# CSC311 Introduction to Machine Learning

## The Backpropagation Algorithm

Alice Gao and Marina Tawfik

# Learning Outcomes

By the end of this lecture, students should be able to

- Explain why the backpropagation algorithm can efficiently compute many gradients for a neural network.

- Derive non-vectorized expressions for the gradients with respect to weights, biases, pre-activations, and activations in a small neural network.

- Explain the role of the forward pass in preparing quantities required for backpropagation.

# Outline

- Challenge of Training a Neural Network

- Gradients for a 3-Layer Neural Network

- Computing Gradients Efficiently

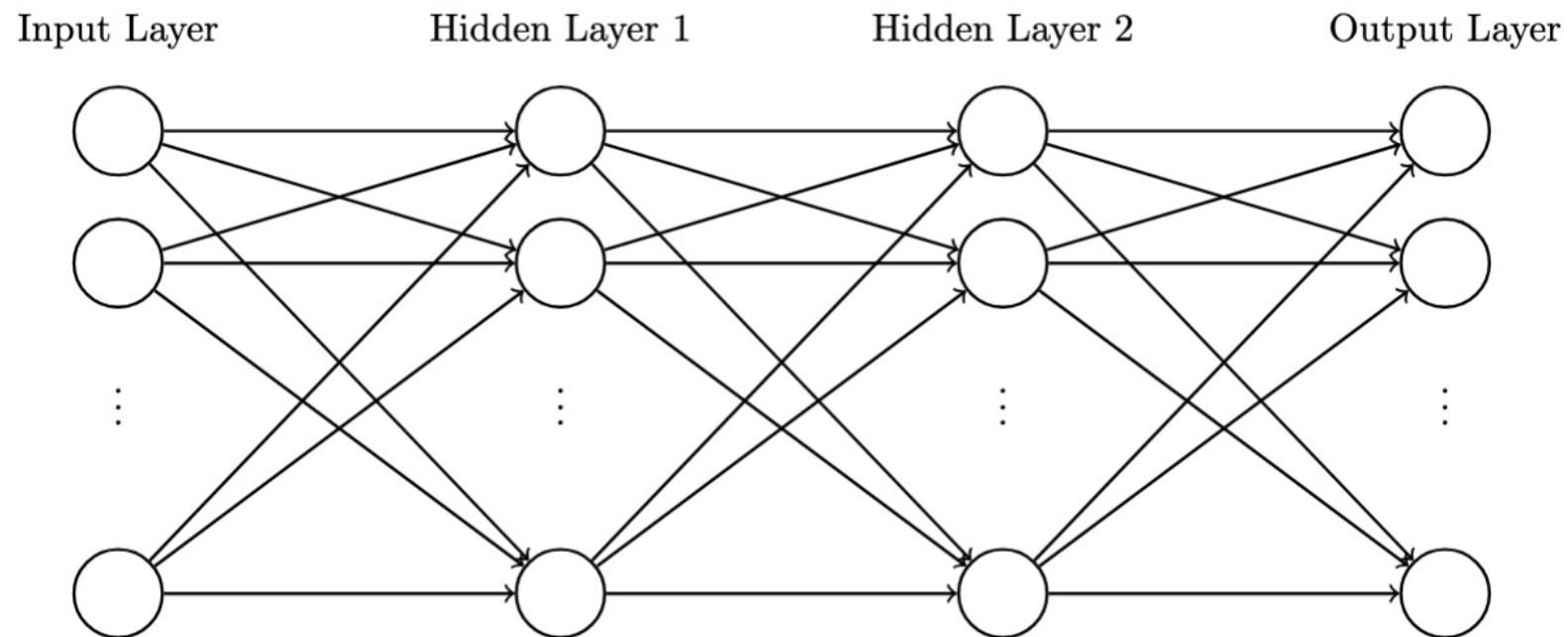- Backpropagation for a 3-Layer Neural Network

- Forward Pass

UNIVERSITY OF TORONTO

# Challenge of Training a Neural Network

# Optimizing a Neural Network

How do we find good weights for a neural network?

Use gradient descent to adjust the weights to reduce the loss.

# Gradient Descent for a Neural Network

1. Initialize the weights **w**.

2. Perform the updates.

This is what we called $\alpha$ in previous lectures.

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \, \nabla_\mathbf{w} \, \mathcal{E}(\mathbf{w})$$

$\alpha$

where

- $\eta$ is the learning rate

- $\nabla_\mathbf{w} \, \mathcal{E}(\mathbf{w})$ is the gradient of cost function with respect to the weights **w**

# The Challenge of Computing All the Gradients

To train a neural network with gradient descent

Need the partial derivative of loss function with respect to every weight.

But a neural network can have millions of weights.

How can we compute the gradient of the cost function efficiently?

The Backpropagation Algorithm

UNIVERSITY OF
TORONTO

# Gradients for 3-Layer Neural Network

# 3-Layer Neural Network

Here, we have a single output i.e. this would work for regression & binary classification.

Input layer

Hidden layer 1
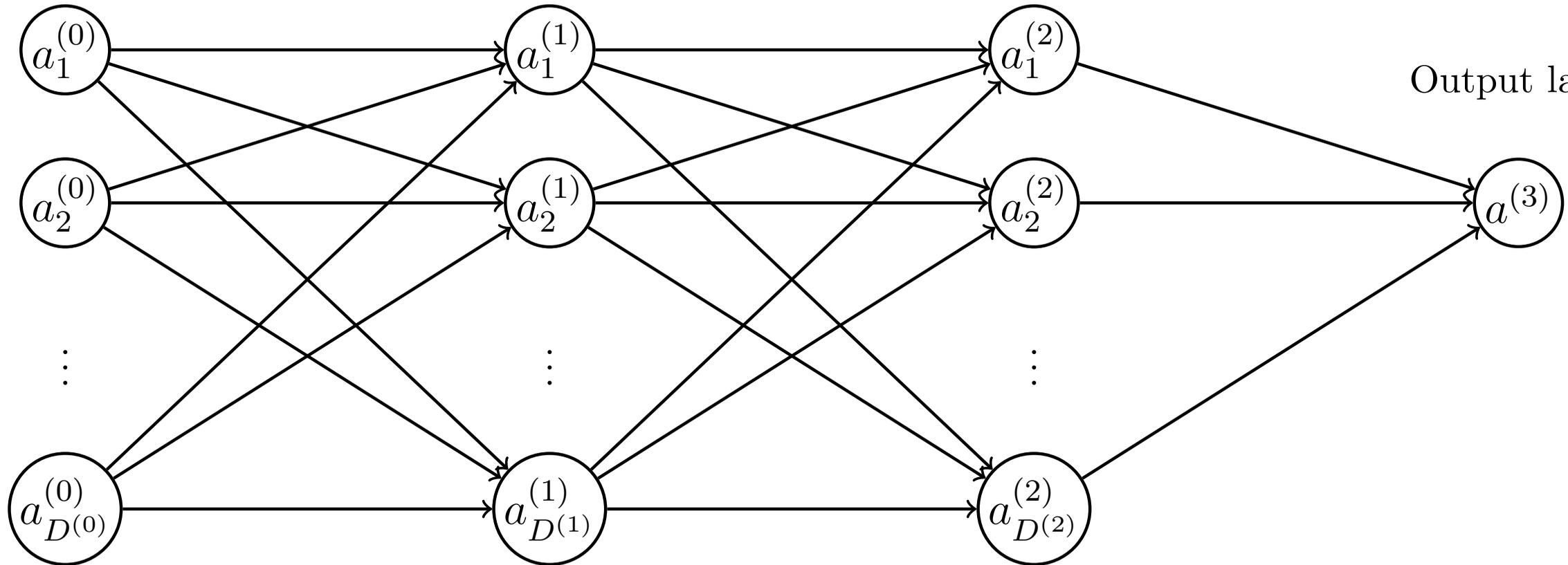
Hidden layer 2



Output layer

UNIVERSITY OF TORONTO

# Notation and Definitions

| | | | |
|---|---|---|---|
| $L$ | | $m$ | |
| | | $D^{(m)}$ | |
| $\mathbf{a}^{(0)}$ | | $\mathbf{W}^{(m)}$ | |
| | | $\mathbf{b}^{(m)}$ | |
| $t$ | | $\sigma^{(m)}$ | |
| $C(\mathbf{a}^{(L)}, y)$ | | $\mathbf{z}^{(m)}$ | |
| | | $\mathbf{a}^{(m)}$ | |

