Topics:
- Gradient Descent
- Neural Networks

# CS 4803-DL / 7643-A
# ZSOLT KIRA

- **Assignment 1 out!**
  - **Due Feb 7th**
  - Start now, start now, start now!
  - Start now, start now, start now!
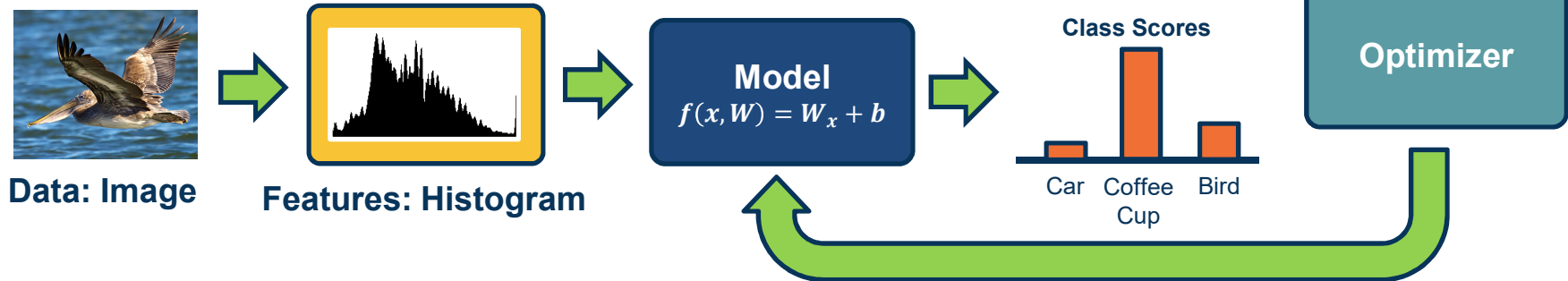  - Start now, start now, start now!

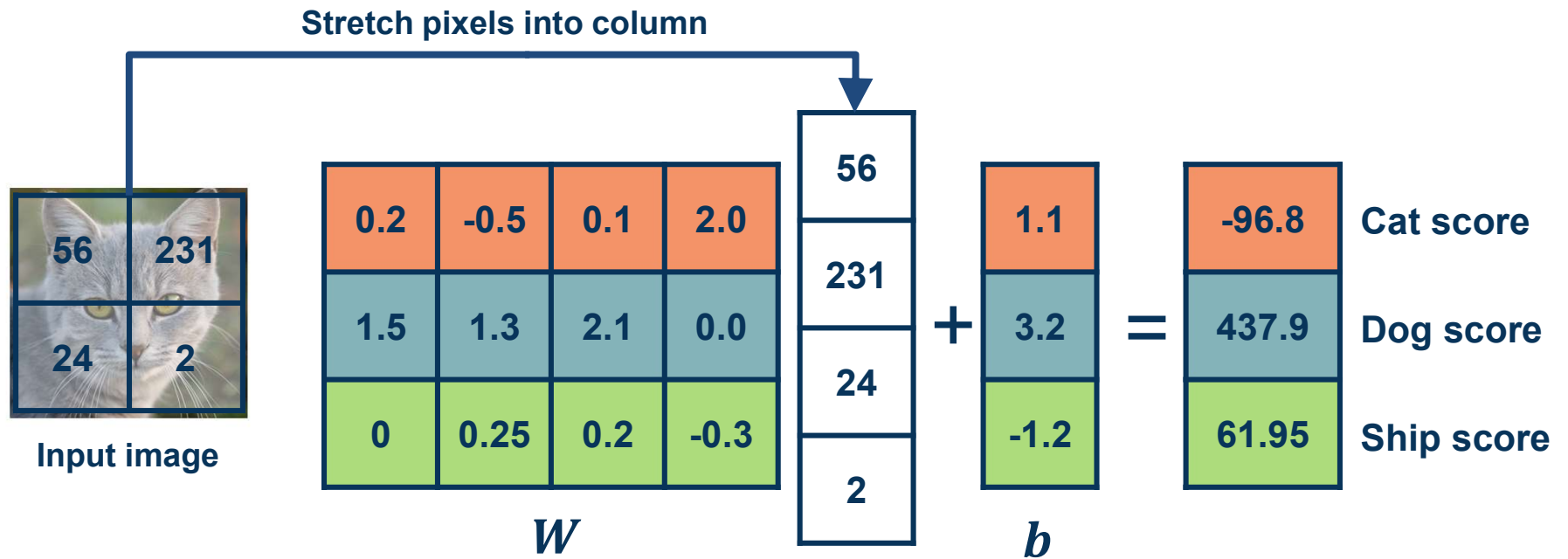- **Piazza**
  - Be active!!!

