

# **Neural Networks from Scratch**

## Neural Networks from Scratch Pt 2

Caleb Hallinan

01/07/2026

# Announcements

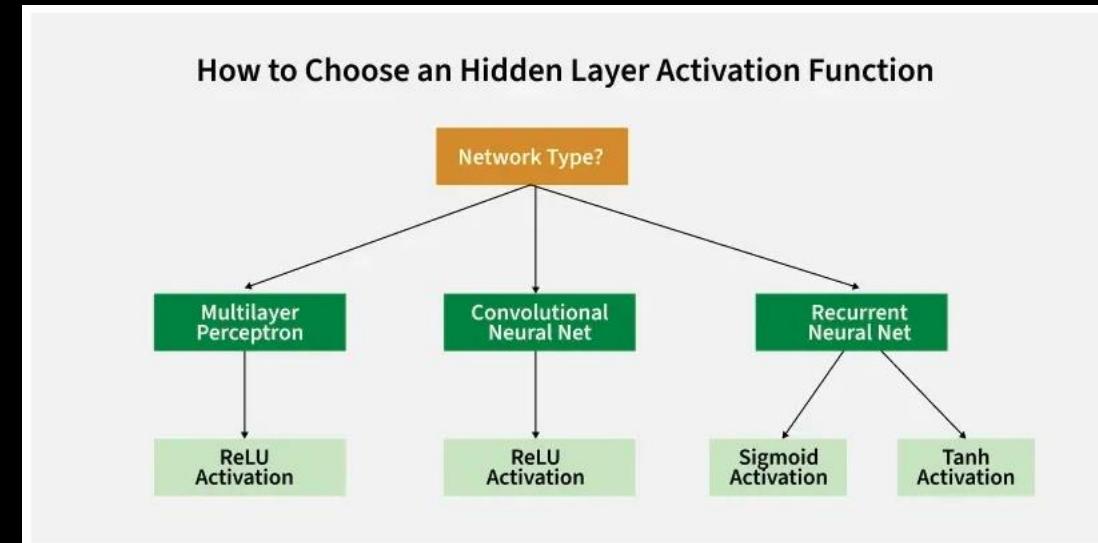
- Lecture and code from yesterday is on GitHub
- Today's code also on GitHub if you want to follow along
  - Going to try and make CNN and PyTorch code more interactive!



A few questions from  
reflection cards

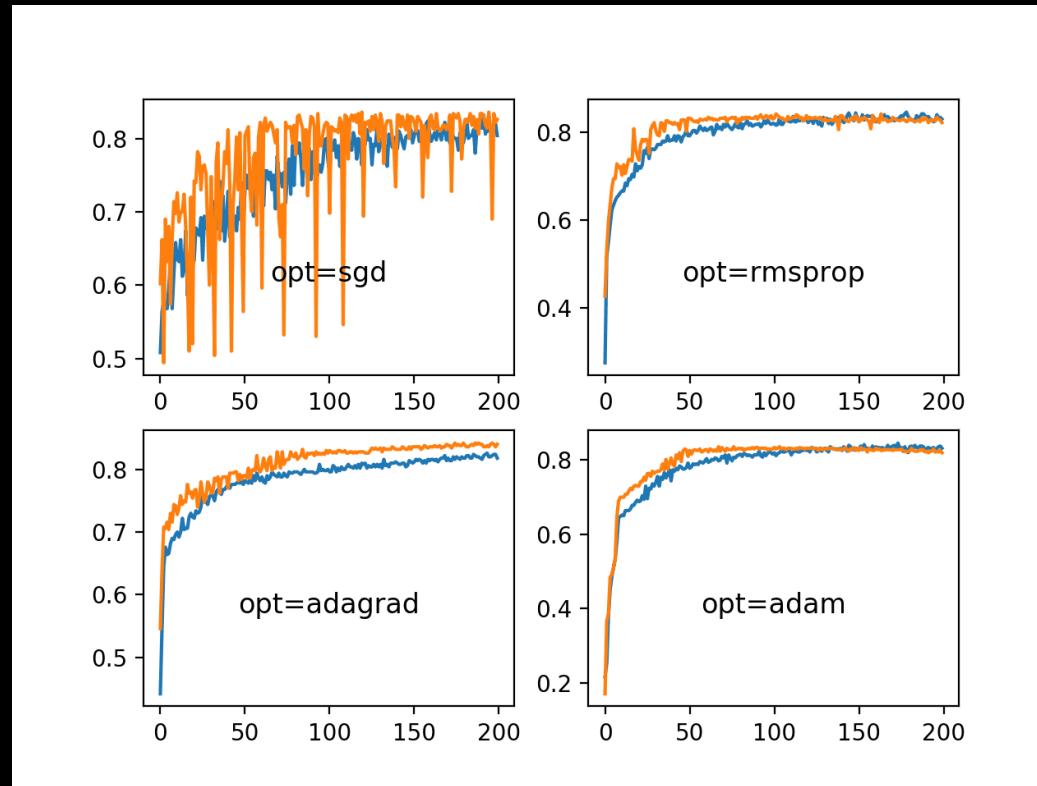
# How do you choose the best activation function?

