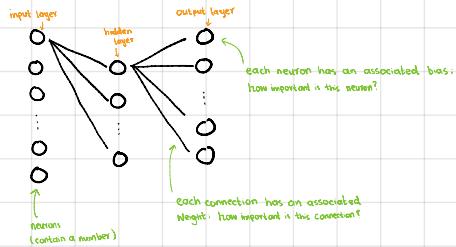


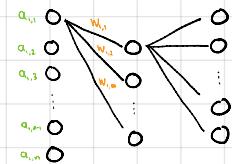

Introduction

A network of neurons; analogous to the brain.

Neurons and Layers



Forward Propagation



The activation of each neuron is given by:

$$1. z = \sum_{i=1}^n W_i a_i + b$$

sum of weighted activations of previous layer individual bias

$$= \begin{bmatrix} W_{11} & W_{12} & \dots & W_{1n} \\ W_{21} & W_{22} & \dots & W_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ W_{m1} & W_{m2} & \dots & W_{mn} \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}$$

outputs $a' = \begin{bmatrix} W_{11}a_1 + W_{12}a_2 + \dots + W_{1n}a_n \\ W_{21}a_1 + W_{22}a_2 + \dots + W_{2n}a_n \\ \vdots \\ W_{m1}a_1 + W_{m2}a_2 + \dots + W_{mn}a_n \end{bmatrix}$

$$2. z' = \text{act}(z')$$

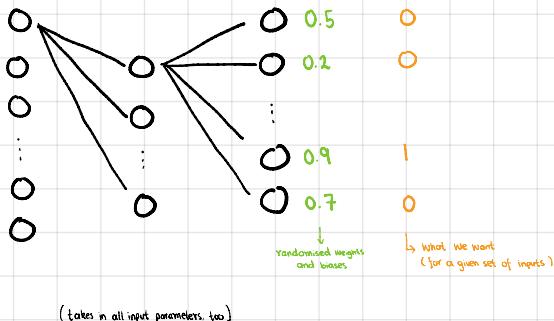
perform some activation function to turn z' into an activation. E.g. ReLU = $\begin{cases} x & x > 0 \\ 0 & \text{otherwise} \end{cases}$

$$3. a = \sigma(z')$$

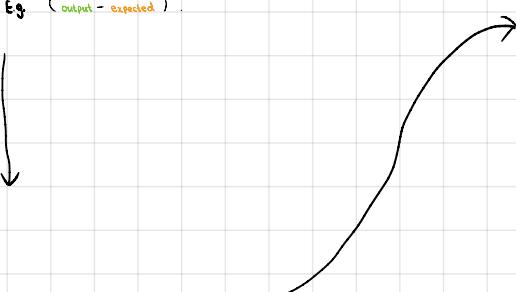
squish down z' into $[0, 1]$. E.g. $\sigma(z) = \frac{1}{1+e^{-z}}$

Training the Network: Cost Function

We train a neural network until the weights and biases are calibrated such that given input activations, the right output neurons will fire. We start by randomising.



A cost function describes how off the current weights and biases are by how off the final output is. E.g. $(\text{output} - \text{expected})^2$.



Gradient Descent

It becomes intuitive that the process of a neural network learning is really just minimising this cost function (by adjusting weights and biases).

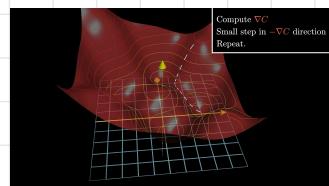
How does one do this? Consider the graph of this multivariable function. If we start at any random point, then constantly step in a downward direction, we will eventually reach a local minimum.

Note how these activation functions and squashing functions make the network a function, not just a linear computation.

Backpropagation

Backpropagation is an algorithm for computing that negative gradient. Recall that $-\nabla C$ is an n-dimensional vector that tells us how to adjust all the weights and biases to decrease the cost most efficiently. [...]