- **Office hours**
  - Let us know special topic requests (e.g. PS0, Assignment 1, research paper discussion, etc. )

- Input (and representation)
- Functional form of the model
  - Including parameters
- Performance measure to improve
  - Loss or objective function
- Algorithm for finding best parameters
  - Optimization algorithm

Class Scores

Car   Coffee   Bird
      Cup

Loss Function

Optimizer

Data: Image

Features: Histogram

Model
$f(x, W) = W_x + b$

Class Scores

Car   Coffee   Bird
      Cup

**Components of a Parametric Model**

Georgia Tech

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

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

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

**Example**

Georgia Tech

# Softmax Classifier (Multinomial Logistic Regression)

Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

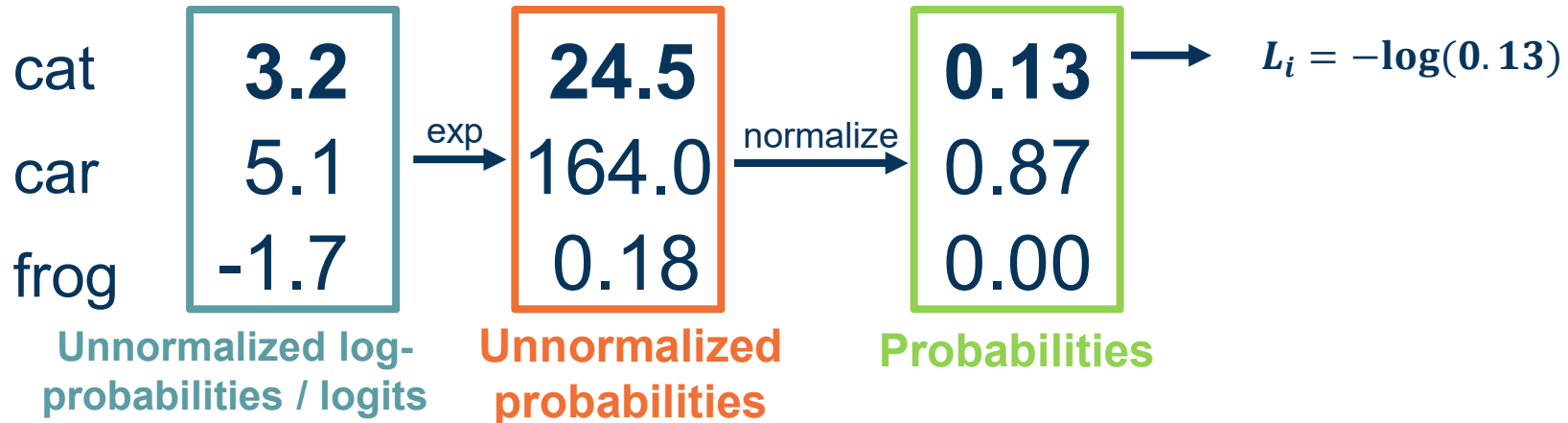$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

**Softmax Function**

$$L_i = -\log P(Y = y_i | X = x_i)$$

Probabilities must be >= 0

Probabilities must sum to 1

| | cat | car | frog |
|---|---|---|---|
| Unnormalized log-probabilities / logits | **3.2** | 5.1 | -1.7 |

exp →

| Unnormalized probabilities | **24.5** | 164.0 | 0.18 |
|---|---|---|---|

normalize →

| Probabilities | **0.13** | 0.87 | 0.00 |
|---|---|---|---|

→ $L_i = -\log(0.13)$

**Cross-Entropy Loss Example**

Georgia Tech

Often, we add a **regularization term** to the loss function

**L1 Regularization**

$$L_i = |y - W x_i|^2 + |W|$$

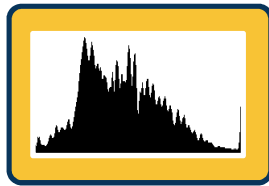**Example regularizations:**

⬡ L1/L2 on weights (encourage small values)

