

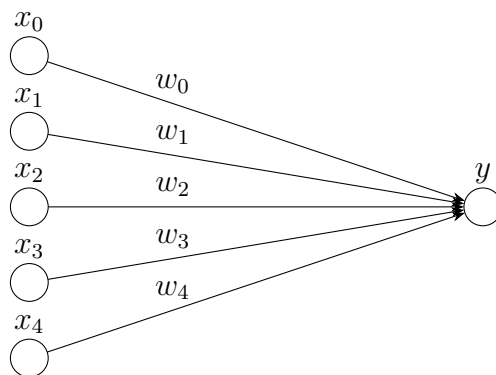
Chapter 1

Introduction

1.1 Neural Networks

Neural networks represent a mathematical tool in machine learning that is useful for performing function approximation; they can be thought of as a generalization of linear regression. The power of a neural network stems from three main properties that we will go over: nonliterary, differentiability, and hidden layers. These key properties allow neural networks to be able to improve its approximation of a set of values via “training.”

To kick things off we will start with a simple example of a neural network. In the simplest form, a neural network consists of a **vector input** $\vec{x} \in \mathbb{R}^n$, a set of **weights** $w_i \in \mathbb{R}$, and a final **output** $y \in \mathbb{R}$. Below is a graphical representation of a such network.



In such a graphical representation, the set of x_i nodes are called the **input layer** (or simply the first layer) and the y node is called the **output layer**. In the above diagram, the output layer consists of one node, but as we will see it can also consist of multiple nodes.

The output is obtained as a function of the input and weights as below.

$$y = \sum_i w_i x_i$$

From a statistics perspective, this is actually just a **linear model**. In statistics, one “trains” such a linear model via linear regression on some dataset. If you have taken a statistics course, you might remember that this strategy works on simple examples (e.g., a suspiciously-already-linear weather

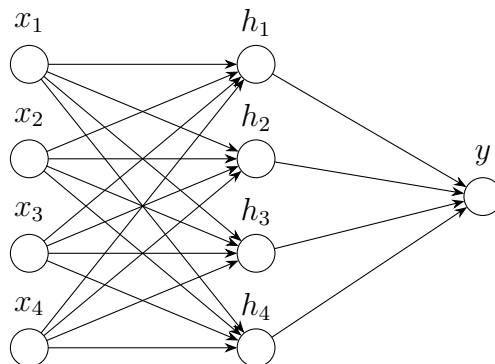
dataset in a Pearson textbook), but linear models do not generalize and often fail to capture complex behavior.

As we will see, neural networks are different from linear models since they add properties of hidden layers and nonlinearity.

1.2 Hidden Layers

Neural networks extend our previous notion of a linear model via **hidden layers**, which can be defined as one or more layers between the input and output layer. Below, we have a neural network which has one hidden layer. The hidden layer can have a variable number of nodes, but in our example below we have five.

In this example, each input node x_i connects to each node h_j in the hidden layer via a weight w_{ij} . These weights are illustrated by the arrows, although we are choosing to suppress the w_{ij} notation in the diagram below to not over complicate the figure.



The calculation of a hidden layer is simply

$$h_i = \sum_j w_{ij} x_j$$

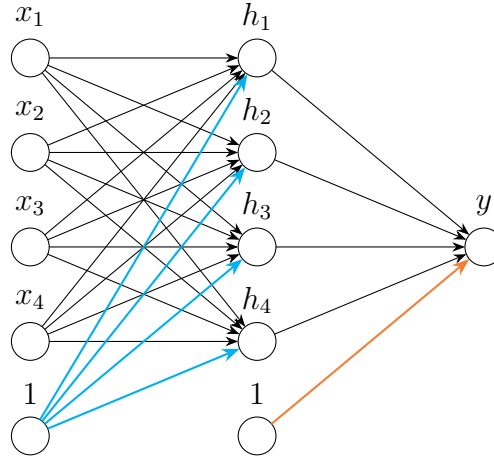
Intuitively, this means that each input value x_i makes a weighted contribution of w_{ij} to the value h_j . Something to observe at this point is that we can summarize the entire hidden layer calculation as a matrix equation:

$$\vec{h} = \begin{bmatrix} h_1 \\ h_2 \\ h_3 \\ h_4 \end{bmatrix} = \begin{bmatrix} \sum_i w_{1i} x_i \\ \sum_i w_{2i} x_i \\ \sum_i w_{3i} x_i \\ \sum_i w_{4i} x_i \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \\ w_{41} & w_{42} & w_{43} & w_{44} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = W\vec{x} \quad (1.1)$$

This suggests the concept of a **weight matrix** W , which is the key to calculating the hidden layer \vec{h} from the input \vec{x} . For a neural network that has multiple hidden layers, there are multiple corresponding weight matrices. In fact, a neural network with N layers will have $(N - 1)$ -many weight matrices.

