

CS 446/ECE 449: Machine Learning

Shenlong Wang

University of Illinois at Urbana-Champaign, 2024

L9: Deep Neural Networks

Goals of this lecture

- Understanding forward and backward pass
- Learning about backpropagation

Reading material

- I. Goodfellow et al.; Deep Learning; Chapters 6-9

Goals of this lecture

- Stochastic Gradient Decent (SGD)
- Understanding forward and backward pass
- Learning about backpropagation

Reading material

- I. Goodfellow et al.; Deep Learning; Chapters 6-9

Recall: Multi-Layer Neural Networks

- The function computed by the network is a composition of the functions computed by individual layer (e.g., linear layers and nonlinearities):



Neural networks as functions.

A linear predictor (**one layer network**) has the form

$$\mathbf{x} \mapsto \mathbf{w}^\top \mathbf{x}.$$

A **two layer network** has the form

$$\mathbf{x} \mapsto \mathbf{w}^\top \sigma_1(\mathbf{A}_1 \mathbf{x} + \mathbf{b}_1).$$

Iterating, a **multi-layer network** has the form

$$\mathbf{x} \mapsto \mathbf{w}^\top \sigma_1 \left(\mathbf{A}_1 \sigma_2 \left(\cdots \mathbf{A}_{L-2} \sigma_{L-1} \left(\mathbf{A}_{L-1} \mathbf{x} + \mathbf{b}_{L-1} \right) + \mathbf{b}_{L-2} \cdots \right) + \mathbf{b}_1 \right).$$

ERM now takes the form

$$\arg \min_{\mathbf{w}, \mathbf{A}_1, \dots, \mathbf{A}_{L-1}, \mathbf{b}_1, \dots, \mathbf{b}_{L-1}} \frac{1}{n} \sum_{i=1}^n \ell \left(y^{(i)} \mathbf{w}^\top \sigma_1 \left(\cdots \sigma_{L-1} (\mathbf{A}_{L-1} \mathbf{x}^{(i)} + \mathbf{b}_{L-1}) \cdots \right) \right).$$

How?

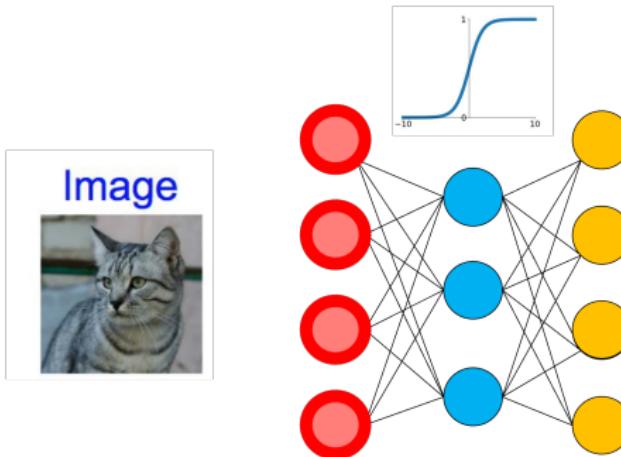
Back-propagation

Consider function: $y = f(g(x))$, what is the gradient of $\frac{\partial y}{\partial x}$?

Chain rule: Let $y = f(z)$, $z = g(x)$, then using chain rule we have:

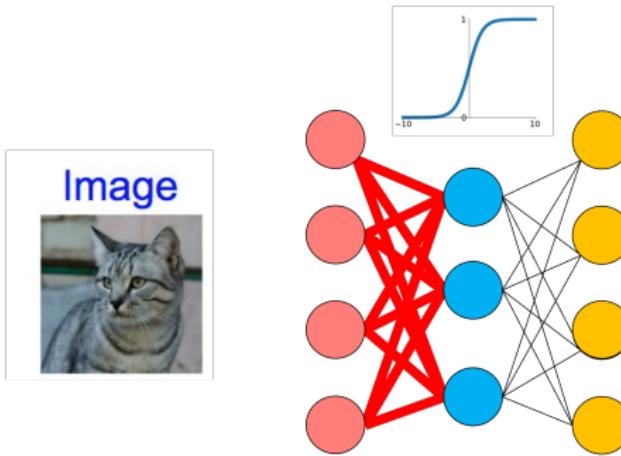
$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial z} \frac{\partial z}{\partial x}$$

Forward pass



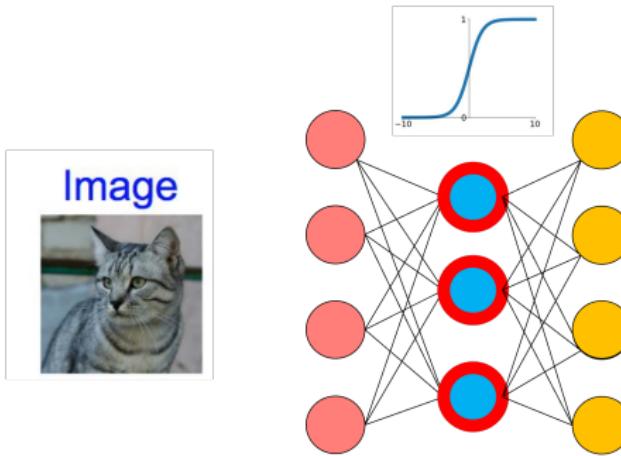
$$y = W_2 \sigma(W_1 x + b_1) + b_2$$

Forward pass



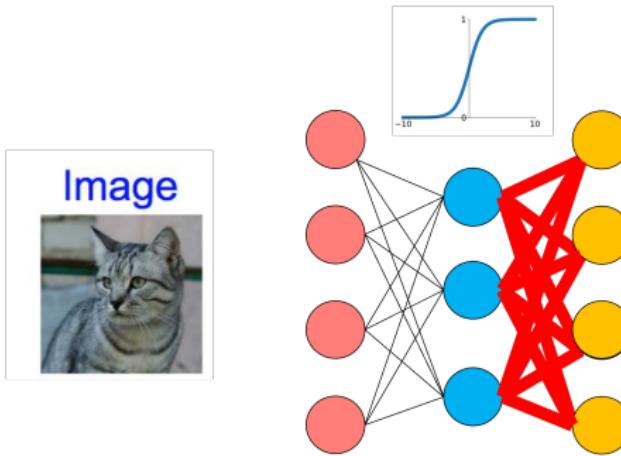
$$y = W_2 \sigma(W_1 x + b_1) + b_2$$

Forward pass



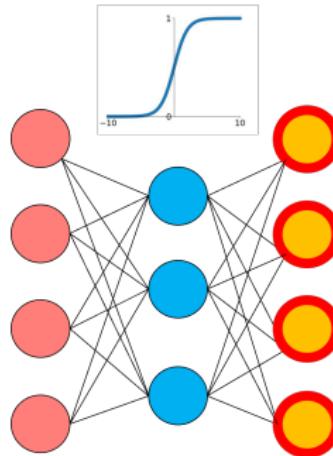
$$y = W_2 \sigma(W_1 x + b_1) + b_2$$

Forward pass



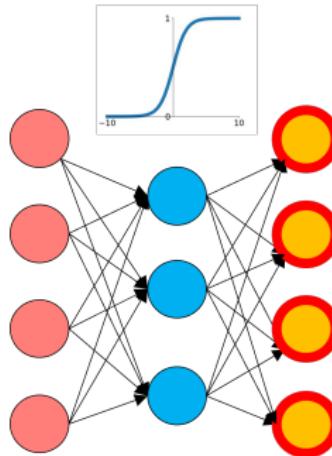
$$y = W_2 \sigma(W_1 x + b_1) + b_2$$

Forward pass



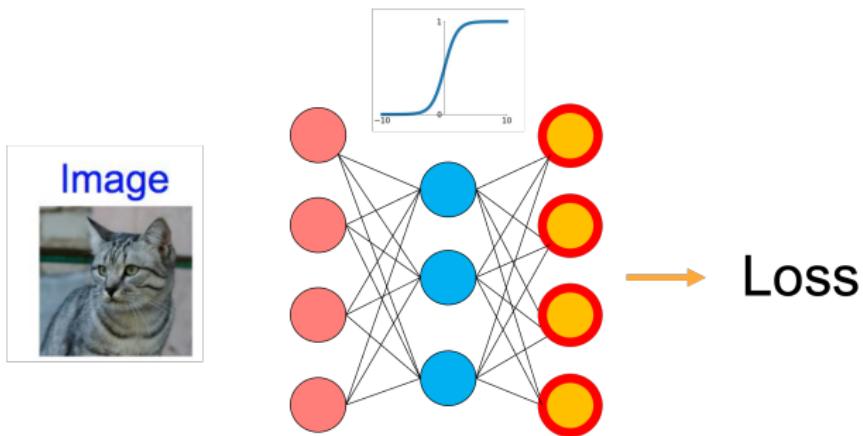
$$y = W_2 \sigma(W_1 x + b_1) + b_2$$

Forward pass



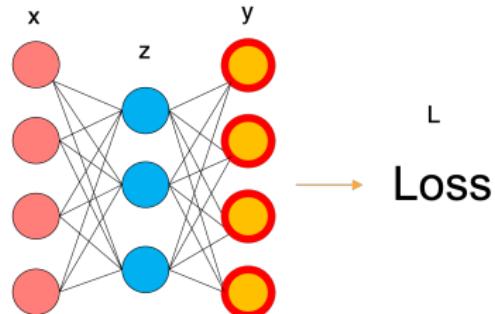
$$y = W_2 \sigma(W_1 x + b_1) + b_2$$

Backward pass



$$y = W_2 \sigma(W_1 x + b_1) + b_2$$

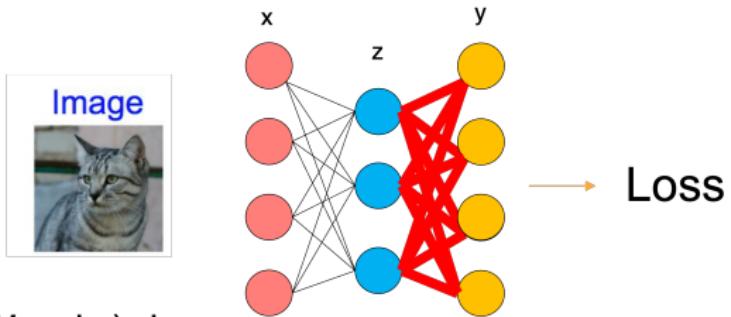
Backward pass



$$y = W_2 \sigma(W_1 x + b_1) + b_2$$

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

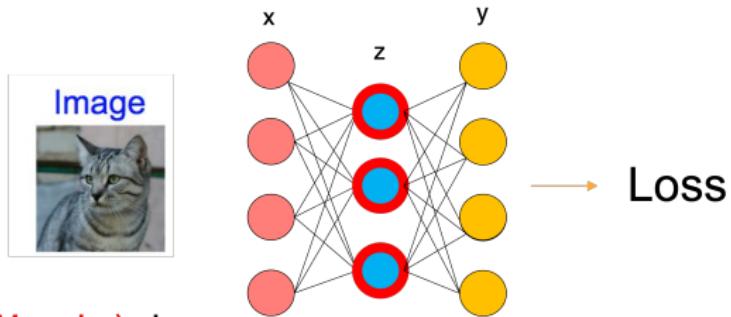
Backward pass



$$y = W_2 \sigma(W_1 x + b_1) + b_2$$

$$\frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial W_2}$$

Backward pass

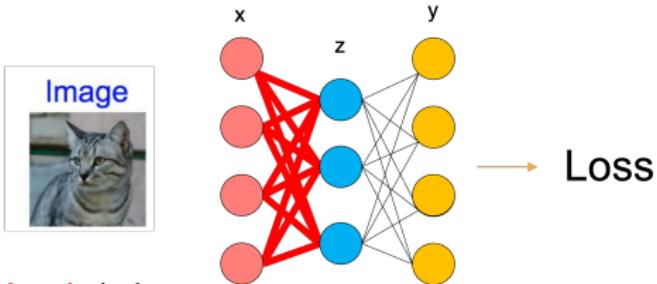


$$y = W_2 \sigma(W_1 x + b_1) + b_2$$

$$\frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial W_2}$$

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

Backward pass



$$y = W_2 \sigma(W_1 x + b_1) + b_2$$

$$\frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial W_2}$$

$$\boxed{\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial W_1}}$$

Recall: Multi-Layer Neural Networks

- The function computed by the network is a composition of the functions computed by individual layer (e.g., linear layers and nonlinearities):



How to implement this gradient computation?

