

# Lecture 9

## Artificial Neural Networks, Part 2

Alice Gao

November 2, 2021

### Contents

<b>1</b>	<b>Learning Goals</b>	<b>2</b>
<b>2</b>	<b>Gradient Descent</b>	<b>2</b>
2.1	A motivating example . . . . .	2
2.2	Introduction . . . . .	2
2.3	Steps of the algorithm . . . . .	3
2.4	Deriving some intuition for gradient descent . . . . .	3
2.5	Alternative ways to update the weights . . . . .	4
<b>3</b>	<b>The Back-propagation Algorithm</b>	<b>5</b>
3.1	Introducing the Back-Propagation Algorithm . . . . .	6
3.2	An Intuitive Description of the Back-Propagation Algorithm . . . . .	6
3.3	The Forward Pass . . . . .	7
3.4	The Backward Pass . . . . .	7
3.5	The recursive relationship in the delta values . . . . .	8
3.6	Deriving the gradients for $W_2$ . . . . .	9
3.7	Deriving the gradients for $W_1$ . . . . .	12
3.8	Defining the delta values . . . . .	13
<b>4</b>	<b>Neural Networks vs. Decision Trees</b>	<b>15</b>
4.1	When should we use a neural network? . . . . .	15
4.2	Disadvantages of neural networks . . . . .	15
4.3	Choosing a decision tree or a neural network . . . . .	15

# 1 Learning Goals

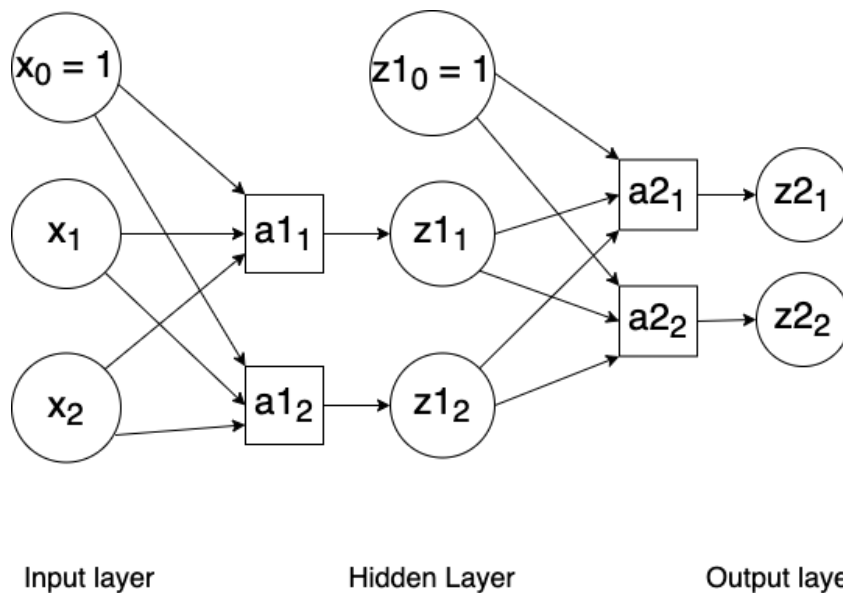
By the end of the lecture, you should be able to

- Explain the steps of the gradient descent algorithm.
- Explain how we can modify gradient descent to speed up learning and ensure convergence.
- Describe the backpropagation algorithm including the forward and backward passes.
- Compute the gradient for a weight in a multi-layer feed-forward neural network.
- Describe situations in which it is appropriate to use a neural network or a decision tree.

## 2 Gradient Descent

### 2.1 A motivating example

Throughout this lecture, we will work with this simple two-layer feed-forward neural network. The input layer has three nodes, one of which is a dummy. The hidden layer also has three nodes, one of which is a dummy. The output layer has two nodes.



Our goal is to learn the weights  $Wk_{ij}$  such that the actual outputs  $z_{21}$  and  $z_{22}$  are close to the expected values given by the training data.

