Topics:

- Neural Networks
- Backpropagation

**CS 4644-DL / 7643-A**
**ZSOLT KIRA**

- **Assignment 1 out!**
  - **Due Feb 5$^{th}$**
  - Start now, start now, start now!
  - Start now, start now, start now!
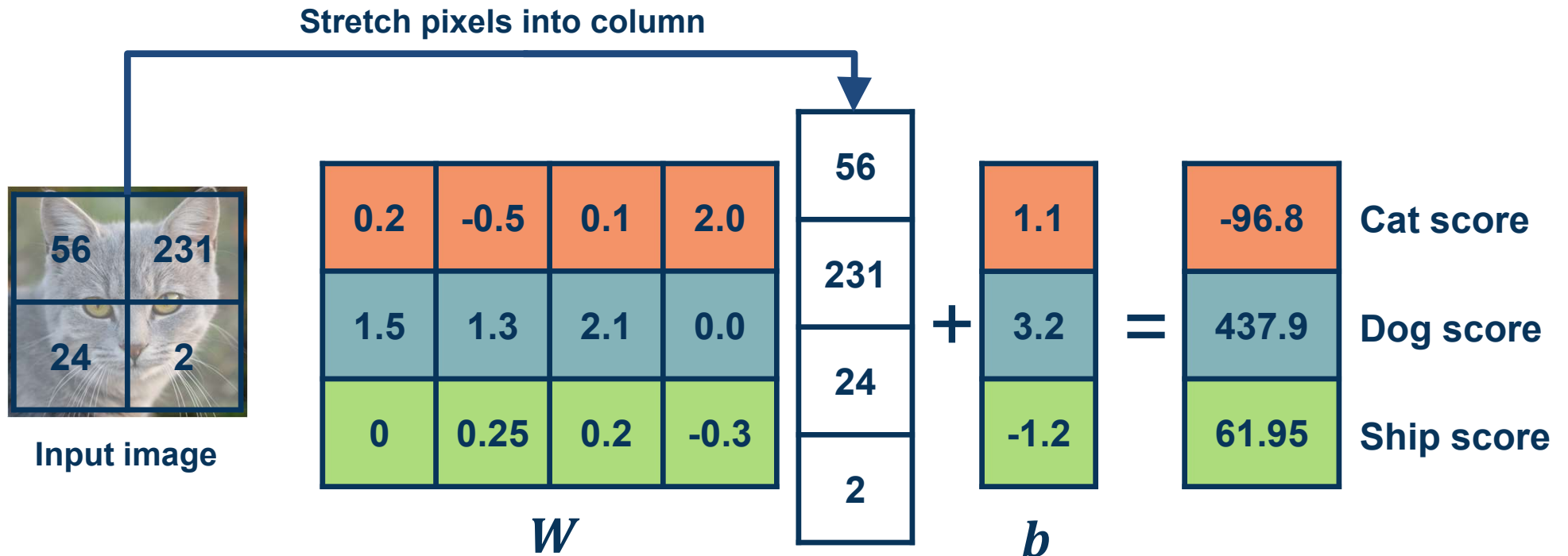  - Start now, start now, start now!

- **Piazza**
  - Be active!!!
  - Extra credit!

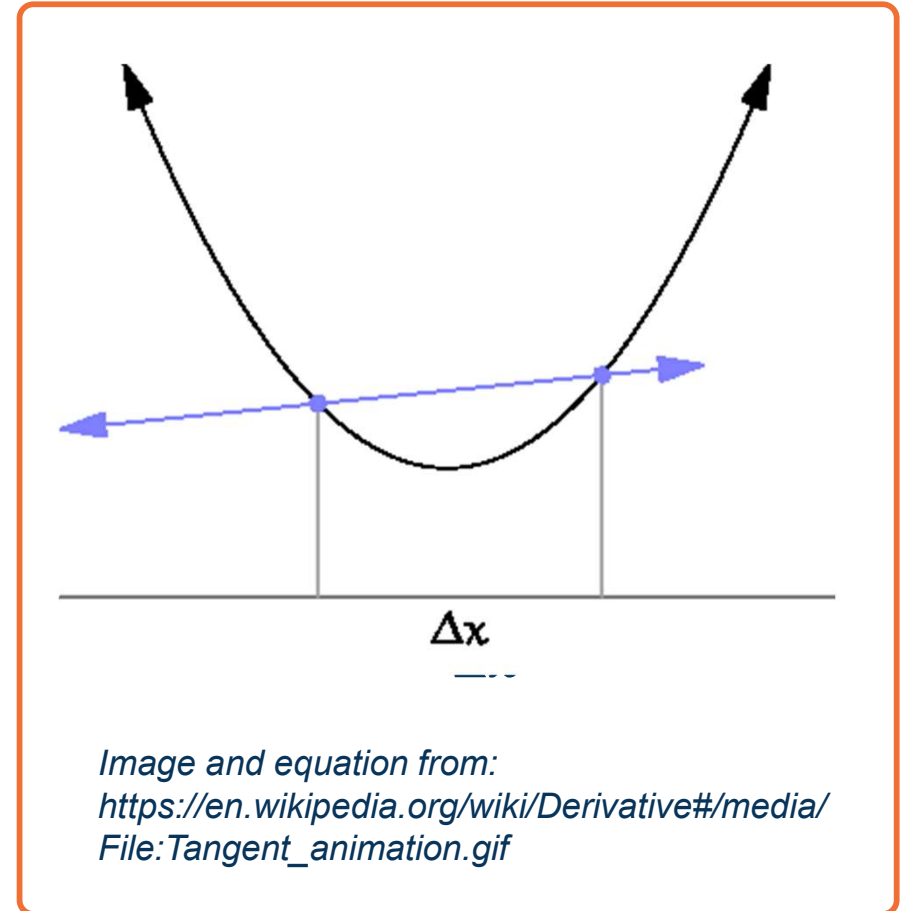- **Office hours**
  - Assignment (@41) and matrix calculus (@46)