For our current example, we will let U denote the matrix of the weights between the output layer and the hidden layer, so that $y = U\vec{h}$.

Often times with neural networks it is necessary to introduce a **bias** in each layer. A bias is an extra node, assigned a value of 1, that we can add to a layer which will be used to contribute to the calculation of the next layer. Below we illustrate the network we'd obtain by adding bias to our previous network.



With this neural network architecture, we have the following weight matrices:

$$W = \begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} & w_{15} \\ w_{21} & w_{22} & w_{23} & w_{24} & w_{25} \\ w_{31} & w_{32} & w_{33} & w_{34} & w_{35} \\ w_{41} & w_{42} & w_{43} & w_{44} & w_{45} \end{bmatrix} \quad U = \begin{bmatrix} u_{11} & u_{12} & u_{13} & u_{14} & u_{15} \end{bmatrix}$$

Thus, adding a bias to a layer is equivalent to adding a new column to the weight matrix. Our input $[x_1, x_2, x_3, x_4]$ is still in \mathbb{R}^4 , but we now instead feed the neural network a value of $[x_1, x_2, x_3, x_4, 1]^T \in \mathbb{R}^5$.

However, note that we're not really doing much mathematically by adding a hidden layer. Observe that we can rewrite y as

$$y = U\vec{h} = U(W\vec{x}) = W'\vec{x}$$

where $W' = UW$. This reduces our above network, with three layers, to a boring network, with two layers (similar to the one we started with), just with a different weight matrix W' . The reason is because our network is linear, which means no matter how many layers we add it will always reduce to the same boring network we started with. As we know, linear patterns cannot adequately capture complex patterns. Thus in order to get something interesting with hidden layers we need to add some nonlinearity.

1.3 Nonlinearity

Neural networks achieve our desired property of nonlinearity via usage of **activation functions**. An activation function is a function f that is called on the summation of a given node in a network,

allowing us to modify the calculation for a node h_i as below.

$$h_i = f \left(\sum_j w_{ij} x_j \right)$$

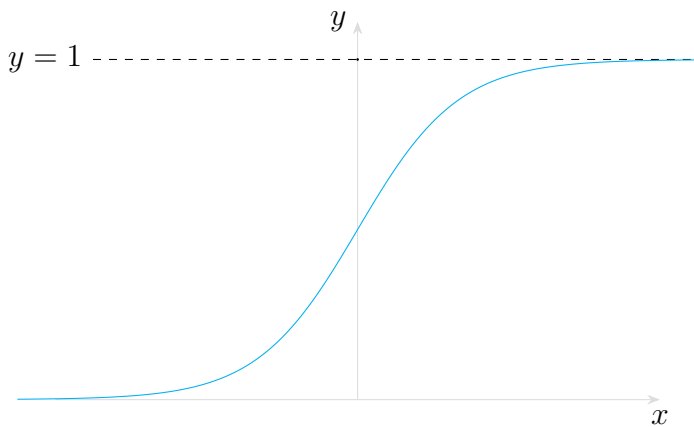
We can introduce nonlinearity into our system if we design f to be nonlinear.

A few common examples of such functions are the sigmoid (also known as logistic), tanh or RELU functions. The machine learning community has gone through several iterations of what is considered to be best practice for an activation function. In the 1990s, the sigmoid function was used very widely. In the later 90s and early 2000s, it was discovered that the tanh function lead to be better training performance. Later, it was the discovered that the ReLU function lead to even better training performance. As a result, most modern neural networks will use this for the activation function.