- There is no universal rule for choosing the best activation functions – comes down to trial and error
  - However, there are some ideas for when it is best to use each one
  - <https://www.geeksforgeeks.org/deep-learning/choosing-the-right-activation-function-for-your-neural-network/>
- General rules:
  - ReLU or ReLU variant in hidden layers
  - Sigmoid or Softmax for last layer (why?)



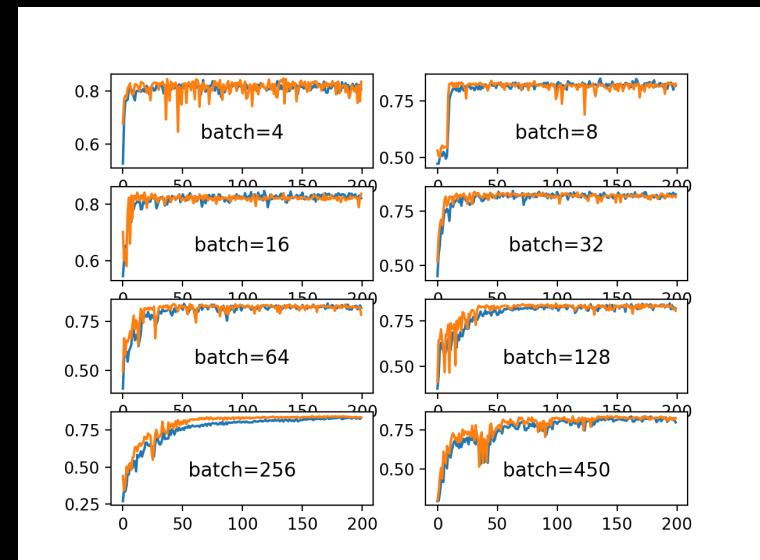
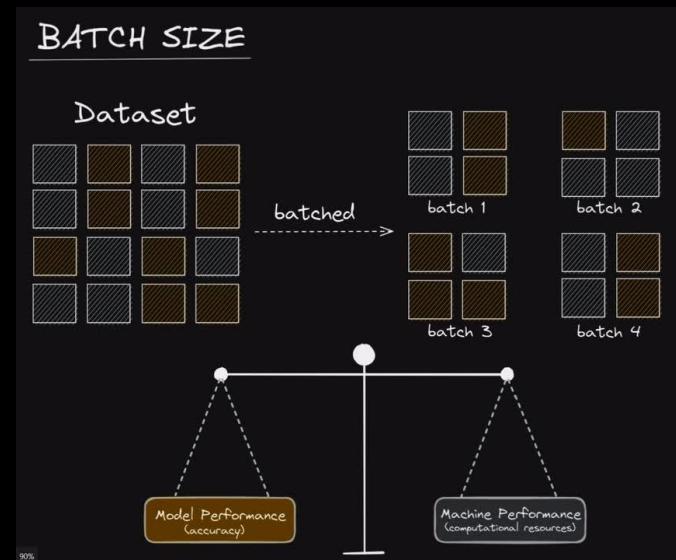
# How do adaptive learning rates work?

