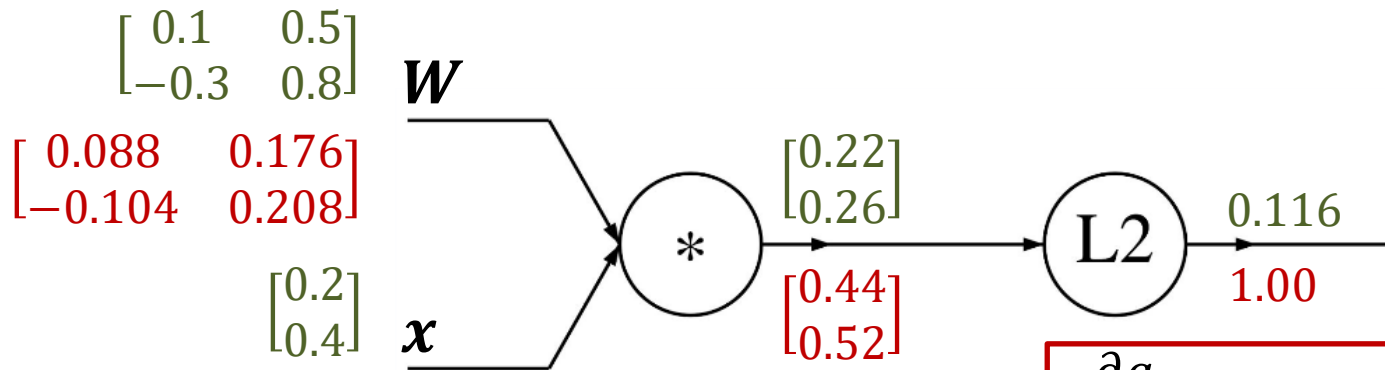
$$\mathbf{q} = \mathbf{W}\mathbf{x} = \begin{pmatrix} W_{1,1}x_1 & \cdots & W_{1,n}x_n \\ \vdots & \ddots & \vdots \\ W_{n,1}x_1 & \cdots & W_{n,n}x_n \end{pmatrix}$$

$$f(\mathbf{q}) = \|\mathbf{q}\|^2 = q_1^2 + \cdots + q_n^2$$

$$\begin{aligned} \frac{\partial q_k}{\partial W_{i,j}} &= \mathbf{1}_{k=i} x_j \\ \frac{\partial f}{\partial W_{i,j}} &= \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial W_{i,j}} \\ &= \sum_k (2q_k) (\mathbf{1}_{k=i} x_j) \\ &= 2q_i x_j \end{aligned}$$

Vectorized Example

$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W}\mathbf{x}\|^2 = \sum_i^n (\mathbf{W}\mathbf{x})_i^2$$



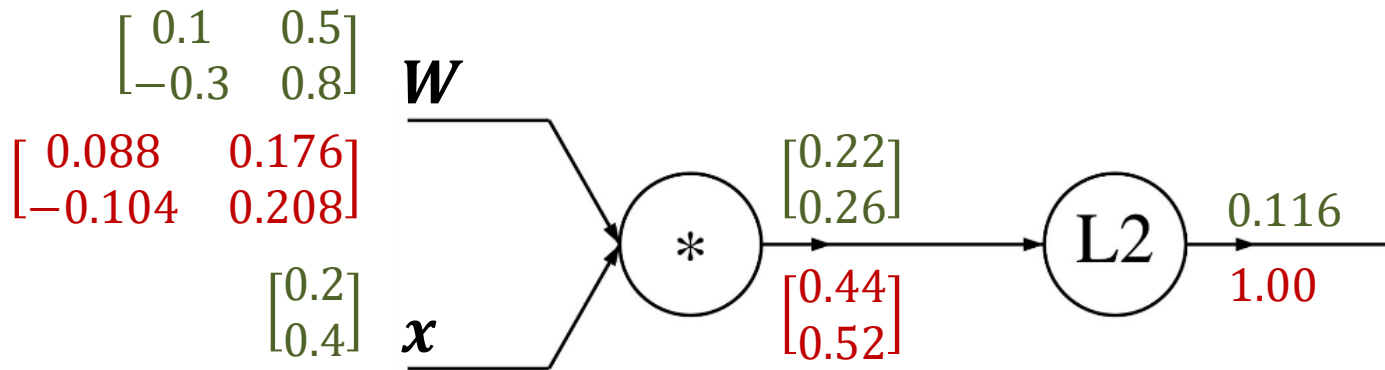
$$\mathbf{q} = \mathbf{W}\mathbf{x} = \begin{pmatrix} W_{1,1}x_1 & \cdots & W_{1,n}x_n \\ \vdots & \ddots & \vdots \\ W_{n,1}x_1 & \cdots & W_{n,n}x_n \end{pmatrix}$$