Let's introduce the activation functions we just discussed. Below, we have the sigmoid function

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

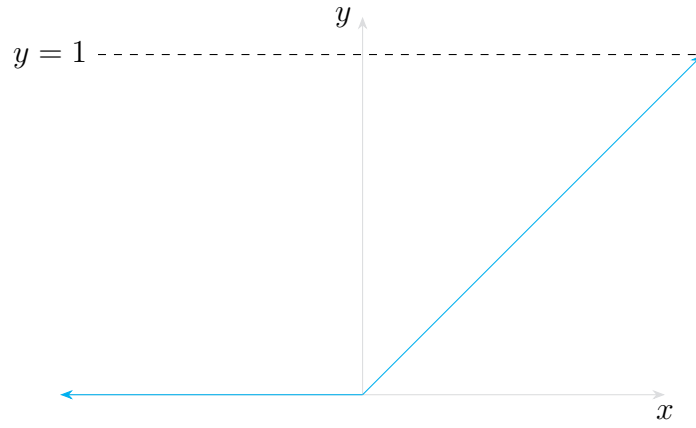
and the graph of the sigmoid function is given below.



Finally, ReLU itself is given by

$$r(x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

Below we have ReLU, which we will use in our examples as it generally results in better training performance.



Using our previous network, we can add nonlinearity by defining the computation of a hidden unit to be

$$h_j = \sigma \left(\sum_i w_{ij} x_i + b_i \right)$$

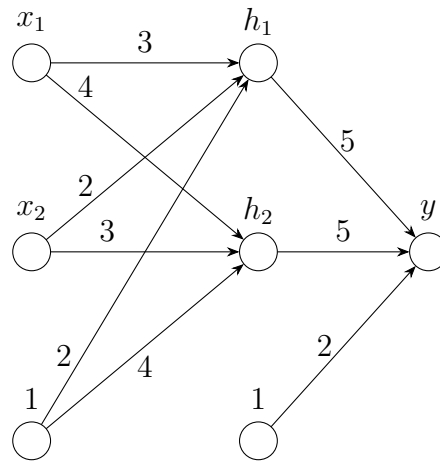
where σ is the activation function of choice.

1.4 Backpropagation: A first stab

At this point, as we have discussed the basic properties of a neural network, we will introduce a concrete example of a neural network and attempt to train it to approximate the XOR function. The XOR function performs the following mapping on these two bit inputs:

$$\begin{bmatrix} 1 \\ 0 \end{bmatrix} \rightarrow 1 \quad \begin{bmatrix} 0 \\ 0 \end{bmatrix} \rightarrow 0 \quad \begin{bmatrix} 0 \\ 1 \end{bmatrix} \rightarrow 1 \quad \begin{bmatrix} 1 \\ 1 \end{bmatrix} \rightarrow 0 \quad (1.2)$$

Below is our proposed network. We'll use ReLU, denoted as $r(x)$, for the activation function on our nodes.



Note that we have a bias inbetween each layer. The first one has weights 2 and -4, and the second one has weight -2. Thus for this network, the weight matrices are defined to be

$$W = \begin{bmatrix} 3 & 4 & 2 \\ 2 & 3 & 4 \end{bmatrix} \quad U = \begin{bmatrix} 5 & 5 & 2 \end{bmatrix}$$

allowing us to write $\vec{h} = r(W\vec{x})$ and $y = r(U\vec{h})$, or more explicitly, for a given input (x_0, x_1)

$$h_1 = r(3x_1 + 4x_2 + 2) \quad (1.3)$$

$$h_2 = r(2x_1 + 3x_2 + 4) \quad (1.4)$$

$$y = r(5h_1 + 5h_2 + 2) \quad (1.5)$$

which we can use to calculate the network. Below is a table of how this neural network currently computes the XOR values.

(x_0, x_1)	h_0	h_1	y	target
(1, 0)	$\sigma(1) = 0.731$	$\sigma(-2) = 0.119$	$\sigma(1.060) = 0.743$	1
(0, 0)	$\sigma(-2) = 0.119$	$\sigma(-4) = 0.018$	$\sigma(-1.495) = 0.183$	0
(0, 1)	$\sigma(2) = 0.881$	$\sigma(-1) = 0.269$	$\sigma(1.060) = 0.743$	1
(1, 1)	$\sigma(5) = 0.993$	$\sigma(1) = 0.731$	$\sigma(-0.690) = 0.334$	0

Based on the above table, we can see that so far it's not performing that well, but after all it is a first stab. This now begs the question: What does it mean for a model to perform well, and how do we know when it is improving? The answer to this is to introduce a **cost function** which can give a measure of error. There are many possible choices for a cost function, but for simplicity we will use the **least squares** cost function. If we have a dataset of target values t_i , and our model currently approximates this data with a set of values y_i , then the measured loss is

$$L = \frac{1}{2} \sum_i (t_i - y_i)^2$$

In our case, since the output of our function is in \mathbb{R} , we have $L(t, y) = \frac{1}{2}(t - y)^2$.