### 2.2 Introduction

Gradient descent is a well-known optimization algorithm; you may have learned about it in another math class.

One way to think about gradient descent is as a local search algorithm to find the minimum of a function. Consider a search space whose dimension is given by the number of weights in the neural network. We want to find a combination of weights which minimizes the total error made on the training examples.

Gradient descent is very similar to greedy descent, but differs in that we are now dealing with continuous rather than discrete values. Accordingly, we will use the gradient, or partial derivative, to determine the direction and size of each step we want to take. The function in this case will be the error function: the difference between the expected outputs based on the training data and the actual outputs based on the training data and the network.

This quotation might give you an intuitive way to understand the algorithm:

*“Walking downhill and always taking a step in the direction that goes down the most.”*

We descend to the minimum by following the steepest direction at each step.

## 2.3 Steps of the algorithm

- Initialize the weights randomly.
- For each training example, change each weight in proportion to the negative of the partial derivative of the error with respect to that weight. Overall, we will update a weight  $w$  by

$$w := w - \sum_{\text{examples}} \eta \frac{\partial \text{error}}{\partial w}.$$

- Here,  $\eta$  is a constant called the learning rate.
- Terminate after some number of steps when the error is small enough or when the changes made get small enough.

## 2.4 Deriving some intuition for gradient descent

You may find gradient descent to be unintuitive at first, but here are some of the intuition that may help.

Consider a simple function:  $y = x^2$ . We want to minimize this function given some starting point  $x_0$ ; the minimum is at  $x = 0$ . To do so, we need to consider direction and magnitude.

First suppose we start with  $x_0$  to the right of the minimum. Then  $\frac{dy}{dx} > 0$ , but to decrease our value, we should move in the direction of the negative of the gradient.

Instead, suppose we start with  $x_0$  to the left of the minimum. Then  $\frac{dy}{dx} < 0$ , but to increase our value, we should move in the direction in the negative of the gradient.

As for direction, first note that the gradient is flat (zero) at the minimum. Then consider a region where the gradient is large in magnitude. In such a region, it is certain that we must be far from the minimum (since the function is continuous), so we can afford to take a larger step. On the other hand, in a region where the gradient is small in magnitude, we are likely

to be close to the minimum. Thus, we should take a smaller step to avoid overshooting the minimum.

This discussion is not rigorous by any means, but should provide some intuition behind the ideas and derivation of gradient descent.

## 2.5 Alternative ways to update the weights

In the steps given in a previous section, we said that gradient descent updates the weights after sweeping through all the training examples. But when the training data set is large, wouldn't the weight updates be very infrequent?

To speed up learning, we can instead update weights after each example. This variant is called **incremental gradient descent**. A related variant is **stochastic gradient descent**, where we randomly choose an example at each step instead of proceeding sequentially.

The advantage of incremental gradient descent is increased learning speed, since we update the weights much more frequently. The disadvantage is that we may not converge to the local minimum—a single example may take us away from rather than towards it.

To move the trade-off in our favour, we have another variant called **batched gradient descent**. As the name suggests, we will update the weights after a batch of examples.

We can choose the batch size. On one extreme, a batch size of one is identical to incremental gradient descent. On the other extreme, putting all examples into the batch is identical to the original gradient descent. By moving between the extremes, we can work in our favour. Often, people will start with small batches to learn quickly, then gradually increase the batch size until the weights eventually converge.

### 3 The Back-propagation Algorithm

There's a huge hype around neural networks right now. If you haven't learned the back-propagation algorithm already, you are probably pretty excited about this video.

What is the back-propagation algorithm? It is essentially an algorithm to learn the weights in a neural network by using the gradient descent optimization algorithm.