# Example with an image with **4 pixels**, and **3 classes (cat/dog/ship)**

Stretch pixels into column



| | | | |
|---|---|---|---|
| 0.2 | -0.5 | 0.1 | 2.0 |
| 1.5 | 1.3 | 2.1 | 0.0 |
| 0 | 0.25 | 0.2 | -0.3 |

$W$

| |
|---|
| 56 |
| 231 |
| 24 |
| 2 |

**+**

| |
|---|
| 1.1 |
| 3.2 |
| -1.2 |

$b$

**=**

| | |
|---|---|
| -96.8 | Cat score |
| 437.9 | Dog score |
| 61.95 | Ship score |

Input image

*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*

**Example**

Georgia Tech

- We can find the steepest descent direction by computing the **derivative (gradient):**

$$f'(a) = \lim_{h \to 0} \frac{f(a+h) - f(a)}{h}$$

- Steepest descent direction is the **negative gradient**

- **Intuitively:** Measures how the function changes as the argument a changes by a small step size

  - As step size goes to zero

- **In Machine Learning:** Want to know how the **loss function** changes **as weights** are varied

  - Can consider each parameter separately by taking **partial derivative** of loss function with respect to that parameter



*Image and equation from:*
*https://en.wikipedia.org/wiki/Derivative#/media/*
*File:Tangent_animation.gif*

This idea can be turned into an **algorithm (gradient descent)**

- Choose a model: $f(x, W) = \mathrm{W}x$

- Choose loss function: $L_i = |y - Wx_i|^2$

- Calculate partial derivative for each parameter: $\frac{\partial L}{\partial w_i}$

- Update the parameters: $w_i = w_i - \frac{\partial L}{\partial w_i}$

- Add learning rate to prevent too big of a step: $w_i = w_i - \alpha \frac{\partial L}{\partial w_i}$

- Repeat (from Step 3)

**Gradient Descent**

Georgia Tech

Often, we only compute the gradients across a small subset of data

- Full Batch Gradient Descent $\quad L = \dfrac{1}{N} \sum L\left(f(x_i, W), y_i\right)$

- Mini-Batch Gradient Descent $\quad L = \dfrac{1}{M} \sum L\left(f(x_i, W), y_i\right)$

  - Where M is a *subset* of data

- We iterate over mini-batches:

  - Get mini-batch, compute loss, compute derivatives, and take a set

**Mini-Batch Gradient Descent**

Georgia Tech

http://demonstrations.wolfram.com/VisualizingTheGradientVector/

$w_2$

original W

negative gradient direction

$w_1$

**Gradient Descent**

Georgia Tech

# For some functions, we can analytically derive the partial derivative

## Example:

### Derivation of Update Rule

### Function

$$f(w, x_i) = w^T x_i$$

### Loss

$$(y_i - w^T x_i)^2$$

(Assume $w$ and $x_i$ are column vectors, so same as $w \cdot x_i$)

**Dataset:** N examples (indexed by $k$)

### Update Rule

$$w_j \leftarrow w_j + 2\eta \sum_{k=1}^{N} \delta_k x_{kj}$$

$L = \sum_{k=1}^{N} (y_k - w^T x_k)^2$

Gradient descent tells us we should update $w$ as follows to minimize $L$:

$$w_j \leftarrow w_j - \eta \frac{\partial L}{\partial w_j}$$

So what's $\frac{\partial L}{\partial w_j}$?

$$\frac{\partial L}{\partial w_j} = \sum_{k=1}^{N} \frac{\partial}{\partial w_j} (y_k - w^T x_k)^2$$

$$= \sum_{k=1}^{N} 2(y_k - w^T x_k) \frac{\partial}{\partial w_j} (y_k - w^T x_k)$$

$$= -2 \sum_{k=1}^{N} \delta_k \frac{\partial}{\partial w_j} w^T x_k$$

…where…
$$\delta_k = y_k - w^T x_k$$

$$= -2 \sum_{k=1}^{N} \delta_k \frac{\partial}{\partial w_j} \sum_{i=1}^{m} w_i x_{ki}$$

$$= -2 \sum_{k=1}^{N} \delta_k x_{kj}$$

If we add a **non-linearity (sigmoid),** derivation is more complex
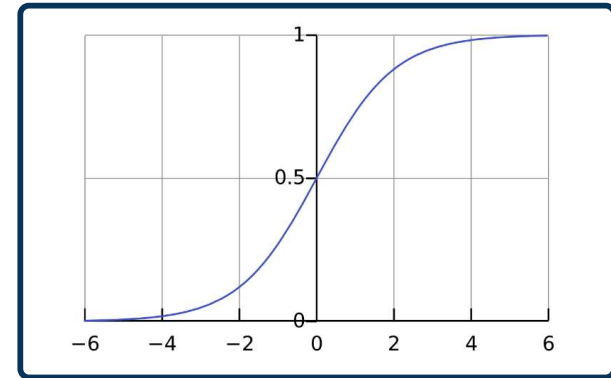
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