The concept of a cost function now leads to the strategy of back propagation: we seek an optimal set of weights u_i and w_{ij} such that our cost function is minimized with respect to our dataset. Intuitively this makes sense, but what does this mean mathematically?

Let us take an abstract perspective. Up until this point we have defined $y : \mathbb{R}^2 \rightarrow \mathbb{R}$ to be a function that takes in two inputs. By abuse of notation, we can just as well view y to be a function of two inputs x_1, x_2 and the weights $u_1, u_2, u_3, w_{11}, w_{12}, w_{13}, w_{21}, w_{22}, w_{23}$. For our specific model architecture, the exact expression would be

$$y = r(u_1 h_1 + u_2 h_2 + u_3) \quad (1.6)$$

$$= r(u_1 r(w_{11} x_1 + w_{12} x_2 + w_{13}) + u_2 r(w_{21} x_1 + w_{22} x_2 + w_{23}) + u_3) \quad (1.7)$$

and we can declare that $y : \mathbb{R}^{11} \rightarrow \mathbb{R}$. In this case, we can write

$$y(u_1, u_2, u_3, w_{11}, w_{12}, w_{13}, w_{21}, w_{22}, w_{23}, x_1, x_2)$$

to emphasize that y is capable of taking in 11 inputs. From this perspective, our initially proposed model is a special case of the above expression, just with appropriate values for the weights filled in for u_i and w_{ij} . In fact, our initial model is just $y(5, 5, 2, 3, 4, 2, 2, 3, 4, x_1, x_2)$.

Next, consider our cost function $L : \mathbb{R}^2 \rightarrow \mathbb{R}$. While we have considered it to be a function of

two inputs, we can similarly abuse notation and view it as a function of many inputs.

$$L(t, y) = L(t, y(u_1, u_2, u_3, w_{11}, w_{12}, w_{13}, w_{21}, w_{22}, w_{23}, x_1, x_2))$$

In this case, we have that $L : \mathbb{R}^{12} \rightarrow \mathbb{R}$.

Let $(t_1, \vec{x}_1), (t_2, \vec{x}_2), \dots, (t_n, \vec{x}_n)$ be a dataset of interests, where $t_i \in \mathbb{R}$ denotes the target values and $\vec{x}_i = (x_1^{(i)}, x_2^{(i)}) \in \mathbb{R}^2$ denotes the input values to our model.

For one dataset pair (t_k, \vec{x}_k) , we obtain a function $L_k : \mathbb{R}^9 \rightarrow \mathbb{R}$ by taking $L : \mathbb{R}^{12} \rightarrow \mathbb{R}$ and plugging in the value t_k for the variable t and the values $x_1^{(k)}, x_2^{(k)}$ into the variables x_1, x_2 .

$$L_k = L(t_k, u_1, u_2, u_3, w_{11}, w_{12}, w_{13}, w_{21}, w_{22}, w_{23}, x_1^{(k)}, x_2^{(k)})$$

The above expression is exactly what we need to minimize. We desire a set of weights such that L_k is minimized for all $k = 1, 2, \dots, n$. This would let us know the weights we need in our neural network to fit our dataset. If only there was an algorithm to do this!

If the reader may allow us to abuse notation once more, denote L_k as L , keeping in mind we are discussing this in the context of the training example (t_k, \vec{x}_k)

We can minimize L with respect to the training example (t_k, \vec{x}_k) by calculating the gradient with respect to each of the weights. In terms of calculus, this means we are interested in the quantities

$$\frac{\partial L}{\partial u_i} \quad \frac{\partial L}{\partial w_{ij}}$$

Once we obtain these quantities, we can update our weights after reviewing one training example via

$$u'_i = u_i - \frac{\partial L}{\partial u_i} \tag{1.8}$$

$$w'_{ij} = w_{ij} - \frac{\partial L}{\partial w_{ij}} \tag{1.9}$$

First, let us calculate how the loss is affected by the weights in the final layer (i.e. the entries of the matrix U). Since $y = r(\sum_k u_k h_k)$, we have that

$$\frac{\partial L}{\partial u_i} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial u_i}$$