- Input (and representation)
- Functional form of the model
  - Including parameters
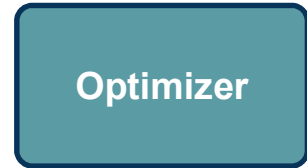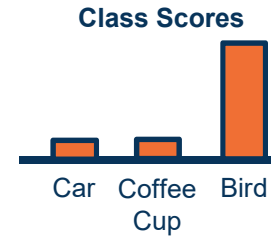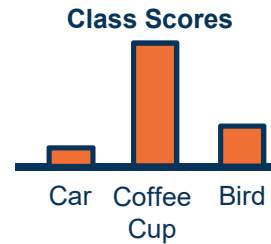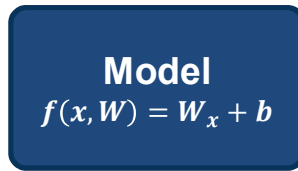- Performance measure to improve
  - Loss or objective function
- **Algorithm for finding best parameters**
  - **Optimization algorithm**

Class Scores

Car    Coffee    Bird
        Cup

**Loss Function**

**Optimizer**

Class Scores

Car    Coffee    Bird
        Cup

**Model**
$f(x, W) = W_x + b$

**Data: Image**

**Features: Histogram**

**Components of a Parametric Model**

Georgia Tech

Given a model and loss function, finding the best set of weights is a **search problem**

- Find the best combination of weights that minimizes our loss function

**Several classes of methods:**

- Random search
- Genetic algorithms (population-based search)
- Gradient-based optimization

In deep learning, **gradient-based methods are dominant** although not the only approach possible

$$\begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1m} & b1 \\ w_{21} & w_{22} & \cdots & w_{2m} & b2 \\ w_{21} & w_{22} & \cdots & w_{3m} & b3 \end{bmatrix}$$

**Loss**

**As weights change, the loss changes as well**

- This is often somewhat-smooth locally, so small changes in weights produce small changes in the loss

We can therefore think about **iterative algorithms** that take **current values of weights and modify them** a bit

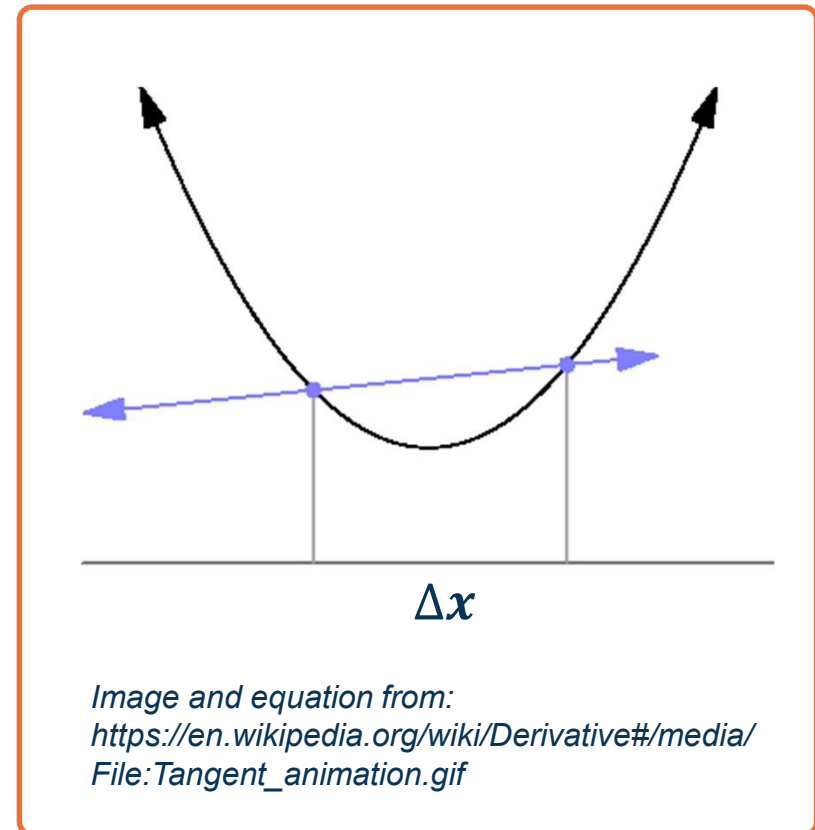Strategy: Follow the Slope!