$$f(\mathbf{q}) = \|\mathbf{q}\|^2 = q_1^2 + \cdots + q_n^2$$

$$\begin{aligned} \frac{\partial q_k}{\partial W_{i,j}} &= \mathbf{1}_{k=i} x_j \\ \frac{\partial f}{\partial W} &= \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial W} \\ \nabla_{\mathbf{W}} f &= 2\mathbf{q} \cdot \mathbf{x}^T \\ &= \sum_k (2q_k) (\mathbf{1}_{k=i} x_j) \\ &= 2q_i x_j \end{aligned}$$

Vectorized Example

$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W}\mathbf{x}\|^2 = \sum_i^n (\mathbf{W}\mathbf{x})_i^2$$



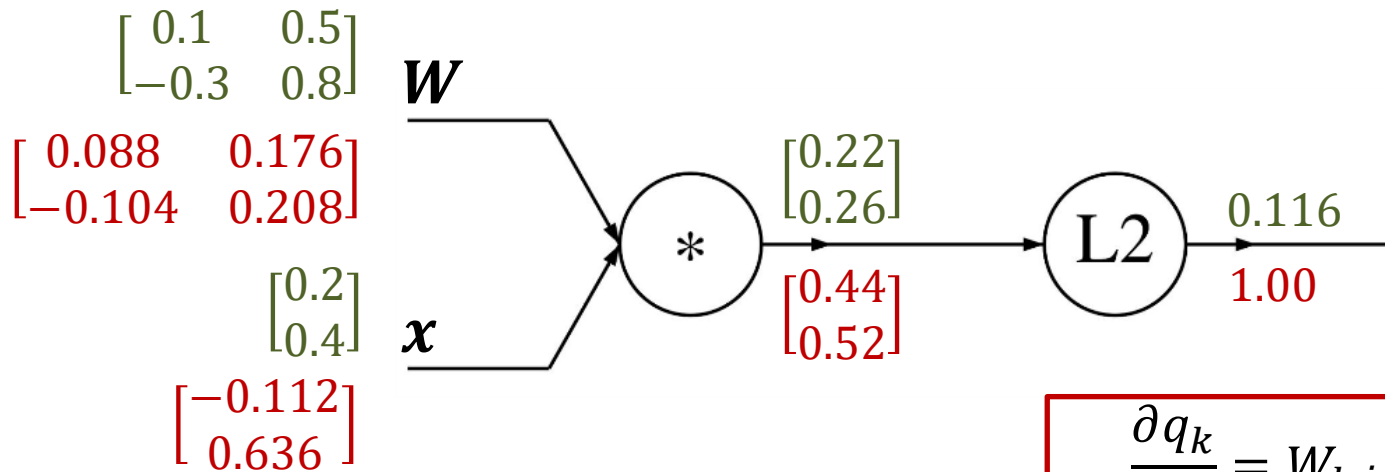
$$\mathbf{q} = \mathbf{W}\mathbf{x} = \begin{pmatrix} W_{1,1}x_1 & \cdots & W_{1,n}x_n \\ \vdots & \ddots & \vdots \\ W_{n,1}x_1 & \cdots & W_{n,n}x_n \end{pmatrix}$$

$$f(\mathbf{q}) = \|\mathbf{q}\|^2 = q_1^2 + \cdots + q_n^2$$

$$\begin{aligned} \frac{\partial q_k}{\partial x_i} &= W_{k,i} \\ \frac{\partial f}{\partial x_i} &= \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial x_i} \\ &= \sum_k 2q_k W_{k,i} \end{aligned}$$

Vectorized Example

$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W}\mathbf{x}\|^2 = \sum_i^n (\mathbf{W}\mathbf{x})_i^2$$



$$\mathbf{q} = \mathbf{W}\mathbf{x} = \begin{pmatrix} W_{1,1}x_1 & \cdots & W_{1,n}x_n \\ \vdots & \ddots & \vdots \\ W_{n,1}x_1 & \cdots & W_{n,n}x_n \end{pmatrix}$$

$$f(\mathbf{q}) = \|\mathbf{q}\|^2 = q_1^2 + \cdots + q_n^2$$

$$\frac{\partial q_k}{\partial x_i} = W_{k,i}$$

$$\frac{\partial f}{\partial \mathbf{x}} \sqsubset \frac{\partial f}{\partial \mathbf{q}_k}$$

$$\nabla_{\mathbf{x}} f = 2\mathbf{W}^T \mathbf{q}$$

$$= \sum_k 2q_k W_{k,i}$$

Two approaches to backpropagation

1. Symbol-to-number differentiation

- Take a computational graph and a set of numerical values for inputs to the graph
- Return a set of numerical values describing gradient at those input values
- Used by libraries: Torch and Caffe

