

The Mathematics of Dense Neural Networks

Trevor Dick

January 2021

1 Notation

Let L be the number of layers, not including the input layer, which will be designated layer 0. Let $g(z)$ be an activation function for some node in the network. The input layer will accept a vector $\vec{x} \in \mathbb{R}^n$ and the target will be a vector $\vec{y} \in \mathbb{R}^m$. Keep in mind that a dense net is designed to accept one sample at a time. The process of training a dense net can be thought of in time steps, and these time steps are coming from the set of samples in the training set. For example in batch training, each batch uses a subset of the training samples, and if the batch size is not one, then the number of batches is necessarily less than the number of total samples. In this example we have one set of time steps corresponding to each batch and another set of time steps for the samples used within a given batch.

2 Constructing the Network

For each layer $\ell \in \{0, 1, \dots, L\}$, we have the number of nodes per layer $n_\ell \in \{n_0, n_1, \dots, n_L\}$. We also have activation functions $g^\ell \in \{g^1, \dots, g^L\}$ for each layer excluding the input layer. A single node in one layer is connected to all the nodes in the previous layer by a weight matrix $W_{[n_\ell \times n_{\ell-1}]}^\ell$ and each node may add a bias $b_{\nu_\ell}^\ell$ before the activation g^ℓ is applied. In total the network may be expressed by two indices, the layer ℓ and node in that layer ν_ℓ . For each pair of indices (ℓ, ν_ℓ) we can define the following functions that describe the values any node in the network may have:

$$z_{\nu_\ell}^\ell = \left(\sum_{\nu_{\ell-1}=1}^{n_{\ell-1}} w_{\nu_\ell \nu_{\ell-1}}^\ell a_{\nu_{\ell-1}}^{\ell-1} \right) + b_{\nu_\ell}^\ell$$
$$a_{\nu_\ell}^\ell = g^\ell(z_{\nu_\ell}^\ell)$$

In vector notation the equations become the following:

$$\vec{z}^\ell = W^\ell \cdot \overrightarrow{a^{\ell-1}} + \vec{b}^\ell$$

$$\overrightarrow{a}^\ell = g^\ell(\vec{z}^\ell)$$

3 Derivatives w.r.t W and \vec{b}

Consider the final activation layer, $\overrightarrow{a}^L = \overrightarrow{y_{pred}}$, since it is used to compare against the target \vec{y} . We will consider the index representation for the following calculations.

Differentiate any node in the final layer L , with respect to the p, q -th weights in any layer ℓ . Note: δ_{ij} is the Kronecker delta function.

