# CSC311 Introduction to Machine Learning

# Vectorizing The Backpropagation Algorithm

Alice Gao and Marina Tawfik

# Learning Outcomes

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

- Derive vectorized expressions for quantities in the forward pass of the backpropagation algorithm for a small neural network.

- Derive vectorized expressions for gradients with respect to activations, pre-activations, and weights in the backward pass of the backpropagation algorithm for a small neural network.
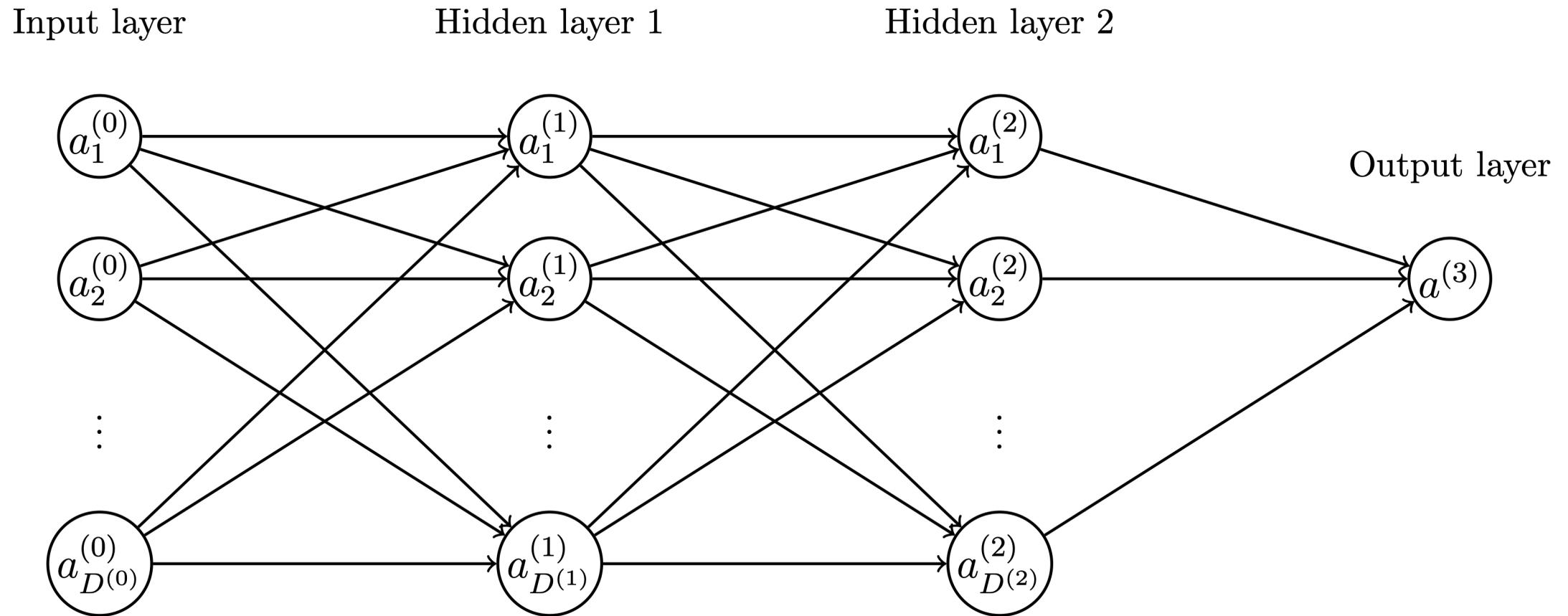
# Outline

- Review of 3-Layer Neural Network

- Vectorizing Forward Pass

- Vectorizing Backward Pass

# Review: 3-Layer Neural Network

# 3-Layer Neural Network



Input layer     Hidden layer 1     Hidden layer 2

Output layer

# Notation and Definitions

$L$: number of layers

$\mathbf{a}^{(0)}$: input to the model ($\mathbf{x}$)

$t$: target output

$C$: loss function $C(a^{(L)}, t)$

For each layer $m$ where $m \in [1, L]$

$D^{(m)}$: dimensionality of layer $m$

$\mathbf{W}^{(m)}$: weight matrix for layer $m$, $\mathbb{R}^{D^{(m)} \times D^{(m-1)}}$