Example:

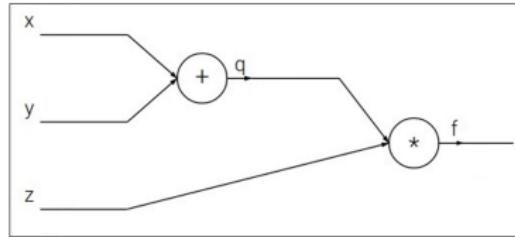
Derive each $\frac{\partial L}{\partial w}$ using chain rule on paper? —Bad idea! Why?

- Very tedious; lots of matrix calculus
- What if we want to change loss? Need to re-derive from scratch.
Not modular

Better idea: Computation graph.

Backpropagation: Simple Example

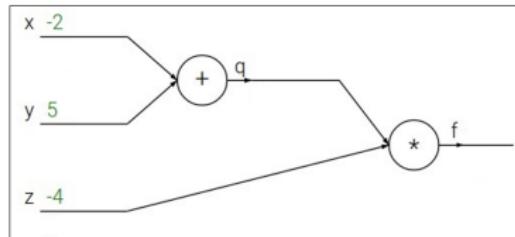
$$f(x, y, z) = (x + y)z$$



Backpropagation: Simple Example

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

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



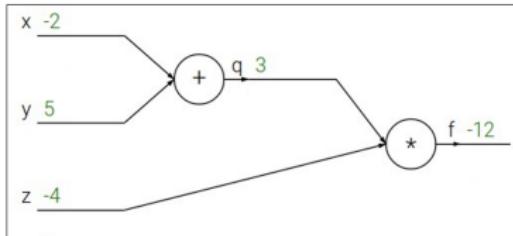
Backpropagation: Simple Example

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

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

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$



Backpropagation: Simple Example

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

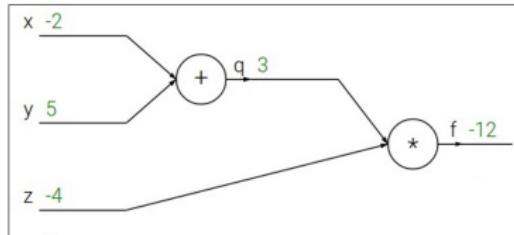
e.g. $x = -2, y = 5, z = -4$

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

2. Backward pass: Compute derivatives

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



Backpropagation: Simple Example

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

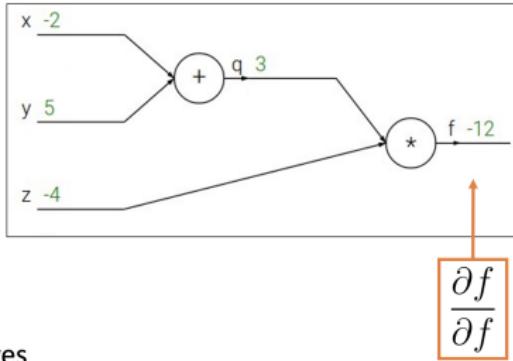
e.g. $x = -2, y = 5, z = -4$

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

2. Backward pass: Compute derivatives

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



Backpropagation: Simple Example

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

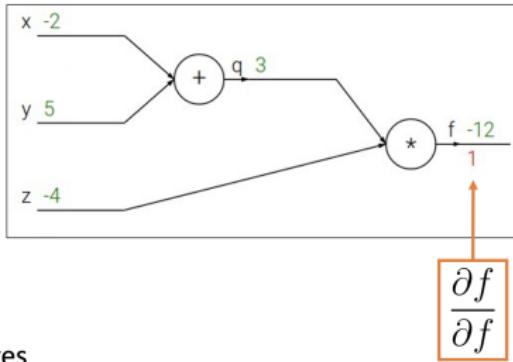
e.g. $x = -2, y = 5, z = -4$

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

2. Backward pass: Compute derivatives

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



Backpropagation: Simple Example

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

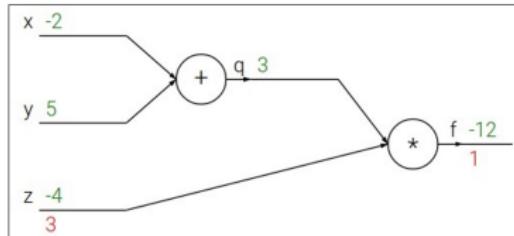
e.g. $x = -2, y = 5, z = -4$

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

2. Backward pass: Compute derivatives

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



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

Backpropagation: Simple Example

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

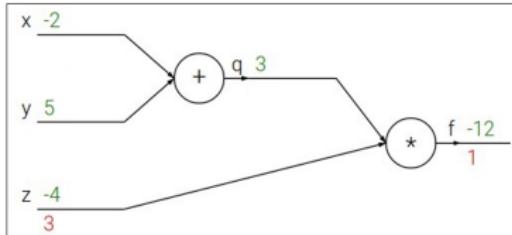
e.g. $x = -2, y = 5, z = -4$

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

2. Backward pass: Compute derivatives

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



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

Backpropagation: Simple Example

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

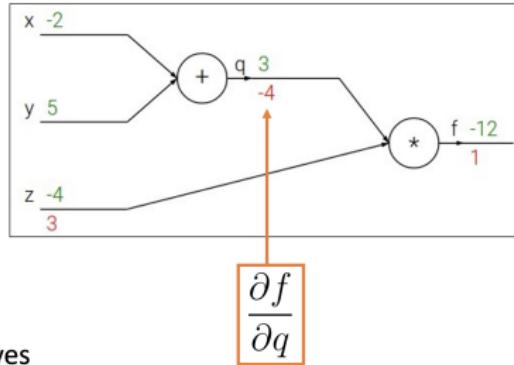
e.g. $x = -2, y = 5, z = -4$

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

2. Backward pass: Compute derivatives

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



Backpropagation: Simple Example

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

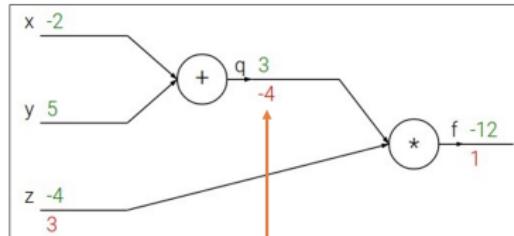
e.g. $x = -2, y = 5, z = -4$

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

2. Backward pass: Compute derivatives

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



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

Backpropagation: Simple Example

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

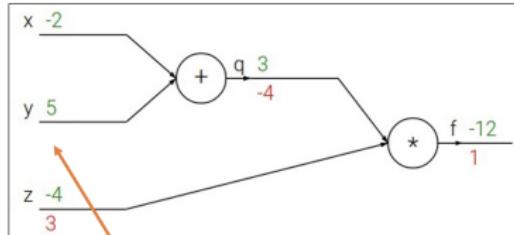