- We will go over this in more detail tomorrow!
  - Specifically, we will discuss learning rate decay and momentum



# Why use batches?

- Computing gradients on a small set of samples at once is much faster and more memory-efficient than processing the entire dataset every step, especially on GPUs
- Give a “noisy” but useful gradient estimate that often helps optimization and generalization compared to purely full-dataset updates.



# Why is $\text{np.dot}(AxB)$ the same as $\text{np.dot}(AxB^T)$

Depends on how the weights are stored

```
# Dense layer
class Layer_Dense:
    # Layer Initialization
    def __init__(self, n_inputs, n_neurons):
        # Initialize weights and biases
        self.weights = 0.01 * np.random.randn(n_inputs, n_neurons)
        self.biases = np.zeros((1, n_neurons))

    # Forward Pass
    def forward(self, inputs):
        # Calculate output values from inputs, weights and biases
        self.output = np.dot(inputs, self.weights) + self.biases
```

Our code from yesterday

If weights are initialized as  $(n_{\text{inputs}}, n_{\text{neurons}})$  then it is already ready to be multiplied to the inputs because  $(\text{batch}, n_{\text{inputs}}) \times (n_{\text{inputs}}, n_{\text{neurons}})$  is good!

If it was instead  $(n_{\text{neurons}}, n_{\text{inputs}})$ , then you would need to transpose them