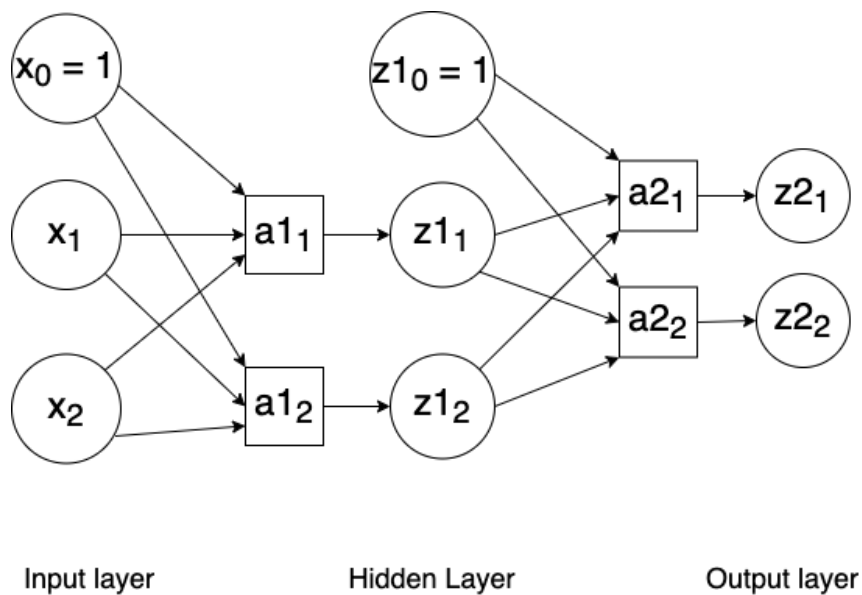


Figure 1: 2-layer feed-forward neural network

I will use a 2-layer feed-forward neural network as an example. Let me explain some notation. Let's start from the inputs  $x_i$ .  $W_{1ij}$  denotes the weights between the input layer and the hidden layer.  $a_{1j}$  denotes the  $j$ -th weighted sum of the input units.  $z_{1j}$  denotes the  $j$ -th hidden unit, and it's a result of applying the activation function  $g$  to  $a_{1j}$ .

The notation for the next layer is very similar.  $z_{1j}$  is the value of the hidden unit.  $W_{2jk}$  denotes the weights between the hidden layer and the output layer.  $a_{2k}$  denotes the  $k$ -th weighted sum of the hidden units. Finally,  $z_{2k}$  denotes the  $k$ -th output unit, and it's a result of applying the activation function  $g$  to  $a_{2j}$ .

You can also think of the network structure as a chain as shown below.

$$x_i \xrightarrow{W_{1ij}} a_{1j} \xrightarrow{g} z_{1j} \xrightarrow{W_{2jk}} a_{2k} \xrightarrow{g} z_{2k}$$

Figure 2: Chain to represent a neural network

### 3.1 Introducing the Back-Propagation Algorithm

Let's discuss the back-propagation algorithm. The most challenging step in the gradient descent algorithm is to calculate the gradient for each weight. We can write down an expression for the gradient for each weight and calculate it using the expression. Unfortunately, this approach is quite in-efficient. Nowadays, neural networks tend to be quite large, and there could be thousands or more weights that we need to learn. The back-propagation algorithm is an efficient way of calculating the gradients for the weights.

Consider a setting where we have  $n$  training examples. For each training example,  $x$  consists of the input feature values, and  $y$  is the label. In other words,  $x$  is the input to our neural network and  $y$  is the expected output.

To evaluate how close the actual output values  $z_2$  are to the expected output values  $y$ , we will use an error/loss function  $E$ . When we execute gradient descent, our goal is to minimize the error/loss  $E$  by adjusting the weights in the neural network.

Given the training examples, we will calculate the gradients by performing 2 passes in the network: a forward pass and a backward pass.

**The forward pass** takes the input values  $x$ , the current weights  $W_1$  and  $W_2$ , and calculates  $E(z_2, y)$  the error/loss  $E$  between the actual output values  $z_2$  and the expected output values  $y$ .