e.g. $x = -2, y = 5, z = -4$

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

2. Backward pass: Compute derivatives

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



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

Backpropagation: Simple Example

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

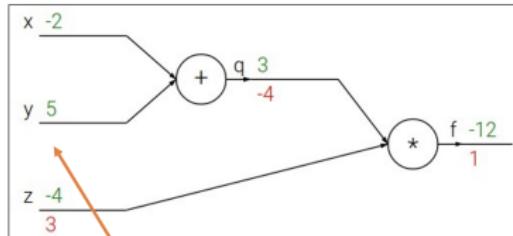
e.g. $x = -2, y = 5, z = -4$

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

2. Backward pass: Compute derivatives

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



Chain Rule

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

Backpropagation: Simple Example

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

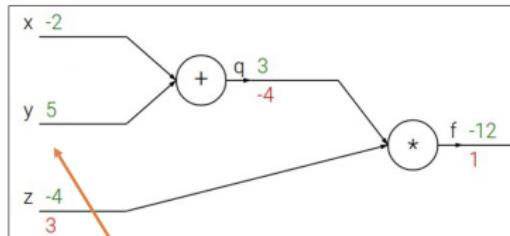
e.g. $x = -2, y = 5, z = -4$

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

2. Backward pass: Compute derivatives

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



Chain Rule

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

Downstream Gradient Local Gradient Upstream Gradient

Backpropagation: Simple Example

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

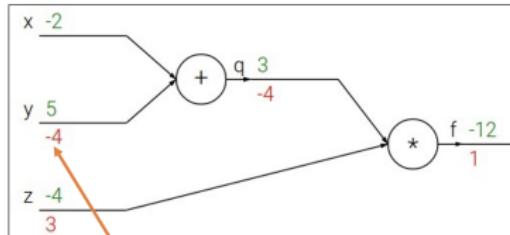
e.g. $x = -2, y = 5, z = -4$

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

2. Backward pass: Compute derivatives

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



Chain Rule

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

Downstream Gradient Local Gradient Upstream Gradient

Backpropagation: Simple Example

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

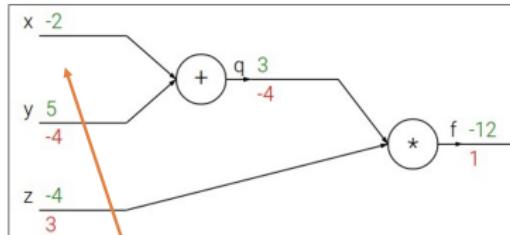
e.g. $x = -2, y = 5, z = -4$

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

2. Backward pass: Compute derivatives

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



Chain Rule

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

$$\frac{\partial q}{\partial x} = 1$$

Downstream Gradient Local Gradient Upstream Gradient

Backpropagation: Simple Example

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

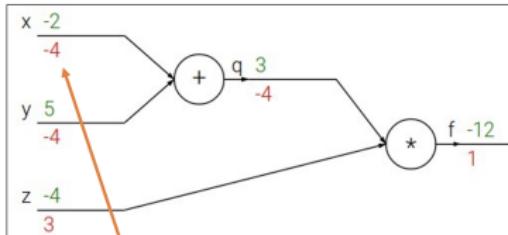
e.g. $x = -2, y = 5, z = -4$

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

2. Backward pass: Compute derivatives

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

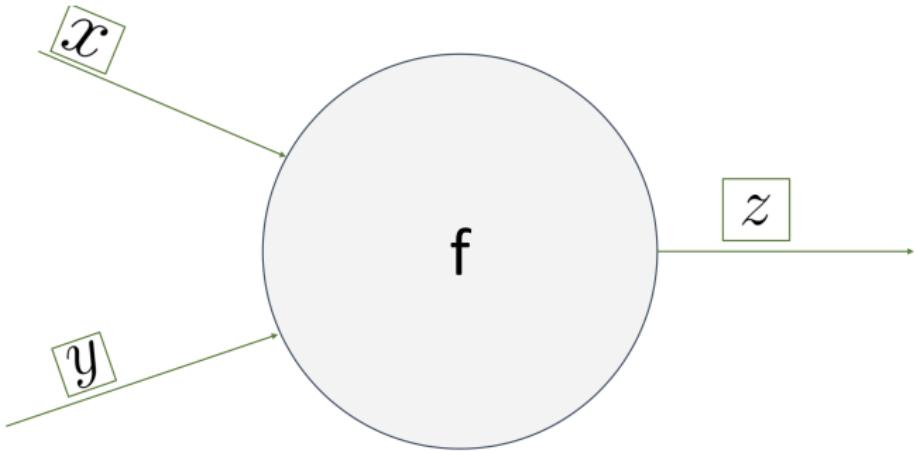


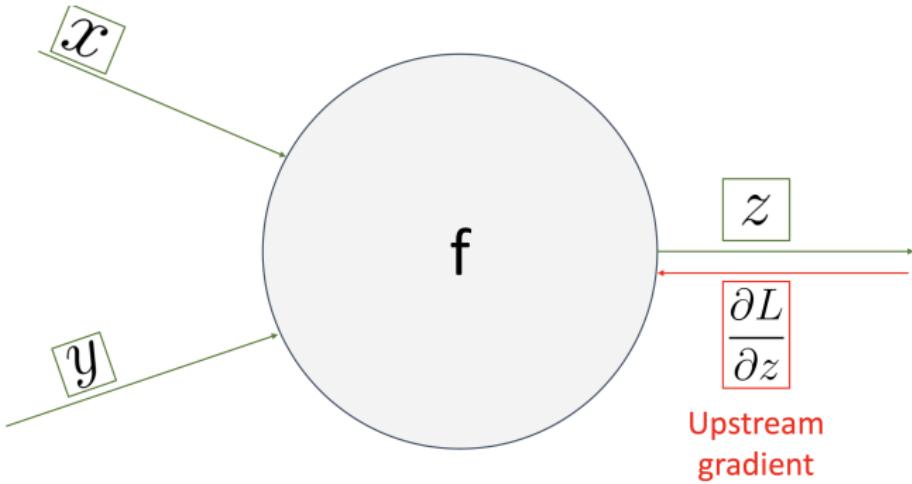
Chain Rule

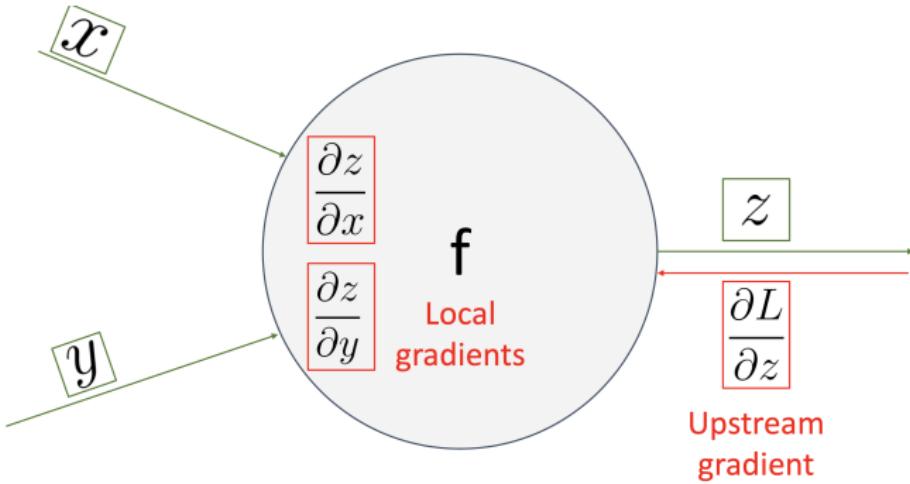
$$\frac{\partial f}{\partial x} = \frac{\partial q}{\partial x} \frac{\partial f}{\partial q}$$

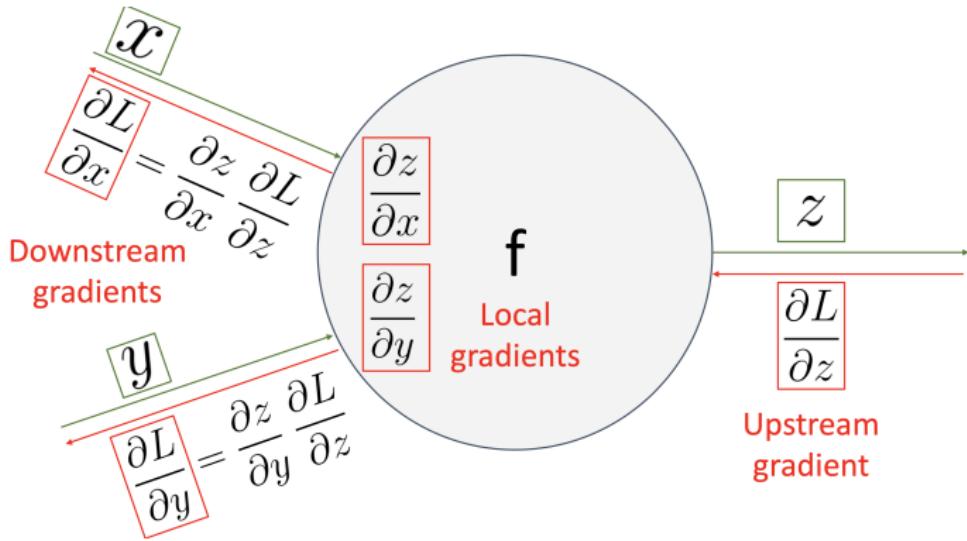
$$\frac{\partial q}{\partial x} = 1$$

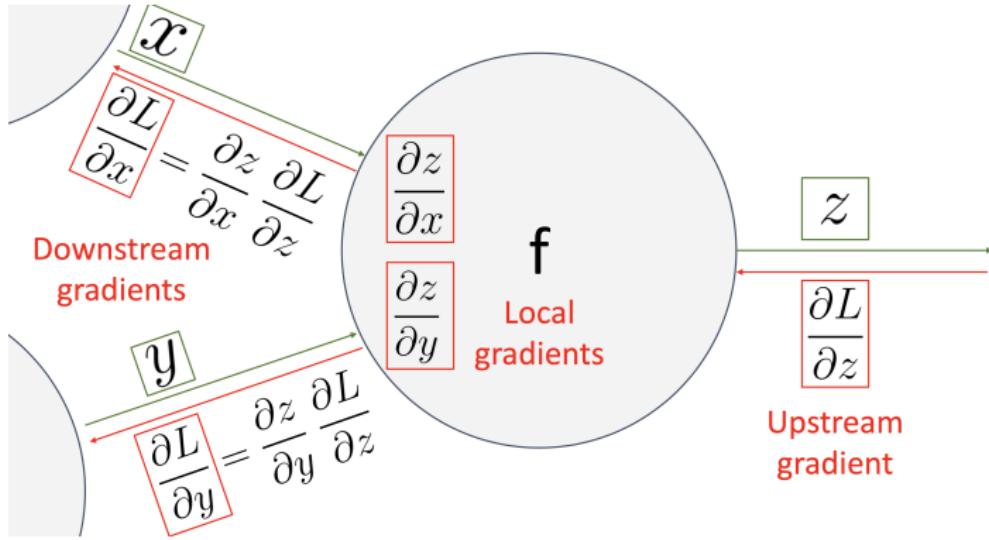
Downstream Gradient Local Gradient Upstream Gradient



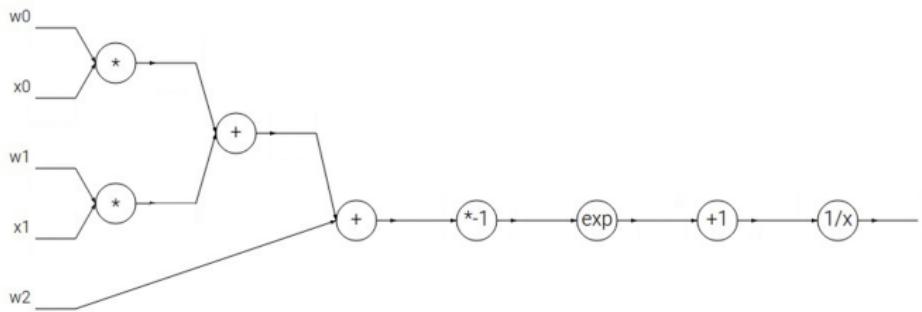








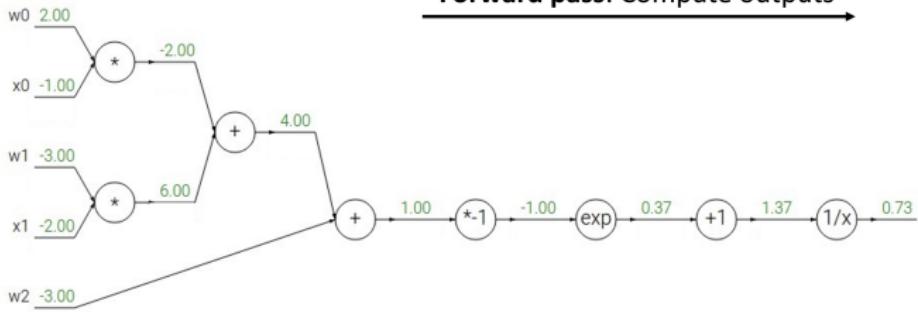
Another Example $f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



Another Example

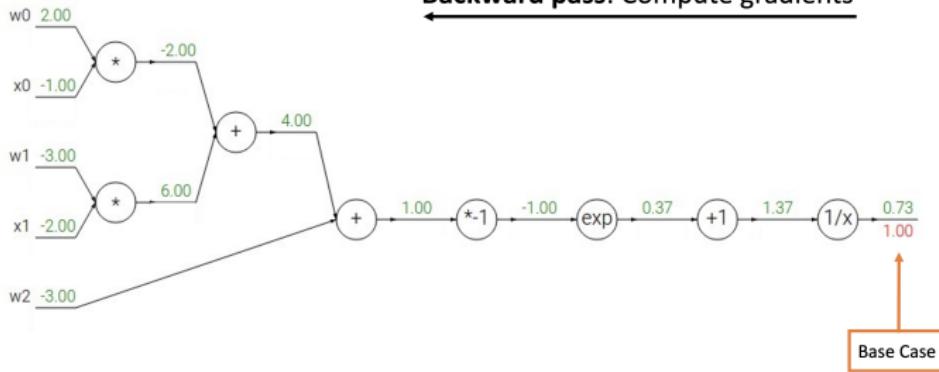
$$f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

Forward pass: Compute outputs



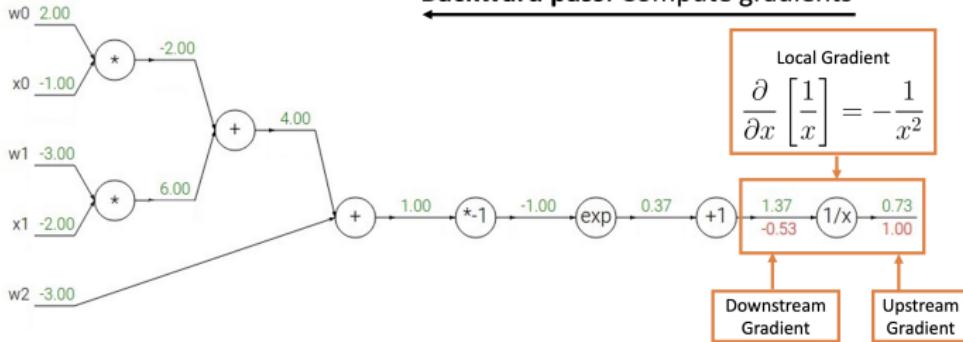
Another Example $f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$

Backward pass: Compute gradients



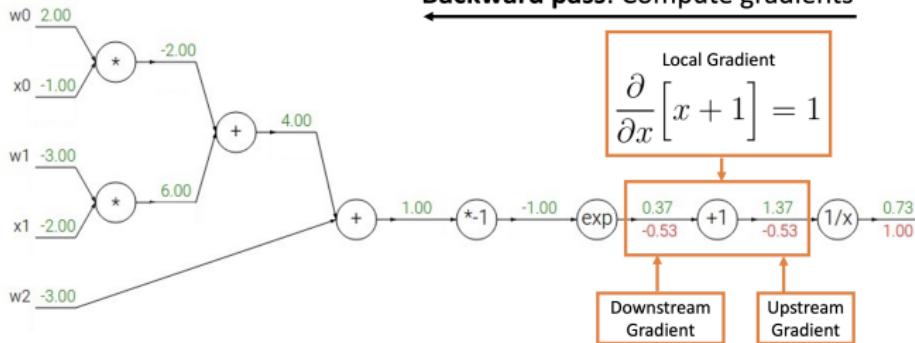
Another Example $f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$