*output* *input*

*weights going into the $m^{th}$ layer*

$\mathbf{b}^{(m)}$: bias vector for layer $m$, $\mathbb{R}^{D^{(m)}}$

$\sigma^{(m)}$: non-linearity for layer $m$

$\mathbf{z}^{(m)}$: pre-activations for layer $m$  $\mathbb{R}^{D^{(m)}}$
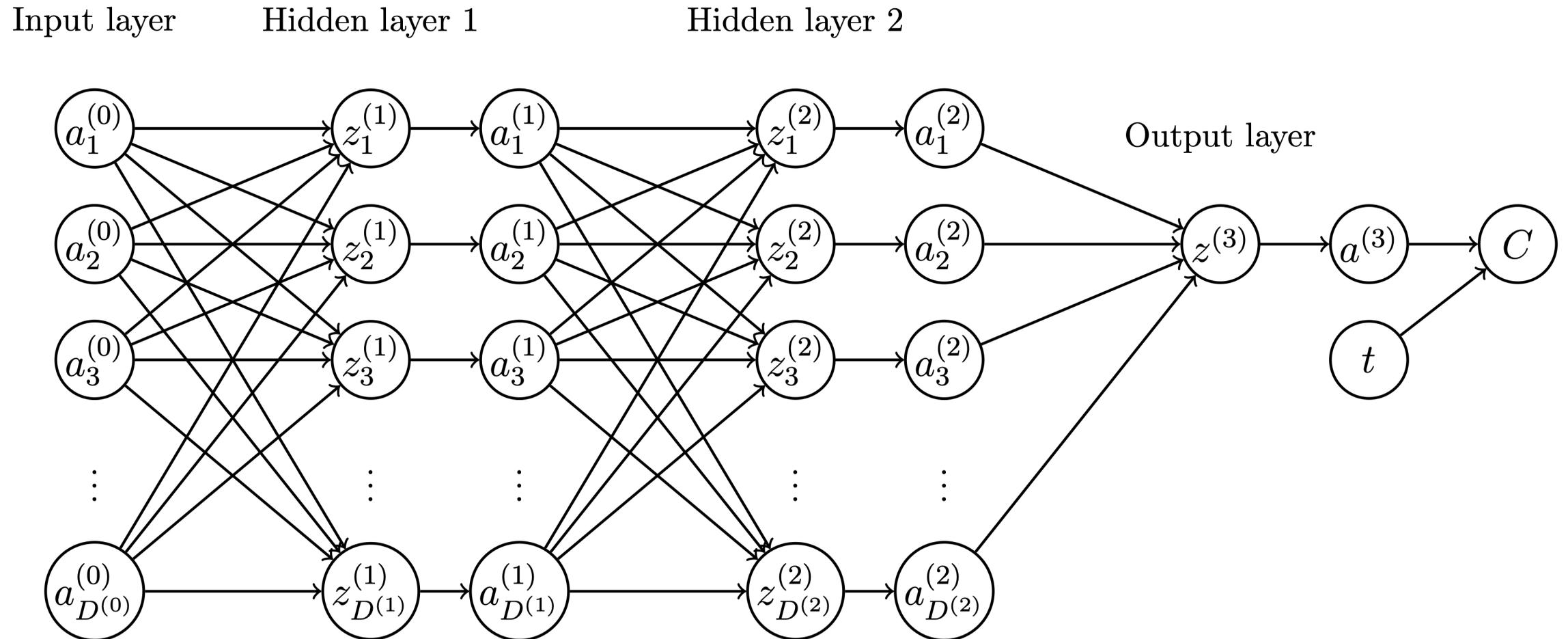
$\mathbf{a}^{(m)}$: activations for layer $m$  $\mathbb{R}^{D^{(m)}}$

UNIVERSITY OF
TORONTO

# Computation Graph

# Vectorizing Forward Pass

# Forward Pass Computations

layer 1:

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

$$a_i^{(1)} = \sigma^{(1)}\left(z_i^{(1)}\right)$$

layer 2:

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

$$a_i^{(2)} = \sigma^{(2)}\left(z_i^{(2)}\right)$$

layer 3:

