

Basic Operations of Neural Networks

Sourya Dey

1 Notation

- Scalars are written as lower case letters.
- Vectors are written as lower case bold letters, such as \mathbf{x} , and can be either row (dimensions $1 \times n$) or column (dimensions $n \times 1$). Column vectors are the default choice, unless otherwise mentioned. Individual elements are indexed by subscripts, such as x_i ($i \in \{1, \dots, n\}$).
- Matrices are written as upper case bold letters, such as \mathbf{X} , and have dimensions $m \times n$ corresponding to m rows and n columns. Individual elements are indexed by double subscripts for row and column, such as X_{ij} ($i \in \{1, \dots, m\}, j \in \{1, \dots, n\}$).
- Bracketed superscripts are used to denote layers, for example, $W_{ij}^{(l)}$ denotes the (i, j) th element of the weight matrix of layer l .

The **derivative** of f with respect to x is $\frac{\partial f}{\partial x}$. Both x and f can be a scalar, vector, or matrix. The **gradient** of f w.r.t x is $\nabla_x f = \left(\frac{\partial f}{\partial x}\right)^T$, i.e. **gradient is transpose of derivative**. The gradient at any point x_0 in the domain has a physical interpretation, its direction is the direction of maximum increase of the function f at the point x_0 , and its magnitude is the rate of increase in that direction.

2 Neural Network Operations

A neural network has $(L + 1)$ layers having $(N^{(0)}, N^{(1)}, \dots, N^{(L)})$ neurons respectively, i.e. there are $N^{(0)}$ input and $N^{(L)}$ output neurons. Any layer l ($l \neq 0$) has a bias vector $\mathbf{b}^{(l)}$, an activation vector $\mathbf{a}^{(l)}$, and a delta (error) vector $\boldsymbol{\delta}^{(l)}$, each of dimensions $N^{(l)} \times 1$, and a weight matrix $\mathbf{W}^{(l)}$ preceding it of dimensions $N^{(l)} \times N^{(l-1)}$.

2.1 Feedforward

The input layer is fed an input sample vector $\mathbf{a}^{(0)}$ of dimensions $N^{(0)} \times 1$. Then the feedforward operation for all layers $l \in \{1, \dots, L\}$ proceeds as:

$$\mathbf{s}^{(l)} = \mathbf{W}^{(l)}\mathbf{a}^{(l-1)} + \mathbf{b}^{(l)} \quad (1a)$$

$$\mathbf{a}^{(l)} = \mathbf{h}\left(\mathbf{s}^{(l)}\right) \quad (1b)$$

where $\mathbf{h}(\cdot)$ is the activation function. Activation functions are discussed further in Sec. 3.

The final layer output activation vector $\mathbf{a}^{(L)}$ is compared with the ground truth output vector for that input sample, $\mathbf{y}^{(L)}$, to compute a scalar-valued cost C . A popular cost function is cross-entropy:

$$C = - \sum_{i=1}^{N^{(L)}} y_i^{(L)} \ln a_i^{(L)} \quad (2)$$

2.2 Backpropagation

The goal of backpropagation is to compute the gradients of the cost w.r.t all the network parameters, i.e. $\nabla_{\mathbf{W}^{(l)}} C$ and $\nabla_{\mathbf{b}^{(l)}} C$ for all $l \in \{1, \dots, L\}$. These gradients are used for updating the parameter values to make the network learn.