Backward pass: Compute gradients



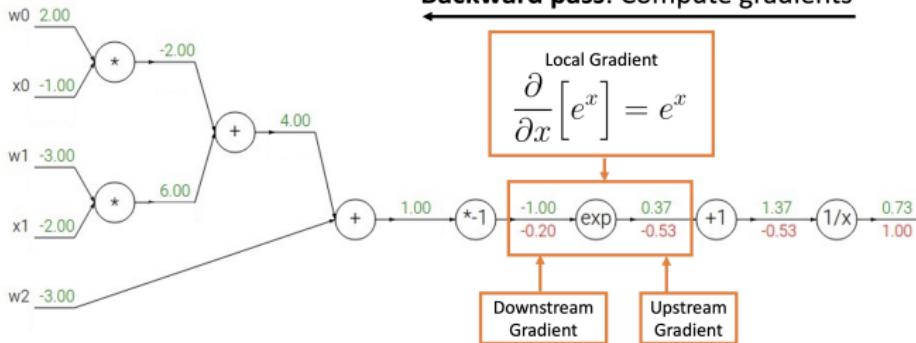
Another Example $f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$

Backward pass: Compute gradients



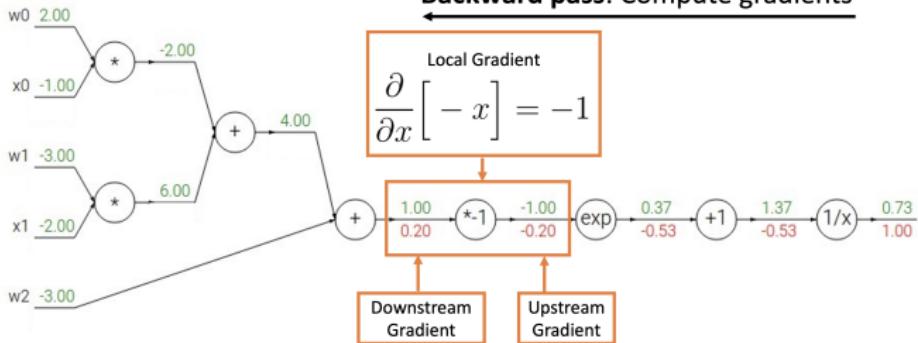
Another Example $f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$

Backward pass: Compute gradients



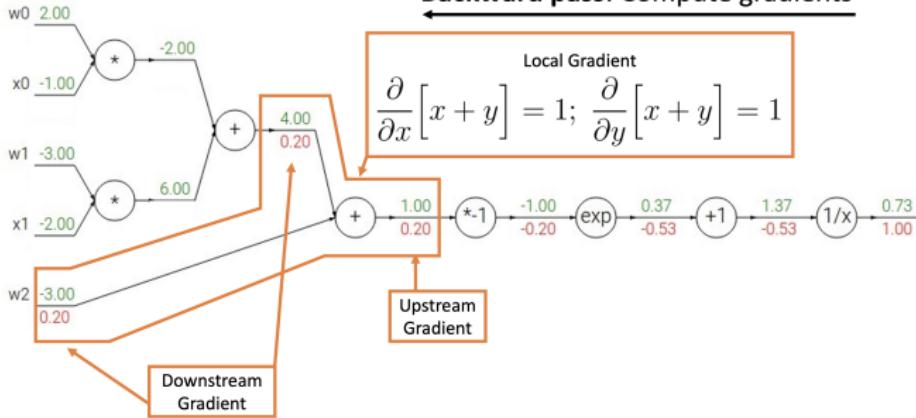
Another Example $f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$

Backward pass: Compute gradients



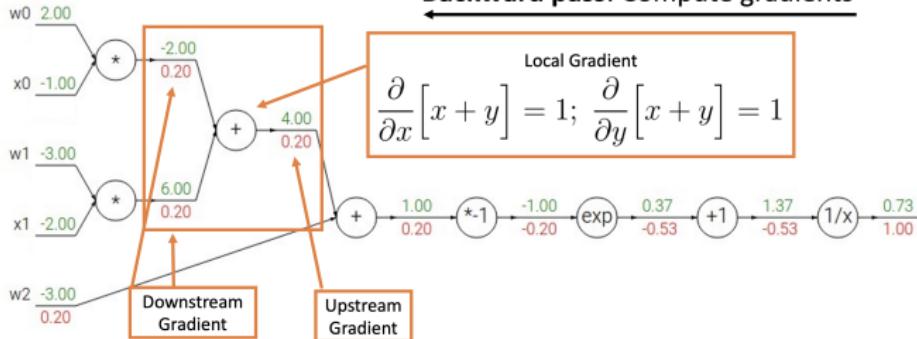
Another Example $f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$

Backward pass: Compute gradients



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

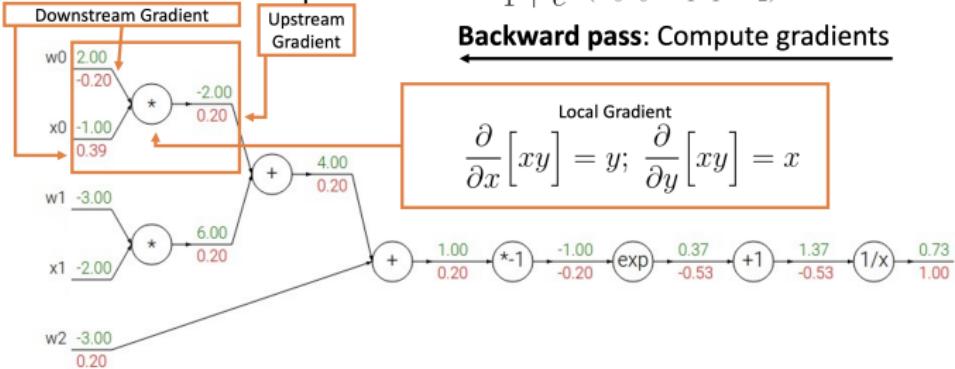
Backward pass: Compute gradients



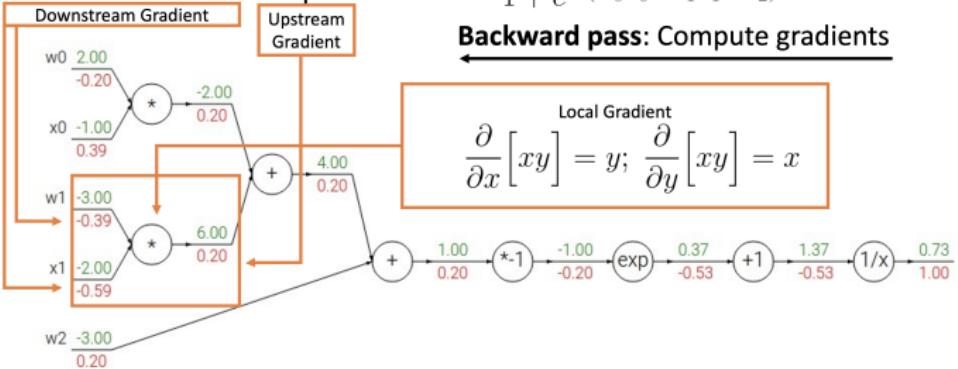
Another Example

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

Backward pass: Compute gradients



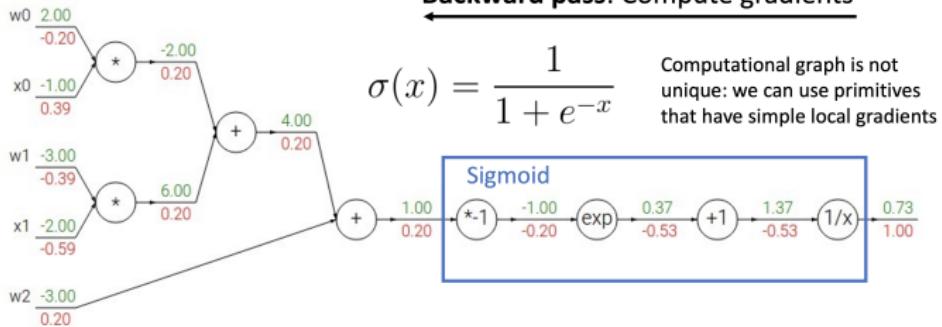
Another Example



Another Example

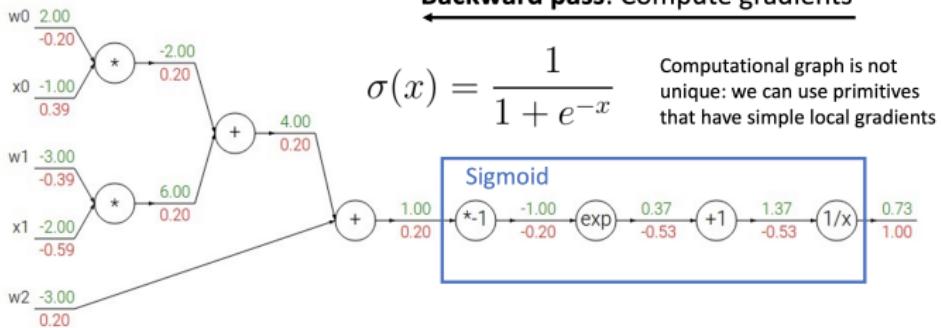
$$f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}} = \sigma(w_0x_0 + w_1x_1 + w_2)$$

Backward pass: Compute gradients



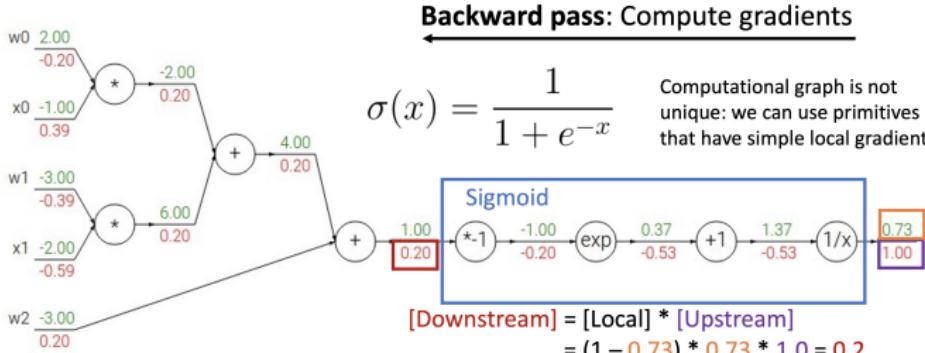
Another Example $f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}} = \sigma(w_0x_0 + w_1x_1 + w_2)$

Backward pass: Compute gradients



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

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



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

So far: backprop with scalars What about vector-valued functions?

Recap: Vector Derivatives

$x \in \mathbb{R}, y \in \mathbb{R}$

Regular derivative:

$$\frac{\partial y}{\partial x} \in \mathbb{R}$$

If x changes by a small amount, how much will y change?

Recap: Vector Derivatives

$$x \in \mathbb{R}, y \in \mathbb{R}$$

$$x \in \mathbb{R}^N, y \in \mathbb{R}$$

Regular derivative:

$$\frac{\partial y}{\partial x} \in \mathbb{R}$$

If x changes by a small amount, how much will y change?

Derivative is **Gradient**:

$$\frac{\partial y}{\partial x} \in \mathbb{R}^N \quad \left(\frac{\partial y}{\partial x} \right)_n = \frac{\partial y}{\partial x_n}$$

For each element of x , if it changes by a small amount then how much will y change?

Recap: Vector Derivatives

$$x \in \mathbb{R}, y \in \mathbb{R}$$

Regular derivative:

$$\frac{\partial y}{\partial x} \in \mathbb{R}$$

If x changes by a small amount, how much will y change?

$$x \in \mathbb{R}^N, y \in \mathbb{R}$$

Derivative is **Gradient**:

$$\frac{\partial y}{\partial x} \in \mathbb{R}^N \quad \left(\frac{\partial y}{\partial x} \right)_n = \frac{\partial y}{\partial x_n}$$

For each element of x , if it changes by a small amount then how much will y change?