# Notation and Definitions

*[handwritten: $W$ $\vec{x}$ — output × input, input ×1]*

| | | | |
|---|---|---|---|
| $L$ | Number of layers in the model | $m$ | Index of layer, $m \in [1, L]$ |
| | | $D^{(m)}$ | Dimensionality of layer $m$  *[handwritten: # of units in the $m^{th}$ layer]* |
| $\mathbf{a}^{(0)}$ | Input to the model ($\mathbf{x}$) | $\mathbf{W}^{(m)} \in \mathbb{R}^{D^{(m)} \times D^{(m-1)}}$  *[handwritten: output dim, input dim, going to layer m]* | Weight matrix for layer $m$ |
| | | $\mathbf{b}^{(m)} \in \mathbb{R}^{D^{(m)}}$ | Bias vector for layer $m$ |
| $t$ | Target output | $\sigma^{(m)}$ | Non-linearity for layer $m$ |
| $C(\mathbf{a}^{(L)}, y)$ | Cost (or loss) function | $\mathbf{z}^{(m)} \in \mathbb{R}^{D^{(m)}}$ | Pre-activations for layer $m$ |
| | | $\mathbf{a}^{(m)} \in \mathbb{R}^{D^{(m)}}$ | Activations for layer $m$ |

UNIVERSITY OF TORONTO

# Computations in the 3-Layer Neural Network

$$z_i^{(2)} = \sum_{j=1}^{D^{(1)}} W_{ij}^{(2)} a_j^{(1)} + b_i^{(2)}, i = 1, \ldots, D^{(2)}$$

$$a_i^{(2)} = \sigma^{(2)}\left(z_i^{(2)}\right), i = 1, \ldots, D^{(2)}$$

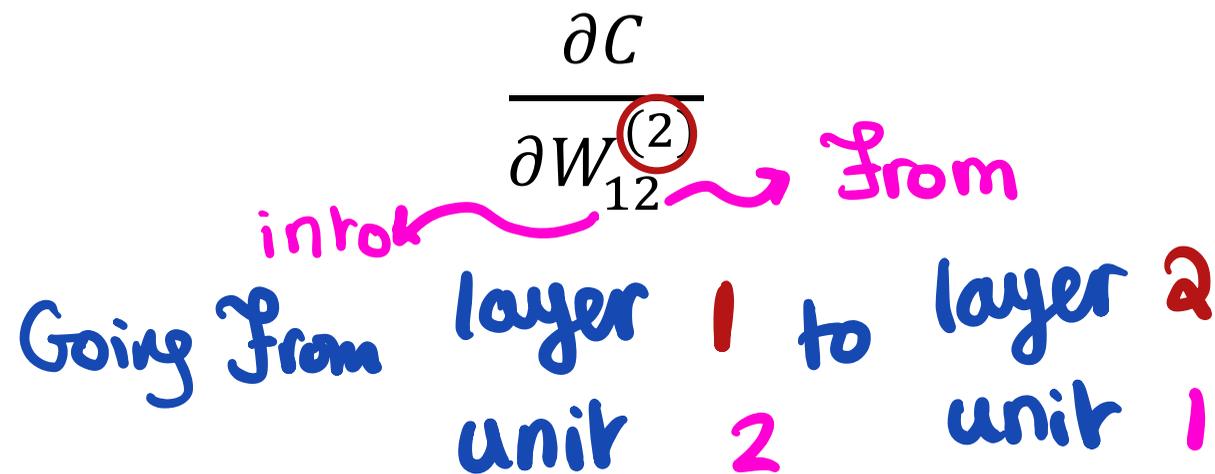$$z^{(3)} = \sum_{j=1}^{D^{(2)}} w_j^{(3)} a_j^{(2)} + b^{(3)}$$

$$a^{(3)} = \sigma^{(3)}\left(z^{(3)}\right)$$

Input layer      Hidden layer 1      Hidden layer 2



Output layer

$\text{Cost: } C = C\left(a^{(3)}, t\right)$

UNIVERSITY OF
TORONTO

# Two Exercises of Computing Derivatives

Compute

$$\frac{\partial C}{\partial W_{12}^{(2)}}$$

intro → From

Going From layer 1 to layer 2
unit 2        unit 1

Compute

$$\frac{\partial C}{\partial W_{23}^{(1)}}$$

$a^{(0)}$ input

Going From layer 0 to layer 1
unit 3        unit 2

UNIVERSITY OF
TORONTO

# Computation Graph

$\omega_{12}^{(2)}$

$\omega_{23}^{(1)}$

we didn't add the weights and biases to the computation graph since it is already busy.

Input layer    Hidden layer 1    Hidden layer 2

$\omega_{12}^{(2)}$



Output layer

but they are usually part of the graph.
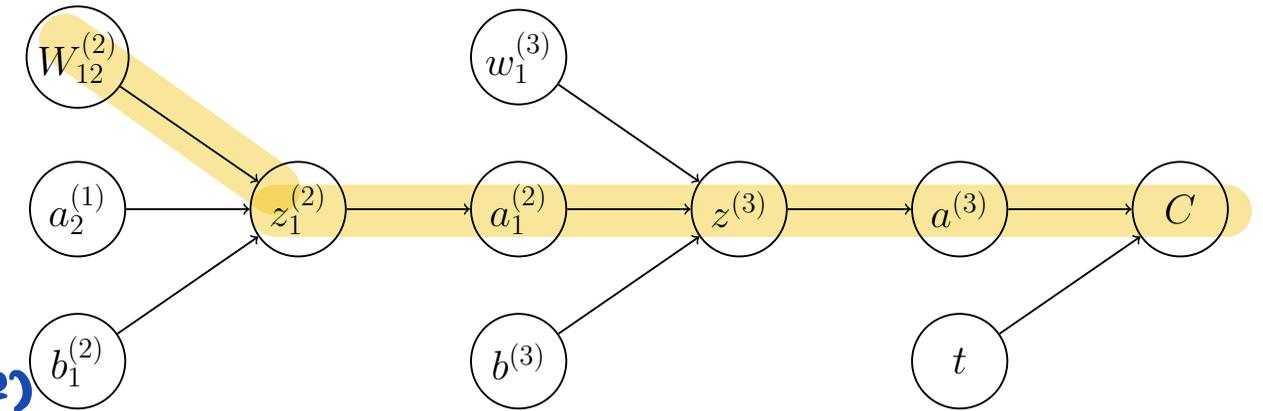
UNIVERSITY OF TORONTO