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

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

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

UNIVERSITY OF
TORONTO

# Vectorizing Forward Pass Computations

## Non-Vectorized

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

$$z_i^{(m+1)} = \left( \sum_j W_{ij}^{(m+1)} a_j^{(m)} \right) + b_i^{(m+1)}$$

$$a_i^{(m+1)} = \sigma^{(m+1)}\left( z_i^{(m+1)} \right)$$

For output layer $L$
$$C = C(a^{(L)}, t)$$

## Vectorized

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

For output layer $L$
$$C = C(a^{(L)}, t)$$

UNIVERSITY OF TORONTO

# Result of Vectorizing Weighted Sum

$$W^{(m+1)} \qquad \vec{a}^{(m)}$$
$$D^{(m+1)} \times D^{(m)} \qquad D^{(m)} \times 1$$

all the weights going into the 1st unit in the (m+1)th layer

What is the vectorized expression for the following?

$$W^{(m+1)} = \begin{bmatrix} \omega_{11} & \omega_{12} & \cdots & \cdot \\ & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}$$

$$\begin{bmatrix} z_1^{(m+1)} \\ z_2^{(m+1)} \\ \vdots \\ z_{D^{(m+1)}}^{(m+1)} \end{bmatrix} \longleftarrow z_i^{(m+1)} = \sum_{j=1}^{D^{(m)}} W_{ij}^{(m+1)} a_j^{(m)}, \qquad i = 1, ..., D^{(m+1)}$$

dim: $D^{(m+1)} \times 1$

output $\times$ input
$D^{(m+1)} \times D^{(m)}$

(A) $\left(\mathbf{W}^{(m+1)}\right)^{\top} \mathbf{a}^{(m)}$

(C) $\mathbf{W}^{(m+1)} \mathbf{a}^{(m)}$

(B) $\left(\mathbf{a}^{(m)}\right)^{\top} \mathbf{W}^{(m+1)}$

(D) $\mathbf{a}^{(m)} \mathbf{W}^{(m+1)}$

UNIVERSITY OF
TORONTO

# Solution: Result of Vectorizing Weighted Sum

What is the vectorized expression for the following?

$$\sum_{j=1}^{D^{(m)}} W_{ij}^{(m+1)} a_j^{(m)}, \qquad i = 1, ..., D^{(m+1)}$$

(A) $\left(\mathbf{W}^{(m+1)}\right)^{\top} \mathbf{a}^{(m)}$

(C) $\mathbf{W}^{(m+1)} \mathbf{a}^{(m)}$ **(Correct)**

(B) $\left(\mathbf{a}^{(m)}\right)^{\top} \mathbf{W}^{(m+1)}$

(D) $\mathbf{a}^{(m)} \mathbf{W}^{(m+1)}$

Recall the dimensions $\mathbf{W}^{(m+1)} \in \mathbb{R}^{D^{(m+1)} \times D^{(m)}}$, $\mathbf{a}^{(m)} \in \mathbb{R}^{D^{(m)} \times 1}$

UNIVERSITY OF
TORONTO

# Vectorizing Forward Pass Computations

## Non-Vectorized

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

$$z_i^{(m+1)} = \sum_j W_{ij}^{(m+1)} a_j^{(m)} + b_i^{(m+1)}$$

$$a_i^{(m+1)} = \sigma^{(m+1)}\left(z_i^{(m+1)}\right)$$

For output layer $L$

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

## Vectorized

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

$$\mathbf{z}^{(m+1)} = \mathbf{W}^{(m+1)} \mathbf{a}^{(m)} + \mathbf{b}^{(m+1)}$$

$$\mathbf{a}^{(m+1)} = \sigma^{(m+1)}\left(\mathbf{z}^{(m+1)}\right)$$

For output layer $L$

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

# Vectorizing Backward Pass

# Backward Pass Computations

1. Compute gradients for output layer $\longrightarrow$ this assumes a scalar output

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

2. Compute gradients for each hidden layer recursively

$$\frac{\partial C}{\partial z_j^{(m)}} = \left(\sum_i \frac{\partial C}{\partial z_i^{(m+1)}} \cdot W_{ij}^{(m+1)}\right) \cdot \sigma^{(m)'}\left(z_j^{(m)}\right), \qquad j = 1, \dots, D^{(m)}$$
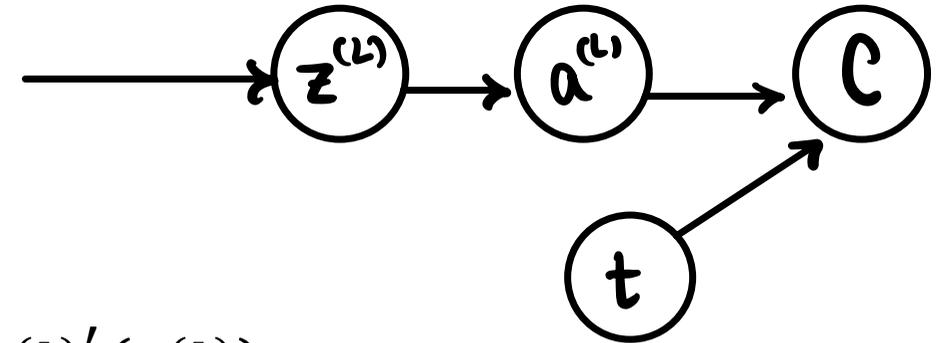
3. Compute gradients for the weights