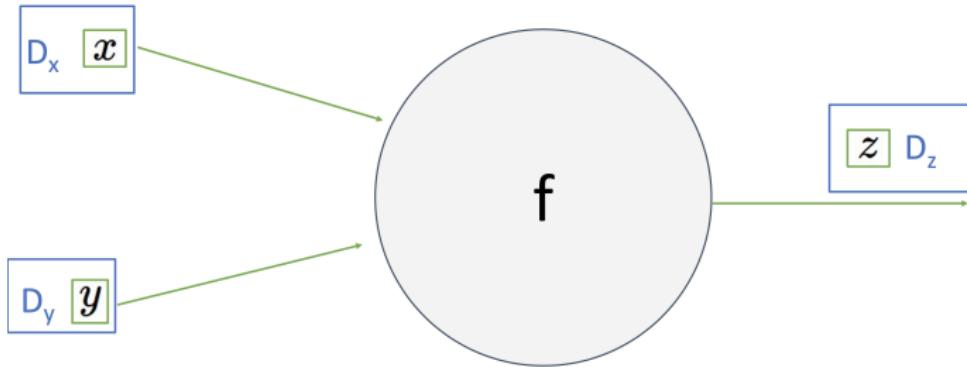
$$x \in \mathbb{R}^N, y \in \mathbb{R}^M$$

Derivative is **Jacobian**:

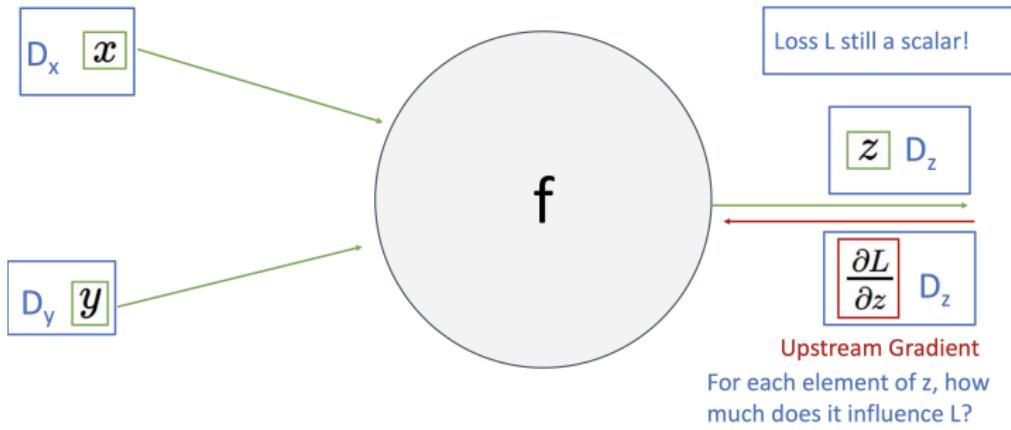
$$\frac{\partial y}{\partial x} \in \mathbb{R}^{N \times M} \quad \left(\frac{\partial y}{\partial x} \right)_{n,m} = \frac{\partial y_m}{\partial x_n}$$

For each element of x , if it changes by a small amount then how much will each element of y change?

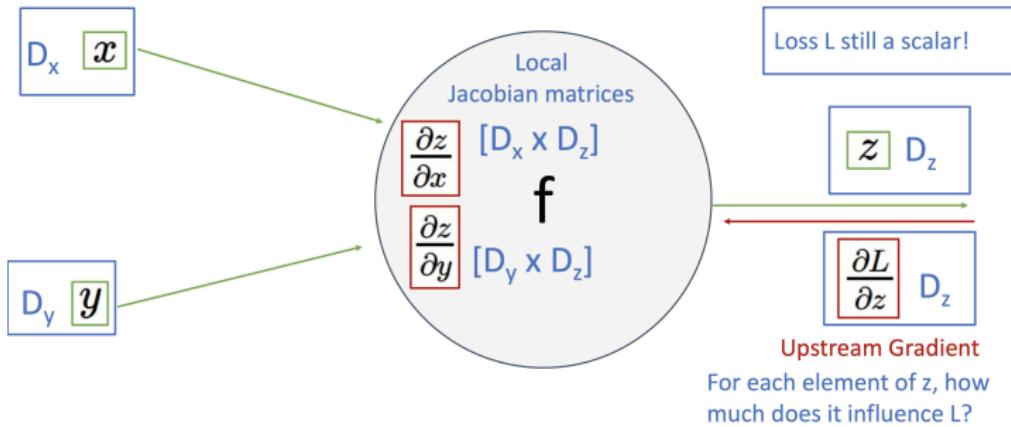
Backprop with Vectors



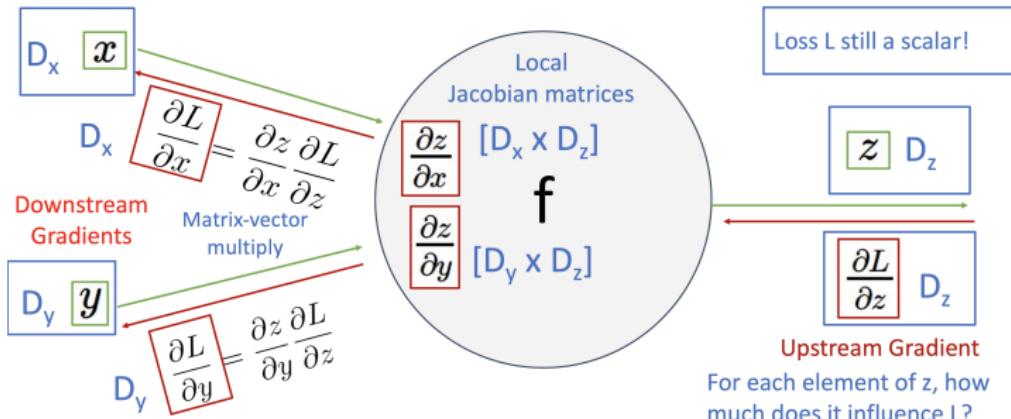
Backprop with Vectors



Backprop with Vectors



Backprop with Vectors



Backprop with Vectors

4D input x:

$$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix} \longrightarrow$$

4D output y:

$$\begin{array}{l} \longrightarrow [1] \\ \longrightarrow [0] \\ \longrightarrow [3] \\ \longrightarrow [0] \end{array}$$

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

(elementwise)

Backprop with Vectors

4D input x:

$$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix} \longrightarrow$$

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

(elementwise)

4D output y:

$$\longrightarrow \begin{bmatrix} 1 \\ 0 \\ 3 \\ 0 \end{bmatrix}$$

4D dL/dy :

$$\begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix} \longleftarrow$$

Upstream
gradient

Backprop with Vectors

4D input x:

$$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix} \longrightarrow$$

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

(elementwise)

4D output y:

$$\longrightarrow \begin{bmatrix} 1 \\ 0 \\ 3 \\ 0 \end{bmatrix}$$

Jacobian dy/dx

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

4D dL/dy :

$$\begin{array}{c} \longleftarrow [4] \longleftarrow \\ \longleftarrow [-1] \longleftarrow \\ \longleftarrow [5] \longleftarrow \\ \longleftarrow [9] \longleftarrow \end{array}$$

Upstream
gradient

Backprop with Vectors

4D input x:

$$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix} \longrightarrow$$

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

(elementwise)

4D output y:

$$\longrightarrow \begin{bmatrix} 1 \\ 0 \\ 3 \\ 0 \end{bmatrix}$$

[dy/dx] [dL/dy]

$$\begin{bmatrix} 1 & 0 & 0 & 0 \end{bmatrix} [4]$$
$$\begin{bmatrix} 0 & 0 & 0 & 0 \end{bmatrix} [-1]$$
$$\begin{bmatrix} 0 & 0 & 1 & 0 \end{bmatrix} [5]$$
$$\begin{bmatrix} 0 & 0 & 0 & 0 \end{bmatrix} [9]$$

4D dL/dy:

$$\begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix} \longleftarrow$$

Upstream
gradient

Backprop with Vectors

4D input x:

$$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix} \longrightarrow$$

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

(elementwise)

4D output y:

$$\longrightarrow \begin{bmatrix} 1 \\ 0 \\ 3 \\ 0 \end{bmatrix}$$

4D dL/dx :

$$\begin{bmatrix} 4 \\ 0 \\ 5 \\ 0 \end{bmatrix} \longleftarrow$$

$[dy/dx] [dL/dy]$

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix}$$

4D dL/dy :

$$\begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix} \longleftarrow$$

Upstream
gradient

Backprop with Vectors

Jacobian is **sparse**: off-diagonal entries all zero! Never explicitly form Jacobian; instead use **implicit** multiplication

4D input x:

$$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix}$$

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

(elementwise)

4D output y:

$$\begin{bmatrix} 1 \\ 0 \\ 3 \\ 0 \end{bmatrix}$$

4D dL/dx :

$$\begin{bmatrix} 4 \\ 0 \\ 5 \\ 0 \end{bmatrix} \leftarrow \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix}$$

4D dL/dy :

$$\begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix} \leftarrow \begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix} \quad \text{Upstream gradient}$$

Backprop with Vectors

Jacobian is **sparse**: off-diagonal entries all zero! Never **explicitly** form Jacobian; instead use **implicit** multiplication

4D input x:

$$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix} \longrightarrow$$

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

(elementwise)

4D output y:

$$\begin{array}{l} \longrightarrow [1] \\ \longrightarrow [0] \\ \longrightarrow [3] \\ \longrightarrow [0] \end{array}$$

4D dL/dx :

$$\begin{bmatrix} 4 \\ 0 \\ 5 \\ 0 \end{bmatrix} \leftarrow$$

$[dy/dx]$ $[dL/dy]$

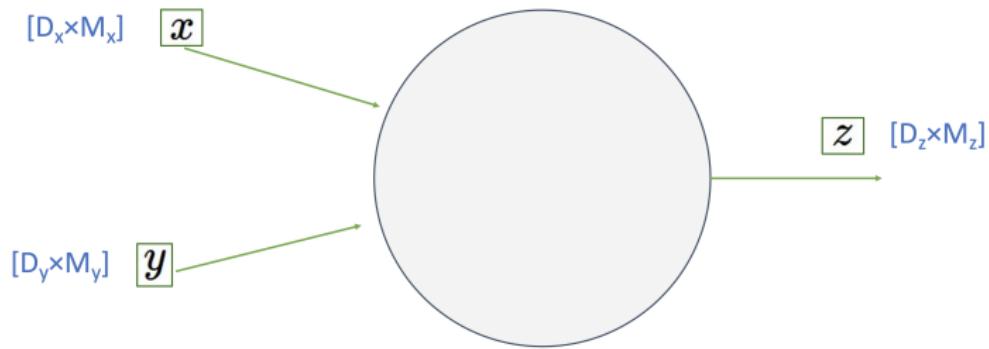
$$\left(\frac{\partial L}{\partial x} \right)_i = \begin{cases} \left(\frac{\partial L}{\partial y} \right)_i & \text{if } x_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

4D dL/dy :

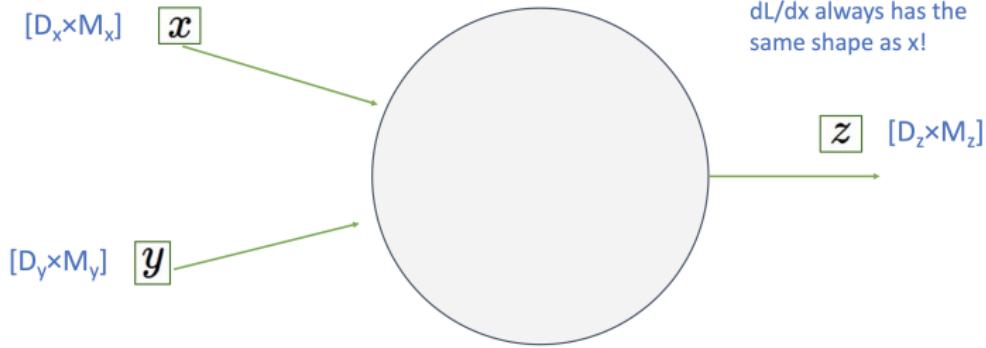
$$\begin{array}{l} \leftarrow [4] \leftarrow \\ \leftarrow [-1] \leftarrow \\ \leftarrow [5] \leftarrow \\ \leftarrow [9] \leftarrow \end{array}$$