# Review: Univariate Chain Rule

$$y = f(x) \qquad z = g(y) = g(f(x))$$

$$\frac{dz}{dx} = \frac{dz}{dy} * \frac{dy}{dx}$$

$$x \longrightarrow y \longrightarrow z$$

A. Gao and M. Tawfik, CSC311, Introduction to Machine Learning.

# Loss Derivative with respect to $W_{12}^{(2)}$

$$\frac{\partial C}{\partial W_{12}^{(2)}}$$



$$= \frac{\partial C}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial a_1^{(2)}} \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \frac{\partial z_1^{(2)}}{\partial w_{12}^{(2)}}$$

$$=$$

$$z_i^{(2)} = \sum_{j=1}^{D^{(1)}} W_{ij}^{(2)} a_j^{(1)} + b_i^{(2)}, i = 1, \ldots, D^{(2)} \quad z^{(3)} = \sum_{j=1}^{D^{(2)}} w_j^{(3)} a_j^{(2)} + b^{(3)}$$
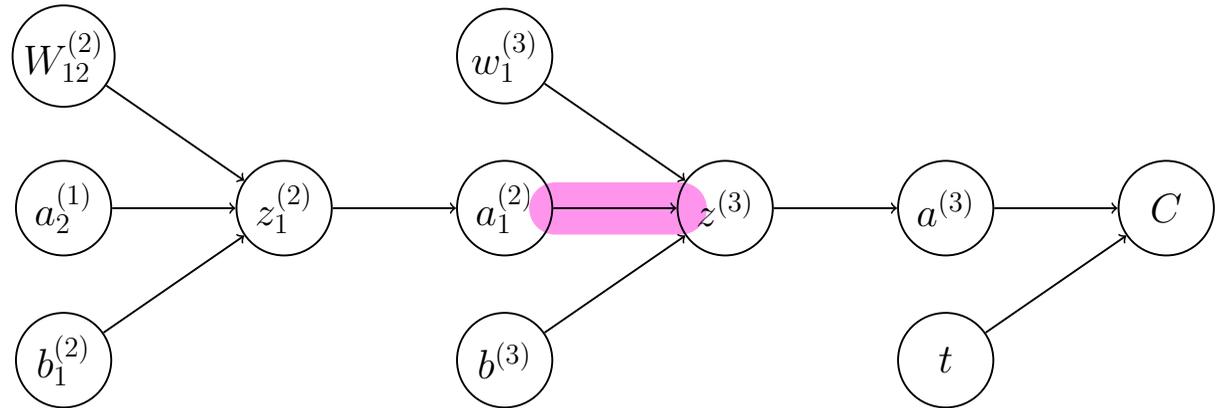
$$a_i^{(2)} = \sigma^{(2)}\left(z_i^{(2)}\right), i = 1, \ldots, D^{(2)}$$

$$a^{(3)} = \sigma^{(3)}\left(z^{(3)}\right)$$

$$C = C(a^{(3)}, t)$$

UNIVERSITY OF TORONTO

# Solution: Loss Derivative with respect to $W_{12}^{(2)}$

$$\frac{\partial C}{\partial W_{12}^{(2)}} = \frac{\partial C}{\partial a^{(3)}} \cdot \frac{\partial a^{(3)}}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial W_{12}^{(2)}}$$

$$= C' \cdot \sigma^{(3)'}\left(z^{(3)}\right) \cdot \left(w_1^{(3)}\right) \cdot \sigma^{(2)'}\left(z_1^{(2)}\right) \cdot a_2^{(1)}$$



$$z_i^{(2)} = \sum_{j=1}^{D^{(1)}} W_{ij}^{(2)} a_j^{(1)} + b_i^{(2)}, i = 1, \ldots, D^{(2)}$$

$$a_i^{(2)} = \sigma^{(2)}\left(z_i^{(2)}\right), i = 1, \ldots, D^{(2)}$$

$$z^{(3)} = \sum_{j=1}^{D^{(2)}} w_j^{(3)} a_j^{(2)} + b^{(3)}$$

$$a^{(3)} = \sigma^{(3)}\left(z^{(3)}\right)$$

$$C = C(a^{(3)}, t)$$

UNIVERSITY OF TORONTO

# Computing Loss Derivative with respect to $W_{23}^{(1)}$

$$\frac{\partial C}{\partial \omega_{23}^{(1)}}$$

# Review: Multi-Variate Chain Rule

$$y_1 = f_1(x), \qquad y_2 = f_2(x), \qquad C = g(y_1, y_2)$$



$$\frac{dC}{dx} = \frac{\partial C}{\partial y_1} * \frac{dy_1}{dx} + \frac{\partial C}{\partial y_2} * \frac{dy_2}{dx}$$

UNIVERSITY OF TORONTO

# Computing the Loss

layer 1

layer 2

layer 3

$$z_i^{(1)} = \sum_{j=1}^{D^{(0)}} W_{ij}^{(1)} a_j^{(0)} + b_i^{(1)},$$
$$i = 1, \dots, D^{(1)}$$

$$z_i^{(2)} = \sum_{j=1}^{D^{(1)}} W_{ij}^{(2)} a_j^{(1)} + b_i^{(2)},$$
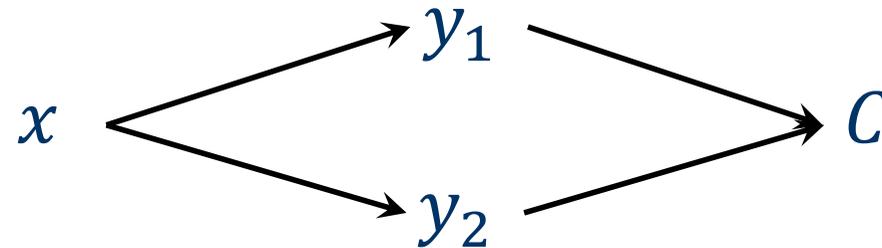$$i = 1, \dots, D^{(2)}$$

$$z^{(3)} = \sum_{j=1}^{D^{(2)}} W_j^{(3)} a_j^{(2)} + b^{(3)}$$

$$a^{(3)} = \sigma^{(3)}\left(z^{(3)}\right)$$

$$a_i^{(1)} = \sigma^{(1)}\left(z_i^{(1)}\right),$$
$$i = 1, \dots, D^{(1)}$$

$$a_i^{(2)} = \sigma^{(2)}\left(z_i^{(2)}\right),$$
$$i = 1, \dots, D^{(2)}$$

$$C = C(a^{(3)}, t)$$

UNIVERSITY OF TORONTO

# Loss Derivative with respect to $W_{23}^{(1)}$

$$\frac{\partial C}{\partial z_2^{(1)}} = \left( \sum_{i=1}^{D^{(2)}} \frac{\partial C}{\partial z_i^{(2)}} \frac{\partial z_i^{(2)}}{\partial a_2^{(1)}} \right) \frac{\partial a_2^{(1)}}{\partial z_2^{(1)}}$$

$$\frac{\partial C}{\partial z^{(3)}} = \frac{\partial C}{\partial a^{(3)}} \cdot \frac{\partial a^{(3)}}{\partial z^{(3)}}$$

$$\frac{\partial C}{\partial W_{23}^{(1)}} = \frac{\partial C}{\partial z_2^{(1)}} \frac{\partial z_2^{(1)}}{\partial w_{23}^{(1)}}$$