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

UNIVERSITY OF
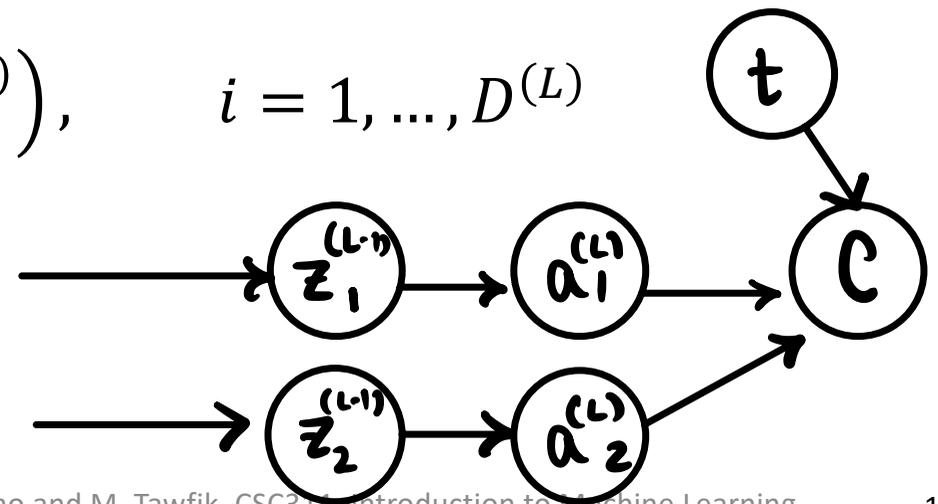TORONTO

# Gradients for the Output Layer



In our network, $z^{(L)}, a^{(L)}, t \in \mathbb{R}$ are scalars.

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

In general, $\mathbf{z}^{(L)}, \mathbf{a}^{(L)}, \mathbf{t} \in \mathbb{R}^{D^{(L)} \times 1}$ can be vectors. e.g. multi-class classification

$$\frac{\partial C}{\partial z_i^{(L)}} = \frac{\partial C}{\partial a_i^{(L)}} \cdot \frac{\partial a_i^{(L)}}{\partial z_i^{(L)}} = \frac{\partial C}{\partial a_i^{(L)}} \cdot {\sigma^{(L)}}'\left(z_i^{(L)}\right), \qquad i = 1, \ldots, D^{(L)}$$

UNIVERSITY OF
TORONTO

# Exercise 1: Gradients for the Output Layer

$\nabla_{\mathbf{a}^{(L)}} C$
$D^{(L)} \times 1$

What is the vectorized expression for the following?

$$\frac{\partial C}{\partial z_i^{(L)}} = \frac{\partial C}{\partial a_i^{(L)}} \cdot \sigma^{(L)'}\left(z_i^{(L)}\right), \qquad i = 1, \ldots, D^{(L)}$$

$D^{(L)} \times 1$

(A) $\left(\nabla_{\mathbf{a}^{(L)}} C\right)^{\top} \left(\sigma^{(L)'}(\mathbf{z}^{(L)})\right)$

(D) $\left(\sigma^{(L)'}(\mathbf{z}^{(L)})\right) \left(\nabla_{\mathbf{a}^{(L)}} C\right)^{\top}$

(B) $\left(\sigma^{(L)'}(\mathbf{z}^{(L)})\right)^{\top} \left(\nabla_{\mathbf{a}^{(L)}} C\right)$

(E) $\left(\nabla_{\mathbf{a}^{(L)}} C\right) \left(\sigma^{(L)'}(\mathbf{z}^{(L)})\right)^{\top}$