$$\begin{aligned}
\frac{\partial a_{\nu_L}^L}{\partial w_{pq}^\ell} &= \frac{\partial g_{\nu_L}^L}{\partial z_{\nu_L}^L} \frac{\partial z_{\nu_L}^L}{\partial w_{pq}^\ell} \\
&= \frac{\partial g_{\nu_L}^L}{\partial z_{\nu_L}^L} \frac{\partial}{\partial w_{pq}^\ell} \left[\left(\sum_{\nu_{L-1}=1}^{n_{L-1}} w_{\nu_L \nu_{L-1}}^L a_{\nu_{L-1}}^{L-1} \right) + b_{\nu_L}^L \right] \\
&= \frac{\partial g_{\nu_L}^L}{\partial z_{\nu_L}^L} \sum_{\nu_{L-1}=1}^{n_{L-1}} w_{\nu_L \nu_{L-1}}^L \frac{\partial a_{\nu_{L-1}}^{L-1}}{\partial w_{pq}^\ell} \\
&= \frac{\partial g_{\nu_L}^L}{\partial z_{\nu_L}^L} \sum_{\nu_{L-1}=1}^{n_{L-1}} w_{\nu_L \nu_{L-1}}^L \frac{\partial g_{\nu_{L-1}}^{L-1}}{\partial z_{\nu_{L-1}}^{L-1}} \frac{\partial z_{\nu_{L-1}}^{L-1}}{\partial w_{pq}^\ell} \\
&= \frac{\partial g_{\nu_L}^L}{\partial z_{\nu_L}^L} \sum_{\nu_{L-1}=1}^{n_{L-1}} w_{\nu_L \nu_{L-1}}^L \frac{\partial g_{\nu_{L-1}}^{L-1}}{\partial z_{\nu_{L-1}}^{L-1}} \frac{\partial}{\partial w_{pq}^\ell} \left[\left(\sum_{\nu_{L-2}=1}^{n_{L-2}} w_{\nu_{L-1} \nu_{L-2}}^{L-1} a_{\nu_{L-2}}^{L-2} \right) + b_{\nu_{L-1}}^{L-1} \right] \\
&= \frac{\partial g_{\nu_L}^L}{\partial z_{\nu_L}^L} \sum_{\nu_{L-1}=1}^{n_{L-1}} w_{\nu_L \nu_{L-1}}^L \frac{\partial g_{\nu_{L-1}}^{L-1}}{\partial z_{\nu_{L-1}}^{L-1}} \sum_{\nu_{L-2}=1}^{n_{L-2}} w_{\nu_{L-1} \nu_{L-2}}^{L-1} \frac{\partial a_{\nu_{L-2}}^{L-2}}{\partial w_{pq}^\ell} \\
&= \vdots \\
&= \frac{\partial g_{\nu_L}^L}{\partial z_{\nu_L}^L} \sum_{\nu_{L-1}=1}^{n_{L-1}} w_{\nu_L \nu_{L-1}}^L \frac{\partial g_{\nu_{L-1}}^{L-1}}{\partial z_{\nu_{L-1}}^{L-1}} \sum_{\nu_{L-2}=1}^{n_{L-2}} w_{\nu_{L-1} \nu_{L-2}}^{L-1} \frac{\partial g_{\nu_{L-2}}^{L-2}}{\partial z_{\nu_{L-2}}^{L-2}} \cdots \frac{\partial g_{\nu_{\ell+1}}^{\ell+1}}{\partial z_{\nu_{\ell+1}}^{\ell+1}} \sum_{\nu_\ell=1}^{n_\ell} w_{\nu_{\ell+1} \nu_\ell}^{\ell+1} \frac{\partial a_{\nu_\ell}^\ell}{\partial w_{pq}^\ell} \\
&= \frac{\partial g_{\nu_L}^L}{\partial z_{\nu_L}^L} \sum_{\nu_{L-1}=1}^{n_{L-1}} w_{\nu_L \nu_{L-1}}^L \frac{\partial g_{\nu_{L-1}}^{L-1}}{\partial z_{\nu_{L-1}}^{L-1}} \sum_{\nu_{L-2}=1}^{n_{L-2}} w_{\nu_{L-1} \nu_{L-2}}^{L-1} \frac{\partial g_{\nu_{L-2}}^{L-2}}{\partial z_{\nu_{L-2}}^{L-2}} \cdots \frac{\partial g_{\nu_{\ell+1}}^{\ell+1}}{\partial z_{\nu_{\ell+1}}^{\ell+1}} \sum_{\nu_\ell=1}^{n_\ell} w_{\nu_{\ell+1} \nu_\ell}^{\ell+1} \frac{\partial g_{\nu_\ell}^\ell}{\partial z_{\nu_\ell}^\ell} \frac{\partial z_{\nu_\ell}^\ell}{\partial w_{pq}^\ell} \\
&= \frac{\partial g_{\nu_L}^L}{\partial z_{\nu_L}^L} \sum_{\nu_{L-1}=1}^{n_{L-1}} w_{\nu_L \nu_{L-1}}^L \frac{\partial g_{\nu_{L-1}}^{L-1}}{\partial z_{\nu_{L-1}}^{L-1}} \sum_{\nu_{L-2}=1}^{n_{L-2}} w_{\nu_{L-1} \nu_{L-2}}^{L-1} \frac{\partial g_{\nu_{L-2}}^{L-2}}{\partial z_{\nu_{L-2}}^{L-2}} \cdots \frac{\partial g_{\nu_{\ell+1}}^{\ell+1}}{\partial z_{\nu_{\ell+1}}^{\ell+1}} \sum_{\nu_\ell=1}^{n_\ell} w_{\nu_{\ell+1} \nu_\ell}^{\ell+1} \frac{\partial g_{\nu_\ell}^\ell}{\partial z_{\nu_\ell}^\ell} \sum_{\nu_{\ell-1}=1}^{n_{\ell-1}} \frac{\partial w_{\nu_\ell \nu_{\ell-1}}^\ell}{\partial w_{pq}^\ell} a_{\nu_{\ell-1}}^{\ell-1} \\
&= \frac{\partial g_{\nu_L}^L}{\partial z_{\nu_L}^L} \sum_{\nu_{L-1}=1}^{n_{L-1}} w_{\nu_L \nu_{L-1}}^L \frac{\partial g_{\nu_{L-1}}^{L-1}}{\partial z_{\nu_{L-1}}^{L-1}} \sum_{\nu_{L-2}=1}^{n_{L-2}} w_{\nu_{L-1} \nu_{L-2}}^{L-1} \frac{\partial g_{\nu_{L-2}}^{L-2}}{\partial z_{\nu_{L-2}}^{L-2}} \cdots \frac{\partial g_{\nu_{\ell+1}}^{\ell+1}}{\partial z_{\nu_{\ell+1}}^{\ell+1}} \sum_{\nu_\ell=1}^{n_\ell} w_{\nu_{\ell+1} \nu_\ell}^{\ell+1} \frac{\partial g_{\nu_\ell}^\ell}{\partial z_{\nu_\ell}^\ell} \sum_{\nu_{\ell-1}=1}^{n_{\ell-1}} \delta_{\nu_\ell p} \delta_{\nu_{\ell-1} q} a_{\nu_{\ell-1}}^{\ell-1}
\end{aligned}$$