$$\frac{\partial C}{\partial z_i^{(2)}} = \frac{\partial C}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a_i^{(2)}} \cdot \frac{\partial a_i^{(2)}}{\partial z_i^{(2)}}, i = 1, \dots D^{(2)}$$
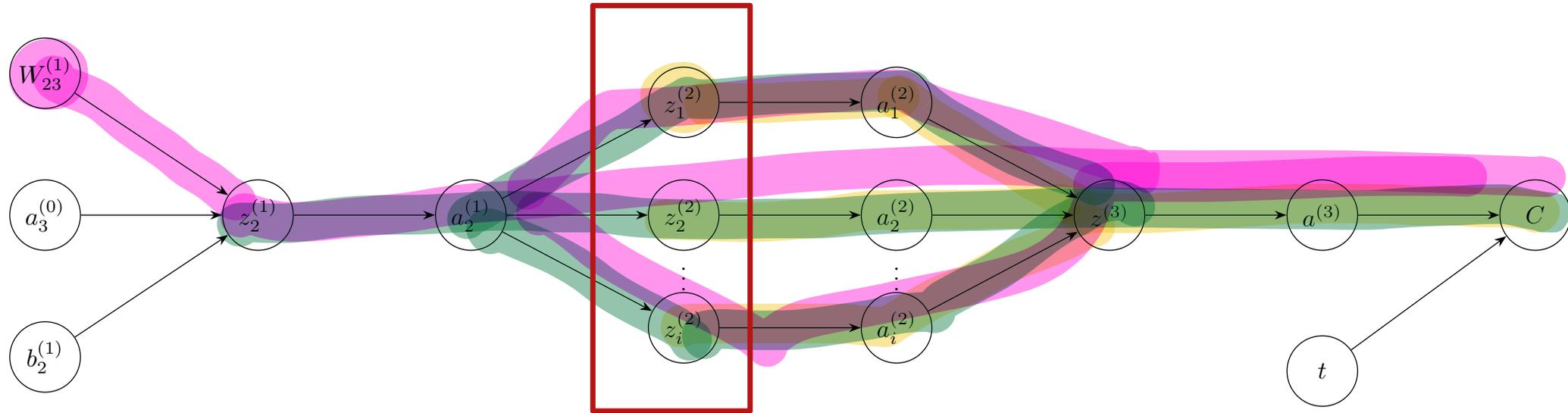
UNIVERSITY OF TORONTO

# Solution: Loss Derivative with respect to $W_{23}^{(1)}$

$$\frac{\partial C}{\partial z_2^{(1)}} = \left( \sum_{i=1}^{D^{(2)}} \frac{\partial C}{\partial z_i^{(2)}} \cdot \frac{\partial z_i^{(2)}}{\partial a_2^{(1)}} \right) \frac{\partial a_2^{(1)}}{\partial z_2^{(1)}}$$

$$\frac{\partial C}{\partial W_{23}^{(1)}} = \frac{\partial C}{\partial z_2^{(1)}} \cdot \frac{\partial z_2^{(1)}}{\partial W_{23}^{(1)}}$$

$$\frac{\partial C}{\partial z^{(3)}} = \frac{\partial C}{\partial a^{(3)}} \cdot \frac{\partial a^{(3)}}{\partial z^{(3)}}$$

$$\frac{\partial C}{\partial z_i^{(2)}} = \frac{\partial C}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a_i^{(2)}} \cdot \frac{\partial a_i^{(2)}}{\partial z_i^{(2)}}, i = 1, \dots D^{(2)}$$

UNIVERSITY OF
TORONTO

# Computing Gradients Efficiently

# Notice Any Inefficiency in These Gradient Computations?

$$\frac{\partial C}{\partial z^{(3)}} = \frac{\partial C}{\partial a^{(3)}} \cdot \frac{\partial a^{(3)}}{\partial z^{(3)}}$$

$$\frac{\partial C}{\partial z_1^{(2)}} = \frac{\partial C}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}}$$

$$\frac{\partial C}{\partial W_{12}^{(2)}} = \frac{\partial C}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial W_{12}^{(2)}}$$

$$\frac{\partial C}{\partial z^{(3)}} = \frac{\partial C}{\partial a^{(3)}} \cdot \frac{\partial a^{(3)}}{\partial z^{(3)}}$$

$$\frac{\partial C}{\partial z_i^{(2)}} = \frac{\partial C}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a_i^{(2)}} \cdot \frac{\partial a_i^{(2)}}{\partial z_i^{(2)}}, i = 1, \dots D^{(2)}$$
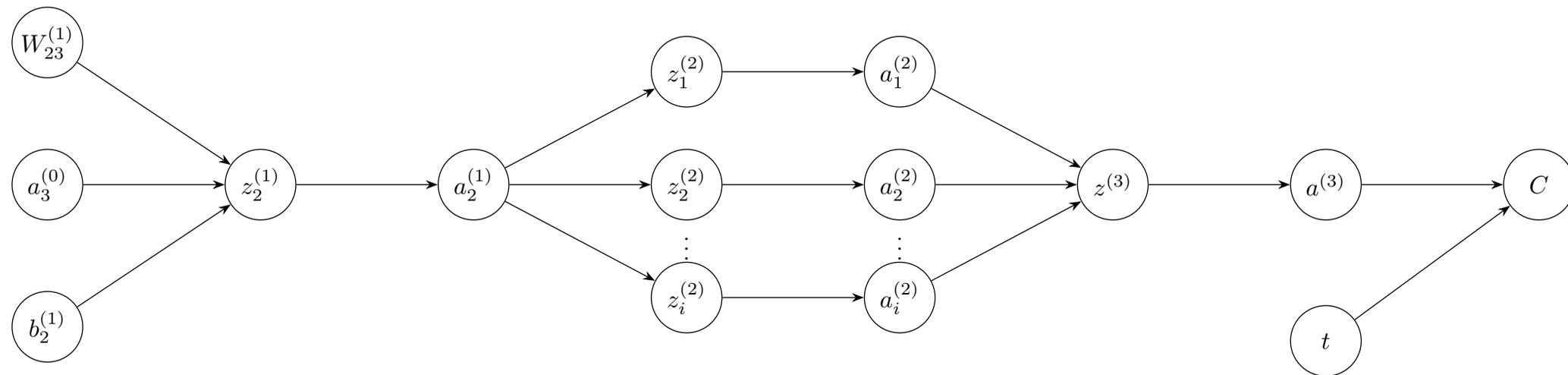
$$\frac{\partial C}{\partial z_2^{(1)}} = \left( \sum_{i=1}^{D^{(2)}} \frac{\partial C}{\partial z_i^{(2)}} \cdot \frac{\partial z_i^{(2)}}{\partial a_5^{(1)}} \right) \frac{\partial a_2^{(1)}}{\partial z_2^{(1)}}$$

$$\frac{\partial C}{\partial W_{23}^{(1)}} = \frac{\partial C}{\partial z_2^{(1)}} \cdot \frac{\partial z_2^{(1)}}{\partial W_{23}^{(1)}}$$

*highlighted computation are repeated between the 2 derivatives.*

A. Gao and M. Tawhik, CSC311, Introduction to Machine Learning.  24

# Reuse Computations for All Gradients

*This is a vectorized computation graph.*

Many repeated computations to compute each derivative

Compute all derivatives in one pass by reusing intermediate results (dynamic programming).



*base case*

$$\mathbf{a}^{(0)} \longrightarrow \mathbf{z}^{(1)} \longrightarrow \mathbf{a}^{(1)} \longrightarrow \mathbf{z}^{(2)} \longrightarrow \mathbf{a}^{(2)} \longrightarrow z^{(3)} \longrightarrow a^{(3)} \longrightarrow C$$

$\mathbf{W}^{(1)} \quad \mathbf{b}^{(1)} \quad \mathbf{W}^{(2)} \quad \mathbf{b}^{(2)} \quad \mathbf{w}^{(3)} \quad b^{(3)} \quad \mathbf{t}$

UNIVERSITY OF TORONTO

# Backpropagation for a 3-Layer Neural Network

# Backpropagation for 3-Layer Network

Step 1: Gradients for the output layer

$$\frac{\partial C}{\partial z^{(3)}} = ? \qquad \frac{\partial C}{\partial w^{(3)}} = ? \qquad \frac{\partial C}{\partial b^{(3)}} = ?$$

Step 2: Gradients for hidden layer 2

$$\frac{\partial C}{\partial z_i^{(2)}} = ? \qquad \frac{\partial C}{\partial W_{ij}^{(2)}} = ? \qquad \frac{\partial C}{\partial b_i^{(2)}} = ?$$