(C) $\left(\nabla_{\mathbf{a}^{(L)}} C\right) \odot \left(\sigma^{(L)'}(\mathbf{z}^{(L)})\right)$

$$\left[ \frac{\partial C}{\partial z_1^{(L)}} \quad \frac{\partial C}{\partial z_2^{(L)}} \quad \cdots \quad \frac{\partial C}{\partial z_{D^{(L)}}^{(L)}} \right]^{\top}$$

$\nabla_{\mathbf{z}^{(L)}} C$

$D^{(L)} \times 1$

UNIVERSITY OF TORONTO

# Solution 1: Gradients for the Output Layer

What is the vectorized expression for the following?

$$\frac{\partial C}{\partial z_i^{(L)}} = \frac{\partial C}{\partial a_i^{(L)}} \cdot \sigma^{(L)'}\left(z_i^{(L)}\right), \qquad i = 1, \ldots, D^{(L)}$$

(A) $\left(\nabla_{\mathbf{a}^{(L)}}C\right)^{\top}\left(\sigma^{(L)'}(\mathbf{z}^{(L)})\right)$

(D) $\left(\sigma^{(L)'}(\mathbf{z}^{(L)})\right)\left(\nabla_{\mathbf{a}^{(L)}}C\right)^{\top}$

(B) $\left(\sigma^{(L)'}(\mathbf{z}^{(L)})\right)^{\top}\left(\nabla_{\mathbf{a}^{(L)}}C\right)$

(E) $\left(\nabla_{\mathbf{a}^{(L)}}C\right)\left(\sigma^{(L)'}(\mathbf{z}^{(L)})\right)^{\top}$

(C) $\left(\nabla_{\mathbf{a}^{(L)}}C\right)\odot\left(\sigma^{(L)'}(\mathbf{z}^{(L)})\right)$ **(Correct)**

Recall the dimensions: $\nabla_{\mathbf{a}^{(L)}}C, \sigma^{(L)'}(\mathbf{z}^{(L)}), \nabla_{\mathbf{z}^{(L)}}C \in \mathbb{R}^{D^{(L)}\times 1}$

# Vectorizing Gradients for the Output Layer

Non-vectorized:

$$\frac{\partial C}{\partial z_i^{(L)}} = \frac{\partial C}{\partial a_i^{(L)}} \cdot \frac{\partial a_i^{(L)}}{\partial z_i^{(L)}} = \frac{\partial C}{\partial a_i^{(L)}} \cdot \sigma^{(L)'}\left(z_i^{(L)}\right), \qquad i = 1, \dots, D^{(L)}$$

Vectorized:

UNIVERSITY OF
TORONTO

# Solution: Vectorizing Gradients for the Output Layer

Non-vectorized:

$$\frac{\partial C}{\partial z_i^{(L)}} = \frac{\partial C}{\partial a_i^{(L)}} \cdot \frac{\partial a_i^{(L)}}{\partial z_i^{(L)}} = \frac{\partial C}{\partial a_i^{(L)}} \cdot \sigma^{(L)'}\left(z_i^{(L)}\right), \qquad i = 1, \dots, D^{(L)}$$

Vectorized:

$$\nabla_{\mathbf{z}^{(L)}} C = \left(\nabla_{\mathbf{a}^{(L)}} C\right) \odot \left(\sigma^{(L)'}\left(\mathbf{z}^{(L)}\right)\right)$$

# Backward Pass Computations

✓ 1. Compute gradients for output layer

$$\nabla_{\mathbf{z}^{(L)}} C = \left(\nabla_{\mathbf{a}^{(L)}} C\right) \odot \left(\sigma^{(L)'}\left(\mathbf{z}^{(L)}\right)\right)$$

2. Compute gradients for each hidden layer recursively

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

3. Compute gradients for the weights

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

UNIVERSITY OF
TORONTO

# Vectorizing the Recursive Step

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

1. Vectorizing

$$\frac{\partial C}{\partial a_j^{(m)}} = \sum_i \left( \frac{\partial C}{\partial z_i^{(m+1)}} \cdot W_{ij}^{(m+1)} \right), \qquad j = 1, \ldots, D^{(m)}$$

2. Vectorizing

$$\frac{\partial C}{\partial z_j^{(m)}} = \frac{\partial C}{\partial a_j^{(m)}} \cdot \sigma^{(m)'} \left( z_j^{(m)} \right), \qquad j = 1, \ldots, D^{(m)}$$

UNIVERSITY OF
TORONTO

# Exercise 2: Loss Derivative with respect to Activations

What is the vectorized expression for the following?