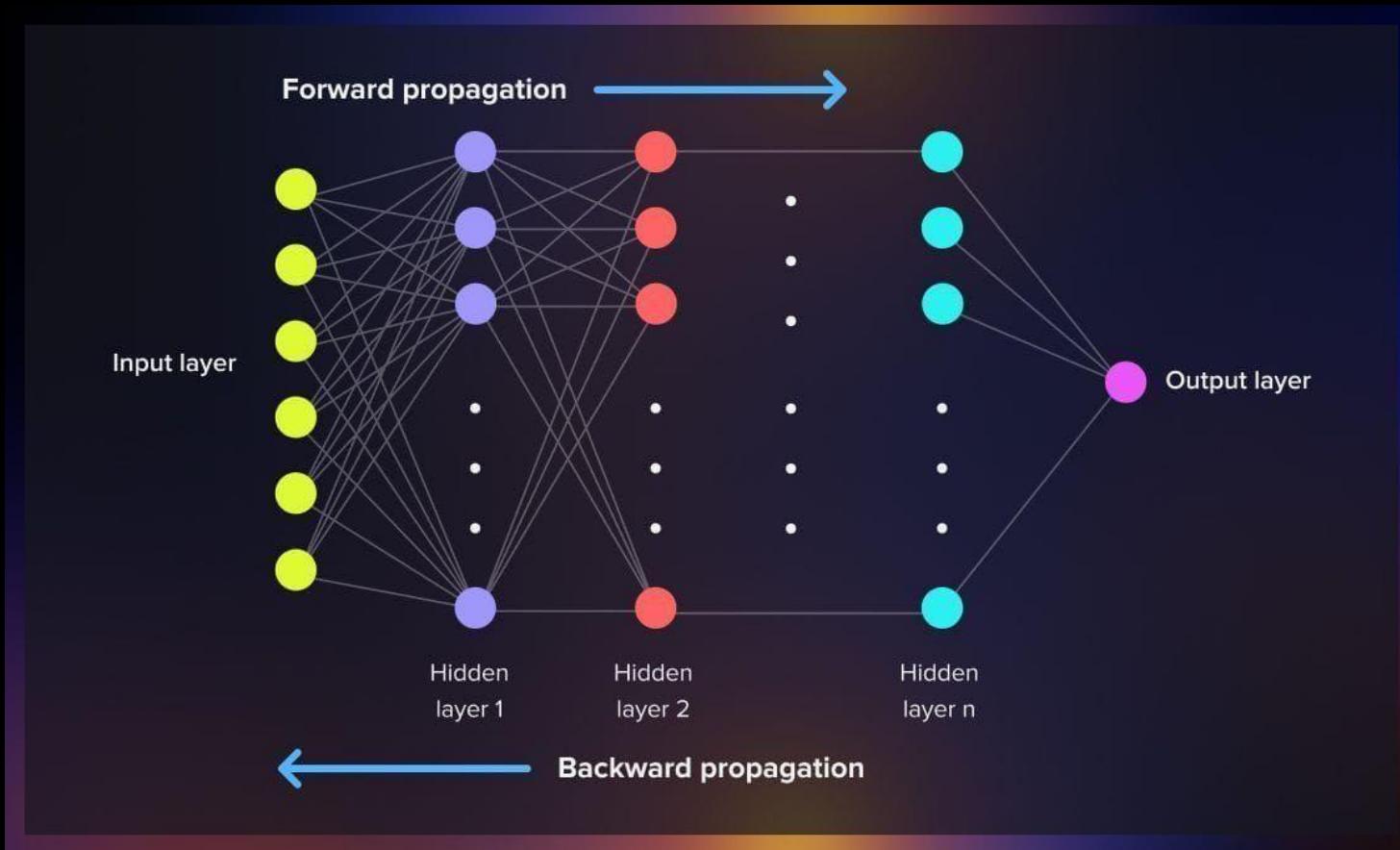
# Outline for today

- Continue going through neural networks from scratch – building from the ground up



# Today's notes and coding

Last Class: How do NN do forward propagation?



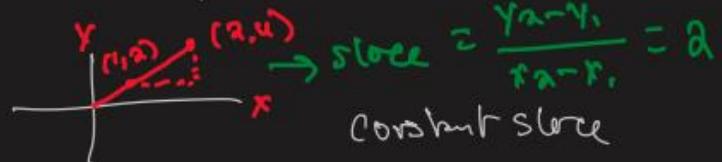
This class: How do NN do backward propagation?

## Derivatives

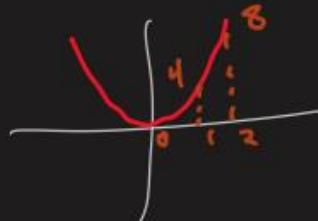
- shows impact of  $w, b$  on loss
- need to know effect of  $x$  on  $y$   
or  $\Delta x$  on  $\Delta y$

# Re impact of a param on output

$$\rightarrow f(x) = 2x$$



$$\rightarrow f(x) = 2x^2$$



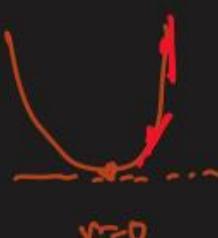
- we see as  $x$  increases  
 $y$  increases **linearly**

- we measure the  
slope of the tangent  
point

$$x \rightarrow 0 \quad \Delta x \cdot \Delta y$$

$$x \rightarrow 1 \quad \Delta x \rightarrow \Delta y \uparrow$$

$$x \rightarrow 5 \quad \Delta x \rightarrow \Delta y \uparrow\uparrow$$



So, to measure the impact which  
 $x$  has on  $y$ , we measure the slope  
of the tangent line at  $x$   
= Derivative!

If we change  $x$ , how does  $y$  change

---

Derivative notation:  $\frac{dy}{dx}$

$$\text{So if } ax^n \rightarrow a \cdot n x^{n-1}$$

etc.

Ask students derivatives

$$f(x) = x$$

$$\left( \begin{array}{l} f(x) = x^2 \\ f(x) = x^3 + 2x^2 + 6 \end{array} \right)$$

$$\frac{df}{dx}(x) = 1 = f'(x), \quad f'(x) = 2x \quad \left( \begin{array}{l} \frac{df}{dx} = 3x^2 + 4x + 0 \end{array} \right)$$

---

Gradients, partial derivatives, chain rule

What about  $f(x, y, z) \rightarrow f'(x, y, z)$ ?