Step 3: Gradients for hidden layer 1

$$\frac{\partial C}{\partial z_j^{(1)}} = ? \qquad \frac{\partial C}{\partial W_{ij}^{(1)}} = ? \qquad \frac{\partial C}{\partial b_i^{(1)}} = ?$$

$C$

$\uparrow$

$a^{(3)}$

$\uparrow$

$\mathbf{w}^{(3)} \longrightarrow z^{(3)} \longleftarrow b^{(3)}$

$\uparrow$

$\mathbf{a}^{(2)}$

$\uparrow$

$\mathbf{W}^{(2)} \longrightarrow \mathbf{z}^{(2)} \longleftarrow \mathbf{b}^{(2)}$

$\uparrow$

$\mathbf{a}^{(1)}$

$\uparrow$

$\mathbf{W}^{(1)} \longrightarrow \mathbf{z}^{(1)} \longleftarrow \mathbf{b}^{(1)}$

$\uparrow$

$\mathbf{a}^{(0)}$

UNIVERSITY OF
TORONTO

output is a scalar

# Backprop Step 1: Gradients for Output Layer

For pre-activation $z^{(3)}$:

$$\frac{\partial C}{\partial z^{(3)}} = \frac{\partial C}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial z^{(3)}}$$

depends on the activation func

For weight:

$$\frac{\partial C}{\partial w_i^{(3)}} = \frac{\partial C}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial w_i^{(3)}} = \frac{\partial C}{\partial z^{(3)}} a_i^{(2)}$$

For bias:

$$\frac{\partial C}{\partial b^{(3)}} = \frac{\partial C}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial b^{(3)}} = \frac{\partial C}{\partial z^{(3)}} \cdot 1$$

$$z^{(3)} = w_1^{(3)} a_1^{(2)} + w_2^{(3)} a_2^{(2)} + \dots + b^{(3)}$$

$C$, $a^{(3)}$ — output is scalar; $\mathbf{w}^{(3)}$ (vector) $\to z^{(3)} \leftarrow b^{(3)}$; $\mathbf{a}^{(2)}$; $\mathbf{W}^{(2)} \to \mathbf{z}^{(2)} \leftarrow \mathbf{b}^{(2)}$; $\mathbf{a}^{(1)}$; $\mathbf{W}^{(1)} \to \mathbf{z}^{(1)} \leftarrow \mathbf{b}^{(1)}$; $\mathbf{a}^{(0)}$

# Backprop Step 1: Gradients for Output Layer

For pre-activation $z^{(3)}$:

$$\frac{\partial C}{\partial z^{(3)}} = \frac{\partial C}{\partial a^{(3)}} \cdot \frac{\partial a^{(3)}}{\partial z^{(3)}} = \frac{\partial C}{\partial a^{(3)}} \cdot \sigma^{(3)'}\left(z^{(3)}\right)$$
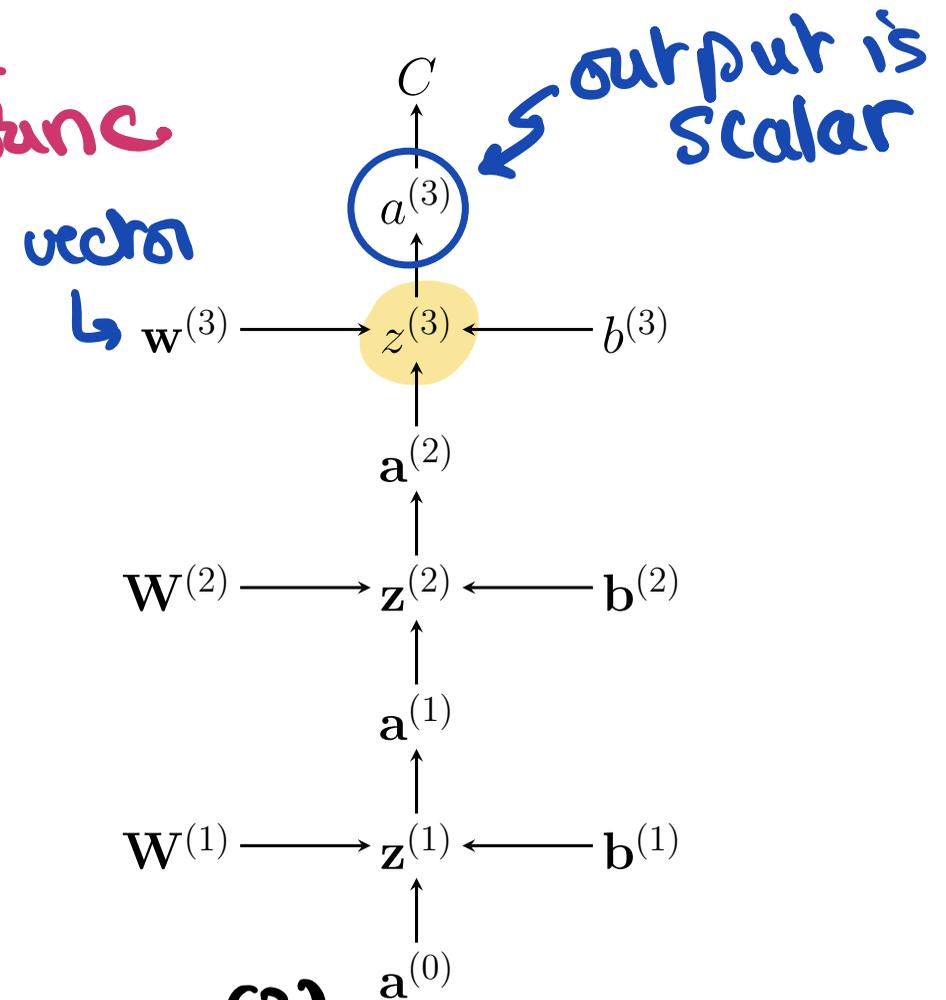
For weight:

$$\frac{\partial C}{\partial w_i^{(3)}} = \frac{\partial C}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial w_i^{(3)}} = \frac{\partial C}{\partial z^{(3)}} \cdot a_i^{(2)}$$
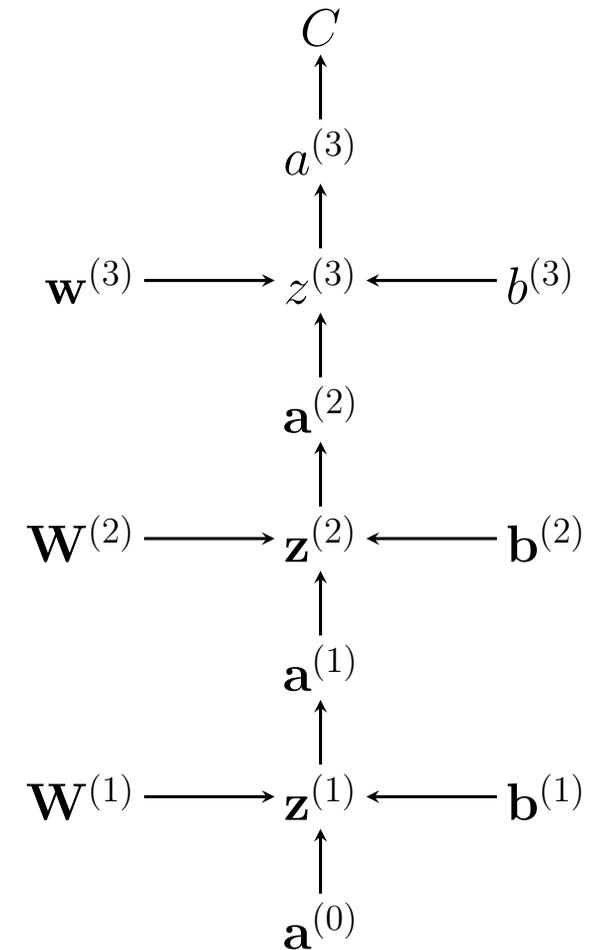
For bias:

$$\frac{\partial C}{\partial b^{(3)}} = \frac{\partial C}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial b^{(3)}} = \frac{\partial C}{\partial z^{(3)}}$$