The result is a $[n_L \times n_\ell \times n_{\ell-1}]$ array, arising from the free indices (ν_L, p, q) .

This shape matches our expectations that there should be enough information to encode changes of each node in the final layer w.r.t. changes of any element in the W^ℓ matrix. Somewhat more precisely, during training the loss function is aggregating over the index ν_L leaving a matrix with shape $[n_\ell \times n_{\ell-1}]$, the same as W^ℓ . This new matrix must be of the same shape as W^ℓ so that it can update W^ℓ during learning, typically through back propagation.

We can reuse much of the above work in deriving the equations for derivatives w.r.t b^ℓ . Since w_{pq}^ℓ is just another variable, we can easily replace it with b_p^ℓ , for every place where any function does not explicitly depend on it as we apply the chain rule.

$$\begin{aligned}
\frac{\partial a_{\nu_L}^L}{\partial b_p^\ell} &= \frac{\partial g_{\nu_L}^L}{\partial z_{\nu_L}^L} \frac{\partial z_{\nu_L}^L}{\partial b_p^\ell} \\
&= \vdots \\
&= \frac{\partial g_{\nu_L}^L}{\partial z_{\nu_L}^L} \sum_{\nu_{L-1}=1}^{n_{L-1}} w_{\nu_L \nu_{L-1}}^L \frac{\partial g_{\nu_{L-1}}^{L-1}}{\partial z_{\nu_{L-1}}^{L-1}} \sum_{\nu_{L-2}=1}^{n_{L-2}} w_{\nu_{L-1} \nu_{L-2}}^{L-1} \frac{\partial g_{\nu_{L-2}}^{L-2}}{\partial z_{\nu_{L-2}}^{L-2}} \cdots \frac{\partial g_{\nu_{\ell+1}}^{\ell+1}}{\partial z_{\nu_{\ell+1}}^{\ell+1}} \sum_{\nu_\ell=1}^{n_\ell} w_{\nu_{\ell+1} \nu_\ell}^{\ell+1} \frac{\partial g_{\nu_\ell}^\ell}{\partial z_{\nu_\ell}^\ell} \frac{\partial z_{\nu_\ell}^\ell}{\partial b_p^\ell} \\
&= \frac{\partial g_{\nu_L}^L}{\partial z_{\nu_L}^L} \sum_{\nu_{L-1}=1}^{n_{L-1}} w_{\nu_L \nu_{L-1}}^L \frac{\partial g_{\nu_{L-1}}^{L-1}}{\partial z_{\nu_{L-1}}^{L-1}} \sum_{\nu_{L-2}=1}^{n_{L-2}} w_{\nu_{L-1} \nu_{L-2}}^{L-1} \frac{\partial g_{\nu_{L-2}}^{L-2}}{\partial z_{\nu_{L-2}}^{L-2}} \cdots \\
&\quad \cdots \frac{\partial g_{\nu_{\ell+1}}^{\ell+1}}{\partial z_{\nu_{\ell+1}}^{\ell+1}} \sum_{\nu_\ell=1}^{n_\ell} w_{\nu_{\ell+1} \nu_\ell}^{\ell+1} \frac{\partial g_{\nu_\ell}^\ell}{\partial z_{\nu_\ell}^\ell} \frac{\partial}{\partial b_p^\ell} \left[\left(\sum_{\nu_{\ell-1}=1}^{n_{\ell-1}} w_{\nu_\ell \nu_{\ell-1}}^\ell a_{\nu_{\ell-1}}^{\ell-1} \right) + b_{\nu_\ell}^\ell \right] \\
&= \frac{\partial g_{\nu_L}^L}{\partial z_{\nu_L}^L} \sum_{\nu_{L-1}=1}^{n_{L-1}} w_{\nu_L \nu_{L-1}}^L \frac{\partial g_{\nu_{L-1}}^{L-1}}{\partial z_{\nu_{L-1}}^{L-1}} \sum_{\nu_{L-2}=1}^{n_{L-2}} w_{\nu_{L-1} \nu_{L-2}}^{L-1} \frac{\partial g_{\nu_{L-2}}^{L-2}}{\partial z_{\nu_{L-2}}^{L-2}} \cdots \frac{\partial g_{\nu_{\ell+1}}^{\ell+1}}{\partial z_{\nu_{\ell+1}}^{\ell+1}} \sum_{\nu_\ell=1}^{n_\ell} w_{\nu_{\ell+1} \nu_\ell}^{\ell+1} \frac{\partial g_{\nu_\ell}^\ell}{\partial z_{\nu_\ell}^\ell} \frac{\partial b_{\nu_\ell}^\ell}{\partial b_p^\ell} \\
&= \frac{\partial g_{\nu_L}^L}{\partial z_{\nu_L}^L} \sum_{\nu_{L-1}=1}^{n_{L-1}} w_{\nu_L \nu_{L-1}}^L \frac{\partial g_{\nu_{L-1}}^{L-1}}{\partial z_{\nu_{L-1}}^{L-1}} \sum_{\nu_{L-2}=1}^{n_{L-2}} w_{\nu_{L-1} \nu_{L-2}}^{L-1} \frac{\partial g_{\nu_{L-2}}^{L-2}}{\partial z_{\nu_{L-2}}^{L-2}} \cdots \frac{\partial g_{\nu_{\ell+1}}^{\ell+1}}{\partial z_{\nu_{\ell+1}}^{\ell+1}} \sum_{\nu_\ell=1}^{n_\ell} w_{\nu_{\ell+1} \nu_\ell}^{\ell+1} \frac{\partial g_{\nu_\ell}^\ell}{\partial z_{\nu_\ell}^\ell} \delta_{\nu_\ell p}
\end{aligned}$$