PD notation:  $f(x, y, z) \rightarrow \frac{\partial}{\partial x} f(x, y, z)$  PD of  $f$   
 $\downarrow \frac{\partial}{\partial y} f(x, y, z)$  PD of  $f$  wrt  $x$   
 $\downarrow \frac{\partial}{\partial z} f(x, y, z)$  PD of  $f$  wrt  $y$





Chain rule = backbone of backprop

When you have function built from other functions, you can use the chain rule to take the derivative by:

differentiating the outside function (keeping the inside the same), then multiplying by the derivative of the inside.

Does anyone know the formula for the chain rule?

$$\frac{dy}{dx} = f'(g(x)) \cdot g'(x)$$

$$\text{where } y = f(g(x))$$

To solve:

$$\text{let } u = g(x), v = f(u)$$

$$\text{so } \frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx}$$

$$\text{Ex. } y = (3x^2 - 5x + 1)^7$$

What would  $u$  be?

$$u = 3x^2 - 5x + 1 \text{ s.t.}$$

$$y = u^7$$

$$\text{so } \frac{dy}{du} = 7u^6 \text{ and}$$

$$\frac{du}{dx} = 6x - 5 \text{ so}$$

$$\frac{dy}{dx} = 7u^6 (6x - 5) \quad \checkmark$$

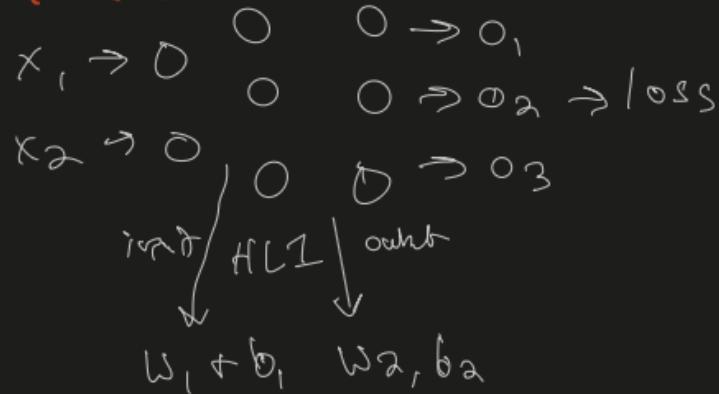
$$\frac{dy}{dx} = 7(3x^2 - 5x + 1)^6 (6x - 5)$$

A ANN is a chain of chain rules.  
Would look like

$$F(g(h(z(x))))$$

## Optimization

here is where we are



We know how to calc loss, but how do we use it to learn?

So we have loss function



How do you think we can adjust weights, biases to decrease loss?

This is optimization!

### Strategy 1

- ① pick random  $w, b$
- ② find loss
- ③ iterate 100000 times
- ④ find  $w, b$  w/o constraints



### Strategy 2

- ① init random  $w, b$
- ② randomly adjust  $w, b$
- ③ if loss decreases, then adjust next  $w, b$  to be closer to old  $w, b$ 
  - if loss increases, don't update weights/biases

- ④ keep track of lowest loss



### Best Method

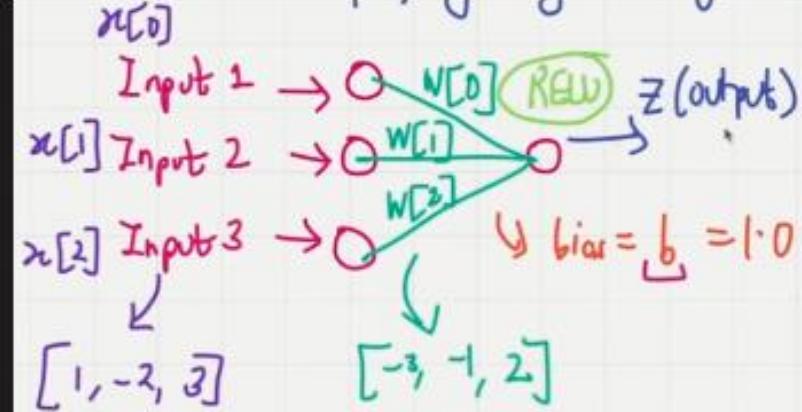
Go along negative gradient of score



Negative gradient shows the steepest descent

## Backpropagation → engine of NN

single neuron



So how do we optimize the weights & biases of neurons so that the output is close to 0?