Since cost is directly a function of $\mathbf{a}^{(L)}$, we start by computing $\frac{\partial C}{\partial \mathbf{a}^{(L)}}$:

$$\frac{\partial C}{\partial \mathbf{a}^{(L)}} = - \begin{bmatrix} \frac{y_1^{(L)}}{a_1^{(L)}} & \frac{y_2^{(L)}}{a_2^{(L)}} & \dots & \frac{y_{N^{(L)}}^{(L)}}{a_{N^{(L)}}^{(L)}} \end{bmatrix} \quad (3)$$

Then we work backwards:

$$\frac{\partial C}{\partial \mathbf{s}^{(L)}} = \frac{\partial C}{\partial \mathbf{a}^{(L)}} \frac{\partial \mathbf{a}^{(L)}}{\partial \mathbf{s}^{(L)}} = \frac{\partial C}{\partial \mathbf{a}^{(L)}} \mathbf{H}'^{(L)} \quad (4)$$

where $\mathbf{H}'^{(L)} = \frac{\partial \mathbf{a}^{(L)}}{\partial \mathbf{s}^{(L)}}$ is the matrix of activation function derivatives (described further in Sec. 3).

The delta vector of any layer is defined as $\delta^{(l)} = \nabla_{\mathbf{s}^{(l)}} C$. Since gradient is transpose of the derivative, we get:

$$\delta^{(L)} = \nabla_{\mathbf{s}^{(L)}} C = \mathbf{H}'^{(L)T} \nabla_{\mathbf{a}^{(L)}} C \quad (5)$$

Continuing the chain rule, we get:

$$\frac{\partial C}{\partial \mathbf{b}^{(L)}} = \frac{\partial C}{\partial \mathbf{s}^{(L)}} \frac{\partial \mathbf{s}^{(L)}}{\partial \mathbf{b}^{(L)}} \quad (6)$$

From (1a), $\frac{\partial \mathbf{s}^{(L)}}{\partial \mathbf{b}^{(L)}} = I$, the identity matrix. Thus, $\frac{\partial C}{\partial \mathbf{b}^{(L)}} = \frac{\partial C}{\partial \mathbf{s}^{(L)}}$, transposing which we get our first important result:

$$\nabla_{\mathbf{b}^{(L)}} C = \delta^{(L)} \quad (7)$$

which is a vector of the same dimensions as $\mathbf{b}^{(L)}$, i.e. $(N^{(L)} \times 1)$.

Again from the chain rule, we get:

$$\frac{\partial C}{\partial \mathbf{W}^{(L)}} = \frac{\partial C}{\partial \mathbf{s}^{(L)}} \frac{\partial \mathbf{s}^{(L)}}{\partial \mathbf{W}^{(L)}} \quad (8)$$

The quantity $\frac{\partial \mathbf{s}^{(L)}}{\partial \mathbf{W}^{(L)}}$ is a 3rd order tensor of dimensions $(N^{(L)} \times N^{(L-1)} \times N^{(L)})$, i.e. an outer matrix of dimensions $(N^{(L-1)} \times N^{(L)})$ with each element being an $N^{(L)}$ -dimensional column vector. Using (1a) for $l = L$, note that:

$$s_i^{(L)} = W_{i,1}^{(L)} a_1^{(L-1)} + W_{i,2}^{(L)} a_2^{(L-1)} + \dots + W_{i,N^{(L-1)}}^{(L)} a_{N^{(L-1)}}^{(L-1)} \quad (9)$$

So the $N^{(L)}$ -dimensional column vector at $\left(\frac{\partial \mathbf{s}^{(L)}}{\partial \mathbf{W}^{(L)}}\right)_{jk}$ has zeroes everywhere except for location k having value $a_j^{(L-1)}$. This is visualized below for the simple case $N^{(L)} = 2$, $N^{(L-1)} = 3$:

$$\frac{\partial \mathbf{s}^{(L)}}{\partial \mathbf{W}^{(L)}} = \begin{bmatrix} \begin{bmatrix} a_1^{(L-1)} \\ 0 \end{bmatrix} & \begin{bmatrix} 0 \\ a_1^{(L-1)} \end{bmatrix} \\ \begin{bmatrix} a_2^{(L-1)} \\ 0 \end{bmatrix} & \begin{bmatrix} 0 \\ a_2^{(L-1)} \end{bmatrix} \\ \begin{bmatrix} a_3^{(L-1)} \\ 0 \end{bmatrix} & \begin{bmatrix} 0 \\ a_3^{(L-1)} \end{bmatrix} \end{bmatrix} \quad (10)$$