Upstream
gradient

Backprop with Matrices (or Tensors):



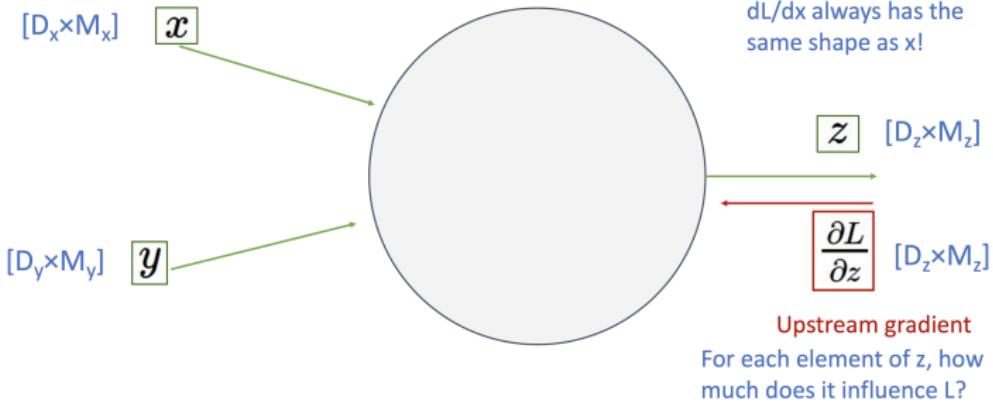
Backprop with Matrices (or Tensors):



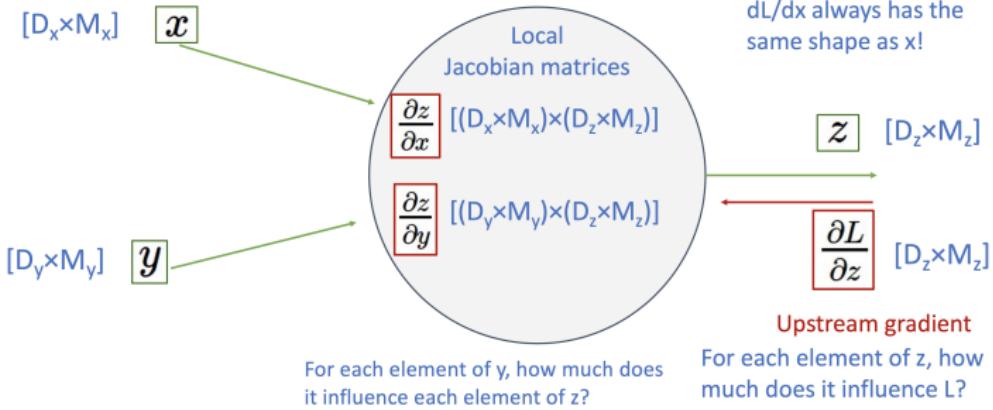
Loss L still a scalar!

dL/dx always has the same shape as x !

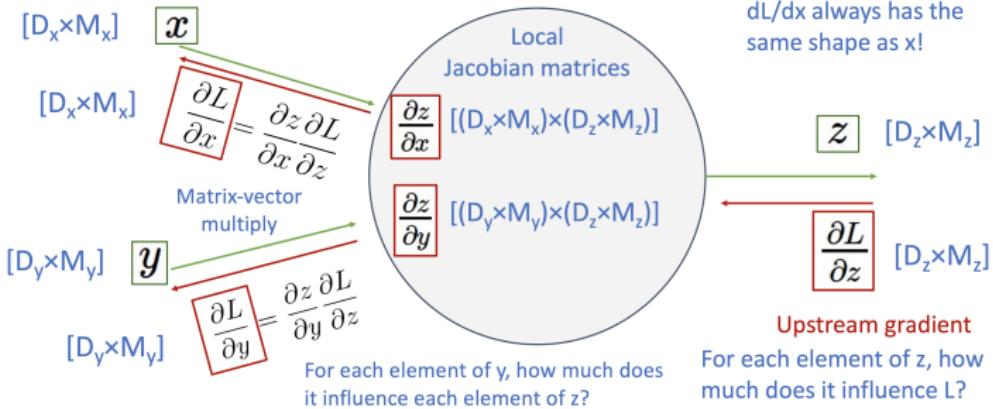
Backprop with Matrices (or Tensors):



Backprop with Matrices (or Tensors):



Backprop with Matrices (or Tensors):



Example: Matrix Multiplication

$$\begin{array}{l} \mathbf{x}: [\mathbf{N} \times \mathbf{D}] \\ \begin{bmatrix} 2 & 1 & -3 \\ -3 & 4 & 2 \\ 3 & 2 & 1-2 \end{bmatrix} \end{array} \quad \begin{array}{l} \mathbf{w}: [\mathbf{D} \times \mathbf{M}] \\ \begin{bmatrix} 3 & 2 & 1 & -1 \\ 2 & 1 & 3 & 2 \end{bmatrix} \end{array}$$

Matrix Multiply $y = \mathbf{x}\mathbf{w}$

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

$$\begin{array}{l} \mathbf{y}: [\mathbf{N} \times \mathbf{M}] \\ \begin{bmatrix} -1 & -1 & 2 & 6 \\ 5 & 2 & 11 & 7 \end{bmatrix} \end{array}$$

Example: Matrix Multiplication

$$x: [N \times D] \quad w: [D \times M]$$
$$\begin{bmatrix} 2 & 1 & -3 \\ -3 & 4 & 2 \\ 3 & 2 & 1-2 \end{bmatrix} \quad \begin{bmatrix} 3 & 2 & 1 & -1 \\ 2 & 1 & 3 & 2 \\ 3 & 2 & 1 & -2 \end{bmatrix}$$

$$dL/dx: [N \times D]$$
$$\begin{bmatrix} ? & ? & ? \\ ? & ? & ? \end{bmatrix}$$

Matrix Multiply $y = xw$

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

$$y: [N \times M]$$
$$\begin{bmatrix} -1 & -1 & 2 & 6 \\ 5 & 2 & 11 & 7 \end{bmatrix}$$
$$dL/dy: [N \times M]$$
$$\begin{bmatrix} 2 & 3 & -3 & 9 \\ -8 & 1 & 4 & 6 \end{bmatrix}$$

Example: Matrix Multiplication

$$\begin{array}{l} x: [N \times D] \\ \begin{bmatrix} 2 & 1 & -3 \\ -3 & 4 & 2 \\ 3 & 2 & 1-2 \end{bmatrix} \end{array} \quad \begin{array}{l} w: [D \times M] \\ \begin{bmatrix} 3 & 2 & 1 & -1 \\ 2 & 1 & 3 & 2 \\ 3 & 2 & 1 & -2 \end{bmatrix} \end{array}$$

$$\begin{array}{l} dL/dx: [N \times D] \\ \begin{bmatrix} ? & ? & ? \\ ? & ? & ? \end{bmatrix} \end{array}$$

Matrix Multiply $y = xw$

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

$$\begin{array}{l} y: [N \times M] \\ \begin{bmatrix} -1 & -1 & 2 & 6 \\ 5 & 2 & 11 & 7 \end{bmatrix} \end{array}$$
$$\begin{array}{l} dL/dy: [N \times M] \\ \begin{bmatrix} 2 & 3 & -3 & 9 \\ -8 & 1 & 4 & 6 \end{bmatrix} \end{array}$$

Jacobians:

$$\begin{array}{l} dy/dx: [(N \times D) \times (N \times M)] \\ dy/dw: [(D \times M) \times (N \times M)] \end{array}$$

For a neural net we may have

$$N=64, D=M=4096$$

Each Jacobian takes 256 GB of memory! Must work with them implicitly!

Example: Matrix Multiplication

$$x: [N \times D] \quad w: [D \times M]$$
$$\begin{bmatrix} 2 & 1 & -3 \\ -3 & 4 & 2 \\ 3 & 2 & 1-2 \end{bmatrix} \quad \begin{bmatrix} 3 & 2 & 1-1 \\ 2 & 1 & 3 & 2 \\ 3 & 2 & 1-2 \end{bmatrix}$$

$$dL/dx: [N \times D]$$
$$\begin{bmatrix} ? & ? & ? \\ ? & ? & ? \end{bmatrix}$$

$$dL/dx_{1,1}$$
$$= (dy/dx_{1,1}) \cdot (dL/dy)$$

Matrix Multiply $y = xw$

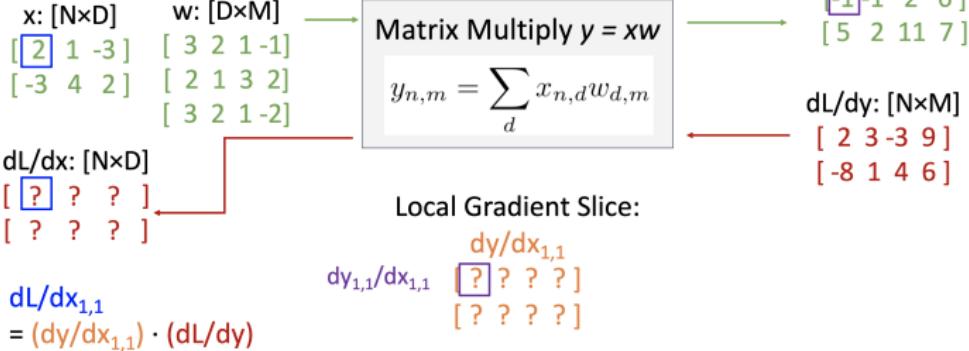
$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

$$y: [N \times M]$$
$$\begin{bmatrix} -1 & -1 & 2 & 6 \\ 5 & 2 & 11 & 7 \end{bmatrix}$$
$$dL/dy: [N \times M]$$
$$\begin{bmatrix} 2 & 3 & -3 & 9 \\ -8 & 1 & 4 & 6 \end{bmatrix}$$