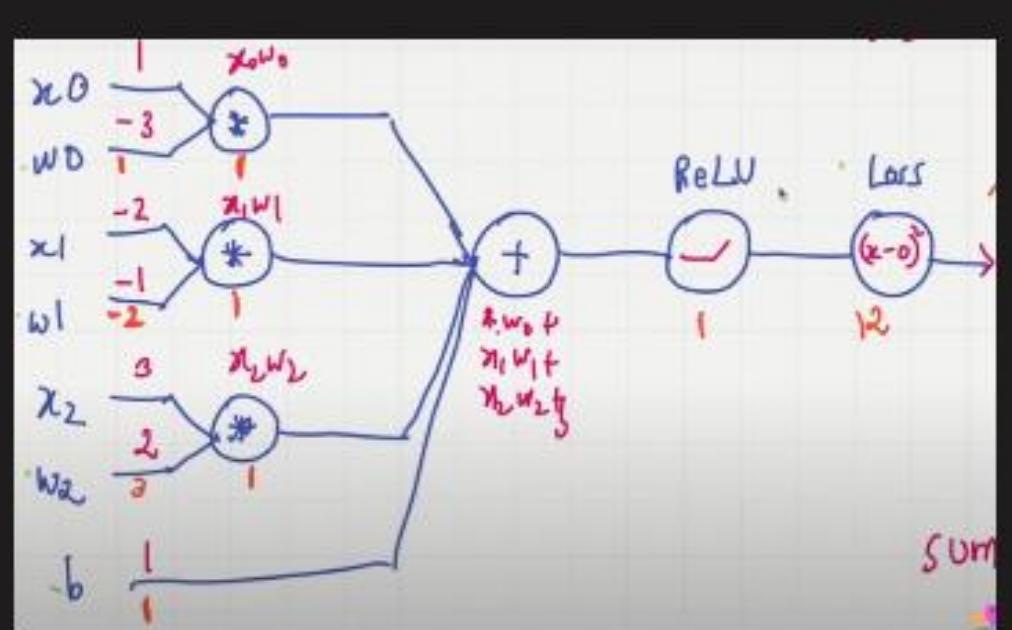
\*\* Let us minimize the loss function \*\*

→ We need to find how to update  $w[0], w[1], w[2], b$  so that loss is minimized.

→ For this, we have to move in negative gradient direction

We have to find derivative of loss with respect to:  $w[0], w[1], w[2], b$

$$\begin{aligned} z &= \text{ReLU}(x[0]w[0] + x[1]w[1] + x[2]w[2] + b) \\ l &= \max(xw_0 + x_1w_1 + x_2w_2 + b, 0) \\ l &\leftarrow \text{loss} = (z - o)^2 = z^2 \\ l &\rightarrow \text{minimize, use GD} \end{aligned}$$



Want to find what changes these

$$\begin{pmatrix} \frac{\partial L}{\partial w_0}, \frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2} \\ \frac{\partial L}{\partial b} \end{pmatrix} \xrightarrow{\text{so we can move by}} w_0 \rightarrow w_0 - \eta \frac{\partial L}{\partial w_0}$$

Let us consider  $w_0$ . We need to find  $\frac{\partial \text{Loss}}{\partial w_0}$

$$\text{Loss} = (\text{ReLU}(\text{sum}(\text{mul}(x_0, w_0), \text{mul}(x_1, w_0), \text{mul}(x_2, w_0), b)))^2$$

$$\frac{\partial \text{Loss}}{\partial w_0} = \frac{\partial \text{Loss}}{\partial \text{ReLU}()} * \frac{\partial \text{ReLU}()}{\partial \text{sum}()} * \frac{\partial \text{sum}()}{\partial \text{mul}(x_0, w_0)} * \frac{\partial \text{mul}(x_0, w_0)}{\partial w_0}$$