The resultant object is a $[n_L \times n_\ell]$ array, or matrix in this case. This is also according to the free indices (ν_L, p) .

A few things to note. The Kronecker delta function when summed over does behave as expected by selecting out only the one index such that $\delta = 1$. I have chosen not to apply this operation in the formulas, since there is a situation where one might attempt to convert all the above notation over to index (Einstein) notation and arrive at incorrect array shapes. Specifically, the common index notation convention that treats repeated indices as summed over is not compatible with these equations. There are repeated indices that correctly appear more than twice, barring the use of Einstein notation.

3.1 Vector Representation

Let us introduce some minor notation changes to aid in simplifying expressions.

Let $\gamma_{\nu_\ell}^\ell = \frac{\partial g^\ell}{\partial z_{\nu_\ell}^\ell}$, and we will begin by rearranging the two derivative equations by bringing all terms within the summations.

$$\frac{\partial a_{\nu_L}^L}{\partial w_{pq}^\ell} = \sum_{\nu_{L-1}=1}^{n_{L-1}} \cdots \sum_{\nu_\ell=1}^{n_\ell} \sum_{\nu_{\ell-1}=1}^{n_{\ell-1}} \gamma_{\nu_L}^L w_{\nu_L \nu_{L-1}}^L \gamma_{\nu_{L-1}}^{L-1} w_{\nu_{L-1} \nu_{L-2}}^{L-1} \cdots \gamma_{\nu_{\ell+1}}^{\ell+1} w_{\nu_{\ell+1} \nu_\ell}^{\ell+1} \gamma_{\nu_\ell}^\ell \delta_{\nu_\ell p} \delta_{\nu_{\ell-1} q} a_{\nu_{\ell-1}}^{\ell-1}$$

$$\frac{\partial a_{\nu_L}^L}{\partial b_p^\ell} = \sum_{\nu_{L-1}=1}^{n_{L-1}} \cdots \sum_{\nu_\ell=1}^{n_\ell} \gamma_{\nu_L}^L w_{\nu_L \nu_{L-1}}^L \gamma_{\nu_{L-1}}^{L-1} w_{\nu_{L-1} \nu_{L-2}}^{L-1} \cdots \gamma_{\nu_{\ell+1}}^{\ell+1} w_{\nu_{\ell+1} \nu_\ell}^{\ell+1} \gamma_{\nu_\ell}^\ell \delta_{\nu_\ell p}$$

