edX
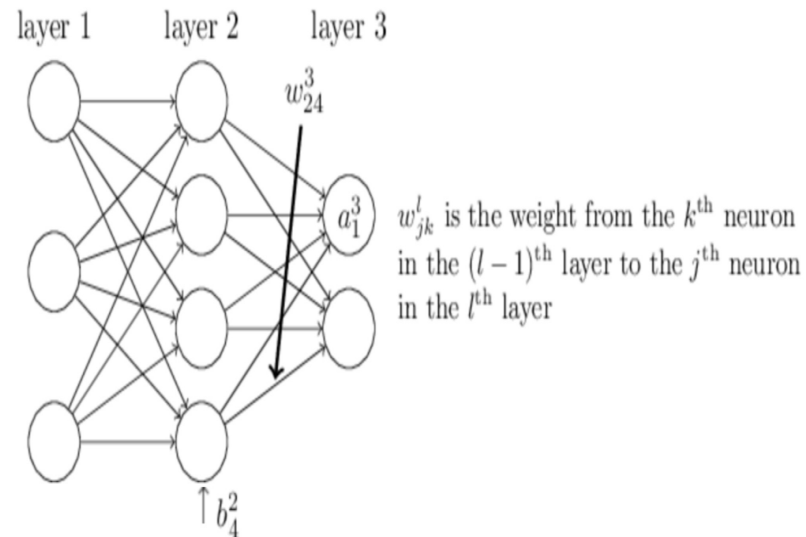
Course  >  Unit 3 Neural networks (2.5 weeks)  >  Homework 4  >  3. Backpropagation

# 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:

layer 1  layer 2  layer 3

$w_{jk}^l$ is the weight from the $k^{\text{th}}$ neuron in the $(l-1)^{\text{th}}$ layer to the $j^{\text{th}}$ neuron in the $l^{\text{th}}$ layer

In the following problems, we will 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:

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

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

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

# 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?

- ○ $\frac{\partial C}{\partial a_j^L} f'\left(z_j^L\right)$ ✔

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

- ○ $\frac{\partial C}{\partial a_j^L}$

- ○ $f'\left(z_j^L\right)$

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

○ $\sum_k w_{kj}^{l+1} \delta_k^{l+1} f'\left(z_j^l\right)$ ✔

○ $\delta_k^{l+1} f'\left(z_j^l\right)$

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

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

**Solution:**

We make use of the chain rule.

1. By definition, $\delta_j^L = \frac{\partial C}{\partial a_j^L}\frac{\partial a_j^L}{\partial z_j^L} = \frac{\partial C}{\partial a_j^L} f'\left(z_j^L\right)$.

2. We have:

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

$$= \sum_k \frac{\partial z_k^{l+1}}{\partial z_j^l} \delta_k^{l+1}$$

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:

$delta\_j\backslash \hat{\ } l$ for $\delta_j^l$

$w\_\{jk\}\backslash\hat{}\,l$ for $w^l_{jk}$

$\dfrac{\partial C}{\partial w^l_{jk}} =$ | delta_j^l*a_k^(l-1) | ✔ **Answer:** a_k^(l-1)*delta_j^l

$\delta^l_j \cdot a^{l-1}_k$

What is $\dfrac{\partial C}{\partial b^l_j}$? Write in terms of the variables $a^{l-1}_k$, $w^l_j$, $b^l_j$, and $\delta^l_j$ if necessary.

$\dfrac{\partial C}{\partial b^l_j} =$ | delta_j^l | ✔ **Answer:** delta_j^l

$\delta^l_j$

[ STANDARD NOTATION ]

**Solution:**

1. $\dfrac{\partial C}{\partial w^l_{jk}} = \dfrac{\partial C}{\partial z^l_j}\dfrac{\partial z^l_j}{\partial w^l_{jk}} = a^{l-1}_k \delta^l_j$

2. $\dfrac{\partial C}{\partial b^l_j} = \dfrac{\partial C}{\partial z^l_j}\dfrac{\partial z^l_j}{\partial b^l_j} = 1 * \delta^l_j$

[ Submit ]　　You have used 1 of 5 attempts

ℹ   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, $\sigma\left(z\right) = \frac{1}{1+e^{-z}}$? Please write your answer in terms of $e$ and $z$:

| e^(-z)/(1+e^(-z))^2 |

✔ **Answer:** e^(-z) / (1 + e^(-z))^2

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

Which of the following is true of $\sigma'\left(z\right)$ 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, $\mathrm{ReLU}\left(z\right) = \max\left(0, z\right)$ for $z > 0$?

1

✔ **Answer:** 1

1

For $z < 0$?

0

✔ **Answer:** 0

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:

- $z_1 = w_1 x$

- $a_1 = \mathrm{ReLU}(z_1)$

- $z_2 = w_2 a_1 + b$

- $y = \sigma(z_2)$

- $C = \frac{1}{2}(y - t)^2$

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?

$\frac{\partial C}{\partial w_1} =$ 2.0809165621704553    ✔ **Answer:** 2.0809165621704553

$\frac{\partial C}{\partial w_2} =$ | -0.00416183312434091 |     ✔ **Answer:** -0.00416183312434091

$\frac{\partial C}{\partial b} =$ | -0.13872777081136367 |     ✔ **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

ℹ   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)/(part |  ✔ **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

ℹ   Answers are displayed within the problem

# Discussion

**Hide Discussion**

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

**Add a Post**

Show all posts    ▼                                                          by recent activity ▼

?   Simple Network, Derivative                                                            1

☑ [Staff]Coarse grader?

4

💬 Simple Network

I'm taking the derivatives w.r.t each parameter, and the expressions are very long. Are we perhaps supposed to plug-in the values we're given, or...

10

💬 Simple Network

any clues? I'm at the third try. I got the Loss and ∂C/∂w2 but didn't figure out ∂C/∂w1 and ∂C/∂b

8

? [STAFF] Parameter Derivatives

6

💬 [staff] Simple network, why is my second answer wrong? All others are correct? Is this a notation problem again?

Why is my second answer wrong? I calculated several times. I want to know if this is a typo or a calculation mistake.

4

? computing the error part 2 reset

3

💬 Can the course follow the same notation convention please?

May I suggest a little improvement to the course? As stated in the title, please state a set of notation convention at the beginning of the course, a...

3

? [Simple Network]: To staff, can you look into my entered answer

since it takes the derivatives with respective to w1, w2 and b respectively but the system states that there is w1 and w2, b are invalid input

2

💬 Learning Source - Hint

for those who still struggle to understand our course material, watch this video. Took some time to understand the mechanics behind by this vid...

8

💬 Backpropagation -- a little history, and a cautionary tale

Backpropagation is one of those ideas like Bayes Rule that's "whoa, why didn't I think of that"...*after* you see it. It's just the plain old ordinary c...

👤 Community TA

4

💬 What the σ(x)function in y=σ(z2)?

What the σ(x)function in y=σ(z2)?

3

[Staff] Online Derivatives

☑ <u>Is there any objection to us looking up the requested derivatives online? i.e. RelU and Sigmoid. Additional note: Well, I just looked them up online...</u>　**3**

2　Simple Network - puzzled with rejection of my submission