

Back propagation for Deep Feedforward Neural Networks

Aaron Avram

June 10 2025

Introduction

In this write up I will go through the derivation of the back propagation algorithm for deep feed-forward neural networks and discuss how I will implement it using NumPy.

Set up

We are working with a classic neural net set up where we have $n + 1$ layers $\{x^{(0)}, \dots, x^{(n)}\}$ with the $i - th$ layer having l_i neurons. We then define a weight matrix and bias vector for each of the n non-input neurons. $W_{jk}^{(i)}$ is the weight of the k -th neuron in the $i - 1$ th on the j th neuron in the i th layer. And, $b_j^{(i)}$ is the bias associated with the j th neuron in the i th layer. Denote the activation function by σ (which could be tanh for example). Note that when σ is applied to a vector we do so element wise. Now we can represent the forward pass between each layer succinctly where

$$x^{(i)} = \sigma(W^{(i)}x^{(i-1)} + b^{(i)})$$

It will be convenient later to denote $z^{(i)} := W^{(i)}x^{(i-1)} + b^{(i)}$, so that $x^{(i)} = \sigma(z^{(i)})$. Note, that there may be some transformation on the output layer so that it represents a probability distribution (I.e. applying the softmax), but this is trivial in the context of back propagation calculations so it won't be included here.

Now when training this model on some input, we optimize its performance by minimizing some loss function, which measures some numerical representation of how well the model performs on the training data. Formally: $L : \mathbb{R}^{(l_n)} \rightarrow \mathbb{R}$. In the following calculations we will take gradients and derivatives of the loss function, and so we will assume that all of the mathematical functions involved are differentiable, which is practically always the case. Now, to train our model we want to find a method for computing the gradient of the weights and biases with respect to our loss function then perform gradient descent. Back propagation is such a method.

The Plan

The back propagation algorithm is essentially a fancy application of the chain rule from calculus. Namely the chain rule states that for a composition of functions f, g, h we have:

$$\frac{df(g(h))}{dx} = \frac{df}{dg} \frac{dg}{dh} \frac{dh}{dx}$$

In the multivariable case we apply this to partial derivatives for each component of input and output respectively.

Clearly our neural network is just a very large and complicated function and to get the gradients of the weights and biases, we can apply the above concept. This looks similar to dynamic programming where we are using memoization by caching the gradients of the preceeding layers as we differentiate through the network.

Back Prop

Now, we will build up our gradients recursively working backwards from the output layer.

Since our loss function has scalar output, the gradient of the output layer with respect to the loss is straight forward:

$$\nabla_{x^{(n)}}(L) = \begin{bmatrix} \frac{\partial L}{\partial x_1^{(n)}} \\ \vdots \\ \frac{\partial L}{\partial x_{l_n}^{(n)}} \end{bmatrix}$$

Then for any layer $x^{(i)}$ we can consider the Jacobian of $z^{(i)}$ with respect to the output. However since $z_j^{(i)}$ only effects $x_j^{(i)}$ the Jacobian will be diagonal, so we can just write it as a gradient vector (abusing notation slightly).

$$\nabla_{z^{(i)}}(x^{(i)}) = \begin{bmatrix} \frac{\partial x_1^{(i)}}{\partial z_1^{(i)}} \\ \vdots \\ \frac{\partial x_{l_i}^{(i)}}{\partial z_{l_i}^{(i)}} \end{bmatrix}$$

Additionally we know that $x^{(i)} = \sigma(z^{(i)})$, so $\nabla_{z^{(i)}}(x^{(i)}) = \sigma'(z^{(i)})$. Next we want to consider the Jacobian of $x^{(i-1)}$ with respect to $z^{(i)}$, but notice that $z^{(i)} = W^{(i)}x^{(i-1)} + b^{(i)}$ so $\frac{\partial z_j^{(i)}}{\partial x_k^{(i-1)}} = W_{jk}^{(i)}$. Thus:

$$\begin{aligned} J(z^{(i)})(x^{(i-1)}) &= \begin{bmatrix} \frac{\partial z_1^{(i)}}{\partial x_1^{(i-1)}} & \cdots & \frac{\partial z_{l_i}^{(i)}}{\partial x_1^{(i-1)}} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_1^{(i)}}{\partial x_{l_{i-1}}^{(i-1)}} & \cdots & \frac{\partial z_{l_i}^{(i)}}{\partial x_{l_{i-1}}^{(i-1)}} \end{bmatrix} \\ &= (W^{(i)})^T \end{aligned}$$

Abusing notation again we will associate this with a gradient vector $\nabla_{x^{(i-1)}}(z^{(i)})$. Now we can combine the following results to find the Jacobian of $x^{(i-1)}$ with respect to $x^{(i)}$. Observe that $\frac{\partial x_j^{(i)}}{\partial x_k^{(i-1)}} = \frac{\partial x_j^{(i)}}{\partial z_j^{(i)}} \frac{\partial z_j^{(i)}}{\partial x_k^{(i-1)}}$ and we know both of these values. Thus, The Jacobian of $x^{(i-1)}$ with respect to $x^{(i)}$ is just the Jacobian of $x^{(i-1)}$ with respect to $z^{(i)}$ with each row scaled by $\nabla_{z^{(i)}}(x^{(i)})$. So:

$$\nabla_{x^{(i-1)}}(x^{(i)}) = \nabla_{x^{(i-1)}}(z^{(i)}) \odot \nabla_{z^{(i)}}(x^{(i)})^T$$

Note we assume that the row vector in the equation is broadcasted by copying the rows on top of each other.