- We can find the steepest descent direction by computing the **derivative (gradient):**
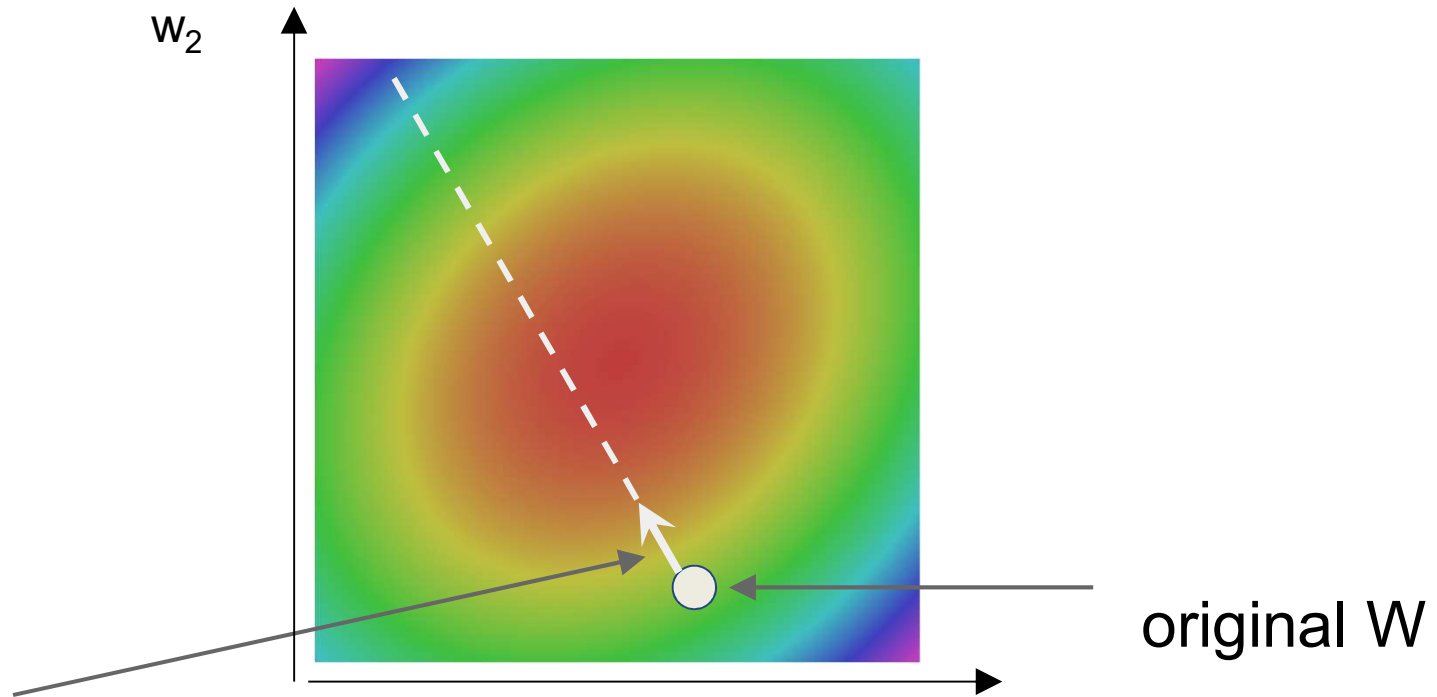
$$f'(a) = \lim_{h \to 0} \frac{f(a+h) - f(a)}{h}$$