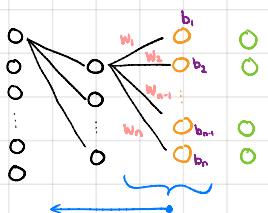
What's the "cost" of this difference?	$(0.43 - 0.00)^2 + (0.28 - 0.00)^2 + (0.19 - 0.00)^2 + (0.88 - 1.00)^2 + (0.72 - 0.00)^2 + (0.01 - 0.00)^2 + (0.64 - 0.00)^2 + (0.86 - 0.00)^2 + (0.99 - 0.00)^2 + (0.63 - 0.00)^2$
Cost of 3	Uter trash



From calculus III, this 'downward direction' is given by the negative of the gradient of the cost function ($-\nabla C$). Further, the magnitude of this gradient vector tells us how much to proportionally step in each direction: as we get closer to the local minimum, this value becomes smaller and smaller so as to not overshoot. We can further control each step through a scalar factor — the learning rate.

We compute the value of the gradient vector (at each step) through a process known as backpropagation.

Backpropagation cont.



First recall that the **cost** is directly related to the output layer's activation. This, itself, is influenced by each neuron's **individual bias** and each connection **weight**. This applies **iteratively** for each layer.

So, there are three (3) avenues that influence an output neuron's activation (and hence the cost):

1. Change the **bias**.
2. Change the **weights**.
3. Change the **activations** from the previous layer.

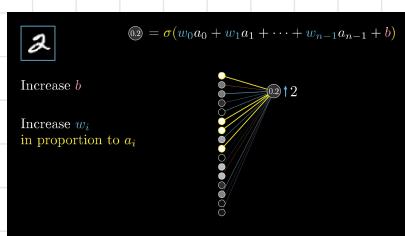
Changing Bias

The effect of a change to the bias on the weighted sum is constant and predictable. Increase if output < desired activation and vice versa.

Changing the Weights

How should we change the weights? First note that changing the weights associated with large activation values will have a stronger effect than changing the weights associated with small activation values.

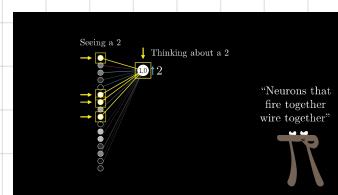
For $-\nabla C$, we don't just care about what should be moved up or down, but also by how much. (what gives you most bang for your buck).



This is somewhat similar to biological neural networks under Hebbian theory: "neurons that fire together wire together". Here, the biggest changes to weights — the biggest strengthening of connections — happens between neurons that are the most active and the ones which wish to become more active.

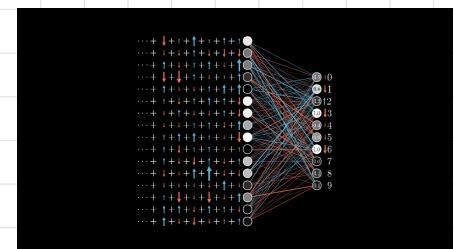
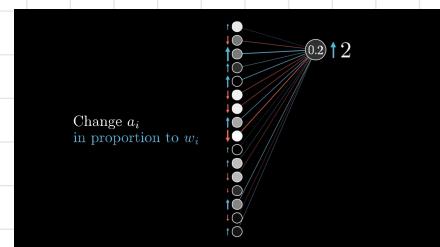
(Strong activations \Rightarrow more strong activations).

Pictorially:



Changing the Activations

We cannot change the activations in a previous layer directly. So, we change the weights of the preceding layer (and biases of the previous layer). This will likely cause competing requests. Again, these changes are proportional to the weights.



Here is where the idea of propagating backwards comes in. By adding all these desired effects, we can get a list of the nudges we want to happen to this 2nd last layer. From there, we recursively apply this process to the relevant weights and biases determining those values, repeating as you move backward through the network.

Repeating for All Training Examples

We take an average of all versions/suggestions of weights and biases.

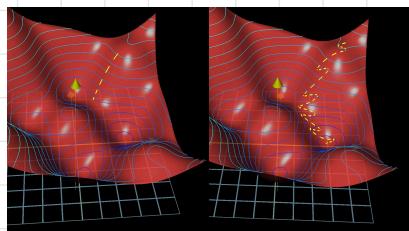
	2	5	0	4	7	9	...	Average over all training data
w_0	-0.08	+0.02	-0.07	+0.11	-0.05	-0.14	...	-0.08
w_1	-0.11	+0.11	+0.07	+0.02	+0.09	+0.05	...	+0.02
w_2	-0.07	-0.04	-0.01	+0.02	+0.13	-0.15	...	-0.06
w_3
$w_{13,000}$	+0.13	+0.08	-0.06	-0.09	-0.02	+0.04	...	+0.04

Stochastic Gradient Descent

It is too exhaustive to train every sample, so we train in 'mini-batches'.