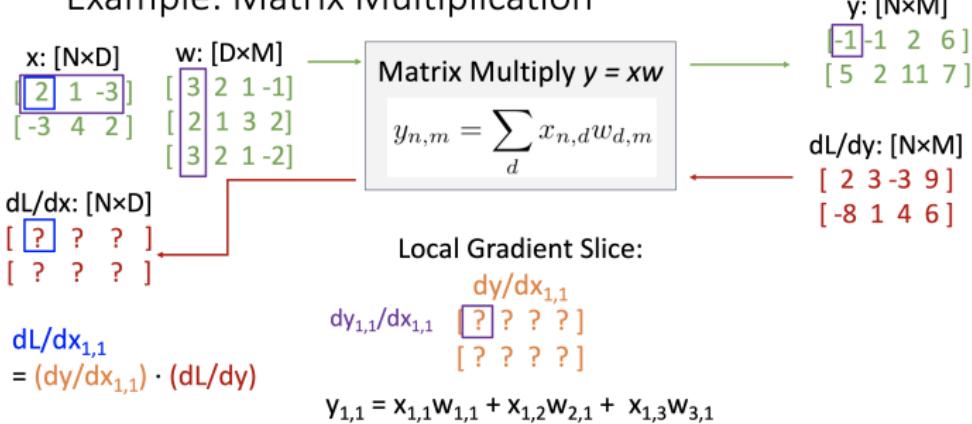
Local Gradient Slice:

$$dy/dx_{1,1}$$
$$\begin{bmatrix} ? & ? & ? & ? \\ ? & ? & ? & ? \end{bmatrix}$$

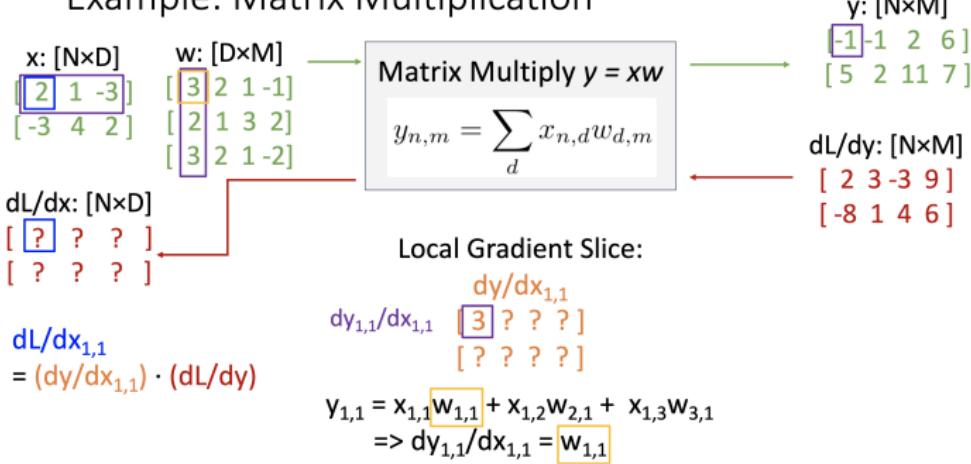
Example: Matrix Multiplication



Example: Matrix Multiplication



Example: Matrix Multiplication



Example: Matrix Multiplication

$$x: [N \times D] \quad w: [D \times M]$$

| | | |
|----|---|-----|
| 2 | 1 | -3 |
| -3 | 4 | 2 |
| 3 | 2 | 1-2 |

Matrix Multiply $y = xw$

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

$$y: [N \times M]$$

| | | | |
|----|---|----|---|
| -1 | 1 | 2 | 6 |
| 5 | 2 | 11 | 7 |

$$dL/dy: [N \times M]$$

| | | | |
|----|---|----|---|
| 2 | 3 | -3 | 9 |
| -8 | 1 | 4 | 6 |

$$dL/dx: [N \times D]$$

| | | |
|---|---|---|
| ? | ? | ? |
| ? | ? | ? |

$$dL/dx_{1,1}$$

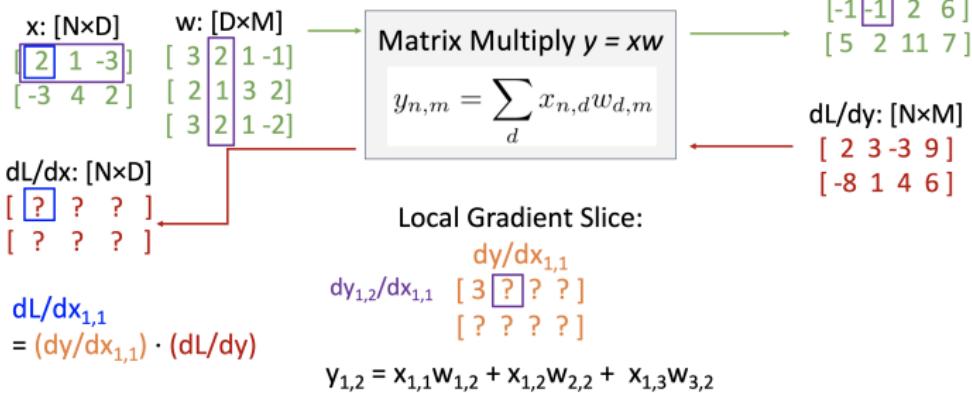
$$= (dy/dx_{1,1}) \cdot (dL/dy)$$

Local Gradient Slice:

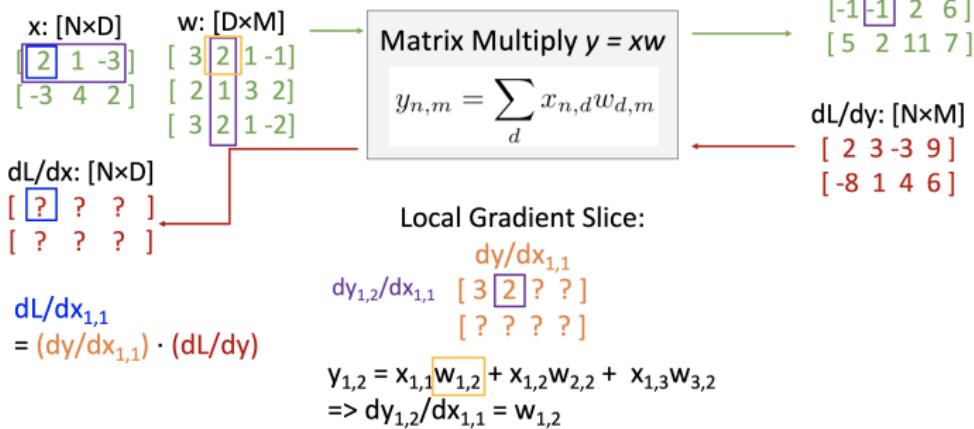
$$\frac{dy}{dx}_{1,1} \quad [3 \quad ? \quad ?]$$