2. Symbol-to-symbol differentiation

- Take a computational graph
- Add additional nodes to the graph that provide a symbolic description of desired derivatives
- Used by libraries: Theano and Tensorflow

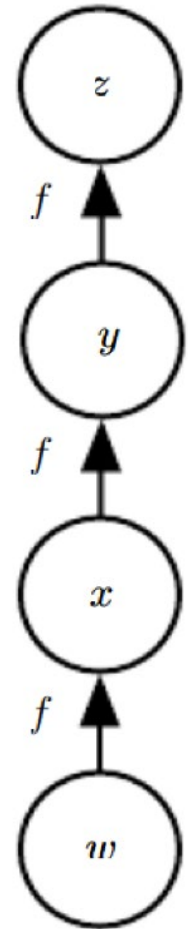
Symbol-to-symbol Derivatives

- To compute derivative using this approach, backpropagation does not need to ever access any actual numerical values
- Instead it adds nodes to a computational graph describing how to compute the derivatives
- A generic graph evaluation engine can later compute derivatives for any specific numerical values

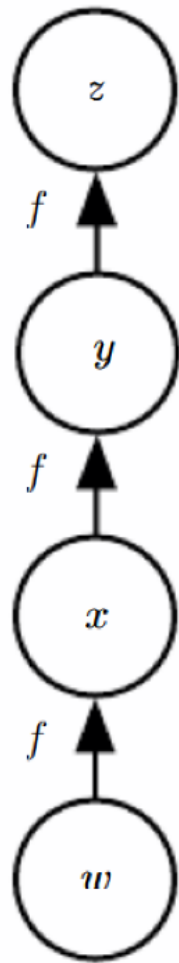
Example

- Consider the following function:

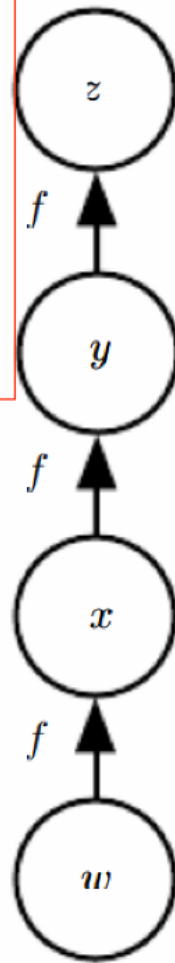
$$z = f(f(f(w)))$$



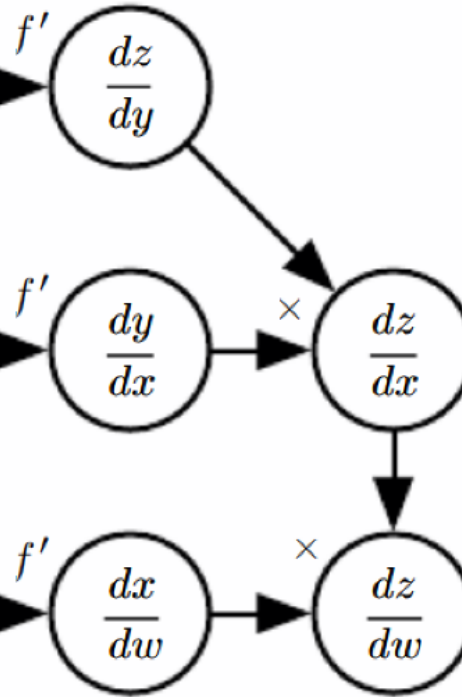
Symbol-to-Symbol Derivative Computation



We run BP instructing it to construct graph for expression corresponding to dz/dw



Result is a computational graph with a symbolic description of the derivative



Advantages of Approach

- Derivatives are described in the same language as the original expression
- Because the derivatives are just another computational graph, it is possible to run back-propagation again
 - Differentiating the derivatives
 - Yields higher-order derivatives

What is Back-Propagation and what not!

- Often simply called backprop
 - Allows information from the cost to flow back through network to compute gradient
- The backpropagation algorithm does this using a simple and inexpensive procedure (and some optimizations, like dynamic programming to avoid evaluating the same expression twice)
- Backpropagation **is not Learning**
 - Only refers to the method for computing gradients
 - Needs to be coupled with a learning algorithm, e.g., stochastic gradient descent
 - Backprop is NOT specific to Deep Learning