Observe that

$$\frac{\partial L}{\partial y} = -(t - y)$$

and

$$\frac{\partial y}{\partial u_i} = r' \left(\sum_k u_k h_k \right) \cdot \frac{\partial}{\partial u_i} \left(\sum_k u_k h_k \right) = r' \left(\sum_k u_k h_k \right) \cdot h_i$$

If we write $s_y = \sum_k u_k h_k$ (i.e. the summation before applying the activation r which calculates the value of y) then we can write

$$\frac{\partial L}{\partial u_i} = -(t - y) r'(s_y) h_i.$$

Next, let us calculate how the loss is affected by weights in the first layer (i.e. the entries of the

matrix W). Once again we have

$$\frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_{ij}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial h_i} \frac{\partial h_i}{\partial w_{ij}}$$

We already know $\frac{\partial L}{\partial y}$. Thus we calculate

$$\frac{\partial y}{\partial h_i} = r' \left(\sum_k u_k h_k \right) \cdot \frac{\partial}{\partial h_i} \left(\sum_k u_k h_k \right) \quad (1.10)$$

$$= r' \left(\sum_k u_k h_k \right) \cdot u_i \quad (1.11)$$

$$= r'(s_y) u_i \quad (1.12)$$

and since $h_i = r(\sum_k w_{ik} x_k)$ we have that

$$\frac{\partial h_i}{\partial w_{ij}} = r' \left(\sum_k w_{ik} x_k \right) \cdot \frac{\partial}{\partial w_{ij}} \left(\sum_k w_{ik} x_k \right) \quad (1.13)$$

$$= r' \left(\sum_k w_{ik} x_k \right) \cdot x_j \quad (1.14)$$

If we write $s_{h_i} = \sum_k w_{ik} h_k$, then we can write

$$\frac{\partial h_i}{\partial w_{ij}} = r'(s_{h_i}) x_j$$

Combining all of our calculations, we then get that

$$\frac{\partial L}{\partial w_{ij}} = -(t - y) \cdot r'(s_y) u_i \cdot r'(s_{h_i}) x_j$$

If we define $\delta = -(t - y)r'(s_y)$, interpreting it as an error term, we obtain the following explicit weight update formulas.

$$u'_1 = u_1 - \delta h_1 \quad (1.15)$$

$$u'_2 = u_2 - \delta h_2 \quad (1.16)$$

$$w'_{11} = w_{11} - \delta u_1 r'(s_{h_1}) x_1 \quad (1.17)$$

$$w'_{12} = w_{12} - \delta u_1 r'(s_{h_1}) x_2 \quad (1.18)$$

$$w'_{21} = w_{21} - \delta u_2 r'(s_{h_2}) x_1 \quad (1.19)$$

$$w'_{22} = w_{22} - \delta u_2 r'(s_{h_2}) x_2 \quad (1.20)$$

$$(1.21)$$

Recall that the weights u_3 (in U) and w_{13}, w_{23} (in W) correspond to the bias parameters in our model. Their weight update formulas are much simpler, since the value of their origin node is automatically

fixed to 1.

$$u'_3 = u_3 - \delta \quad (1.22)$$

$$w'_{13} = w_{13} - \delta u_1 r'(s_{h_1}) \quad (1.23)$$

$$w'_{23} = w_{23} - \delta u_2 r'(s_{h_2}) \quad (1.24)$$

$$(1.25)$$

This completes our description for how we update our model's weights given one training example. In practice, however, we tend to have many training examples t_1, t_2, \dots, t_N that we can use to update our model. There are three main approaches as to exactly how we can update our model's weights using all of the training examples.

- **Stochastic Gradient Descent.** Update the model's weights after each training example t_i .
- **Batch Gradient Descent.** Update the model's weights after showing it every training example. In this case, we could achieve this by collecting all of the gradients calculated from each training example and then averaging them.

$$w'_{ij} = w_{ij} - \frac{1}{N} \sum_{t_k} \frac{\partial L(t_k)}{\partial w_{ij}}$$

- **Mini-Batch Gradient Descent.** Update the model's weights after showing it $n < N$ many training examples; we'd call n our **batch size**. This can be thought of as a compromise between stochastic and batch gradient descent.