With this we have a way of recursively calculating the gradient of each layer with respect to the loss function. The final step is to now compute the gradients with respect to the weights and biases of each layer and store them.

First consider $z^{(i)}$ and $W^{(i)}$, note that we can compute the Jacobian of the flattened version of $W^{(i)}$ with respect to $z^{(i)}$ but if we observe that the j th component of $z^{(i)}$ is only affected by the j th row of $W^{(i)}$ we can construct our jacobian matrix as containing the partial derivatives of each component of the j th row of $W^{(i)}$ with respect to the j th component of $z^{(i)}$ as its columns. Now since $z^{(i)} = W^{(i)}x^{(i-1)} + b^{(i)}$, $\frac{\partial z_j^{(i)}}{\partial W_{jk}^{(i)}} = x_k^{(i-1)}$. Thus:

$$\begin{aligned} J(z^{(i)})(W^{(i)}) &= \begin{bmatrix} \frac{\partial z_1^{(i)}}{\partial W_{11}^{(i)}} & \cdots & \frac{\partial z_{l_i}^{(i)}}{\partial W_{l_i 1}^{(i)}} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_1^{(i)}}{\partial W_{1 l_{i-1}}^{(i)}} & \cdots & \frac{\partial z_{l_i}^{(i)}}{\partial W_{l_i l_{i-1}}^{(i)}} \end{bmatrix} \\ &= \begin{bmatrix} x_1^{(i-1)} & \cdots & x_{l_{i-1}}^{(i-1)} \\ \vdots & \ddots & \vdots \\ x_1^{(i-1)} & \cdots & x_{l_{i-1}}^{(i-1)} \end{bmatrix} \end{aligned}$$

Thus we can succinctly write the combined gradient as a product:

$$\begin{aligned} \nabla_{W^{(i)}}(x^{(i)}) &= J(z^{(i)})(W^{(i)}) \odot \nabla_{z^{(i)}}(x^{(i)}) \\ &= x^{(i-1)} \nabla_{z^{(i)}}(x^{(i)})^T \end{aligned}$$

Next for the $b^{(i)}$ term since it is just a constant term and only its j th component affects the j -th component of $z^{(i)}$, $\nabla_{b^{(i)}}(z^{(i)}) = 1_{l_i}$ and $\nabla_{b^{(i)}}(x^{(i)}) = \nabla_{z^{(i)}}(x^{(i)})$. With this, we have fully detailed the recursive algorithm steps.

Implementation

Since I am building this model in pure NumPy and Python, I want to leverage NumPy BLAS calls as much as possible, thus I will opt for storing the necessary gradients as arrays and using vectorized operations to compute each back prop step. Thus, I won't pass in individual examples during training but instead I will pass in an array of them. The backprop step will look similar to what has been derived above, however I will store a matrix for each layer gradient. The weight and bias gradients will be the summed across training examples which happens naturally with matrix multiplication of the layer arrays. I.e. notice that, if the input batch has m training examples with x now having its i th row being the i th training example. With this setup there is an elegant way to express the jacobian of the loss with respect to the last unactivated layer if the output activation is a softmax, and the loss is cross entropy:

$$J(L)(z^{(n)}) = x^{(n)} - y_{\text{onehot}}$$

And we can express the gradients of the weights by considering the same construction as in the single example case, but with the gradients per example summed. Thus, for $i < n$:

$$\begin{aligned} J(L)(W^{(i)})_{jk} &= \sum_{p=1}^m \frac{L}{\partial W_{kj}^{(i)}} \\ &= \sum_{p=1}^m \frac{L}{\partial z_{pk}^{(i)}} \frac{\partial z_{pk}^{(i)}}{\partial W_{kj}^{(i)}} \end{aligned}$$

Plugging in values:

$$J(L)(W^{(i)})_{jk} = \sum_{p=1}^m \frac{L}{\partial z_{pk}^{(i)}} x_{pj}^{(i-1)}$$

If we let G be the m, l_i matrix of the gradients of each unactivated neuron in the i th layer we can write this as a matrix multiplication: $J(L)(W^{(i)}) = G^T X^{i-1}$.

Then, $J(L)(B^{(i)})$ is just a sum across training examples for the previous unactivated layer grads.

Next to find $J(L)(X^{(i-1)})$ consider:

$$\begin{aligned} J(L)(X^{(i-1)})_{jk} &= \frac{\partial L}{\partial x_j^{(i-1)} k} \\ &= \sum_{p=0}^{l_i} \frac{\partial L}{\partial z_{jp}^{(i)}} \frac{\partial z_{jp}^{(i)}}{x_{jk}^{(i-1)}} \\ &= \sum_{p=0}^{l_i} \frac{\partial L}{\partial z_{jp}^{(i)}} W_{pk}^{(i)} \end{aligned}$$

If we let the matrix G now represent the m, l_i matrix containing the gradients of the unactivated i th we have:

$$J(L)(X^{(i-1)}) = G W^{(i)}$$

Finally, there is one important note, the data I will be using for this specific implementation is name data, so the model will be given a sequence of characters and must predict the next in the form a probability distribution over the possible next characters. Thus, a popular technique, as seen in the Bengio et. al paper at the word level, is to embed the individual tokens, characters for us, into a higher dimensional feature space with an embedding matrix. Then these embeddings are trained with the weights and biases of the network, so that characters that are similar have similar embeddings. The embedding allows us to capture more fine properties of each character by tuning the dimension of the feature space. Thus, our $X^{(0)}$ is actually the flattened embedded inputs. Thus, the gradient $X^{(0)}$ with respect to the loss captures the gradients of certain embedding rows with respect to the loss, so we use this data to update the embedding during gradient descent.