- 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*

**Derivatives**

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

- Choose a model: $f(x, W) = Wx$

- 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: $w_i = w_i - \alpha \frac{\partial L}{\partial w_i}$
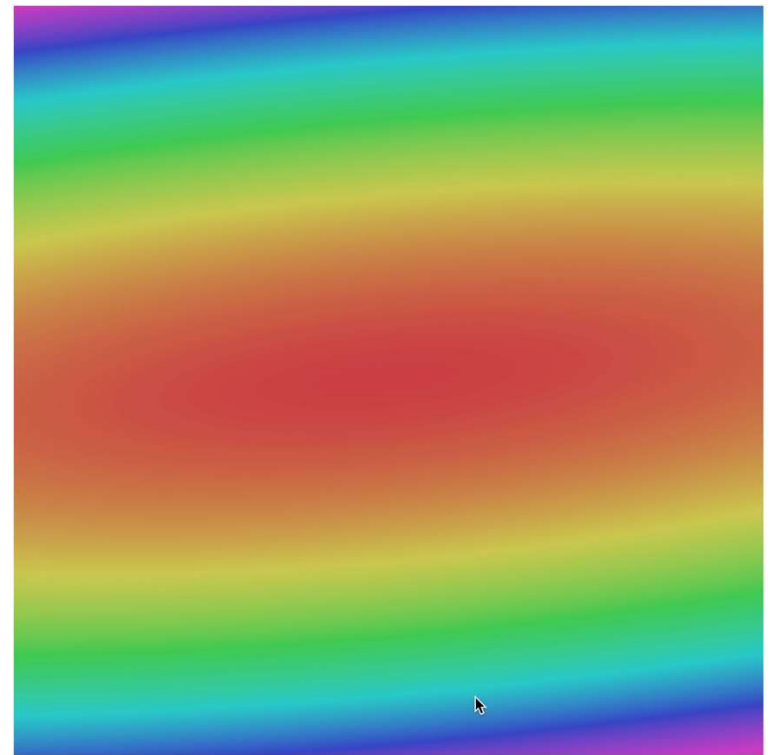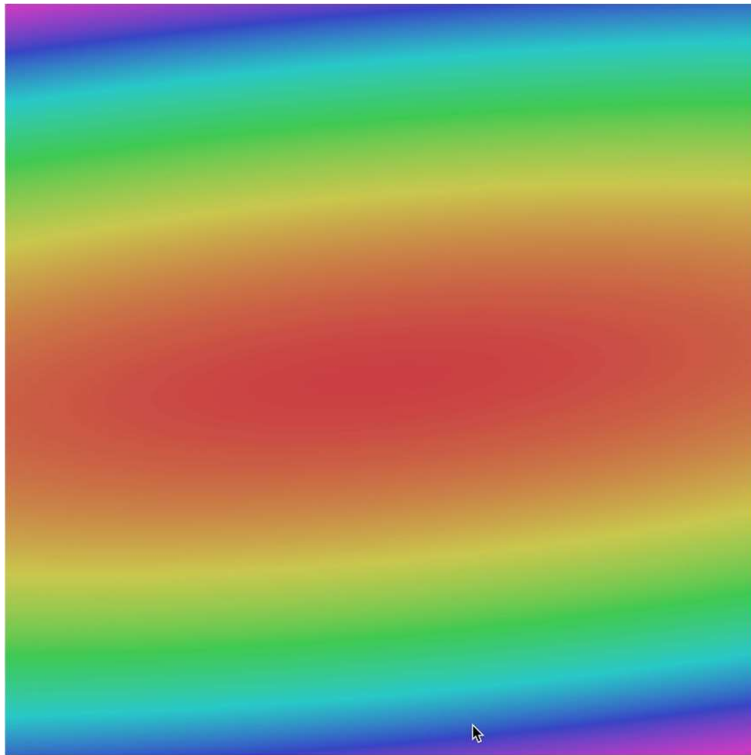
- Repeat (from Step 3)

**Gradient Descent**

Georgia Tech

$w_2$

original W

negative gradient direction

$w_1$

**Gradient Descent**

Georgia
Tech

$$W_1$$

**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

Gradient descent is guaranteed to converge under some conditions

- For example, learning rate has to be appropriately reduced throughout training

- It will converge to a *local* minima

  - Small changes in weights would not decrease the loss

- It turns out that some of the local minima that it finds in practice (if trained well) are still pretty good!

We know how to compute the **model output and loss function**

**Several ways to compute** $\dfrac{\partial L}{\partial w_i}$

- ⬢ Manual differentiation

- ⬢ Symbolic differentiation

- ⬢ Numerical differentiation

- ⬢ Automatic differentiation

Georgia Tech

**current W:**

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25347**

**gradient dW:**

[?,
?,
?,
?,
?,
?,
?,
?,
?,…]

| **current W:** | **W + h** (first dim): | **gradient dW:** |
|---|---|---|
| [0.34, | [0.34 + **0.0001**, | [?, |
| -1.11, | -1.11, | ?, |
| 0.78, | 0.78, | ?, |
| 0.12, | 0.12, | ?, |
| 0.55, | 0.55, | ?, |
| 2.81, | 2.81, | ?, |
| -3.1, | -3.1, | ?, |
| -1.5, | -1.5, | ?, |
| 0.33,…] | 0.33,…] | ?,…] |
| **loss 1.25347** | **loss 1.25322** | |

**current W:**

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25347**

**W + h** (first dim)**:**

[0.34 + **0.0001**,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25322**

**gradient dW:**

[**-2.5**,
?,
?,

(1.25322 - 1.25347)/0.0001
= -2.5

$$\frac{df(x)}{dx} = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$

?,
?,…]

**current W:**

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25347**

**W + h** (second dim)**:**

[0.34,
-1.11 + **0.0001**,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25353**

**gradient dW:**

[-2.5,
?,
?,
?,
?,
?,
?,
?,
?,…]

**current W:**

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25347**

**W + h** (second dim)**:**

[0.34,
-1.11 + **0.0001**,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25353**