$$\frac{dy}{dx}_{1,2} \quad [? \quad ? \quad ?]$$

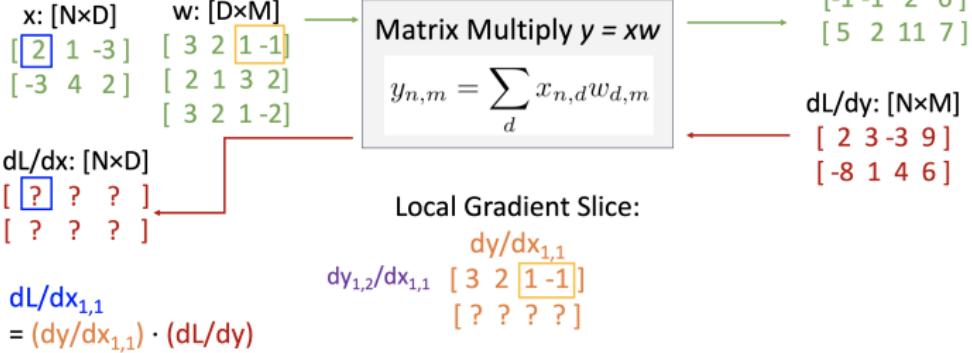
Example: Matrix Multiplication



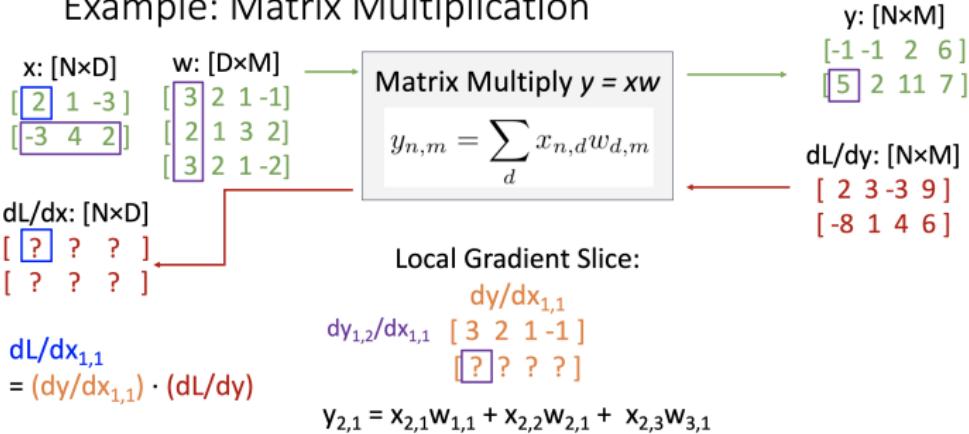
Example: Matrix Multiplication



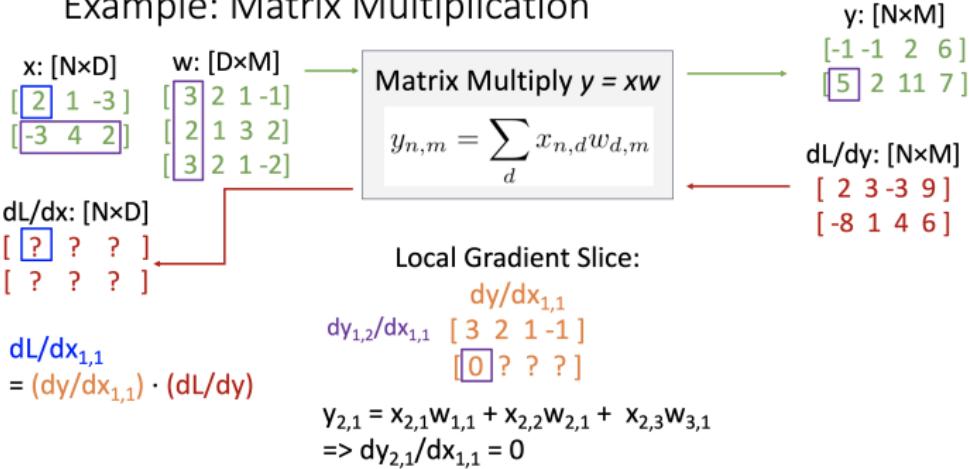
Example: Matrix Multiplication



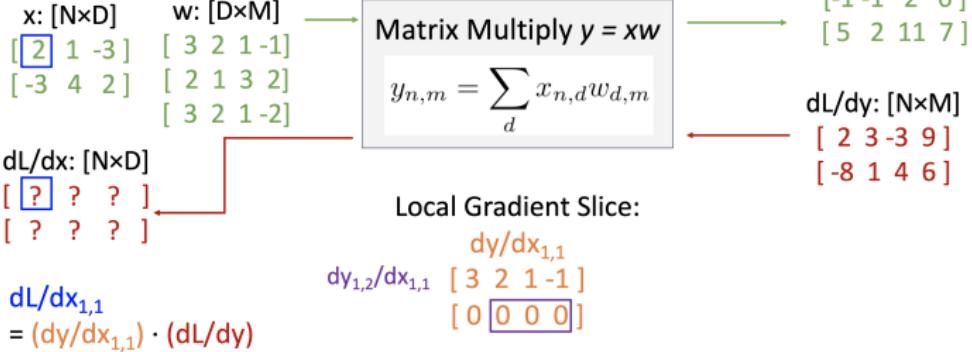
Example: Matrix Multiplication



Example: Matrix Multiplication



Example: Matrix Multiplication



Example: Matrix Multiplication

$$x: [N \times D] \quad w: [D \times M]$$

| | | |
|----|---|-----|
| 2 | 1 | -3 |
| -3 | 4 | 2 |
| 3 | 2 | 1-2 |

Matrix Multiply $y = xw$

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

$$y: [N \times M]$$

| | | | |
|----|----|----|---|
| -1 | -1 | 2 | 6 |
| 5 | 2 | 11 | 7 |

$$dL/dy: [N \times M]$$

| | | | |
|----|---|----|---|
| 2 | 3 | -3 | 9 |
| -8 | 1 | 4 | 6 |

$$dL/dx: [N \times D]$$

| | | |
|---|---|---|
| ? | ? | ? |
| ? | ? | ? |

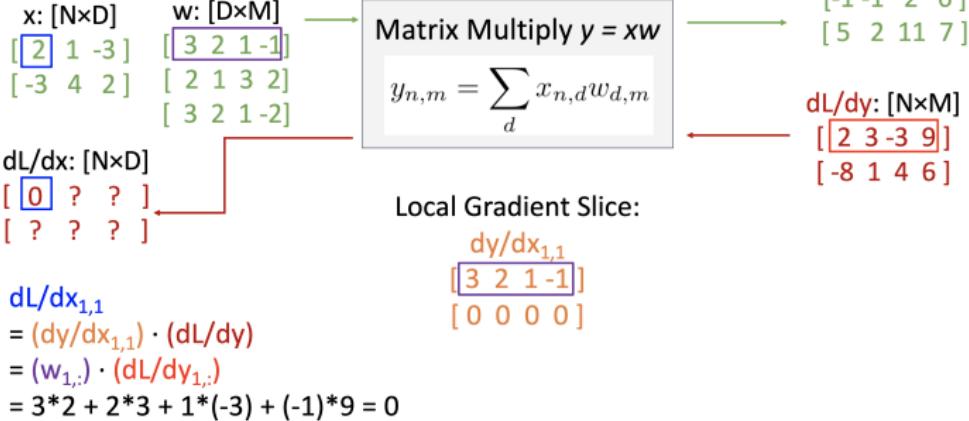
$$dL/dx_{1,1}$$

$$= (dy/dx_{1,1}) \cdot (dL/dy)$$

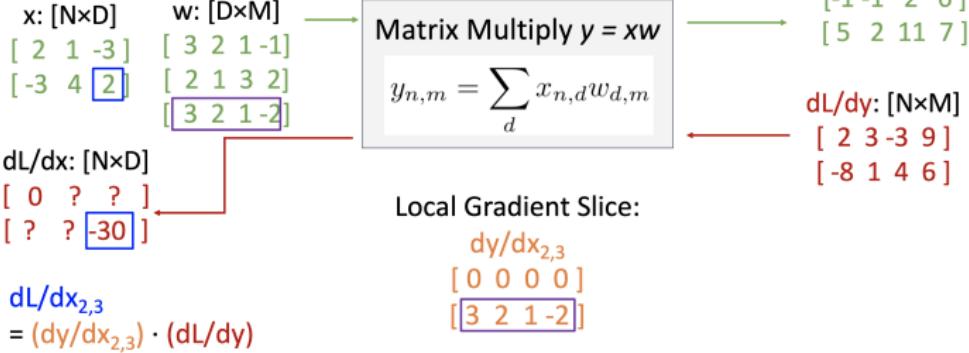
Local Gradient Slice:

$$\begin{aligned} dy/dx_{1,1} \\ [3 & 2 & 1-1] \\ [0 & 0 & 0 & 0] \end{aligned}$$

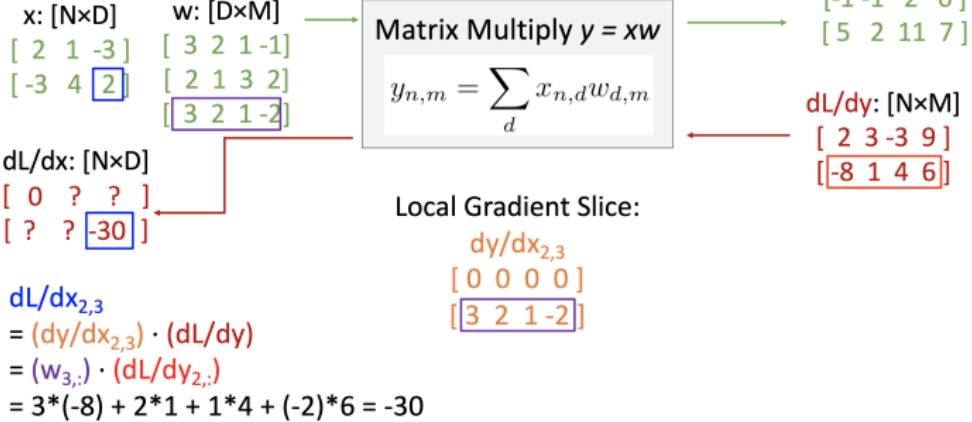
Example: Matrix Multiplication



Example: Matrix Multiplication



Example: Matrix Multiplication



Example: Matrix Multiplication

$$x: [N \times D] \quad w: [D \times M]$$
$$\begin{bmatrix} 2 & 1 & -3 \\ -3 & 4 & 2 \\ 3 & 2 & 1 & -2 \end{bmatrix} \quad \begin{bmatrix} 3 & 2 & 1 & -1 \\ 2 & 1 & 3 & 2 \\ 3 & 2 & 1 & -2 \end{bmatrix}$$

$$dL/dx: [N \times D]$$
$$\begin{bmatrix} 0 & 16 & -9 \\ -24 & 9 & -30 \end{bmatrix}$$

Matrix Multiply $y = xw$

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

$$y: [N \times M]$$
$$\begin{bmatrix} -1 & -1 & 2 & 6 \\ 5 & 2 & 11 & 7 \end{bmatrix}$$
$$dL/dy: [N \times M]$$
$$\begin{bmatrix} 2 & 3 & -3 & 9 \\ -8 & 1 & 4 & 6 \end{bmatrix}$$

$$dL/dx_{i,j}$$
$$= (dy/dx_{i,j}) \cdot (dL/dy)$$
$$= (w_{j,:}) \cdot (dL/dy_{i,:})$$

Example: Matrix Multiplication

$$x: [N \times D] \quad w: [D \times M]$$

$$\begin{bmatrix} 2 & 1 & -3 \\ -3 & 4 & 2 \\ 3 & 2 & 1-2 \end{bmatrix} \quad \begin{bmatrix} 3 & 2 & 1-1 \\ 2 & 1 & 3 & 2 \\ 3 & 2 & 1-2 \end{bmatrix}$$

$$dL/dx: [N \times D]$$

$$\begin{bmatrix} 0 & 16 & -9 \\ -24 & 9 & -30 \end{bmatrix}$$

$$dL/dx_{i,j}$$

$$= (dy/dx_{i,j}) \cdot (dL/dy)$$

$$= (w_{j,:}) \cdot (dL/dy_{i,:})$$

Matrix Multiply $y = xw$

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

$$y: [N \times M]$$

$$\begin{bmatrix} -1 & -1 & 2 & 6 \\ 5 & 2 & 11 & 7 \end{bmatrix}$$

$$dL/dy: [N \times M]$$

$$\begin{bmatrix} 2 & 3 & -3 & 9 \\ -8 & 1 & 4 & 6 \end{bmatrix}$$

$$dL/dx = (dL/dy) w^T$$

$$[N \times D] \quad [N \times M] \quad [M \times D]$$

Easy way to remember:
It's the only way the shapes work out!

Example: Matrix Multiplication

$$x: [N \times D] \quad w: [D \times M]$$

| | |
|--------------|--------------|
| [2 1 -3] | [3 2 1 -1] |
| [-3 4 2] | [2 1 3 2] |
| [3 2 1 -2] | |

$$dL/dx: [N \times D]$$

| |
|---------------|
| [0 16 -9] |
| [-24 9 -30] |

Matrix Multiply $y = xw$

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

$$y: [N \times M]$$

| |
|---------------|
| [-1 -1 2 6] |
| [5 2 11 7] |

$$dL/dy: [N \times M]$$

| |
|--------------|
| [2 3 -3 9] |
| [-8 1 4 6] |

$$dL/dx = (dL/dy) w^T$$

$$[N \times D] \quad [N \times M] \quad [M \times D]$$

$$dL/dw = x^T (dL/dy)$$

$$[D \times M] \quad [D \times N] \quad [N \times M]$$

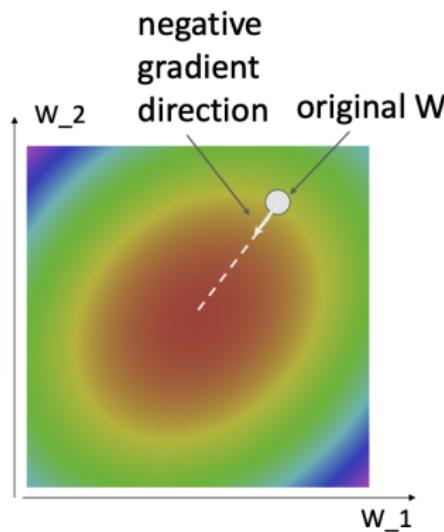
Easy way to remember:
It's the only way the shapes work out!

Recall: Optimization

$$W^* = \arg \min_W L(W)$$

- **Gradient decent:** Iteratively step in the direction of the negative gradient

$$W \leftarrow W - \alpha \nabla_W L(W)$$



Recall: Optimization $L = \frac{1}{N} \sum_{i=1}^N \ell_i(\mathbf{x}^{(i)}, y^{(i)}; W) + \lambda R(W)$

- All the training data is used to perform a single parameter update

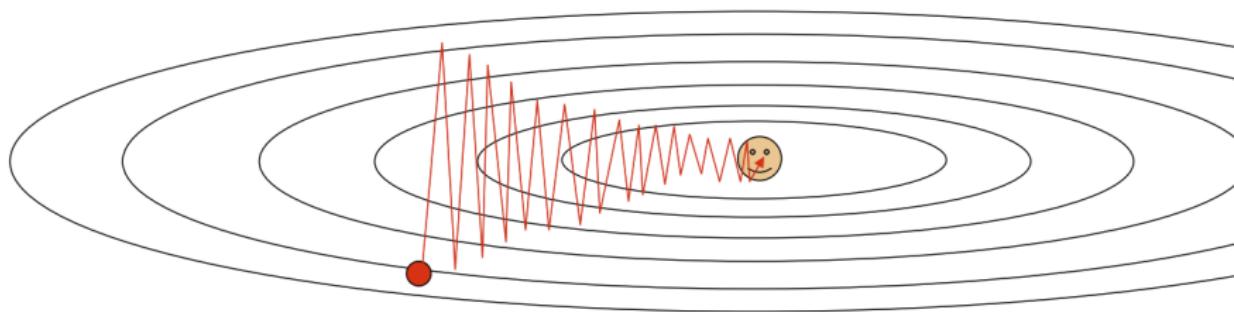
$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W \ell_i(\mathbf{x}^{(i)}, y^{(i)}; W) + \lambda \nabla_W R(W)$$

Full sum is expensive when N is large!

- Mini-batch Gradient Descent: approximate sum using a minibatch of examples; 32/64/128 common
- Stochastic Gradient Descent (SGD): when mini-batch contains only a single example; usually use SGD when referring to mini-batch gradient descent

Problems with SGD

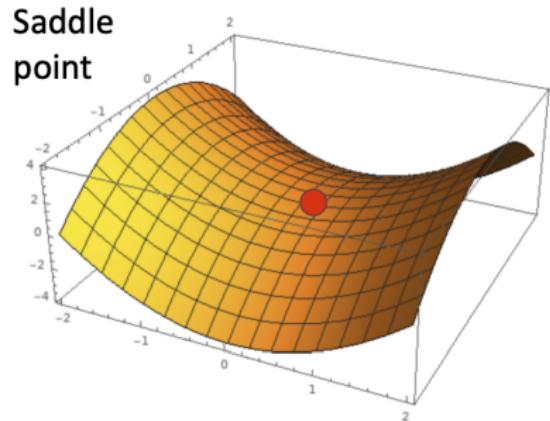
- What if loss changes quickly in one direction and slowly in another? What does gradient decent do?



Very slow progress along shallow direction, jitter along steep direction.

Problems with SGD

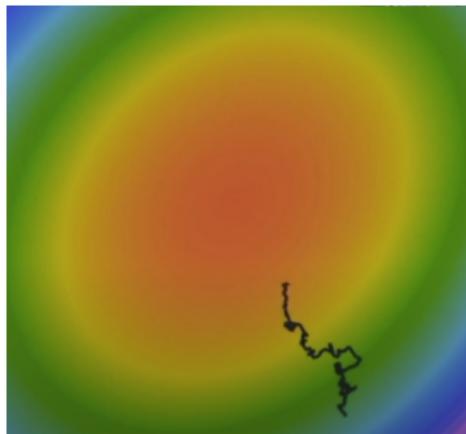
- What if the loss function has a local minimum or saddle point?



Zero gradient, gradient descent gets stuck

Problems with SGD

- Our gradients come from minibatches so they can be noisy!



Use SGD + Momentum instead!

Sutskever et al, "On the importance of initialization and momentum in deep learning", ICML 2013.

Why do you think deep learning became so successful?

Important topics of this lecture

- Deep nets
- Backpropagation

Up next:

- PyTorch Tutorial
- CNN & RNN