$C$

$a^{(3)}$

$\mathbf{w}^{(3)} \longrightarrow z^{(3)} \longleftarrow b^{(3)}$

$\mathbf{a}^{(2)}$

$\mathbf{W}^{(2)} \longrightarrow \mathbf{z}^{(2)} \longleftarrow \mathbf{b}^{(2)}$

$\mathbf{a}^{(1)}$

$\mathbf{W}^{(1)} \longrightarrow \mathbf{z}^{(1)} \longleftarrow \mathbf{b}^{(1)}$

$\mathbf{a}^{(0)}$

UNIVERSITY OF TORONTO

# Backprop Step 2: Gradients for Hidden Layer 2

For pre-activations $z_i^{(2)}$:
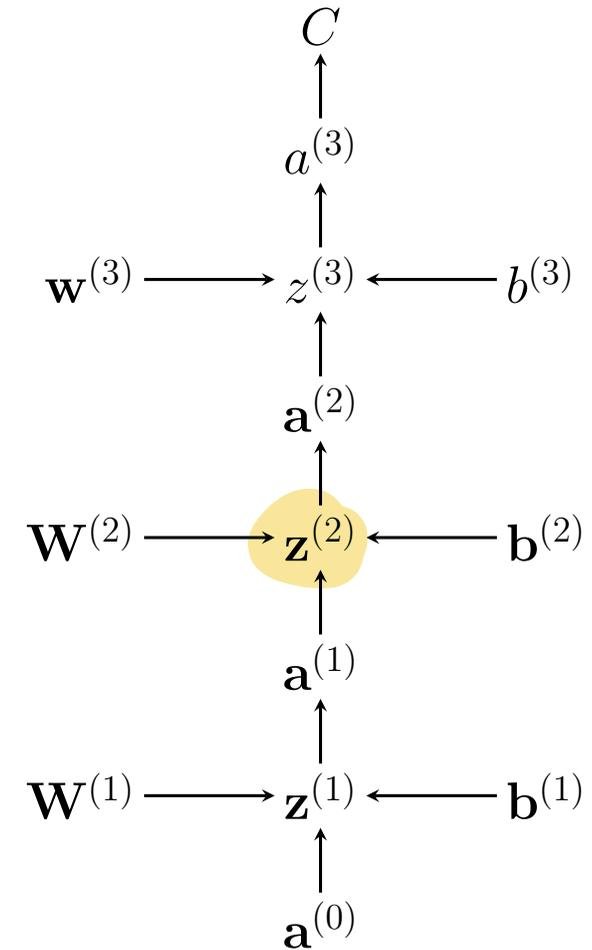
$$\frac{\partial C}{\partial z_i^{(2)}} =$$

For weight:

$$\frac{\partial C}{\partial W_{ij}^{(2)}} =$$

For bias:

$$\frac{\partial C}{\partial b_i^{(2)}} =$$

$C$

$a^{(3)}$

$\mathbf{w}^{(3)} \longrightarrow z^{(3)} \longleftarrow b^{(3)}$

$\mathbf{a}^{(2)}$

$\mathbf{W}^{(2)} \longrightarrow \mathbf{z}^{(2)} \longleftarrow \mathbf{b}^{(2)}$

$\mathbf{a}^{(1)}$

$\mathbf{W}^{(1)} \longrightarrow \mathbf{z}^{(1)} \longleftarrow \mathbf{b}^{(1)}$

$\mathbf{a}^{(0)}$

UNIVERSITY OF
TORONTO

# Backprop Step 2: Gradients for Hidden Layer 2

$$z^{(3)} = \omega_1^{(3)} a_1^{(2)} + \omega_2^{(3)} a_2^{(2)} + \ldots + b^{(3)}$$
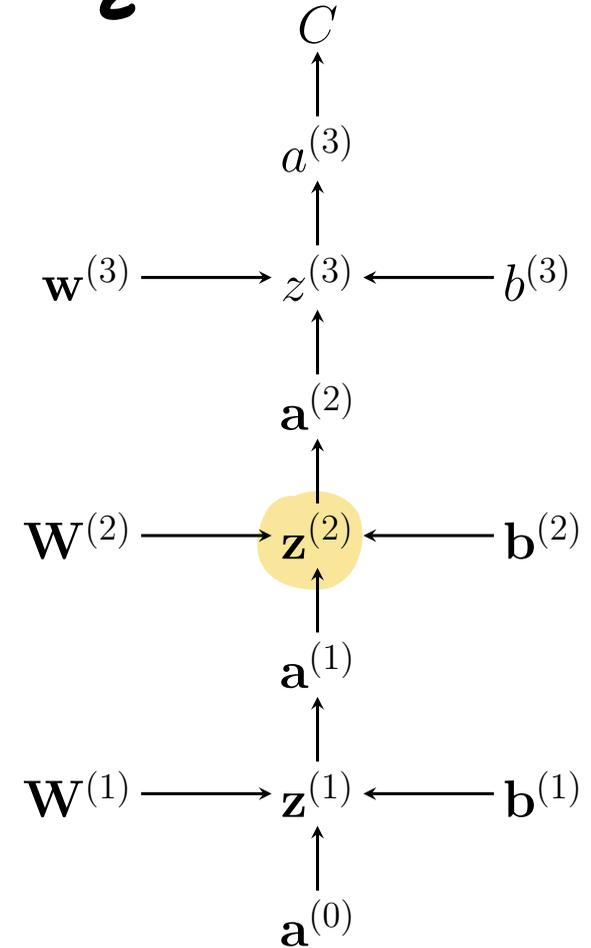
For pre-activations $z_i^{(2)}$:

$$\frac{\partial C}{\partial z_i^{(2)}} = \frac{\partial C}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a_i^{(2)}} \cdot \frac{\partial a_i^{(2)}}{\partial z_i^{(2)}} = \frac{\partial C}{\partial z^{(3)}} \cdot w_i^{(3)} \cdot \sigma^{(2)'}\left(z_i^{(2)}\right)$$

For weight:

$$\frac{\partial C}{\partial W_{ij}^{(2)}} = \frac{\partial C}{\partial z_i^{(2)}} \cdot \frac{\partial z_i^{(2)}}{\partial W_{ij}^{(2)}} = \frac{\partial C}{\partial z_i^{(2)}} \cdot a_j^{(1)}$$

For bias:

$$\frac{\partial C}{\partial b_i^{(2)}} = \frac{\partial C}{\partial z_i^{(2)}} \cdot \frac{\partial z_i^{(2)}}{\partial b_i^{(2)}} = \frac{\partial C}{\partial z_i^{(2)}} \quad 1$$

$$z_i^{(2)} = \omega_{i1}^{(2)} a_1^{(1)} + \omega_{i2}^{(2)} a_2^{(1)} + \ldots + b_i^{(2)}$$

$C$

$a^{(3)}$

$\mathbf{w}^{(3)} \longrightarrow z^{(3)} \longleftarrow b^{(3)}$

$\mathbf{a}^{(2)}$

$\mathbf{W}^{(2)} \longrightarrow \mathbf{z}^{(2)} \longleftarrow \mathbf{b}^{(2)}$

$\mathbf{a}^{(1)}$

$\mathbf{W}^{(1)} \longrightarrow \mathbf{z}^{(1)} \longleftarrow \mathbf{b}^{(1)}$