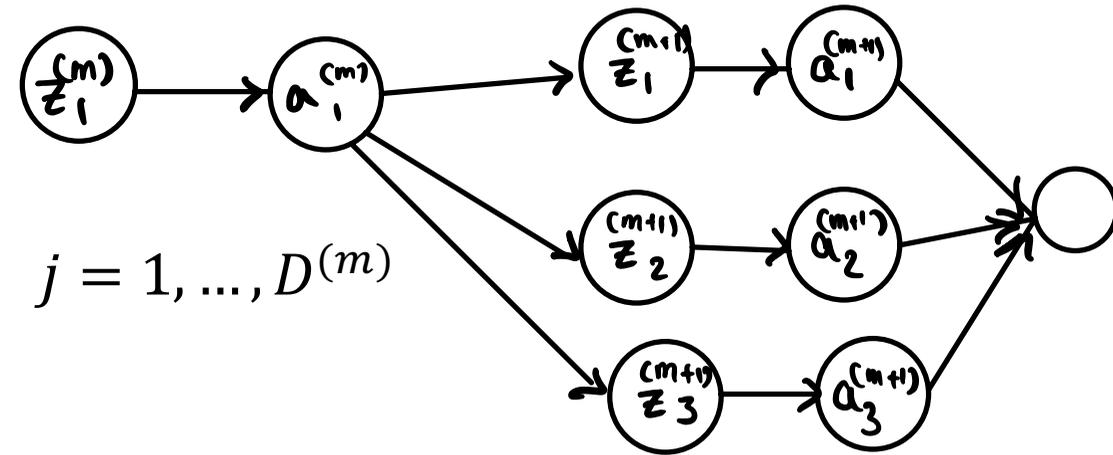$$\nabla_{\vec{a}^{(m)}} C \quad \leftarrow \quad \frac{\partial C}{\partial a_j^{(m)}} = \sum_{i=1}^{D^{(m+1)}} \left( \frac{\partial C}{\partial z_i^{(m+1)}} \cdot W_{ij}^{(m+1)} \right), \quad j = 1, \dots, D^{(m)}$$

$\nabla_{\vec{z}^{(m+1)}} C$

$D^{(m+1)} \times 1$

$\nabla_{\vec{a}^{(m)}} C$

$D^{(m)} \times 1$

$W^{(m+1)}$

$D^{(m+1)} \times D^{(m)}$

| | |
|---|---|
| (A) $\left(\mathbf{W}^{(m+1)}\right)^{\top} \left(\nabla_{\mathbf{z}^{(m+1)}} C\right)$ | (C) $\mathbf{W}^{(m+1)} \left(\nabla_{\mathbf{z}^{(m+1)}} C\right)$ |
| (B) $\left(\nabla_{\mathbf{z}^{(m+1)}} C\right)^{\top} \mathbf{W}^{(m+1)}$ | (D) $\left(\nabla_{\mathbf{z}^{(m+1)}} C\right) \mathbf{W}^{(m+1)}$ |

# Solution 2: Loss Derivative with respect to Activations

What is the vectorized expression for the following?

$$\frac{\partial C}{\partial a_j^{(m)}} = \sum_{i=1}^{D^{(m+1)}} \left( \frac{\partial C}{\partial z_i^{(m+1)}} \cdot W_{ij}^{(m+1)} \right), \qquad j = 1, \ldots, D^{(m)}$$

(A) $\left( \mathbf{W}^{(m+1)} \right)^{\top} \left( \nabla_{\mathbf{z}^{(m+1)}} C \right)$ **(Correct)**

(C) $\mathbf{W}^{(m+1)} \left( \nabla_{\mathbf{z}^{(m+1)}} C \right)$

(B) $\left( \nabla_{\mathbf{z}^{(m+1)}} C \right)^{\top} \mathbf{W}^{(m+1)}$

(D) $\left( \nabla_{\mathbf{z}^{(m+1)}} C \right) \mathbf{W}^{(m+1)}$

Recall the dimensions: $\mathbf{W}^{(m+1)} \in \mathbb{R}^{D^{(m+1)} \times D^{(m)}}$, $\nabla_{\mathbf{z}^{(m+1)}} C \in \mathbb{R}^{D^{(m+1)} \times 1}$

UNIVERSITY OF
TORONTO

# Exercise 3: Loss Derivative with respect to Pre-activations

What is the vectorized expression for the following?