"Mini-batches"											
5	0	4	1	9	2	1	3	1	4	3	5
3	6	1	7	2	8	6	9	4	0	9	1
1	2	4	3	2	7	3	8	6	9	0	5
6	0	7	6	1	8	7	9	3	9	8	5

Then we treat the average weight/bias suggestion of each mini-batch as one training sample, and find the average between these averages.



It looks like larger strides (rather than steps) down the graph.

Backpropagation Calculus

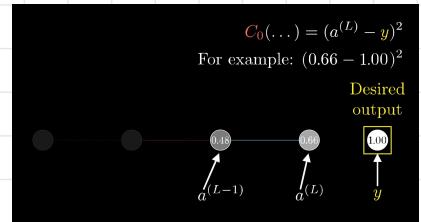
The cost function is a function whose variables are the weights and biases of the neural network. As such, its gradient can be given as

$$\nabla C = \begin{bmatrix} \frac{\partial C}{\partial w_1} \\ \frac{\partial C}{\partial b_1} \\ \vdots \\ \frac{\partial C}{\partial w_n} \\ \frac{\partial C}{\partial b_n} \end{bmatrix}.$$

The goal with this is to find the relationship between a change to C and the necessary change to a given weight and bias to incur that change.

Simple Example

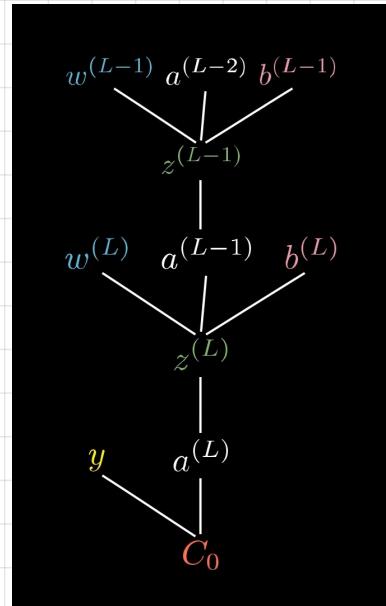
Let's begin by considering a simple network.



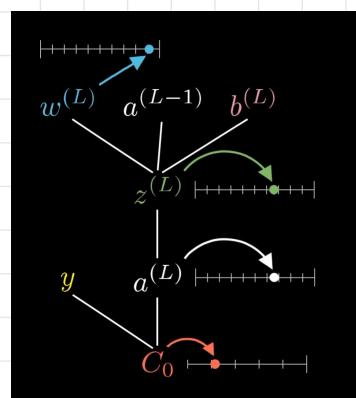
Recall that the last activation ($a^{(L)}$) is determined by a weight, a bias, and the activation of the previous neuron, all pumped through some special non-linear activation function.

$$\begin{aligned} a^{(L)} &= \sigma(w^{(L)} a^{(L-1)} + b^{(L)}) \\ &= \sigma(z^{(L)}), \end{aligned}$$

Where z denotes the weighted sum of the L^{th} layer.
Let's visually keep track of this.



Then to adjust each weight and bias to incur a change in C_0 , we can use the chain rule. (Rate of change with respect to hidden/embedded variables)



Backpropagation Calculus cont.

But this is just one entry in ∇C .

So, we have that:

$$\frac{\partial C_0}{\partial w^{(l)}} = \frac{\partial C_0}{\partial a^{(l)}} \times \frac{\partial a^{(l)}}{\partial z^{(l)}} \times \frac{\partial z^{(l)}}{\partial w^{(l)}}$$

Practically:

$$z^{(l)} = w^{(l)} a^{(l-1)} + b^{(l)}$$

$$\frac{\partial z^{(l)}}{\partial w^{(l)}} = a^{(l-1)}$$

$$a^{(l)} = \sigma(z^{(l)})$$

$$\frac{\partial a^{(l)}}{\partial z^{(l)}} = \sigma'(z^{(l)})$$

$$C_0 = (a^{(l)} - y)^2$$

$$\frac{\partial C_0}{\partial a^{(l)}} = 2(a^{(l)} - y)$$

We can apply this entire process similarly for the bias and find that $\frac{\partial z^{(l)}}{\partial b^{(l)}} = 1$.