**gradient dW:**

[-2.5,
**0.6**,
?,
?,

(1.25353 - 1.25347)/0.0001
= 0.6

$$\frac{df(x)}{dx} = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$

?,…]

| current W: | W + h (third dim): | gradient dW: |
|---|---|---|
| [0.34, | [0.34, | [-2.5, |
| -1.11, | -1.11, | 0.6, |
| 0.78, | 0.78 + **0.0001**, | ?, |
| 0.12, | 0.12, | ?, |
| 0.55, | 0.55, | ?, |
| 2.81, | 2.81, | ?, |
| -3.1, | -3.1, | ?, |
| -1.5, | -1.5, | ?, |
| 0.33,…] | 0.33,…] | ?,…] |
| **loss 1.25347** | **loss 1.25347** | |

**current W:**

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25347**

**W + h** (third dim)**:**

[0.34,
-1.11,
0.78 + **0.0001**,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25347**

**gradient dW:**

[-2.5,
0.6,
**0**,
?,

(1.25347 - 1.25347)/0.0001
= 0

$$\frac{df(x)}{dx} = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$

?,…]

# Numerical vs Analytic Gradients

$$\frac{df(x)}{dx} = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$

**Numerical gradient**: slow :(, approximate :(, easy to write :)
**Analytic gradient**: fast :), exact :), error-prone :(

In practice: Derive analytic gradient, check your
implementation with numerical gradient.
This is called a **gradient check.**

We know how to compute the **model output and loss function**

**Several ways to compute** $\frac{\partial L}{\partial w_i}$

- Manual differentiation

- Symbolic differentiation

- Numerical differentiation

- Automatic differentiation

Georgia Tech

For some functions, we can analytically derive the partial derivative

**Example:**

**Derivation of Update Rule**

**Function**

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

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

**Loss**

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

**Update Rule**

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

**Manual Differentiation**

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 $w \cdot x_i$)

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

### Update Rule

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

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

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

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

...where...

$$\delta_k = y_k - w^T x_k$$

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

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

**Manual Differentiation**

Georgia Tech

If we add a **non-linearity (sigmoid),** derivation is more complex

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

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

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

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

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

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

$$= \sum_i -2\delta_i\sigma(d_i)(1 - \sigma(d_i))x_{ij}$$

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



**The sigmoid perception update rule:**

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

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

$$\delta_i = y_i - \sigma_i$$

**Adding a Non-Linear Function**

Georgia Tech

Neural Network View of a Linear Classifier

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

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

**Input**                **Model**                **Loss Function**

**What Does a Linear Classifier Consist of?**

Georgia Tech

# A simple **neural network** has similar structure as our linear classifier:

- A neuron takes input (firings) from other neurons **(-> input to linear classifier)**
- The inputs are summed in a weighted manner **(-> weighted sum)**
  - Learning is through a modification of the weights
- If it receives enough input, it "fires" (threshold or if weighted sum plus bias is high enough)

**Impulses carried toward cell body**

dendrite

presynaptic terminal

axon

**cell body**

**Impulses carried away from cell body**

$x_0$

$w_0$

axon from a neuron

synapse

$w_0 x_0$

dendrite

cell body

$w_1 x_1$

$\sum_i w_i x_i + b$

$f$

$f\left(\sum_i w_i x_i + b\right)$

output axon

activation function

$w_2 x_2$

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

**Origins of the Term Neural Network**

Georgia Tech

As we did before, the output of a neuron can be modulated by a non-linear function (e.g. sigmoid)



**Sigmoid Activation Function**
$$\frac{1}{1 + e^{-x}}$$

**Impulses carried toward cell body**



dendrite

presynaptic terminal

axon

cell body

**Impulses carried away from cell body**

$x_0$

$w_0$ synapse

axon from a neuron

$w_0 x_0$

dendrite

cell body

$w_1 x_1$

$\sum_i w_i x_i + b$  $f$

$f\left(\sum_i w_i x_i + b\right)$

output axon

activation function

$w_2 x_2$

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

**Adding Non-Linearities**

Georgia Tech

We can have **multiple** neurons connected to the same input

Corresponds to a multi-class classifier

◆ Each output node outputs the score for a class

$$f(x, W) = \sigma(Wx + b) \begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1m} & b1 \\ w_{21} & w_{22} & \cdots & w_{2m} & b2 \\ w_{21} & w_{22} & \cdots & w_{3m} & b3 \end{bmatrix}$$

◆ Often called fully connected layers

  ◆ Also called a linear *projection layer*



input layer

output layer

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

**Connecting Many Neurons**

- Each input/output is a **neuron (node)**

- A linear classifier (+ optional non-linearity) is called a **fully connected** layer

- Connections are represented as **edges**

- Output of a particular neuron is referred to as **activation**

- This will be expanded as we view computation in a neural network as a **graph**



**input layer**

**output layer**

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

We can **stack** multiple layers together

- Input to second layer is output of first layer

Called a **2-layered neural network** (input is not counted)

Because the middle layer is neither input or output, and we don't know what their values represent, we call them **hidden** layers

- We will see that they end up learning effective features

This **increases** the representational power of the function!