These three methods trade off training speed and model accuracy. As mini-batch is a compromise between speed and accuracy, it tends to be preferred in practice. In any case, in practice we tend to train the model on the training set each time; each iteration is called an **epoch**, and so a training procedure may undergo several epochs. Too few epochs will lead to poor accuracy, but too many epochs will lead to poor generalization outside of the training set, a concept known as **overfitting**.

For this example, if we perform batch gradient descent on the model we presented, using our four training examples of the XOR function, training this for tens of thousands of epochs will allow us to converge on these model weights

$$W = \begin{bmatrix} 1.52183263 & 1.55282077 & -1.55282077 \\ 1.0637993 & 1.0637993 & 0.39200322 \end{bmatrix} \quad (1.26)$$

$$U = \begin{bmatrix} -1.31420497 & 0.94002694 & -0.36849359 \end{bmatrix} \quad (1.27)$$

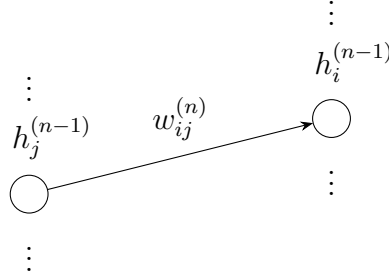
which successfully mimic the XOR function.

1.5 Backpropagation: More generally

At this point we have so far only dealt with neural networks that have at most one hidden layer and only one output node. In general, neural networks can be designed to have many hidden layers and many outputs nodes. We now present the backpropagation algorithm for deep neural networks that

have arbitrarily many output nodes.

First, we introduce some notation. Consider a neural network with N layers (including the input and output layers). Let $h_i^{(n)}$ denote the i -th hidden node in the n -th layer (hence $1 < n \leq N$). We add parentheses to the n superscript to remind the reader that it is not an exponent but is rather an index for notational purposes. Let $w_{ij}^{(n)}$ denote the weight going from node $h_j^{(n-1)}$ to $h_i^{(n)}$. We'll denote the matrix made up of the weights $w_{ij}^{(n)}$ to be W_n . Below is a snapshot of this specific part of the neural network.



With this notation, we can calculate the value of a hidden node as below:

$$h_i^{(n)} = \sigma \left(\sum_t w_{it}^{(n)} h_t^{(n-1)} \right)$$

where σ is our choice of an activation function and N -th layer (i.e. $h_k^{(N)}$) correspond to the output of our neural network.

Suppose we are given a dataset to train the weights of our neural network. Then this means that, given a cost function L , we are interested in the quantity

$$\frac{\partial L}{\partial w_{ij}^{(n)}} = \frac{\partial L}{\partial h_k^{(N)}} \frac{\partial h_k^{(N)}}{\partial w_{ij}^{(n)}}$$

where $1 \leq n \leq N$ and i, j vary as valid indices over the overall weight matrix W_n . Then we have the following result for general backpropagation.

Theorem 1.5.1 (Backpropagation). Let h_k^ℓ denote the k -th hidden unit in layer ℓ where $n \leq \ell \leq N$. Then for the weight $w_{ij}^{(n)}$, we have the recurrence relation

$$\frac{\partial h_k^\ell}{\partial w_{ij}^{(n)}} = h_k'^{(\ell)}(s) \cdot \sum_t w_{kt}^{(\ell)} \frac{\partial h_t^{(\ell-1)}}{\partial w_{ij}^{(n)}}$$

where $s = \sum_t w_{kt}^{(\ell)} h_t^{(\ell-1)}$. When $\ell = n$, we have that

$$\frac{\partial h_k^n}{\partial w_{ij}^{(n)}} = \begin{cases} h_i^n \cdot h_j^{(n-1)} & \text{if } k = i \\ 0 & \text{otherwise} \end{cases}$$

We can then use this recurrence relation to calculate the value of $\frac{\partial L}{\partial w_{ij}^{(n)}}$, allowing us to use an update

formula as below.

$$w_{ij}^{(n)} = w_{ij}^{(n)} - \frac{\partial L}{\partial w_{ij}^{(n)}}$$

The above result can be derived from simple use of the chain rule. The reason for the special case is because the weight $w_{ij}^{(n)}$ only has an affect on node $h_i^{(n)}$ in the n -th layer.