Notice the pattern of terms $\gamma_{\nu_k}^k w_{\nu_k \nu_{k-1}}^k$. Since we are explicitly using summations, this is just another term with indices (ν_k, ν_{k-1}) , so let us define:

$$d_{\nu_k \nu_{k-1}}^k = \gamma_{\nu_k}^k w_{\nu_k \nu_{k-1}}^k$$

From this index notation it is clear that the $\vec{\gamma}^k$ vector is multiplying each column of W^k . Since the product is taken on each element, irrespective of the column index, we can define a new matrix:

$$G^k = \begin{bmatrix} \downarrow & & \downarrow \\ \vec{\gamma}^k & \cdots & \vec{\gamma}^k \\ \uparrow & & \uparrow \end{bmatrix}_{[n_k \times n_{k-1}]} = \begin{bmatrix} \gamma_1^k & \cdots & \gamma_1^k \\ \vdots & \ddots & \vdots \\ \gamma_{\nu_k}^k & \cdots & \gamma_{\nu_k}^k \end{bmatrix}$$

Which is the matrix whose column vectors are made of n_{k-1} copies of $\vec{\gamma}^k$.

If we look at the Hadamard product, $G^k \odot W^k$, we see that its elements are precisely $\gamma_{\nu_k}^k w_{\nu_k \nu_{k-1}}^k$. Therefore, if we define D^k as the matrix with elements $d_{\nu_k \nu_{k-1}}^k$, then:

$$D^k = G^k \odot W^k$$

For completeness and later necessity, let's also define another matrix that is compatible with multiplication along column indices, rather than just row indices as we just saw:

$$\tilde{G}^k = \begin{bmatrix} - & \vec{\gamma}^k & - \\ & \vdots & \\ - & \vec{\gamma}^k & - \end{bmatrix}_{[n_{k+1} \times n_k]}$$

With the tilde, we will say this is the matrix whose row vectors are made of n_{k+1} copies of $\vec{\gamma}^k$.

Going back to our two derivative equations, we replace the terms $\gamma_{\nu_k}^k w_{\nu_k \nu_{k-1}}^k$ with $d_{\nu_k \nu_{k-1}}^k$ and apply the appropriate sums over the Kronecker delta functions.

$$\begin{aligned} \frac{\partial a_{\nu_L}^L}{\partial w_{pq}^\ell} &= \sum_{\nu_{L-1}=1}^{n_{L-1}} \cdots \sum_{\nu_{\ell+1}=1}^{n_{\ell+1}} d_{\nu_L \nu_{L-1}}^L d_{\nu_{L-1} \nu_{L-2}}^{L-1} \cdots d_{\nu_{\ell+2} \nu_{\ell+1}}^{\ell+2} d_{\nu_{\ell+1} p}^{\ell+1} \gamma_p^\ell a_q^{\ell-1} \\ \frac{\partial a_{\nu_L}^L}{\partial b_p^\ell} &= \sum_{\nu_{L-1}=1}^{n_{L-1}} \cdots \sum_{\nu_{\ell+1}=1}^{n_{\ell+1}} d_{\nu_L \nu_{L-1}}^L d_{\nu_{L-1} \nu_{L-2}}^{L-1} \cdots d_{\nu_{\ell+2} \nu_{\ell+1}}^{\ell+2} d_{\nu_{\ell+1} p}^{\ell+1} \gamma_p^\ell \end{aligned}$$

Notice that for each pair $d_{\nu_{k+1} \nu_k}^{k+1}$ and $d_{\nu_k \nu_{k-1}}^k$, we sum over ν_k explicitly like so $\sum_{\nu_k=1}^{n_k} d_{\nu_{k+1} \nu_k}^{k+1} d_{\nu_k \nu_{k-1}}^k$. This is exactly the dot product of the two matrices, D^{k+1} and D^k , and we have one dot product for each summation remaining in each expression.

This now leaves the matter of handling $\gamma_p^\ell a_q^{\ell-1}$ and γ_p^ℓ , which remain unaccounted for, thus far. Let us examine the term $d_{\nu_{\ell+1} p}^{\ell+1} \gamma_p^\ell$ and see that the vector $\vec{\gamma}^\ell$ is multiplying every row in $D^{\ell+1}$. This can be expressed as another Hadamard product of two matrices. We already have a well defined matrix operation representing $d_{\nu_{\ell+1} p}^{\ell+1} \gamma_p^\ell$, which is $D^{\ell+1} \odot \tilde{G}^\ell$. Lastly the term $a_q^{\ell-1}$ has an index independent of all others. In fact it multiplies every element distinctly for every q . The tensor product \otimes is precisely the operation used here.