First, one can derive that: $\sigma'(x) = \sigma(x)(1 - \sigma(x))$

$$\mathbf{f}(\mathbf{x}) = \sigma\left(\sum_k w_k x_k\right)$$

$$L = \sum_i \left(y_i - \sigma\left(\sum_k w_k x_{ik}\right)\right)^2$$

$$\frac{\partial L}{\partial w_j} = \sum_i 2\left(y_i - \sigma\left(\sum_k w_k x_{ik}\right)\right)\left(-\frac{\partial}{\partial w_j}\sigma\left(\sum_k w_k x_{ik}\right)\right)$$

$$= \sum_i -2\left(y_i - \sigma\left(\sum_k w_k x_{ik}\right)\right)\sigma'\left(\sum_k w_k x_{ik}\right)\frac{\partial}{\partial w_j}\sum_k w_k x_{ik}$$

$$= \sum_i -2\delta_i\sigma(\mathbf{d}_i)(1 - \sigma(\mathbf{d}_i))x_{ij}$$

where $\quad \delta_i = y_i - \mathbf{f}(x_i) \qquad d_i = \sum w_k x_{ik}$



**The sigmoid perception update rule:**

$$w_j \leftarrow w_j + 2\eta \sum_{k=1}^{N} \delta_i\sigma_i(1 - \sigma_i)x_{ij}$$

where $\quad \sigma_i = \sigma\left(\sum_{j=1}^{m} w_j x_{ij}\right)$

$$\delta_i = y_i - \sigma_i$$

**Adding a Non-Linear Function**

Georgia Tech

A **linear classifier** can be broken down into:

- Input
- A function of the input
- A loss function

It's all just one function that can be **decomposed** into building blocks

$$X \rightarrow \boxed{w \cdot x} \xrightarrow{u} \boxed{\frac{1}{1 + e^{-u}}} \xrightarrow{p} \boxed{-\log(p)} \xrightarrow{L}$$

Input              Model             Loss Function

Georgia Tech

The same two-layered neural network **corresponds to adding another weight matrix**

- We will prefer the linear algebra view, but use some terminology from neural networks (& biology)

**input layer**

**hidden layer**

**output layer**

$$x \quad\quad W_1 \quad\quad W_2$$

$$=$$

$$f(x, W_1, W_2) = \sigma(W_2 \sigma(W_1 x))$$

*Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n*

Georgia Tech

**Large (deep) networks** can be built by adding more and more layers

==Three-layered neural networks can represent **any function**==

◆ The number of nodes could grow unreasonably (exponential or worse) with respect to the complexity of the function

We will show them **without edges**:



$$f(x, W_1, W_2, W_3) = \sigma(W_2 \sigma(W_1 x))$$

*Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n*

**Adding More Layers!**

Georgia Tech

# Demo

- http://playground.tensorflow.org

# Computation Graphs

Functions can be made **arbitrarily complex** (subject to memory and computational limits), e.g.:

$$f(x, W) = \sigma(W_5 \sigma(W_4 \sigma(W_3 \sigma(W_2 \sigma(W_1 x))))$$

We can use **any type of differentiable function (layer)** we want!

◆   At the end, **add the loss function**

Composition can have **some structure**

The world is **compositional**!

We want our **model** to reflect this

Empirical and theoretical evidence that it makes **learning complex functions easier**

Note that **prior state of art engineered features** often had this compositionality as well



VISION

pixels → edge → texton → motif → part → object

SPEECH

sample → spectral band → formant → motif → phone → word

NLP

character → word → NP/VP/.. → clause → sentence → story

*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

⬡ **Pixels -> edges -> object parts -> objects**

**Compositionality**

Georgia Tech

- We are learning **complex models** with significant amount of parameters (millions or billions)

- How do we compute the gradients of the **loss** (at the end) with respect to **internal** parameters?

- Intuitively, want to understand how **small changes** in weight deep inside **are propagated** to affect the **loss function** at the end



$$\frac{\partial L}{\partial w_i}?$$

Georgia Tech

Given a library of simple functions

$\sin(x)$

$\log(x)$

$\cos(x)$

$x^3$

$\exp(x)$

Compose into a

complicate function

$$-\log\left(\frac{1}{1 + e^{-w \cdot x}}\right)$$

$w \cdot x$ → $u$ → $\dfrac{1}{1 + e^{-u}}$ → $p$ → $-\log(p)$ → $L$

**Decomposing a Function**

Georgia Tech

To develop a general algorithm for this, we will view the function as a **computation graph**

Graph can be any **directed acyclic graph (DAG)**

◆ Modules must be differentiable to support gradient computations for gradient descent

A **training algorithm** will then process this graph, **one module at a time**



*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

**A General Framework**

Georgia Tech

# Directed Acyclic Graphs (DAGs)

- Exactly what the name suggests
  - Directed edges
  - No (directed) cycles
  - Underlying undirected cycles okay

Georgia Tech

# Directed Acyclic Graphs (DAGs)

- Concept
  - Topological Ordering

Georgia Tech

# Directed Acyclic Graphs (DAGs)

$$f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin(x_2)$$

$$-\log\left(\frac{1}{1+e^{-w\cdot x}}\right)$$

$$w\cdot x \xrightarrow{u} \frac{1}{1+e^{-u}} \xrightarrow{p} -\log(p) \xrightarrow{L}$$

*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

**Machine Learning Example**

Georgia Tech

# Backpropagation

**Step 1:** Compute Loss on Mini-Batch: **Forward Pass**

Layer 2

Layer 3

*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

**Neural Network Training**

Georgia Tech

Note that we must store the **intermediate outputs of all layers**!