$$\nabla_{\mathbf{z}^{(m)}} C \leftarrow \frac{\partial C}{\partial z_j^{(m)}} = \frac{\partial C}{\partial a_j^{(m)}} \cdot \sigma^{(m)'}\left(z_j^{(m)}\right), \qquad i = 1 \ldots, D^{(m)}$$

$\nabla_{\mathbf{a}^{(m)}} C$

$D^{(m)} \times 1$

$D^{(m)} \times 1$

$D^{(m)} \times 1$

(A) $\left(\nabla_{\mathbf{a}^{(m)}} C\right)^{\top} \left(\sigma^{(m)'}\left(\mathbf{z}^{(m)}\right)\right)$

(D) $\left(\sigma^{(m)'}\left(\mathbf{z}^{(m)}\right)\right) \left(\nabla_{\mathbf{a}^{(m)}} C\right)^{\top}$

(B) $\left(\sigma^{(m)'}\left(\mathbf{z}^{(m)}\right)\right)^{\top} \left(\nabla_{\mathbf{a}^{(m)}} C\right)$

(E) $\left(\nabla_{\mathbf{a}^{(m)}} C\right) \left(\sigma^{(m)'}\left(\mathbf{z}^{(m)}\right)\right)^{\top}$

(C) $\left(\nabla_{\mathbf{a}^{(m)}} C\right) \odot \left(\sigma^{(m)'}\left(\mathbf{z}^{(m)}\right)\right)$

UNIVERSITY OF TORONTO

# Solution 3: Loss Derivative with respect to Pre-activations

What is the vectorized expression for the following?

$$\frac{\partial C}{\partial z_j^{(m)}} = \frac{\partial C}{\partial a_j^{(m)}} \cdot \sigma^{(m)'}\left(z_j^{(m)}\right), \qquad i = 1 \dots, D^{(m)}$$

(A) $\left(\nabla_{\mathbf{a}^{(m)}} C\right)^\top \left(\sigma^{(m)'}(\mathbf{z}^{(m)})\right)$

(D) $\left(\sigma^{(m)'}(\mathbf{z}^{(m)})\right)\left(\nabla_{\mathbf{a}^{(m)}} C\right)^\top$

(B) $\left(\sigma^{(m)'}(\mathbf{z}^{(m)})\right)^\top \left(\nabla_{\mathbf{a}^{(m)}} C\right)$

(E) $\left(\nabla_{\mathbf{a}^{(m)}} C\right)\left(\sigma^{(m)'}(\mathbf{z}^{(m)})\right)^\top$

(C) $\left(\nabla_{\mathbf{a}^{(m)}} C\right) \odot \left(\sigma^{(m)'}(\mathbf{z}^{(m)})\right)$ **(Correct)**

Recall the dimensions: $\nabla_{\mathbf{a}^{(m)}} C, \sigma^{(m)'}\left(\mathbf{z}^{(m)}\right) \in \mathbb{R}^{D^{(m)} \times 1}$

UNIVERSITY OF
TORONTO

# Vectorizing the Recursive Step

Non-vectorized:

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

Vectorized:

UNIVERSITY OF
TORONTO

# Vectorizing the Recursive Step

Non-vectorized:

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

Vectorized:

$$\nabla_{\mathbf{z}^{(m)}} C = \left( (\mathbf{W}^{(m+1)})^{\top} \left( \nabla_{\mathbf{z}^{(m+1)}} C \right) \right) \odot \left( \sigma^{(m)'}(\mathbf{z}^{(m)}) \right)$$

UNIVERSITY OF
TORONTO

# Backward Pass Computations

✓ 1. Compute gradients for output layer

$$\nabla_{\mathbf{z}^{(L)}} C = \left(\nabla_{\mathbf{a}^{(L)}} C\right) \odot \left(\sigma^{(L)'}\left(\mathbf{z}^{(L)}\right)\right)$$

✓ 2. Compute gradients for each hidden layer recursively

$$\nabla_{\mathbf{z}^{(m)}} C = \left((\mathbf{W}^{(m+1)})^{\top} \left(\nabla_{\mathbf{z}^{(m+1)}} C\right)\right) \odot \left(\sigma^{(m)'}\left(\mathbf{z}^{(m)}\right)\right)$$

3. Compute gradients for the weights

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

UNIVERSITY OF
TORONTO

# Exercise 4: Loss Derivative with respect to Weights

$$\begin{bmatrix} \frac{\partial C}{\partial z_1^{(m)}} \\ \frac{\partial C}{\partial z_2^{(m)}} \\ \vdots \end{bmatrix} \qquad \nabla_{\vec{z}^{(m)}} C \qquad D^{(m)} \times 1$$

What is the vectorized expression for the following?