Now, note that for these functions, the first derivative, $\frac{\partial z^{(l)}}{\partial w^{(l)}} = a^{(l-1)}$, tells us that changing the weight has a

stronger effect on $z^{(l)}$ (and therefore on the cost) when the previous neuron is more active. This is where that biological idea that "neurons that fire together, wire together" is mirrored in our math. (Strong activations \Rightarrow more strong activations).

Further, the third derivative is directly dependent on the loss $a^{(l)} - y$. This reflects the intuitive idea that a more inaccurate set of weights and biases (i.e. greater cost) would incur greater adjustments to the weights and biases.

For All Training Data

So far, this has only been for one training sample. The full cost function for the network is the average all the individual costs for each training example:

$$C = \frac{1}{n} \sum_{k=0}^{n-1} C_k$$

If we want the derivative of C with respect to the weight, we need to take the average of all the individual derivatives.

$$\frac{\partial C}{\partial w^{(l)}} = \frac{1}{n} \sum_{k=0}^{n-1} \frac{\partial C_k}{\partial w^{(l)}}$$

The Derivatives for All Weights And Biases

The full gradient is given as

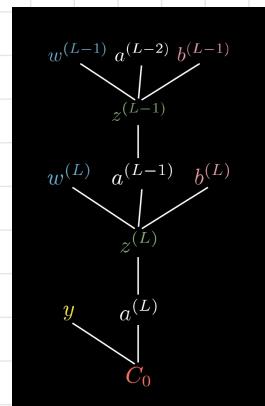
$$\nabla C = \begin{bmatrix} \frac{\partial C}{\partial w_1} \\ \frac{\partial C}{\partial b_1} \\ \vdots \\ \frac{\partial C}{\partial w_n} \\ \frac{\partial C}{\partial b_n} \end{bmatrix}$$

Now, we need to account for the rest of the weights and biases in the network; the 'backward' part. We first think literally:

$$\frac{\partial C_0}{\partial a^{(l-1)}} = \frac{\partial z^{(l)}}{\partial a^{(l-1)}} \times \frac{\partial a^{(l)}}{\partial z^{(l)}} \times \frac{\partial C_0}{\partial a^{(l)}}$$

We can compute $\frac{\partial z^{(l)}}{\partial a^{(l-1)}}$ to be $= w^{(l)}$, but this just means changing the activation in $L-1$ will incur a change in $z^{(l)}$ proportional to $w^{(l)}$. This is somewhat meaningless because we don't care about changing $a^{(l-1)}$ directly.

So, we treat $a^{(l-1)}$ as $a^{(l)}$ and iterate the previously seen chain rule.



Here,

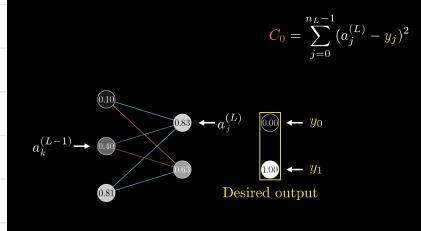
$$\frac{\partial C_0}{\partial w^{(l-1)}} = \frac{\partial C_0}{\partial a^{(l)}} \times \frac{\partial a^{(l)}}{\partial z^{(l)}} \times \frac{\partial z^{(l)}}{\partial a^{(l-1)}} \times \frac{\partial a^{(l-1)}}{\partial z^{(l-1)}} \times \frac{\partial z^{(l-1)}}{\partial w^{(l-1)}}$$

Note, the further back we go, the less we affect the cost.

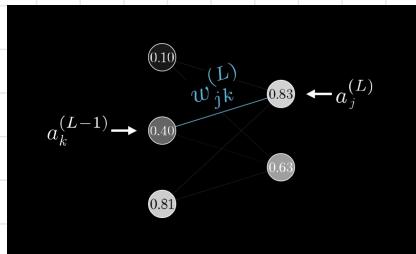
More Complicated Networks

With multiple neurons in each layer, we will need to keep track of the neurons in each layer more closely.

Let's use k to index neurons of layer $(L-1)$ and j to index the layer L . Then we have the cost as a sum:



Additionally, each weight must be made more specific:



Note: we say jk as opposed to kj to align with the weight matrix.