- This is because we will need them to **compute the gradients** (the gradient equations will have terms with the output values in them)

*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

**Neural Network Training**

Georgia Tech
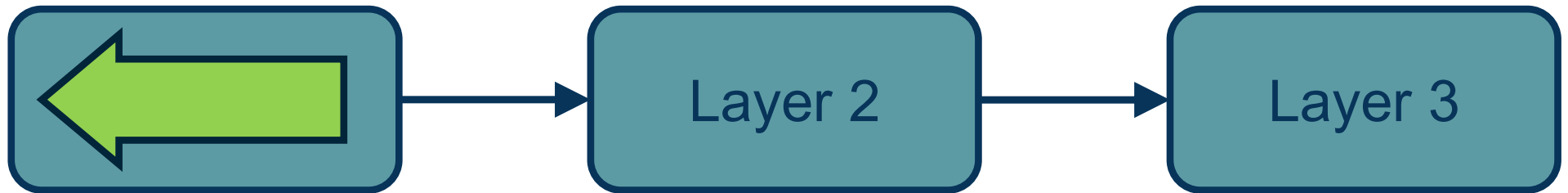
**Step 1:** Compute Loss on Mini-Batch: **Forward Pass**

**Step 2:** Compute Gradients wrt parameters: **Backward Pass**



*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

**Neural Network Training**

**Step 1:** Compute Loss on Mini-Batch: **Forward Pass**

**Step 2:** Compute Gradients wrt parameters: **Backward Pass**



*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

**Step 1:** Compute Loss on Mini-Batch: **Forward Pass**

**Step 2:** Compute Gradients wrt parameters: **Backward Pass**



*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

Georgia Tech

- We want to compute: $\left\{\dfrac{\partial L}{\partial h^{\ell-1}}, \dfrac{\partial L}{\partial W}\right\}$



- We will use the *chain rule* to do this:

**Chain Rule:** $\dfrac{\partial z}{\partial x} = \dfrac{\partial z}{\partial y} \cdot \dfrac{\partial y}{\partial x}$

**Step 1:** Compute Loss on Mini-Batch: **Forward Pass**

**Step 2:** Compute Gradients wrt parameters: **Backward Pass**

**Step 3:** Use **gradient** to update **all parameters** at the end



$$w_i = w_i - \alpha \frac{\partial L}{\partial w_i}$$

**Backpropagation is the application of gradient descent to a computation graph via the chain rule!**

*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

Georgia Tech

Given this computation graph, the training algorithm will:

◆ Calculate the current model's outputs (called the **forward pass**)

◆ Calculate the gradients for each module (called the **backward pass**)

Backward pass is a recursive algorithm that:

◆ Starts at **loss function** where we know how to calculate the gradients

◆ Progresses back through the modules

◆ Ends in the **input layer** where we do not need gradients (no parameters)

This algorithm is called **backpropagation**

Input          Function          Output

$h^{\ell-1}$                          $h^{\ell}$

$W$

Parameters

**Overview of Training**

Georgia Tech

In the **backward pass**, we seek to calculate the gradients of the loss with respect to the module's parameters

- Assume that we have the gradient of the loss with respect to the **module's outputs** (given to us by upstream module)

- We will also pass the gradient of the loss with respect to the **module's inputs**

  - This is not required for update the module's weights, but passes the gradients back to the previous module



$$\frac{\partial L}{\partial h^{\ell-1}} \quad \quad \frac{\partial L}{\partial h^{\ell}}$$

$$\frac{\partial L}{\partial W}$$

**Problem:**

- We can compute local gradients:
  $$\{\frac{\partial h^{\ell}}{\partial h^{\ell-1}}, \frac{\partial h^{\ell}}{\partial W}\}$$

- We are given: $\frac{\partial L}{\partial h^{\ell}}$

- Compute: $\{\frac{\partial L}{\partial h^{\ell-1}}, \frac{\partial L}{\partial W}\}$

*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

**Backward Pass Computations**

Georgia Tech

- We can compute **local gradients**: $\{\frac{\partial h^{\ell}}{\partial h^{\ell-1}}, \frac{\partial h^{\ell}}{\partial W}\}$

- This is just the **derivative of our function** with respect to its parameters and inputs!

**Example:**  If  $h^{\ell} = Wh^{\ell-1}$

then  $\frac{\partial h^{\ell}}{\partial h^{\ell-1}} = W$

and  $\frac{\partial h^{\ell}}{\partial w_i} = h^{\ell-1,T}$  (a sparse matrix with $h^{\ell-1,T}$ in the *i*-th row

Georgia Tech

- We want to to compute: $\left\{\dfrac{\partial L}{\partial h^{\ell-1}}, \dfrac{\partial L}{\partial W}\right\}$



- We will use the *chain rule* to do this:

**Chain Rule:** $\dfrac{\partial z}{\partial x} = \dfrac{\partial z}{\partial y} \cdot \dfrac{\partial y}{\partial x}$

**Computing the Gradients of Loss**

Georgia Tech

We will use the **chain rule** to compute: $\{\frac{\partial L}{\partial h^{\ell-1}}, \frac{\partial L}{\partial W}\}$

**Gradient of loss w.r.t. inputs:** $\frac{\partial L}{\partial h^{\ell-1}} = \frac{\partial L}{\partial h^\ell} \frac{\partial h^\ell}{\partial h^{\ell-1}}$

Given by upstream module **(upstream gradient)**

**Gradient of loss w.r.t. weights:** $\frac{\partial L}{\partial W} = \frac{\partial L}{\partial h^\ell} \frac{\partial h^\ell}{\partial W}$