Note, alternatively, we could have defined a left and right sided Hadamard product, but it seems better to preserve the commutative nature of the Hadamard product and instead use the \tilde{G} transpose notation to the matrix G .

Finally we can reduce our two equations to only matrix dot product, Hadamard product, and tensor product operations.

$$\begin{aligned} \frac{\partial \vec{a}^L}{\partial W^\ell} &= \left(D^L \cdot D^{L-1} \cdots D^{\ell+2} \cdot (D^{\ell+1} \odot \tilde{G}^\ell) \right) \otimes \vec{a}^{\ell-1} \\ \frac{\partial \vec{a}^L}{\partial b^\ell} &= D^L \cdot D^{L-1} \cdots D^{\ell+2} \cdot (D^{\ell+1} \odot \tilde{G}^\ell) \end{aligned}$$

$$D^k = G^k \odot W^k$$

3.2 Examples

1. $L = 5, \ell = 3$

$$\begin{aligned} \frac{\partial a_{\nu_5}^5}{\partial w_{pq}^3} &= \frac{\partial g^5}{\partial z_{\nu_5}^5} \sum_{\nu_4=1}^{n_4} w_{\nu_5 \nu_4}^5 \frac{\partial g^4}{\partial z_{\nu_4}^4} \sum_{\nu_3=1}^{n_3} w_{\nu_4 \nu_3}^4 \frac{\partial g^3}{\partial z_{\nu_3}^3} \sum_{\nu_2=1}^{n_2} \delta_{\nu_3 p} \delta_{\nu_2 q} a_{\nu_2}^2 \\ &= \frac{\partial g^5}{\partial z_{\nu_5}^5} \sum_{\nu_4=1}^{n_4} w_{\nu_5 \nu_4}^5 \frac{\partial g^4}{\partial z_{\nu_4}^4} w_{\nu_4 p}^4 \frac{\partial g^3}{\partial z_p^3} a_q^2 \end{aligned}$$

Array with shape $[n_5 \times n_3 \times n_2]$

2. $L = 2, \ell = 2$

$$\begin{aligned} \frac{\partial a_{pq}^2}{\partial w_{pq}^2} &= \frac{\partial g^2}{\partial z_{\nu_2}^2} \sum_{\nu_1=1}^{n_1} \delta_{\nu_2 p} \delta_{\nu_1 q} a_{\nu_1}^1 \\ &= \frac{\partial g^2}{\partial z_{\nu_2}^2} \delta_{\nu_2 p} a_q^1 \end{aligned}$$

Array with shape $[n_2 \times n_2 \times n_1]$

$$\frac{\partial a_p^2}{\partial b_p^2} = \frac{\partial g^2}{\partial z_{\nu_2}^2} \delta_{\nu_2 p}$$

Array with shape $[n_2 \times n_2]$

3. $L = 2, \ell = 1$

$$\begin{aligned} \frac{\partial a_{\nu_2}^2}{\partial w_{pq}^1} &= \frac{\partial g^2}{\partial z_{\nu_2}^2} \sum_{\nu_1=1}^{n_1} w_{\nu_2 \nu_1}^2 \frac{\partial g^1}{\partial z_{\nu_1}^1} \sum_{\nu_0=1}^{n_0} \delta_{\nu_1 p} \delta_{\nu_0 q} a_{\nu_0}^0 \\ &= \frac{\partial g^2}{\partial z_{\nu_2}^2} w_{\nu_2 p}^2 \frac{\partial g^1}{\partial z_p^1} a_q^0 \end{aligned}$$

Array with shape $[n_2 \times n_1 \times n_0]$

$$\begin{aligned} \frac{\partial a_p^2}{\partial b_p^1} &= \frac{\partial g^2}{\partial z_{\nu_2}^2} \sum_{\nu_1=1}^{n_1} w_{\nu_2 \nu_1}^2 \frac{\partial g^1}{\partial z_{\nu_1}^1} \delta_{\nu_1 p} \\ &= \frac{\partial g^2}{\partial z_{\nu_2}^2} w_{\nu_2 p}^2 \frac{\partial g^1}{\partial z_p^1} \end{aligned}$$

Array with shape $[n_2 \times n_1]$

4 Loss and Cost

4.1 Loss

If we first define the metric function on \mathbb{R}^n ,

$$d(\vec{u}, \vec{v}) : \mathbb{R}^n \times \mathbb{R}^n \longrightarrow \mathbb{R}$$

and a projection function

$$\pi : \mathbb{R}^n \times \mathbb{R}^n \longrightarrow \mathbb{R}^n$$

then loss is the image of the projection of the metric:

$$L : \mathbb{R}^n \longrightarrow \mathbb{R}$$

where $d(\vec{u}, \vec{v}) = (L \circ \pi)(\vec{u}, \vec{v}) = L(\pi(\vec{u}, \vec{v})) = L(\vec{u})$.