$\mathbf{a}^{(0)}$

UNIVERSITY OF TORONTO

# Backprop Step 3: Gradients for Hidden Layer 1

For pre-activations $z_j^{(1)}$:
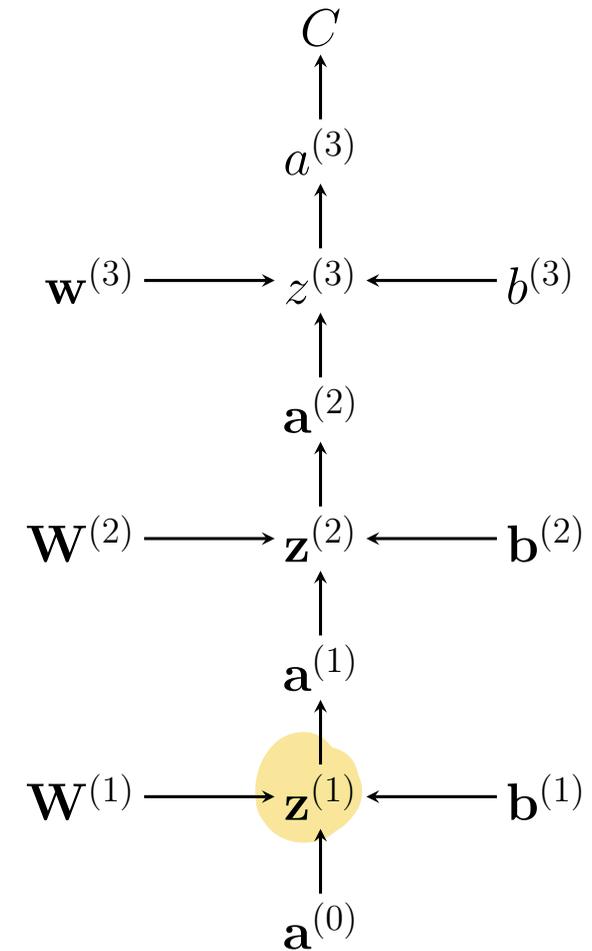
$$\frac{\partial C}{\partial z_j^{(1)}} =$$

For weight:

$$\frac{\partial C}{\partial W_{ij}^{(1)}} =$$

For bias:

$$\frac{\partial C}{\partial b_i^{(1)}} =$$

$C$

$a^{(3)}$

$\mathbf{w}^{(3)} \longrightarrow z^{(3)} \longleftarrow b^{(3)}$

$\mathbf{a}^{(2)}$

$\mathbf{W}^{(2)} \longrightarrow \mathbf{z}^{(2)} \longleftarrow \mathbf{b}^{(2)}$

$\mathbf{a}^{(1)}$

$\mathbf{W}^{(1)} \longrightarrow \mathbf{z}^{(1)} \longleftarrow \mathbf{b}^{(1)}$

$\mathbf{a}^{(0)}$

UNIVERSITY OF TORONTO

# Backprop Step 3: Gradients for Hidden Layer 1

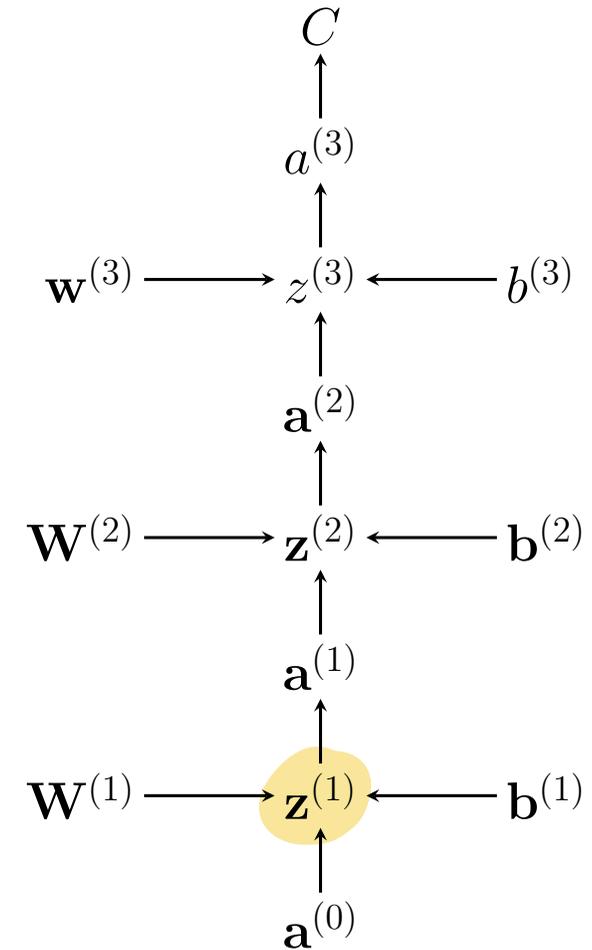For pre-activations $z_j^{(1)}$:

$$\frac{\partial C}{\partial z_j^{(1)}} = \left( \sum_{i=1}^{D^{(2)}} \frac{\partial C}{\partial z_i^{(2)}} \cdot \frac{\partial z_i^{(2)}}{\partial a_j^{(1)}} \right) \cdot \frac{\partial a_j^{(1)}}{\partial z_j^{(1)}} = \left( \sum_{i=1}^{D^{(2)}} \frac{\partial C}{\partial z_i^{(2)}} \cdot W_{ij}^{(2)} \right) \cdot \sigma^{(1)'} \left( z_j^{(1)} \right)$$

For weight:

$$\frac{\partial C}{\partial W_{ij}^{(1)}} = \frac{\partial C}{\partial z_i^{(1)}} \cdot \frac{\partial z_i^{(1)}}{\partial W_{ij}^{(1)}} = \frac{\partial C}{\partial z_i^{(1)}} \cdot a_j^{(0)}$$

For bias:

$$\frac{\partial C}{\partial b_i^{(1)}} = \frac{\partial C}{\partial z_i^{(1)}} \cdot \frac{\partial z_i^{(1)}}{\partial b_i^{(1)}} = \frac{\partial C}{\partial z_i^{(1)}}$$

$C$

$a^{(3)}$

$\mathbf{w}^{(3)} \longrightarrow z^{(3)} \longleftarrow b^{(3)}$

$\mathbf{a}^{(2)}$

$\mathbf{W}^{(2)} \longrightarrow \mathbf{z}^{(2)} \longleftarrow \mathbf{b}^{(2)}$

$\mathbf{a}^{(1)}$

$\mathbf{W}^{(1)} \longrightarrow \mathbf{z}^{(1)} \longleftarrow \mathbf{b}^{(1)}$

$\mathbf{a}^{(0)}$

$$z_i^{(1)} = \omega_{i1}^{(1)} a_1^{(0)} + \omega_{i2}^{(1)} a_2^{(0)} + \cdots + b_i^{(1)}$$