$\frac{\partial L}{\partial h^{\ell-1}}$

$\frac{\partial L}{\partial h^\ell}$

$\frac{\partial L}{\partial W}$

*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

# Backpropagation: a simple example

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

Georgia Tech

# Backpropagation: a simple example

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

# Backpropagation: a simple example

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

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

Georgia Tech

# Backpropagation: a simple example

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

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



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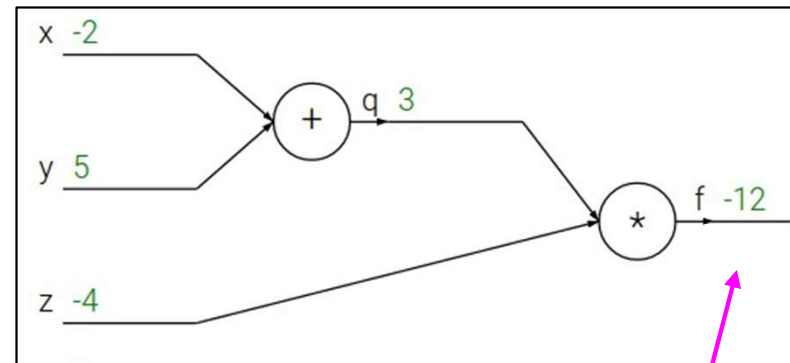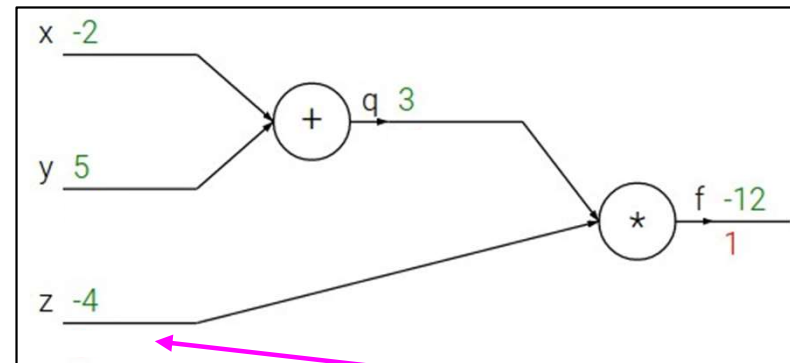
Georgia Tech

# Backpropagation: a simple example

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

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

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

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



x -2
y 5
q 3
z -4
f -12

# Backpropagation: a simple example

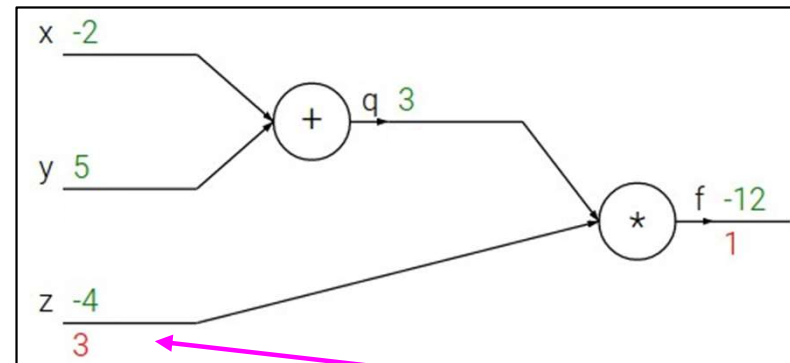

$$f(x, y, z) = (x + y)z$$

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

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

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

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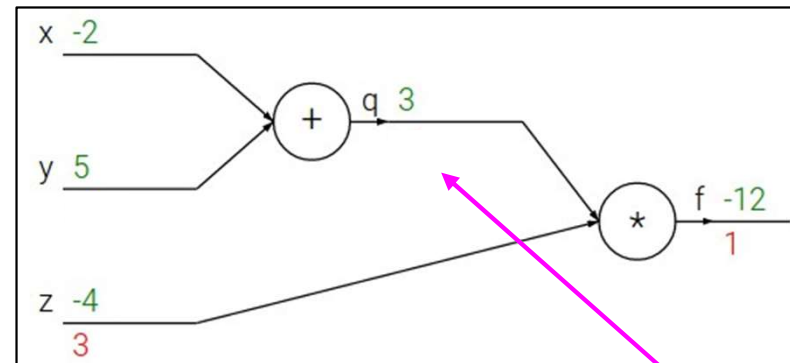
Georgia Tech

# Backpropagation: a simple example

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

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

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

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

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



$$\frac{\partial f}{\partial f}$$
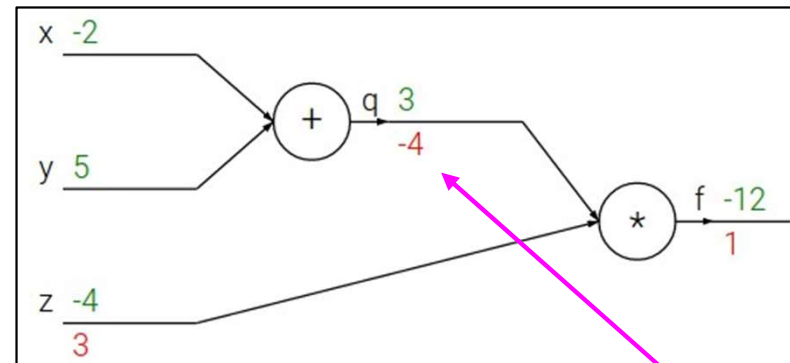
Georgia Tech

