

Explanatory Notes for 6.390

Shauntclair Ruiz (Current TA)

Fall 2022

What do these derivatives equal?

Let's look at each of these derivatives and see if we can't simplify them a bit.

First, every gradient needs

- The **loss derivative**:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{a}^L} \quad (1)$$

This **depends** on our loss function, so we're **stuck** with that one.

Next, within each layer, we have

- The **activation function** - between our activation \mathbf{a} and preactivation \mathbf{z} :

$$\frac{\partial \mathbf{a}^\ell}{\partial \mathbf{z}^\ell} \quad (2)$$

What does the function between these **look** like?

$$\mathbf{a} = f(\mathbf{z}) \quad (3)$$

Well, that's not super interesting: we **don't know** our function. But, at least we can **write** it using f : that way, we know that this term only depends on our **activation function**.

$$\frac{\partial \mathbf{a}^\ell}{\partial \mathbf{z}^\ell} = \left(\overbrace{f^\ell}^{\text{func for layer } \ell} \right)'(\mathbf{z}^\ell) \quad (4)$$

This expression is a bit visually clunky, but it works.

Between layers, we have

- We can also think about the derivative of the **linear function** that **connects two layers**:

$$\frac{\partial \mathbf{z}^\ell}{\partial \mathbf{a}^{\ell-1}} \quad (5)$$

So, we want the function of these two:

$$\mathbf{z}^\ell = \mathbf{w}^\ell \mathbf{a}^{\ell-1} + \mathbf{w}_0^\ell \quad (6)$$

This one is pretty simple! We just take the derivative manually:

Be careful not to get this mixed up with the last one!
They look similar, but one is within the layer, and the other is between layers.

$$\frac{\partial z^\ell}{\partial a^{\ell-1}} = w^\ell \quad (7)$$

Finally, every gradient will end with

- The derivative that directly connects to a **weight**, again using the **linear function**:

$$\frac{\partial z^\ell}{\partial w^\ell} \quad (8)$$

The linear function is the same:

$$z^\ell = w^\ell a^{\ell-1} + w_0^\ell \quad (9)$$

But with a different **variable**, the **derivative** comes out different:

$$\frac{\partial z^\ell}{\partial w^\ell} = a^{\ell-1} \quad (10)$$

Notation 1

Our **derivatives** for the **chain rule** in a **1-D neural network** take the form:

$$\frac{\partial \mathcal{L}}{\partial a^L} \quad (11)$$

$$\frac{\partial a^\ell}{\partial z^\ell} = (f^\ell)'(z^\ell) \quad (12)$$

$$\frac{\partial z^\ell}{\partial a^{\ell-1}} = w^\ell \quad (13)$$

$$\frac{\partial z^\ell}{\partial w^\ell} = a^{\ell-1} \quad (14)$$

Now, we can rewrite our generalized expression for gradient:

$$\frac{\partial \mathcal{L}}{\partial w^\ell} = \overbrace{\left(\frac{\partial \mathcal{L}}{\partial a^L} \right)}^{\text{Loss unit}} \cdot \overbrace{\left((f^L)'(z^L) \cdot w^L \right)}^{\text{Layer L}} \cdot \overbrace{\left((f^{L-1})'(z^{L-1}) \cdot w^L \right)}^{\text{Layer L-1}} \cdot \left(\dots \right) \cdot \overbrace{\left((f^\ell)'(z^\ell) \cdot a^{\ell-1} \right)}^{\text{Layer } \ell} \quad (15)$$

Our expressions are more concrete now. It's still pretty visually messy, though.