$$y = x^2 \quad y = \text{ReLU}(x) \quad y = x+1 \quad y = x_0 w_0$$

$$\frac{\partial y}{\partial x} = 2x \quad \frac{\partial y}{\partial x} = 1 \quad \frac{\partial y}{\partial x} = 1 \quad \frac{\partial y}{\partial x} = 1$$

$$= 2 \cdot \text{ReLU}(\text{sum}(...)) \quad \text{so } \text{sum}(...) = b = 1$$

$$= 2 \cdot \text{ReLU}(b) \quad b = 1$$

$$= -12$$

Ask students what  $y =$  for all these

$x_0 = 1$  here

6 comes from sum; in step 1 above

$$\frac{\partial \text{Loss}}{\partial w_1} = \frac{\partial \text{Loss}}{\partial \text{ReLU}()} * \frac{\partial \text{ReLU}()}{\partial \text{sum}()} * \frac{\partial \text{sum}()}{\partial \text{mul}(x_1, w_1)} * \frac{\partial \text{mul}(x_1, w_1)}{\partial w_1} = 12 * 1 * 1 * -2 = -24$$

✓

$$\frac{\partial \text{Loss}}{\partial w_2} = \frac{\partial \text{Loss}}{\partial \text{ReLU}()} * \frac{\partial \text{ReLU}()}{\partial \text{sum}()} * \frac{\partial \text{sum}()}{\partial \text{mul}(x_2, w_2)} * \frac{\partial \text{mul}(x_2, w_2)}{\partial w_2} = 12 * 1 * 1 * 3 = 36$$

✓

$$\frac{\partial \text{Loss}}{\partial b} = \frac{\partial \text{Loss}}{\partial \text{ReLU}()} * \frac{\partial \text{ReLU}()}{\partial \text{sum}()} * \frac{\partial \text{sum}()}{\partial b} = 12 * 1 * 1 = 12$$

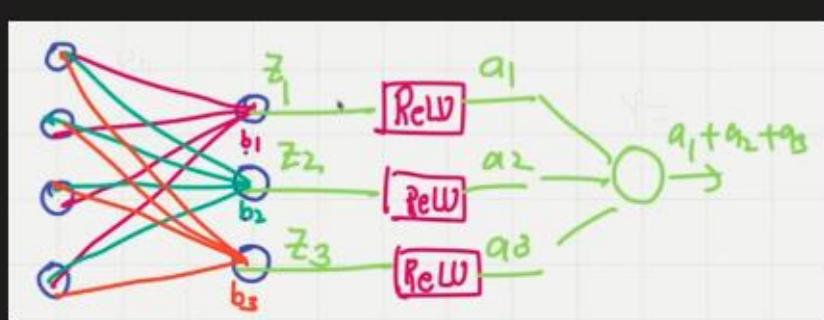
Earlier loss = 36

New loss =  $[\text{ReLU}(\text{sum}(...))]^2$

$$\begin{aligned} \text{current loss} : \quad w_0 &= w_0 - 0.001 * \frac{\partial \text{Loss}}{\partial w_0} \\ &= -3 - 0.001 * (-12) = -3 + 12 = 9 \\ w_1 &= w_1 - 0.001 * \frac{\partial \text{Loss}}{\partial w_1} \\ &= -1 - 0.001 * (-24) = -1 + 24 = 23 \\ w_2 &= w_2 - 0.001 * \frac{\partial \text{Loss}}{\partial w_2} \\ &= 3 - 0.001 * 36 = 3 - 36 = -33 \end{aligned}$$

