

NN Gradients

Archim Jhunjunwala

March 27, 2024

1 Goal

The goal is to efficiently compute the gradient of the cost function - C with respect to the weights and biases of a neural networks - $W_{i,j}^l$ denoting the weight from node j in layer $l - 1$ to node i in layer l . This piece concerns only feed-forward neural networks. Each layer has an activation function $f_l(x)$, as well as the activated a_i^l and unactivated u_i^l values associated with each node i .

2 Single Weights / Biases

Ignoring matrices, let's find $\frac{\partial C}{\partial b_i^l}$, and $\frac{\partial C}{\partial W_{i,j}^l}$ for some layer l in our network of L layers.

2.1 Bias Gradient

$$\begin{aligned}\frac{\partial C}{\partial b_i^l} &= \frac{\partial C}{\partial u_i^l} \frac{\partial u_i^l}{\partial b_i^l} = \frac{\partial C}{\partial a_i^l} \frac{\partial a_i^l}{\partial u_i^l} \frac{\partial u_i^l}{\partial b_i^l} \text{ by the chain rule} \\ \frac{\partial a_i^l}{\partial u_i^l} &= \frac{\partial}{\partial u_i^l} f(u_i^l) = f'(u_i^l)\end{aligned}$$

$$\begin{aligned}\frac{\partial u_i^l}{\partial b_i^l} &= \frac{\partial}{\partial b_i^l}(b_i^l + \dots) = 1 \\ \frac{\partial C}{\partial b_i^l} &= \frac{\partial C}{\partial a_i^l} f'(u_i^l)\end{aligned}$$

2.2 Weight Gradient

$$\begin{aligned}\frac{\partial C}{\partial W_{i,j}^l} &= \frac{\partial C}{\partial u_i^l} \frac{\partial u_i^l}{\partial W_{i,j}^l} = \frac{\partial C}{\partial a_i^l} \frac{\partial a_i^l}{\partial u_i^l} \frac{\partial u_i^l}{\partial W_{i,j}^l} \text{ by the chain rule} \\ \frac{\partial u_i^l}{\partial W_{i,j}^l} &= \frac{\partial}{\partial W_{i,j}^l}(W_{i,j}^l a_{l-1,j}^l + \dots) = a_j^{l-1} \\ \frac{\partial C}{\partial W_{i,j}^l} &= \frac{\partial C}{\partial a_i^l} f'(u_i^l) a_j^{l-1}\end{aligned}$$

2.3 Backpropagation

We must find 2 more values before we are done: $\frac{\partial C}{\partial a_i^l}$ (this value depends on the specific cost function used, and I will ignore it here). $\frac{\partial C}{\partial a_i^{l-1}}$ as an expression of $\frac{\partial C}{\partial a_i^l}$. In this way, we can propagate backwards through the network calculating all the necessary gradients.

$\frac{\partial C}{\partial a_i^{l-1}} = \sum_{j=1}^n \frac{\partial C}{\partial u_j^l} \frac{\partial u_j^l}{\partial a_i^{l-1}} = \sum_{j=1}^n \frac{\partial C}{\partial a_j^l} \frac{\partial a_j^l}{\partial u_j^l} \frac{\partial u_j^l}{\partial a_i^{l-1}}$ by the chain rule (n is the number of nodes in layer l)

$$\begin{aligned}\frac{\partial u_j^l}{\partial a_i^{l-1}} &= \frac{\partial}{\partial a_i^{l-1}}(W_{j,i}^l a_i^{l-1} + \dots) = W_{j,i}^l \\ \frac{\partial C}{\partial a_i^{l-1}} &= \sum_{j=1}^n \frac{\partial C}{\partial a_j^l} f'(u_j^l) W_{j,i}^l\end{aligned}$$

Observe that the bias gradient is present in both the weight gradient, and the derivative cost with respect to the previous layers nodes.

3 Results / Matrix Notation

$$\begin{aligned}\frac{\partial C}{\partial b^l} &= \begin{bmatrix} \frac{\partial C}{\partial a_1^l} f'(u_1^l) \\ \dots \\ \frac{\partial C}{\partial a_n^l} f'(u_n^l) \end{bmatrix} \\ \frac{\partial C}{\partial W_{i,j}^l} &= \frac{\partial C}{\partial b^l} (a^{l-1})^T \\ \frac{\partial C}{\partial a_i^{l-1}} &= \frac{\partial C}{\partial b^l}^T W^l \text{ (this will be a row vector)}\end{aligned}$$