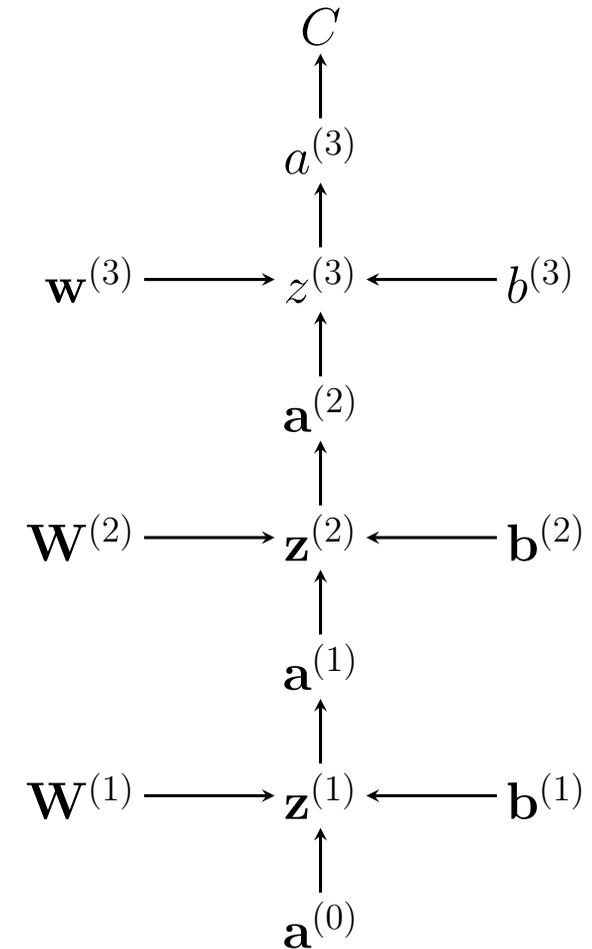# Backpropagation for 3-Layer Network

Step1: Gradients for output layer

$$\frac{\partial C}{\partial z^{(3)}} = \frac{\partial C}{\partial a^{(3)}} \cdot \frac{\partial a^{(3)}}{\partial z^{(3)}} = \frac{\partial C}{\partial a^{(3)}} \cdot \sigma^{(3)'}\left(z^{(3)}\right)$$

Step 2: Gradients for hidden layer 2

$$\frac{\partial C}{\partial z_i^{(2)}} = \frac{\partial C}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a_i^{(2)}} \cdot \frac{\partial a_i^{(2)}}{\partial z_i^{(2)}} = \frac{\partial C}{\partial z^{(3)}} \cdot w_i^{(3)} \cdot \sigma^{(2)'}\left(z_i^{(2)}\right)$$

Step 3: Gradients for hidden layer 1

$$\frac{\partial C}{\partial z_j^{(1)}} = \left(\sum_{i=1}^{D^{(2)}} \frac{\partial C}{\partial z_i^{(2)}} \cdot \frac{\partial z_i^{(2)}}{\partial a_j^{(1)}}\right) \cdot \frac{\partial a_j^{(1)}}{\partial z_j^{(1)}} = \left(\sum_{i=1}^{D^{(2)}} \frac{\partial C}{\partial z_i^{(2)}} \cdot W_{ij}^{(2)}\right) \cdot \sigma^{(1)'}\left(z_j^{(1)}\right)$$

$$
\begin{array}{c}
C \\
\uparrow \\
a^{(3)} \\
\uparrow \\
\mathbf{w}^{(3)} \longrightarrow z^{(3)} \longleftarrow b^{(3)} \\
\uparrow \\
\mathbf{a}^{(2)} \\
\uparrow \\
\mathbf{W}^{(2)} \longrightarrow \mathbf{z}^{(2)} \longleftarrow \mathbf{b}^{(2)} \\
\uparrow \\
\mathbf{a}^{(1)} \\
\uparrow \\
\mathbf{W}^{(1)} \longrightarrow \mathbf{z}^{(1)} \longleftarrow \mathbf{b}^{(1)} \\
\uparrow \\
\mathbf{a}^{(0)}
\end{array}
$$

UNIVERSITY OF TORONTO

# The Backpropagation Algorithm (Almost Complete Version)

1. Compute gradients for output layer

$$\frac{\partial C}{\partial z^{(L)}} = \frac{\partial C}{\partial a^{(L)}} \cdot \frac{\partial a^{(L)}}{\partial z^{(L)}} = \frac{\partial C}{\partial a^{(L)}} \cdot \sigma^{(L)\prime}\left(z^{(L)}\right) \quad \right\} \text{ base case}$$

2. Compute gradients for each hidden layer recursively

$$\frac{\partial C}{\partial z_j^{(m)}} = \left(\sum_{i=1}^{D^{(m+1)}} \frac{\partial C}{\partial z_i^{(m+1)}} \cdot \frac{\partial z_i^{(m+1)}}{\partial a_j^{(m)}}\right) \cdot \frac{\partial a_j^{(m)}}{\partial z_j^{(m)}} = \left(\sum_{i=1}^{D^{(m+1)}} \frac{\partial C}{\partial z_i^{(m+1)}} \cdot W_{ij}^{(m+1)}\right) \cdot \sigma^{(m)\prime}\left(z_j^{(m)}\right) \right\} \text{ recursive case}$$

$$j = 1, \dots, D^{(m)}$$

3. Compute gradients for the weights

This is the backward pass

$$\frac{\partial C}{\partial W_{ij}^{(m)}} = \frac{\partial C}{\partial z_i^{(m)}} \cdot \frac{\partial z_i^{(m)}}{\partial W_{ij}^{(m)}} = \frac{\partial C}{\partial z_i^{(m)}} \cdot a_j^{(m-1)}, i = 1, \dots, D^{(m)}, j = 1, \dots, D^{(m-1)}$$

UNIVERSITY OF TORONTO

# Why The Forward Pass

# Necessary Quantities for Backpropagation

What quantities do we need to carry out backpropagation?

- derivative of loss function: $\frac{\partial C}{\partial a^{(L)}}$

- derivative of non-linearity: $\sigma^{(m)'}\left(z_j^{(m)}\right)$

- weights: $W_{ij}^{(m)}$

- activations: $a_j^{(m)}$

output layer:

$$\frac{\partial C}{\partial z^{(L)}} = \frac{\partial C}{\partial a^{(L)}} \cdot \sigma^{(L)'}\left(z^{(L)}\right)$$

each hidden layer:

$$\frac{\partial C}{\partial z_j^{(m)}} = \left(\sum_{i=1}^{D^{(m+1)}} \frac{\partial C}{\partial z_i^{(m+1)}} \cdot W_{ij}^{(m+1)}\right) \cdot \sigma^{(m)'}\left(z_j^{(m)}\right)$$

weights:

$$\frac{\partial C}{\partial W_{ij}^{(m)}} = \frac{\partial C}{\partial z_i^{(m)}} \cdot a_j^{(m-1)}$$

# The Backpropagation Algorithm (Truly Complete Version)

## Forward Pass

For each $m = 1, \ldots, L$,

$$z_i^{(m)} = \sum_{j=1}^{D^{(m-1)}} W_{ij}^{(m)} a_j^{(m-1)} + b_i^{(m)}, i = 1, \ldots, D^{(m)}$$

$$a_i^{(m)} = \sigma^{(m)}\left(z_i^{(m)}\right), i = 1, \ldots, D^{(m)}$$

*↳ we need these when calculating the derivatives w.r.t the weights*

For output layer $L$,

$$C = C(a^L, t)$$

## Backward Pass

For output layer $L$

$$\frac{\partial C}{\partial z^{(L)}} = \frac{\partial C}{\partial a^{(L)}} \cdot \sigma^{(L)'}\left(z^{(L)}\right)$$

For each $m = 1, \ldots, L$,

$$\frac{\partial C}{\partial z_j^{(m)}} = \left(\sum_{i=1}^{D^{(m+1)}} \frac{\partial C}{\partial z_i^{(m+1)}} \cdot W_{ij}^{(m+1)}\right) \cdot \sigma^{(m)'}\left(z_j^{(m)}\right),$$

$$j = 1, \ldots, D^{(m)}$$

$$\frac{\partial C}{\partial W_{ij}^{(m)}} = \frac{\partial C}{\partial z_i^{(m)}} \cdot a_j^{(m-1)},$$

$$i = 1, \ldots, D^{(m)}, j = 1, \ldots, D^{(m-1)}$$

# Backpropagation Summary

- Efficiently compute gradients for many weights in a neural network

- The forward pass computes and stores activations.

- The backward pass
  - computes all the derivatives in one pass
  - reuses intermediate values (dynamic programming)
  - computes the derivatives w.r.t. pre-activations recursively via the chain rule.

UNIVERSITY OF
TORONTO