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

(A) $\left(\nabla_{\mathbf{z}^{(m)}} C\right)^\top \left(\mathbf{a}^{(m-1)}\right)$  

(D) $\left(\mathbf{a}^{(m-1)}\right)\left(\nabla_{\mathbf{z}^{(m)}} C\right)^\top$

(B) $\left(\mathbf{a}^{(m-1)}\right)^\top \left(\nabla_{\mathbf{z}^{(m)}} C\right)$  

(E) $\left(\nabla_{\mathbf{z}^{(m)}} C\right)\left(\mathbf{a}^{(m-1)}\right)^\top$

(C) $\left(\nabla_{\mathbf{z}^{(m)}} C\right) \odot \left(\mathbf{a}^{(m-1)}\right)$

*outer products*

$$\vec{a}^{(m-1)} \qquad D^{(m-1)} \times 1$$

$$D^{(m)} \times D^{(m-1)}$$

$$\begin{bmatrix} \frac{\partial C}{\partial \omega_{11}} & \frac{\partial C}{\omega_{12}} & \cdots & \frac{\partial C}{\omega_{1D}^{(m-1)}} \\ \vdots & \vdots & & \vdots \end{bmatrix}$$

*output x input*

*weights into the 1st unit in the $m^{th}$ layer*

UNIVERSITY OF TORONTO

# Solution 4: Loss Derivative with respect to Weights

What is the vectorized expression for the following?

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

(A) $\left(\nabla_{\mathbf{z}^{(m)}} C\right)^{\top} \left(\mathbf{a}^{(m-1)}\right)$

(D) $\left(\mathbf{a}^{(m-1)}\right)\left(\nabla_{\mathbf{z}^{(m)}} C\right)^{\top}$

(B) $\left(\mathbf{a}^{(m-1)}\right)^{\top}\left(\nabla_{\mathbf{z}^{(m)}} C\right)$

(E) $\left(\nabla_{\mathbf{z}^{(m)}} C\right)\left(\mathbf{a}^{(m-1)}\right)^{\top}$ **(Correct)**

(C) $\left(\nabla_{\mathbf{z}^{(m)}} C\right) \odot \left(\mathbf{a}^{(m-1)}\right)$

Recall the dimensions: $\nabla_{\mathbf{z}^{(m)}} C \in \mathbb{R}^{D^{(m)} \times 1}, \mathbf{a}^{(m-1)} \in \mathbb{R}^{D^{(m-1)} \times 1}, \nabla_{\mathbf{W}^{(m)}} C \in \mathbb{R}^{D^{(m)} \times D^{(m-1)}}$

UNIVERSITY OF
TORONTO

# Vectorizing the Gradients for the Weights

Non-vectorized:

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

Vectorized:

UNIVERSITY OF
TORONTO

# Solution: Vectorizing the Gradients for the Weights

Non-vectorized:

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

Vectorized:

$$\nabla_{\mathbf{W}^{(m)}} C = \left(\nabla_{\mathbf{z}^{(m)}} C\right)\left(\mathbf{a}^{(m-1)}\right)^{\top}$$

UNIVERSITY OF
TORONTO

# Backward Pass Computations

✓ 1. Compute gradients for output layer

$$\nabla_{\mathbf{z}^{(L)}} C = \left(\nabla_{\mathbf{a}^{(L)}} C\right) \odot \left(\sigma^{(L)'}\left(\mathbf{z}^{(L)}\right)\right)$$

✓ 2. Compute gradients for each hidden layer recursively

$$\nabla_{\mathbf{z}^m} C = \left((\mathbf{W}^{(m+1)})^T \left(\nabla_{\mathbf{z}^{(m+1)}} C\right)\right) \odot \left(\sigma^{(m)'}\left(\mathbf{z}^{(m)}\right)\right)$$

✓ 3. Compute gradients for the weights

$$\nabla_{\mathbf{W}^{(m)}} C = \left(\nabla_{\mathbf{z}^{(m)}} C\right)\left(\mathbf{a}^{(m-1)}\right)^\top$$

UNIVERSITY OF
TORONTO