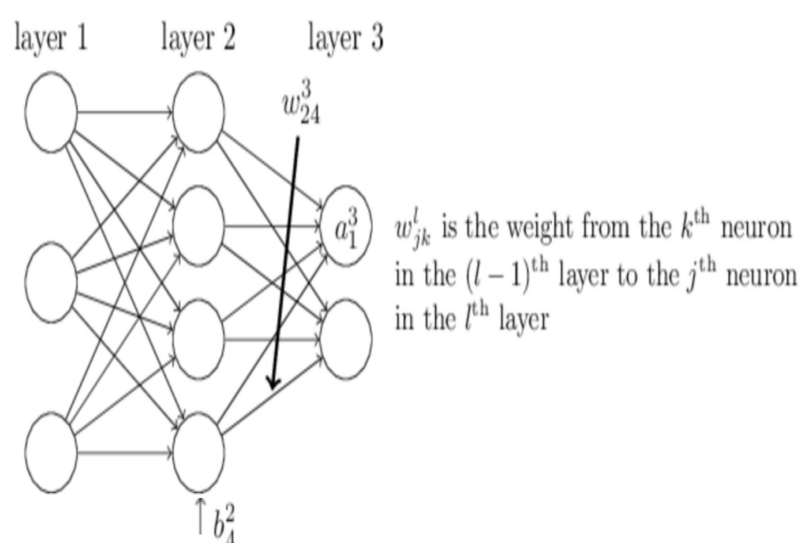


### 3. Backpropagation

*Extension Note:* Homework 4 due date has been extended by 1 day to **July 27 23:59UTC**.

One of the key steps for training multi-layer neural networks is stochastic gradient descent. We will use the back-propagation algorithm to compute the gradient of the loss function with respect to the model parameters.

Consider the  $L$ -layer neural network below:



In the following problems, we will use the following notation:  $b_j^l$  is the bias of the  $j^{th}$  neuron in the  $l^{th}$  layer,  $a_j^l$  is the activation of  $j^{th}$  neuron in the  $l^{th}$  layer, and  $w_{jk}^l$  is the weight for the connection from the  $k^{th}$  neuron in the  $(l-1)^{th}$  layer to the  $j^{th}$  neuron in the  $l^{th}$  layer.

If the activation function is  $f$  and the loss function we are minimizing is  $C$ , then the equations describing the network are:

$$a_j^l = f \left( \sum_k w_{jk}^l a_k^{l-1} + b_j^l \right)$$

$$\text{Loss} = C(a^L)$$

For  $l = 1, \dots, L$ .

### Computing the Error

2/2 points (graded)

Let the weighted inputs to the  $d$  neurons in layer  $l$  be defined as  $z^l \equiv w^l a^{l-1} + b^l$ , where  $z^l \in \mathbb{R}^d$ . As a result, we can also write the activation of layer  $l$  as  $a^l \equiv f(z^l)$ , and the "error" of neuron  $j$  in layer  $l$  as  $\delta_j^l \equiv \frac{\partial C}{\partial z_j^l}$ . Let  $\delta^l \in \mathbb{R}^d$  denote the full vector of errors associated with layer  $l$ .

Back-propagation will give us a way of computing  $\delta^l$  for every layer.

Assume there are  $d$  outputs from the last layer (i.e.  $a^L \in \mathbb{R}^d$ ). What is  $\delta_j^L$  for the last layer?

☒  $\frac{\partial C}{\partial a_j^L} f'(z_j^L)$  ✓

☐  $\sum_{k=1}^d \frac{\partial C}{\partial a_k^L} f'(z_j^L)$

☐  $\frac{\partial C}{\partial a_j^L}$

☐  $f'(z_j^L)$

What is  $\delta_j^l$  for all  $l \neq L$ ?

☒  $\sum_k w_{kj}^{l+1} \delta_k^{l+1} f'(z_j^l)$  ✓

☐  $\delta_k^{l+1} f'(z_j^l)$

☐  $\sum_k w_{jk}^{l-1} \delta_j^{l-1} f'(z_j^l)$

☐  $\sum_k w_{kj}^{l+1} \delta_k^{l+1} f(z_j^l)$

Solution:

We make use of the chain rule.

1. By definition,  $\delta_j^L = \frac{\partial C}{\partial a_j^L} \frac{\partial a_j^L}{\partial z_j^L} = \frac{\partial C}{\partial a_j^L} f'(z_j^L)$ .
2. We have:

$$\begin{aligned} \delta_j^l &= \frac{\partial C}{\partial z_j^l} \\ &= \sum_k \frac{\partial C}{\partial z_k^{l+1}} \frac{\partial z_k^{l+1}}{\partial z_j^l} \\ &= \sum_k \frac{\partial z_k^{l+1}}{\partial z_j^l} \delta_k^{l+1} \end{aligned}$$

Then we have  $z_k^{l+1} = \sum_j w_{kj}^{l+1} a_j^l + b_k^{l+1} = \sum_j w_{kj}^{l+1} f(z_j^l) + b_k^{l+1}$ . Taking the derivative of this with respect to  $z_j^l$  gives  $w_{kj}^{l+1} f'(z_j^l)$ .

Combining the two gives the final answer:  $\delta_j^l = \sum_k w_{kj}^{l+1} \delta_k^{l+1} f'(z_j^l)$ .

Submit

You have used 2 of 2 attempts

📘 Answers are displayed within the problem

## Parameter Derivatives

2/2 points (graded)  
During SGD we are interested in relating the errors computed by back-propagation to the quantities of real interest: the partial derivatives of the loss with respect to our parameters. Here that is  $\frac{\partial C}{\partial w_{jk}^l}$  and  $\frac{\partial C}{\partial b_j^l}$ .