Returning to (8), the (j, k) th element of $\frac{\partial C}{\partial \mathbf{W}^{(L)}}$ is the inner product of the row vector $\frac{\partial C}{\partial \mathbf{s}^{(L)}}$ with the column vector at $\left(\frac{\partial \mathbf{s}^{(L)}}{\partial \mathbf{W}^{(L)}}\right)_{jk}$. This inner product is $a_j^{(L-1)} \delta_k^{(L)}$, so $\frac{\partial C}{\partial \mathbf{W}^{(L)}} = \mathbf{a}^{(L-1)} \boldsymbol{\delta}^{(L)T}$. Transposing this gives our second important result:

$$\nabla_{\mathbf{W}^{(L)}} C = \boldsymbol{\delta}^{(L)} \mathbf{a}^{(L-1)T} \quad (11)$$

which is a matrix of the same dimensions as $\mathbf{W}^{(L)}$, i.e. $(N^{(L)} \times N^{(L-1)})$.

Thus, we have obtained the gradients of cost w.r.t \mathbf{W} and \mathbf{b} of the output layer from its $\boldsymbol{\delta}$ vector. The only remaining task is to obtain $\boldsymbol{\delta}^{(l)}$ from $\boldsymbol{\delta}^{(l+1)}$, i.e. get the delta vector for each layer from that of its next. Using the wonderful chain rule again, we get:

$$\frac{\partial C}{\partial \mathbf{s}^{(l)}} = \frac{\partial C}{\partial \mathbf{s}^{(l+1)}} \frac{\partial \mathbf{s}^{(l+1)}}{\partial \mathbf{a}^{(l)}} \frac{\partial \mathbf{a}^{(l)}}{\partial \mathbf{s}^{(l)}} \quad (12)$$

From (1a), $\frac{\partial \mathbf{s}^{(l+1)}}{\partial \mathbf{a}^{(l)}} = \mathbf{W}^{(l+1)}$, and we have already established that $\frac{\partial \mathbf{a}^{(l)}}{\partial \mathbf{s}^{(l)}} = \mathbf{H}'^{(l)}$, i.e. the matrix of derivatives of the activation function in layer l (see (17)). Combining all this and transposing to get gradient, we get our 3rd and final important result:

$$\boldsymbol{\delta}^{(l)} = \mathbf{H}'^{(l)T} \mathbf{W}^{(l+1)T} \boldsymbol{\delta}^{(l+1)} \quad (13)$$

That's it!

In summary, backpropagation proceeds by computing the gradient of cost w.r.t the following:

1. Output layer: $\boldsymbol{\delta}^{(L)} = \mathbf{H}'^{(L)T} \nabla_{\mathbf{a}^{(L)}} C$
2. Intermediate layers: $\boldsymbol{\delta}^{(l)} = \mathbf{H}'^{(l)T} \mathbf{W}^{(l+1)T} \boldsymbol{\delta}^{(l+1)} \quad \forall l \in \{1, \dots, L-1\}$
3. All biases: $\nabla_{\mathbf{b}^{(l)}} C = \boldsymbol{\delta}^{(l)} \quad \forall l \in \{1, \dots, L\}$
4. All weights: $\nabla_{\mathbf{W}^{(l)}} C = \boldsymbol{\delta}^{(L)} \mathbf{a}^{(L-1)T} \quad \forall l \in \{1, \dots, L\}$

If you do not want to treat the biases separately, you can augment the weight matrix $\mathbf{W}^{(l)}$ with a column for $\mathbf{b}^{(l)}$ and augment the activation vector $\mathbf{a}^{(l-1)}$ with a single 1. Then backprop step 4 takes care of step 3 automatically. Think about it.