If we have a prediction vector, $\overrightarrow{y_{pred}}$ and a truth vector \vec{y} , then the loss is $L(\overrightarrow{y_{pred}}) = d(\overrightarrow{y_{pred}}, \vec{y})$

4.2 Cost

Cost is essentially loss, but now taking into account a set of samples $\{\overrightarrow{a^{L1}}, \dots, \overrightarrow{a^{Lm}}\}$ of the output vectors, $\overrightarrow{a^L} \in \mathbb{R}^n$. The cost function is

$$\mathcal{C} : \Omega \longrightarrow \mathbb{R}$$

where the domain Ω is a finite, $|\Omega| = m$, collection of vectors in \mathbb{R}^n .

We define a few component functions that help to fully describe the various ways in which the cost can be used to aggregate the m vectors while still incorporating the loss.

First define,

$$\phi : \Omega \longrightarrow \mathbb{R}^m$$

such that

$$\phi(\overrightarrow{a^{L1}}, \dots, \overrightarrow{a^{Lm}}) = (L(\overrightarrow{a^{L1}}), \dots, L(\overrightarrow{a^{Lm}}))$$

Notice that a tuple of m vectors in \mathbb{R}^n , $(\overrightarrow{a^{L1}}, \dots, \overrightarrow{a^{Lm}})$ can be made isomorphic to a vector $A \in \mathbb{R}^{mn}$, where $A = (a_1^{L1}, \dots, a_n^{L1}, \dots, a_1^{Lm}, \dots, a_n^{Lm})$. Thus, we can also define a map

$$\tilde{\phi} : \mathbb{R}^{mn} \longrightarrow \mathbb{R}^m$$

such that if $A \in \mathbb{R}^{mn}$ and a projection map,

$$\pi_k : \mathbb{R}^{mn} \longrightarrow \mathbb{R}^n$$

where $\pi_k(A) = (a_1^{Lk}, \dots, a_n^{Lk}) = \overrightarrow{a^{Lk}}$, then $L(\overrightarrow{a^{Lk}}) = (L \circ \pi_k)(A)$. So let

$$L_k = L \circ \pi_k$$

Thus,

$$\begin{aligned}\tilde{\phi}(A) &= (L_1(A), \dots, L_m(A)) \\ &= ((L \circ \pi_1)(A), \dots, (L \circ \pi_m)(A)) \\ &= (L(\overrightarrow{a^{L1}}), \dots, L(\overrightarrow{a^{Lm}})) \\ &= \phi(\overrightarrow{a^{L1}}, \dots, \overrightarrow{a^{Lm}})\end{aligned}$$

Next define the differentiable function (note the similarities to a functional), which gives rise to cost,

$$\varphi : \mathbb{R}^m \longrightarrow \mathbb{R}$$

such that

$$\mathcal{C}(\overrightarrow{a^{L1}}, \dots, \overrightarrow{a^{Lm}}) = (\varphi \circ \phi)(\overrightarrow{a^{L1}}, \dots, \overrightarrow{a^{Lm}})$$

and

$$\mathcal{C}(A) = (\varphi \circ \tilde{\phi})(A)$$

5 Total Derivative of Cost

We can simply apply the chain rule on cost as follows:

$$\begin{aligned}D_A \mathcal{C} &= D_A(\varphi \circ \tilde{\phi}) \\ &= D_{\tilde{\phi}(A)} \varphi \cdot D_A \tilde{\phi} \\ &= \nabla \varphi \cdot D_A \tilde{\phi}\end{aligned}$$

where $\nabla \varphi$ is a $[1 \times m]$ gradient vector.