- Two layered networks can represent any continuous function



**input layer**

**hidden layer**

**output layer**

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

**Connecting Many Layers**

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**

$x$ $W_1$ $W_2$

$$=$$

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

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

**The Linear Algebra View**

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**:



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

**Adding More Layers!**

# 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



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

Given a library of simple functions

$$\sin(x)$$
$$\log(x)$$
$$\cos(x)$$
$$x^3$$
$$\exp(x)$$

Compose into a

complicate function

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

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

*Adapted from slides by: Marc'Aurelio Ranzato, Yann LeCun*

**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*

Georgia Tech

# Directed Acyclic Graphs (DAGs)

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

# Directed Acyclic Graphs (DAGs)

- Concept
  - Topological Ordering

# Directed Acyclic Graphs (DAGs)

Georgia
Tech

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

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



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

**Machine Learning Example**

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

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

**Overview of 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

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

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

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

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

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

and $\frac{\partial h_i^{\ell}}{\partial w_i} = h^{\ell-1,T}$

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

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*

**Computing the Gradients of Loss**

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*

**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*

**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

# 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

# Backpropagation: a simple example

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

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



Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
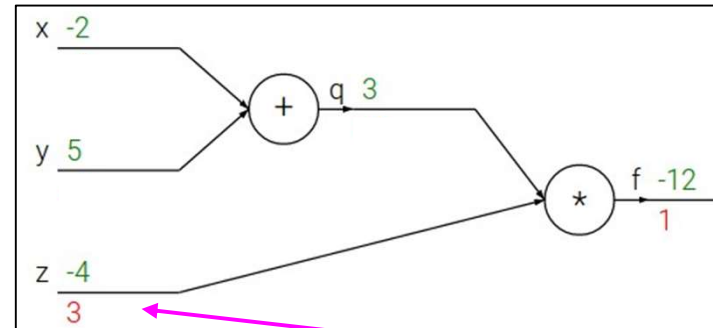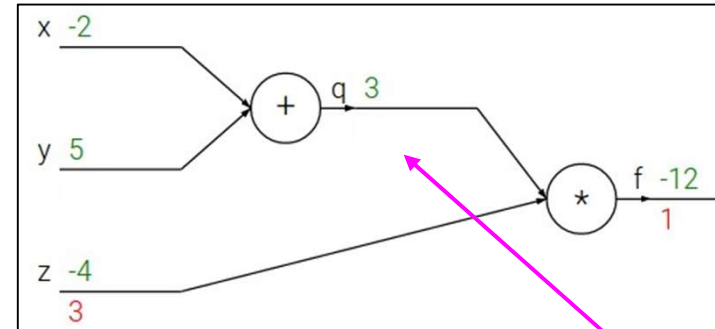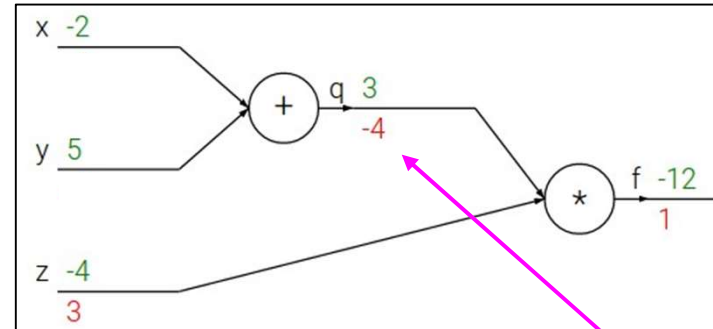
# Backpropagation: a simple example

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

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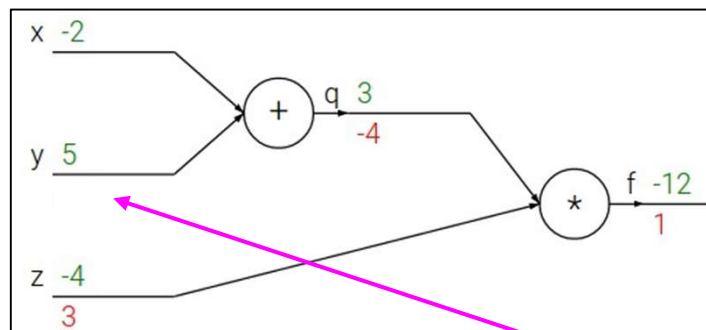
Georgia Tech