# Backpropagation: a simple example

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

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

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

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

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



$$\frac{\partial f}{\partial f}$$
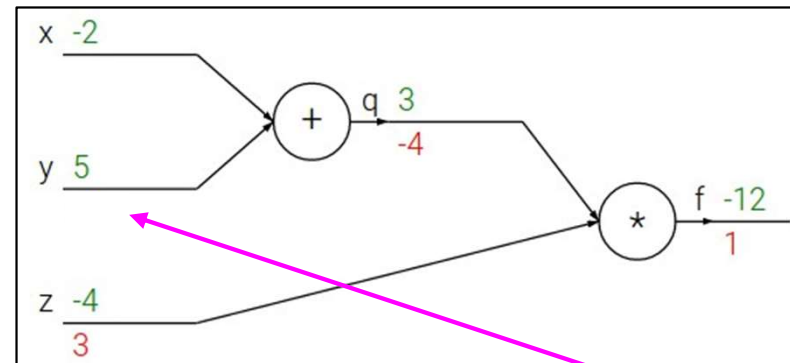
Georgia Tech

# Backpropagation: a simple example

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

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

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

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

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



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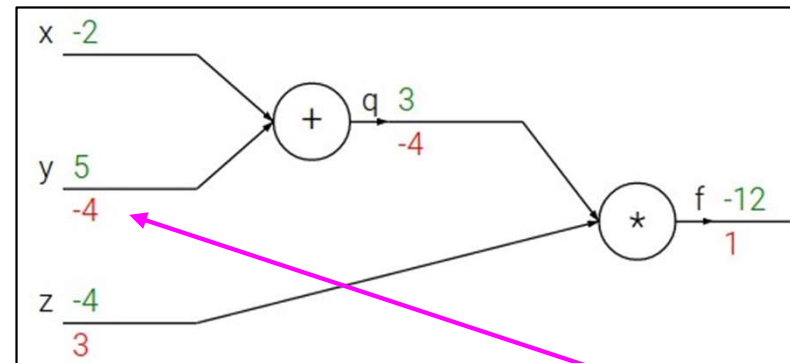
Georgia Tech

# Backpropagation: a simple example



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

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

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

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

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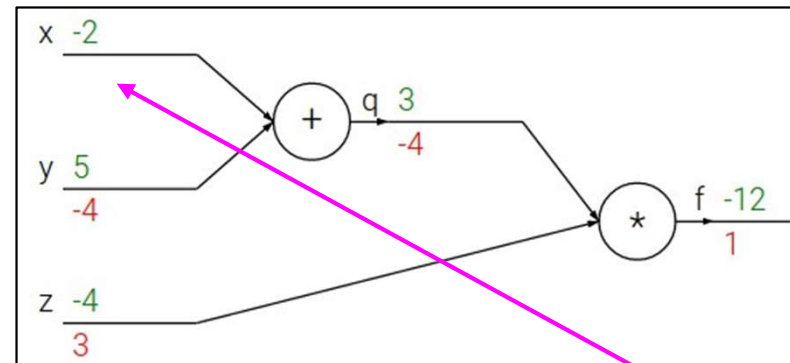
Georgia Tech

# Backpropagation: a simple example

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

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

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

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

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



$$\frac{\partial f}{\partial q}$$
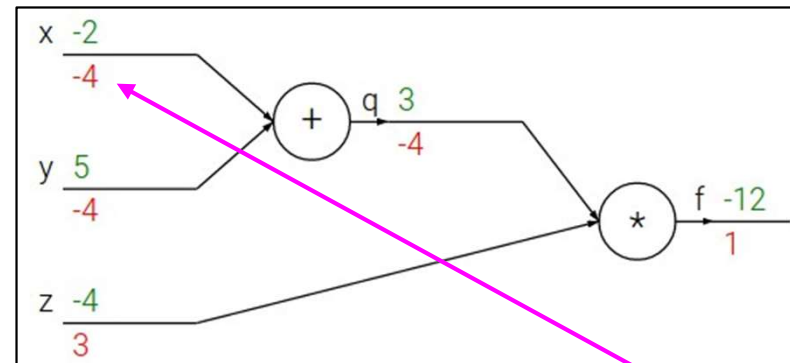
Georgia Tech

# Backpropagation: a simple example



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

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

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

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

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

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

Georgia Tech

# Backpropagation: a simple example

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

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

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

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

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



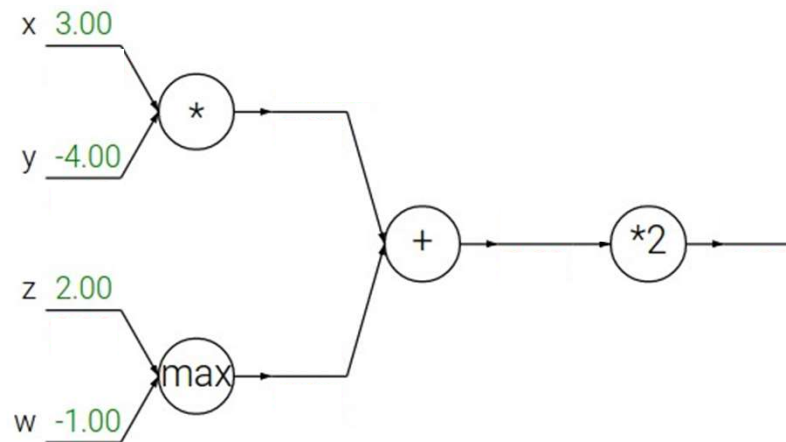$\dfrac{\partial f}{\partial y}$

Chain rule:

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$
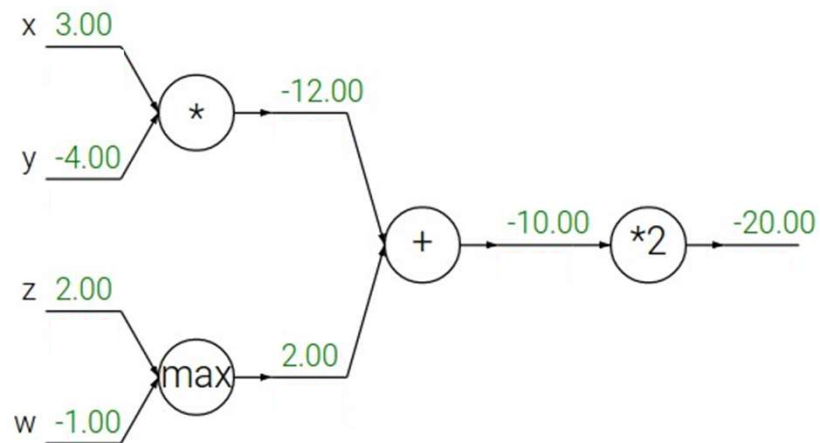
Upstream gradient       Local gradient

Georgia Tech

# Backpropagation: a simple example

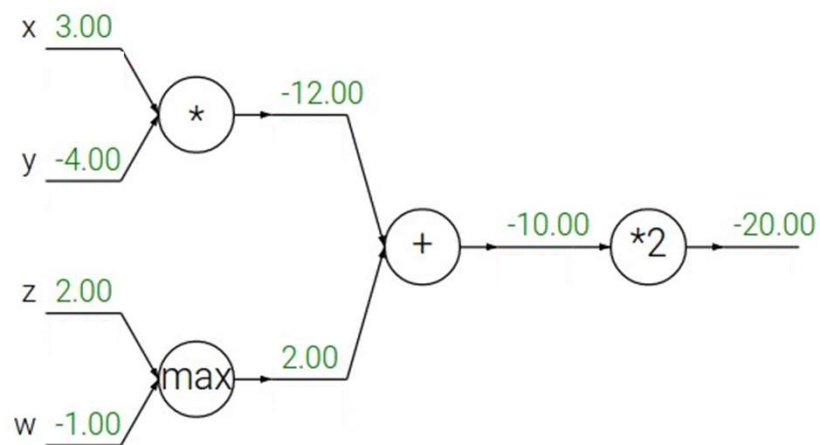$$f(x, y, z) = (x + y)z$$

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

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

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

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



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

Chain rule:

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

Upstream gradient      Local gradient

Georgia Tech

# Backpropagation: a simple example

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

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

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

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

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



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

Chain rule:

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

Upstream gradient    Local gradient

# Backpropagation: a simple example

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

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

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

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

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



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

Chain rule:

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

Upstream gradient    Local gradient

Georgia Tech

# Backpropagation: a simple example

# Backpropagation: a simple example

Georgia
Tech

# Patterns in backward flow

# Patterns in backward flow

Q: What is an **add** gate?

Georgia
Tech

# Patterns in backward flow

**add** gate: gradient distributor

Georgia
Tech

# Patterns in backward flow

**add** gate: gradient distributor

Q: What is a **max** gate?

Georgia Tech

# Patterns in backward flow

**add** gate: gradient distributor

**max** gate: gradient router

# Patterns in backward flow

**add** gate: gradient distributor

**max** gate: gradient router

Q: What is a **mul** gate?

Georgia
Tech

# Patterns in backward flow

**add** gate: gradient distributor

**max** gate: gradient router

**mul** gate: gradient switcher

# Gradients add at branches

Georgia Tech

# Duality in Fprop and Bprop

# Deep Learning = Differentiable Programming

- Computation = Graph
  - Input = Data + Parameters
  - Output = Loss
  - Scheduling = Topological ordering

- What do we need to do?
  - Generic code for representing the graph of modules
  - Specify modules (both forward and backward function)

Georgia Tech

# Modularized implementation: forward / backward API

Graph (or Net) object  *(rough psuedo code)*



```python
class ComputationalGraph(object):
    #...
    def forward(inputs):
        # 1. [pass inputs to input gates...]
        # 2. forward the computational graph:
        for gate in self.graph.nodes_topologically_sorted():
            gate.forward()
        return loss # the final gate in the graph outputs the loss
    def backward():
        for gate in reversed(self.graph.nodes_topologically_sorted()):
            gate.backward() # little piece of backprop (chain rule applied)
        return inputs_gradients
```

# Modularized implementation: forward / backward API



x

z

*

y

(x,y,z are scalars)

```
class MultiplyGate(object):
    def forward(x,y):
        z = x*y
        return z
    def backward(dz):
        # dx = ... #todo
        # dy = ... #todo
        return [dx, dy]
```

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

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

# Modularized implementation: forward / backward API

x

z

*

y

(x,y,z are scalars)

```python
class MultiplyGate(object):
    def forward(x,y):
        z = x*y
        self.x = x # must keep these around!
        self.y = y
        return z
    def backward(dz):
        dx = self.y * dz # [dz/dx * dL/dz]
        dy = self.x * dz # [dz/dy * dL/dz]
        return [dx, dy]
```