What is  $\frac{\partial C}{\partial w_{jk}^l}$ ? Write in terms of the variables  $a_k^{l-1}$ ,  $w_j^l$ ,  $b_j^l$ , and  $\delta_j^l$  if necessary.

Example of writing superscripts and subscripts:

$\delta_j^l$

$w_{jk}^l$

$\frac{\partial C}{\partial w_{jk}^l} =$

$a_k^{l-1} \cdot \delta_j^l$

✔ Answer: a\_k^{l-1} \* delta\_j^l

What is  $\frac{\partial C}{\partial b_j^l}$ ? Write in terms of the variables  $a_k^{l-1}$ ,  $w_j^l$ ,  $b_j^l$ , and  $\delta_j^l$  if necessary.

$\frac{\partial C}{\partial b_j^l} =$

$\delta_j^l$

✔ Answer: delta\_j^l

STANDARD NOTATION

Solution:

1.  $\frac{\partial C}{\partial w_{jk}^l} = \frac{\partial C}{\partial z_j^l} \frac{\partial z_j^l}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$

2.  $\frac{\partial C}{\partial b_j^l} = \frac{\partial C}{\partial z_j^l} \frac{\partial z_j^l}{\partial b_j^l} = 1 * \delta_j^l$

Submit

You have used 1 of 5 attempts

 Answers are displayed within the problem

### Activation Functions: Sigmoid

4/4 points (graded)  
Recall that there are several different possible choices of activation functions  $f$ . Let's get more familiar with them and their gradients.

What is the derivative of the sigmoid function,  $\sigma(z) = \frac{1}{1+e^{-z}}$ ? Please write your answer in terms of  $e$  and  $z$ :

$\frac{e^{-z}}{(1+e^{-z})^2}$

✔ Answer: e^{(-z)} / (1 + e^{(-z)})^2

Which of the following is true of  $\sigma'(z)$  as  $||z||$  gets large?

- ☐ Its magnitude becomes large.
- ☒ Its magnitude becomes small. ✔
- ☐ It suffers from high variance.

What is the derivative of the ReLU function,  $\text{ReLU}(z) = \max(0, z)$  for  $z > 0$ ?

1

✔ Answer: 1

For  $z < 0$ ?

0

✔ Answer: 0

STANDARD NOTATION


Solution:

$\sigma'(z) = \sigma(z)(1 - \sigma(z))$ . As  $z$  gets large in magnitude, the sigmoid function saturates, and the gradient approaches zero.

ReLU is a simple activation function. Above zero, it has a constant gradient of 1. Below zero, it is always zero.

Submit

You have used 1 of 5 attempts

 Answers are displayed within the problem

## Simple Network

4/4 points (graded)

Consider a simple 2-layer neural network with a single neuron in each layer. The loss function is the quadratic loss:  $C = \frac{1}{2}(y - t)^2$ , where  $y$  is the prediction and  $t$  is the target.

Starting with input  $x$  we have:

- $z_1 = w_1 x$
- $a_1 = \text{ReLU}(z_1)$
- $z_2 = w_2 a_1 + b$
- $y = \sigma(z_2)$
- $C = \frac{1}{2}(y - t)^2$

Consider a target value  $t = 1$  and input value  $x = 3$ . The weights and bias are  $w_1 = 0.01$ ,  $w_2 = -5$ , and  $b = -1$ .

Please provide numerical answers accurate to at least three decimal places.

What is the loss?

0.28842841648243966

 Answer: 0.28842841648243966

What are the derivatives with respect to the parameters?

$\frac{\partial C}{\partial w_1} =$

2.0809

 Answer: 2.0809165621704553

$\frac{\partial C}{\partial w_2} =$

-0.0041618

 Answer: -0.00416183312434091

$\frac{\partial C}{\partial b} =$

-0.13872777

 Answer: -0.13872777081136367

STANDARD NOTATION


Solution:

Using the chain rule, we have:

- $\frac{\partial C}{\partial w_1} = \frac{\partial C}{\partial y} \frac{\partial y}{\partial z_2} \frac{\partial z_2}{\partial a_1} \frac{\partial a_1}{\partial z_1} \frac{\partial z_1}{\partial w_1} = (y - t) y (1 - y) w_2 \mathbf{1}\{z_1 > 0\} x$
- $\frac{\partial C}{\partial w_2} = \frac{\partial C}{\partial y} \frac{\partial y}{\partial z_2} \frac{\partial z_2}{\partial w_2} = (y - t) y (1 - y) a_1$
- $\frac{\partial C}{\partial b} = (y - t) y (1 - y)$

Submit

You have used 1 of 5 attempts

 Answers are displayed within the problem

SGD

1/1 point (graded)

Referring to the previous problem, what is the update rule for  $w_1$  in the SGD algorithm with step size  $\eta$ ? Write in terms of  $w_1$ ,  $\eta$ , and  $\frac{\partial C}{\partial w_1}$ ; enter the latter as (partialC)/(partialw\_1), noting the lack of space in the variable names:

Next  $w_1 =$ 

w\_1 - eta\*(partialC)/(pari

 **Answer:** w\_1 - eta \* (partialC)/(partialw\_1)

[STANDARD NOTATION](#)

Solution:

The definition of the simple SGD update rule is new\_parameter = old\_parameter - learning\_rate \* derivative of loss w.r.t old parameter.

Submit

 You have used 1 of 5 attempts

 Answers are displayed within the problem

Discussion

[Show Discussion](#)

Topic: Unit 3 Neural networks (2.5 weeks):Homework 4 / 3. Backpropagation