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

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

# Backpropagation: a simple example
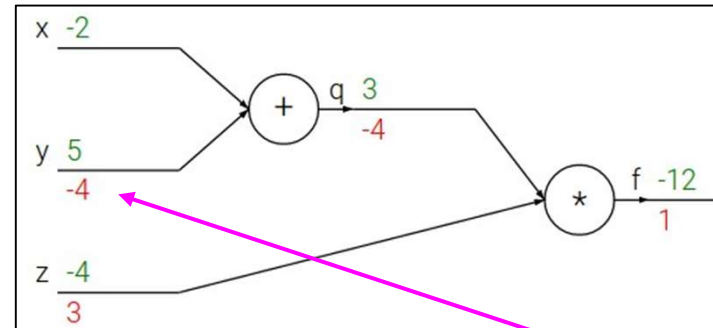
$$f(x, y, z) = (x + y)z$$

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



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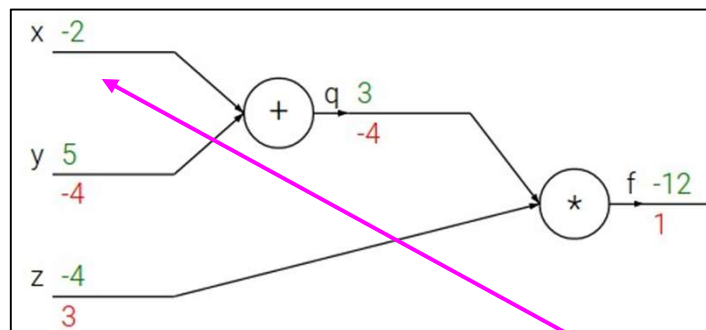
Georgia Tech

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



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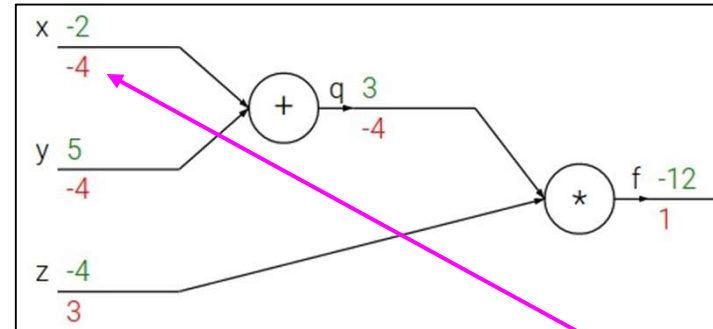
Georgia Tech

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



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

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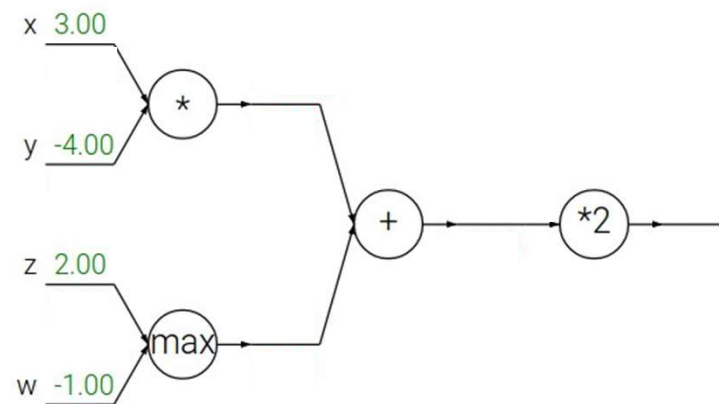


$$\frac{\partial f}{\partial q}$$

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



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

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



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

Chain rule:

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

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

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

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

# Backpropagation: a simple example

Georgia Tech

# Patterns in backward flow

Georgia Tech

# Patterns in backward flow

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

Georgia
Tech

# Patterns in backward flow

**add** gate: gradient distributor

# 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

Georgia
Tech

# Patterns in backward flow

**add** gate: gradient distributor

**max** gate: gradient router

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

# 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

Georgia Tech

# 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)

# 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)

```python
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]
```

# Example: Caffe layers

# Caffe Sigmoid Layer

```cpp
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)