Georgia Tech

# Example: Caffe layers

# Caffe Sigmoid Layer



```
 1  #include <cmath>
 2  #include <vector>
 3
 4  #include "caffe/layers/sigmoid_layer.hpp"
 5
 6  namespace caffe {
 7
 8  template <typename Dtype>
 9  inline Dtype sigmoid(Dtype x) {
10    return 1. / (1. + exp(-x));
11  }
12
13  template <typename Dtype>
14  void SigmoidLayer<Dtype>::Forward_cpu(const vector<Blob<Dtype>*>& bottom,
15      const vector<Blob<Dtype>*>& top) {
16    const Dtype* bottom_data = bottom[0]->cpu_data();
17    Dtype* top_data = top[0]->mutable_cpu_data();
18    const int count = bottom[0]->count();
19    for (int i = 0; i < count; ++i) {
20      top_data[i] = sigmoid(bottom_data[i]);
21    }
22  }
23
24  template <typename Dtype>
25  void SigmoidLayer<Dtype>::Backward_cpu(const vector<Blob<Dtype>*>& top,
26      const vector<bool>& propagate_down,
27      const vector<Blob<Dtype>*>& bottom) {
28    if (propagate_down[0]) {
29      const Dtype* top_data = top[0]->cpu_data();
30      const Dtype* top_diff = top[0]->cpu_diff();
31      Dtype* bottom_diff = bottom[0]->mutable_cpu_diff();
32      const int count = bottom[0]->count();
33      for (int i = 0; i < count; ++i) {
34        const Dtype sigmoid_x = top_data[i];
35        bottom_diff[i] = top_diff[i] * sigmoid_x * (1. - sigmoid_x);
36      }
37    }
38  }
39
40  #ifdef CPU_ONLY
41  STUB_GPU(SigmoidLayer);
42  #endif
43
44  INSTANTIATE_CLASS(SigmoidLayer);
45
46
47  }  // namespace caffe
```

Caffe is licensed under BSD 2-Clause

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

$(1 - \sigma(x))\,\sigma(x)$ * top_diff  (chain rule)

Georgia Tech

Batches of data are **matrices** or **tensors** (multi-dimensional matrices)

**Examples:**

- Each instance is a vector of size $m$, our batch is of size $[B \times m]$
- Each instance is a matrix (e.g. grayscale image) of size $W \times H$, our batch is $[B \times W \times H]$
- Each instance is a multi-channel matrix (e.g. color image with R,B,G channels) of size $C \times W \times H$, our batch is $[B \times C \times W \times H]$

**Jacobians become tensors which is complicated**

- Instead, flatten input to a vector and get a vector of derivatives!
- In practice, figure out Jacobians for simpler items (scalars, vectors), figure out pattern, and *slice* or index appropriate elements to create Jacobians
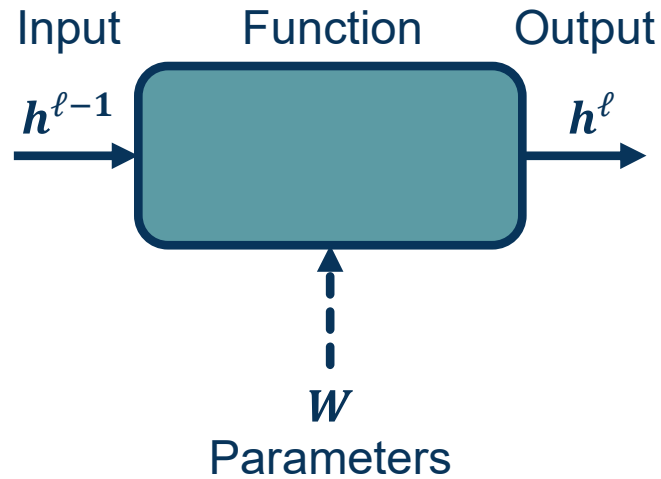
$$\begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nn} \end{bmatrix}$$

**Flatten** ⬇

$$\begin{bmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{21} \\ x_{22} \\ \vdots \\ x_{n1} \\ \vdots \\ x_{nn} \end{bmatrix}$$

**Jacobians of Batches**

Georgia Tech

Input       Function      Output

$h^{\ell-1}$                              $h^{\ell}$

$W$

Parameters

$$h^{\ell} = W h^{\ell-1}$$

$$\leftarrow w_i^T \rightarrow$$

$$|h^{\ell}| \times 1 \qquad |h^{\ell}| \times |h^{\ell-1}| \qquad |h^{\ell-1}| \times 1$$

**Fully Connected (FC) Layer: Forward Function**

Georgia Tech

$$\frac{\partial L}{\partial h^{\ell-1}} \qquad \frac{\partial L}{\partial h^{\ell}}$$

$$\frac{\partial L}{\partial W}$$

Note doing this on full $W$ matrix would result in Jacobian tensor!

But it is *sparse* – each output only affected by corresponding weight row

$$\frac{\partial h^{\ell}}{\partial h^{\ell-1}} = W$$

$$\frac{\partial h^{\ell}}{\partial w_i} = h^{(\ell-1),T}$$

$$\frac{\partial L}{\partial h^{\ell-1}} = \frac{\partial L}{\partial h^{\ell}} \frac{\partial h^{\ell}}{\partial h^{\ell-1}}$$

$$\begin{bmatrix} \ \end{bmatrix} \begin{bmatrix} \ \end{bmatrix} \begin{bmatrix} \ \\ \ \\ \ \end{bmatrix}$$

$$1 \times |h^{\ell-1}| \quad 1 \times |h^{\ell}| \quad |h^{\ell}| \times |h^{\ell-1}|$$

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial h^{\ell}} \frac{\partial h^{\ell}}{\partial w_i}$$

$$\begin{bmatrix} \ \end{bmatrix} \begin{bmatrix} \ \end{bmatrix} \begin{bmatrix} \leftarrow \quad 0 \quad \rightarrow \\ \leftarrow \ \frac{\partial h^{\ell}}{\partial w_i} \ \rightarrow \\ \leftarrow \quad 0 \quad \rightarrow \end{bmatrix}$$

$$1 \times |h^{\ell-1}| \quad 1 \times |h^{\ell}| \quad |h^{\ell}| \times |h^{\ell-1}|$$

**Fully Connected (FC) Layer**

Georgia Tech