Let us examine $D_A \tilde{\phi}$ further. Since we know $\tilde{\phi} : \mathbb{R}^{mn} \longrightarrow \mathbb{R}^m$, then the total derivative is a linear map represented by a $[m \times mn]$ matrix. Now recall $\tilde{\phi}(A) = (L_1(A), \dots, L_m(A))$, therefore:

$$\begin{aligned}D_A \tilde{\phi} &= \begin{bmatrix} \frac{\partial L_1}{\partial a_1^{L1}} & \cdots & \frac{\partial L_1}{\partial a_n^{L1}} & \cdots & \frac{\partial L_1}{\partial a_1^{Lm}} & \cdots & \frac{\partial L_1}{\partial a_n^{Lm}} \\ \vdots & & & \ddots & & & \vdots \\ \frac{\partial L_m}{\partial a_1^{L1}} & \cdots & \frac{\partial L_m}{\partial a_n^{L1}} & \cdots & \frac{\partial L_m}{\partial a_1^{Lm}} & \cdots & \frac{\partial L_m}{\partial a_n^{Lm}} \end{bmatrix} \\ &= \begin{bmatrix} - & D_A L_1 & - \\ & \vdots & \\ - & D_A L_m & - \end{bmatrix}\end{aligned}$$

Next we will compute the total derivatives for L_k ,

$$\begin{aligned}
D_A L_k &= D_A(L \circ \pi_k) \\
&= D_{\pi_k(A)} \vec{L} \cdot D_A \pi_k \\
&= \nabla L(a^{Lk}) \cdot D_A \pi_k
\end{aligned}$$

and since $\pi_k(A) = (a_1^{Lk}, \dots, a_n^{Lk})$,

$$D_A \pi_k = \begin{bmatrix} 0 & | & \dots & | & \mathbf{I}_{n \times n} & | & \dots & | & 0 \end{bmatrix}_{n \times mn}$$

where the identity matrix is k -th block matrix (m total block matrices of size $n \times n$).

Therefore, we can also express,

$$D_A L_k = \nabla L_k = \left(0, \dots, \frac{\partial L}{\partial a_1^{Lk}}, \dots, \frac{\partial L}{\partial a_n^{Lk}}, \dots, 0 \right)$$

where only the k -th, sub-vector of \mathbb{R}^n is non-zero.

Putting these pieces together:

$$\begin{aligned}
D_A \mathcal{C} &= \nabla \varphi \cdot D_A \tilde{\phi} \\
&= \nabla \varphi_{[1 \times m]} \cdot \begin{bmatrix} - & \nabla L_1 & - \\ & \vdots & \\ - & \nabla L_m & - \end{bmatrix}_{[m \times mn]}
\end{aligned}$$

Recall that each element of A is a function of the weights and biases of the network. Using the chain rule once again we compute the total derivatives of cost w.r.t. W^ℓ and \vec{b}^ℓ :

$$\begin{aligned}
D_{W^\ell} \mathcal{C} &= D_{W^\ell}(\varphi \circ \tilde{\phi} \circ A) \\
&= D_{\phi(A)} \varphi \cdot D_{A(W^\ell)} \tilde{\phi} \cdot D_{W^\ell} A \\
&= \nabla \varphi \cdot D_A \tilde{\phi} \cdot D_{W^\ell} A
\end{aligned}$$

where $D_{W^\ell} A$ has shape $mn \times n_\ell \times n_{\ell-1}$.

$$\begin{aligned}
D_{\vec{b}^\ell} \mathcal{C} &= D_{W^\ell}(\varphi \circ \tilde{\phi} \circ A) \\
&= D_{\phi(A)} \varphi \cdot D_{A(W^\ell)} \tilde{\phi} \cdot D_{\vec{b}^\ell} A \\
&= \nabla \varphi \cdot D_A \tilde{\phi} \cdot D_{\vec{b}^\ell} A
\end{aligned}$$

where $D_{\vec{b}^\ell} A$ has shape $mn \times n_\ell$.

The formula for $D_{W^\ell} A$ is exactly the same as that for $\frac{\partial \vec{a}^L}{\partial W^\ell}$, only with a different sized vector. Ignoring a change of notation, where $D_{W^\ell} \vec{a}^L = \frac{\partial \vec{a}^L}{\partial W^\ell}$, the vector $\vec{a}^L \in \mathbb{R}^n$ is simply being replaced with a vector $A \in \mathbb{R}^{mn}$. Thus, we see

in the expression for the shape of the total derivative, $n \times n_\ell \times n_{\ell-1}$, we are only making a substitution of n for mn . The same idea holds for $D_{\vec{b}}\mathcal{C}$, since it is still just a substitution of vectors from a different dimension vector space. The derivatives of the final activation layer w.r.t. the weights or biases are true for any finite dimensional Euclidean vector space.

6 Optimizing Parameters

6.1 Gradient Descent