2.3 Update

Note that the gradient of cost w.r.t any parameter \mathbf{p} is of the same dimensions as \mathbf{p} (\mathbf{p} could be any weight matrix, bias vector, or any other parameter which the network needs to learn). Now we can apply the gradient descent update rule:

$$\mathbf{p} \leftarrow \mathbf{p} - \eta \nabla_{\mathbf{p}} C \quad (14)$$

where η is the learning rate. This takes \mathbf{p} in the direction opposite to the gradient, i.e. in the direction of maximum decrease of C .

In practice, we do not update after every input, but perform feedforward and backpropagation on several inputs before doing a single averaged update. The number of inputs between two updates is the *minibatch size* M , i.e.:

$$\mathbf{p} \leftarrow \mathbf{p} - \frac{\eta}{M} \sum_{i=1}^M (\nabla_{\mathbf{p}} C)^{[i]} \quad (15)$$

where $(\nabla_{\mathbf{p}} C)^{[i]}$ is the gradient for input sample i .

3 Different Activation Functions

For simplicity, we discard the layer superscript (l) in this section. Recall that the activation operation in any layer with N neurons is $\mathbf{a} = \mathbf{h}(\mathbf{s})$, where \mathbf{a} and \mathbf{s} are N -dimensional column vectors.

3.1 Hidden layers

For most hidden layers, \mathbf{h} will be a vectorized scalar function applied element-wise to \mathbf{s} . This is denoted as \underline{h} , i.e.:

$$\mathbf{a} = \underline{h} \left(\begin{bmatrix} s_1 \\ \vdots \\ s_N \end{bmatrix} \right) = \begin{bmatrix} h(s_1) \\ \vdots \\ h(s_N) \end{bmatrix} \quad (16)$$

In this case, the derivative $\mathbf{H}' = \frac{\partial \mathbf{a}}{\partial \mathbf{s}}$ will be a $N \times N$ diagonal matrix:

$$\mathbf{H}' = \begin{bmatrix} h'(s_1) & & 0 \\ & \ddots & \\ 0 & & h'(s_N) \end{bmatrix} \quad (17)$$

3.2 Output layer

For the output layer, we generally use **softmax activation**, where every element in the output vector is influenced by all elements of the input vector:

$$\mathbf{a} = \begin{bmatrix} \frac{e^{s_1}}{\sum_{i=1}^N e^{s_i}} \\ \vdots \\ \frac{e^{s_N}}{\sum_{i=1}^N e^{s_i}} \end{bmatrix} = \begin{bmatrix} \frac{k_1}{k} \\ \vdots \\ \frac{k_N}{k} \end{bmatrix} \quad (18)$$

where $k_i = e^{s_i}$ and $k = \sum_{i=1}^N e^{s_i}$.

Using simple calculus gives the diagonal entries of \mathbf{H}' as $\frac{\partial a_i}{\partial s_i} = \frac{k_i(k - k_i)}{k^2}$, and the off-diagonal entries as $\frac{\partial a_i}{\partial s_j} = \frac{\partial a_j}{\partial s_i} = -\frac{k_i k_j}{k^2}$. Thus, \mathbf{H}' is a symmetric matrix.

As shown in (3), using cross-entropy cost results in $\frac{\partial C}{\partial \mathbf{a}} = -\left[\frac{y_1}{a_1} \quad \dots \quad \frac{y_N}{a_N} \right]$. Then the delta vector of the output layer becomes:

$$\boldsymbol{\delta} = \mathbf{H}' \nabla_{\mathbf{a}} C = \begin{bmatrix} \left(a_1 \sum_{i=1}^N y_i \right) - y_1 \\ \vdots \\ \left(a_N \sum_{i=1}^N y_i \right) - y_N \end{bmatrix} \quad (19)$$

This has a particularly simple interpretation when the outputs are one-hot. This means that if the correct class is n , then $y_n = 1$ and $y_m = 0 \quad \forall m \neq n$. In this case, $\delta_n = a_n - 1$ and $\delta_m = a_m \quad \forall m \neq n$. This means that $\boldsymbol{\delta} = \mathbf{a} - \mathbf{y}$, which makes perfect sense as the ‘error vector’.