$(-3 + 12)^2 + 2(-0.976) + 3(23) + 0.989 = 5.82$

$\text{Loss}_{\text{new}} = (5.82)^2 = 33.8724$



Want loss to be as close to 0 as possible

$$\text{Loss} = (a_1 + a_2 + a_3)^2 = y^2$$

Notation:

$$\begin{aligned} z_1 &= w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + w_{14}x_4 + b_1 \\ z_2 &= w_{21}x_1 + w_{22}x_2 + w_{23}x_3 + w_{24}x_4 + b_2 \\ z_3 &= w_{31}x_1 + w_{32}x_2 + w_{33}x_3 + w_{34}x_4 + b_3 \end{aligned}$$

$$\begin{aligned} a_1 &= \text{ReLU}(z_1) \\ a_2 &= \text{ReLU}(z_2) \\ a_3 &= \text{ReLU}(z_3) \end{aligned}$$

$y = (a_1 + a_2 + a_3)$   
So that's why deriv of  $y$  wrt  $a_1 = 1$

Backward pass:  $\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial y} * \frac{\partial y}{\partial a_1} * \frac{\partial a_1}{\partial z_1} * \frac{\partial z_1}{\partial w_{11}} = 2y * 1_{(z_1 > 0)} * x_1$

How would you write this?

diff is  $\nwarrow$

$$\frac{\partial L}{\partial w_{12}} = \frac{\partial L}{\partial y} * \frac{\partial y}{\partial a_1} * \frac{\partial a_1}{\partial z_1} * \frac{\partial z_1}{\partial w_{12}} = 2y * 1_{(z_1 > 0)} * x_2$$

Forward pass:  $z_1 = 3; z_2 = 7.2; z_3 = 11.4$

$a_1 = 3; a_2 = 7.2; a_3 = 11.4$

$y = a_1 + a_2 + a_3 = 21.6$

$L = y^2 = 466.56$

Inputs = [1, 2, 3, 4]

Did all these partial derivatives to finally update weights

$$\begin{aligned} \frac{\partial L}{\partial w_{11}} &= 2y * x_1 = 43.2; \frac{\partial L}{\partial w_{12}} = 2y * x_2 = 86.4; \frac{\partial L}{\partial w_{13}} = 129.6; \frac{\partial L}{\partial w_{14}} = 172.8 \\ \frac{\partial L}{\partial w_{21}} &= 2y * x_1 = 43.2; \frac{\partial L}{\partial w_{22}} = 2y * x_2 = 86.4; \frac{\partial L}{\partial w_{23}} = 129.6; \frac{\partial L}{\partial w_{24}} = 172.8 \\ \frac{\partial L}{\partial w_{31}} &= 43.2; \frac{\partial L}{\partial w_{32}} = 86.4; \frac{\partial L}{\partial w_{33}} = 129.6; \frac{\partial L}{\partial w_{34}} = 172.8 \\ \frac{\partial L}{\partial b_1} &= 2y = 43.2; \frac{\partial L}{\partial b_2} = 2y = 86.4; \frac{\partial L}{\partial b_3} = 129.6; \end{aligned}$$

Example with  $w_{11}$

Gradient Descent update:

$$w_{11} = w_{11} - \frac{0.001}{\text{step size}} \frac{\partial L}{\partial w_{11}} = 0.1 - 0.001 * 43.2 = 0.0568 \rightarrow \text{Weight}$$

Now can check new loss value based on new weight

	$w_{11}$	$w_{12}$	$w_{13}$	$w_{14}$
$w_1$ :	7	7	7	7
$w_2$ :	0.5	0.6	0.7	0.8
$w_3$ :	0.9	1	1.1	1.2
$b_1$ :	0.1	4.5		
$b_2$ :	0.2			
$b_3$ :	0.3			

# Reflection Cards

## Reflection Card

Please reflect on today's lesson in Neural Networks from Scratch.

Reflection cards are not graded for content. However, the contents of these reflection cards may help identify potential common areas of confusion that can be addressed in the next class along with helping me make the class better :)

Hi, Caleb. When you submit this form, the owner will see your name and email address.

\* Required

- Essentially a means to help me make this class better!
    - Would you be interested in a hybrid version of class?
1. What is something that you learned in today's lecture?
  2. What is something that you are still confused about from today's lecture?
  3. Do you have any other comments/feedback/thoughts/suggestions/concerns?

<https://forms.office.com/r/Kv8LtW4jH>