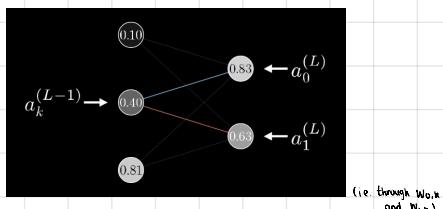
So now we have $\tilde{z}_j^{(L)} = \sigma(w_{j0}^{(L)} a_0^{(L-1)} + \dots + w_{jk}^{(L)} a_k^{(L-1)})$.

$$\text{Then: } \frac{\partial C_0}{\partial w_{jk}^{(L)}} = \frac{\partial C_0}{\partial a_j^{(L)}} \times \frac{\partial a_j^{(L)}}{\partial \tilde{z}_j^{(L)}} \times \frac{\partial \tilde{z}_j^{(L)}}{\partial w_{jk}^{(L)}}$$

Note that these new equations are essentially the same — we just became more specific about which weight (and bias), i.e. letters to index, to account for the many weights.

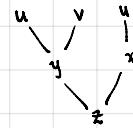
What does change, however, is the derivative of the cost with respect to one of the activations in the layer $(L-1)$, $\frac{\partial C_0}{\partial a_k^{(L-1)}}$.

This neuron now influences the cost through multiple paths.



The total influence of $a_k^{(L-1)}$ is then the sum of the derivatives along multiple paths of influence.

One way to think of it is now, $a_k^{(L-1)}$ is a shared variable:

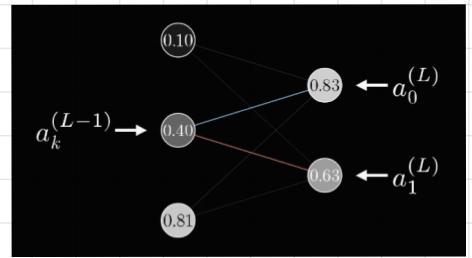


$$\frac{\partial z}{\partial u} = \frac{\partial z}{\partial y} \times \frac{\partial y}{\partial u} + \frac{\partial z}{\partial w} \times \frac{\partial w}{\partial u}$$

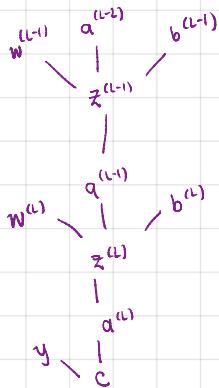
So we finally have that

$$\frac{\partial C_0}{\partial a_k^{(L-1)}} = \sum_{j=0}^{n_L-1} \frac{\partial C_0}{\partial a_j^{(L)}} \times \frac{\partial a_j^{(L)}}{\partial \tilde{z}_j^{(L)}} \times \frac{\partial \tilde{z}_j^{(L)}}{\partial a_k^{(L-1)}}$$

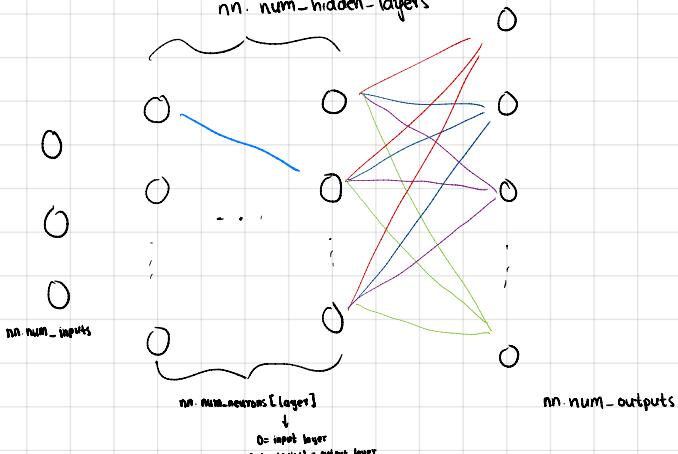
↓
each path (that ends at $a_0^{(L)}, a_1^{(L)}, \dots, a_{n_L-1}^{(L)}$)



Practise Implementation



nn.num-hidden-layers



$$\sigma = \frac{1}{1+e^{-x}} = (1+e^{-x})^{-1}$$

$$\sigma'(x) = -(1+e^{-x})^{-2} (-e^{-x})$$

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