**The backward pass** computes the gradients  $\frac{\partial E}{\partial W_{2jk}}$  and  $\frac{\partial E}{\partial W_{1ij}}$ , which are the partial derivatives of the error function with respect to  $W_2$  and  $W_1$ . For our network, the forward pass flows from left to right, and the backward pass flows from right to left.

For each training example, we will calculate one gradient for each weight. Then, to update each weight, we need to add the gradients for this weight for all the training examples. We will update the weight proportional to this sum of the gradients.

### 3.2 An Intuitive Description of the Back-Propagation Algorithm

Let me give you an intuitive description of the algorithm. Our goal is to set the weights in the network to minimize the error/loss. How do we do that?

Using each training example, we will compute the gradient for each weight. The gradient tells us how much the error/loss changes if I change the weight by a tiny bit. If the gradient is positive, we should decrease the weight and vice versa. The gradient guides us about how we should change the weight locally to minimize the error/loss.

Why do we need to add up the gradients for all the training examples? The reason is that, different training examples may want to change each weight differently. One example may suggest that we should increase the weight, whereas another example may suggest that we decrease the weight. We do not want to only predict one example well. We would like to achieve a high prediction accuracy for all the examples. Therefore, we need to sum up the gradients for all the examples and use the sum to adjust each weight.

### 3.3 The Forward Pass

The forward pass starts from the input values on the left.

First, calculate  $a_1$  as the weighted sum of the input values using the weights  $W_1$ . The hidden unit values  $z_1$  is the activation function  $g$  applied to the weighted sum  $a_1$ .

$$a_{1j} = \sum_i x_i W_{1ij} \qquad z_{1j} = g(a_{1j}) \qquad (1)$$

We will calculate the output values in the same way.  $a_2$  is the weighted sum of the hidden unit values using the weights  $W_2$ . The output value  $z_2$  is the result of applying  $g$  to the weighted sum  $a_2$ .

$$a_{2k} = \sum_j z_{1j} W_{2jk} \qquad z_{2k} = g(a_{2k}) \qquad (2)$$

Finally, the error/loss  $E(z_2, y)$  is a function of the actual output values  $z_2$  and the expected output values  $y$ .

### 3.4 The Backward Pass

The backward pass is the core of the back-propagation algorithm.

Our goal is to calculate the gradients for the weights — the partial derivative of  $E$  with respect to  $W_1$  and  $W_2$ . We will calculate the gradients by going backwards in the network, from the outputs on the right to the inputs on the left.

Starting with the outputs on the right, we will first calculate the gradients for the weights  $W_2$ . I've written the expression as the product of two terms: the partial derivative of  $E$  w.r.t. to  $a_2$  and  $z_1$ . Let me define the first term to be  $\delta_2$ . The second term  $z_1$  is the input going into the edge for the weight  $W_2$ .

$$\frac{\partial E}{\partial W_{2jk}} = \frac{\partial E}{\partial a_{2k}} z_{1j} = \delta_{2k} z_{1j}, \qquad \delta_{2k} = \frac{\partial E}{\partial z_{2k}} g'(a_{2k}) \qquad (3)$$

At this point, it seems unnecessary to define  $\delta_2$ . However, the deltas will be extremely useful for us to understand the back-propagation process. I will define one set of  $\delta$  values for each layer. You will see shortly that the  $\delta$  values for different layers form a recursive relationship, which will allow us to calculate the gradients efficiently.

Next, let's calculate the gradients for the weights  $W_1$ . This expression is similar to that of  $W_2$ . The gradient is the product of two terms: the partial derivative of  $E$  w.r.t. to  $a_1$  and

$x$ . Similarly, I'll define the first term to be  $\delta_1$ . The second term  $x_1$  is the input going into the edge for the weight  $W_1$ .

$$\frac{\partial E}{\partial W_{1ij}} = \frac{\partial E}{\partial a_{1j}} x_i = \delta_{1j} x_i, \quad \delta_{1j} = \left( \sum_k \delta_{2k} W_{2jk} \right) g'(a_{1j}) \quad (4)$$

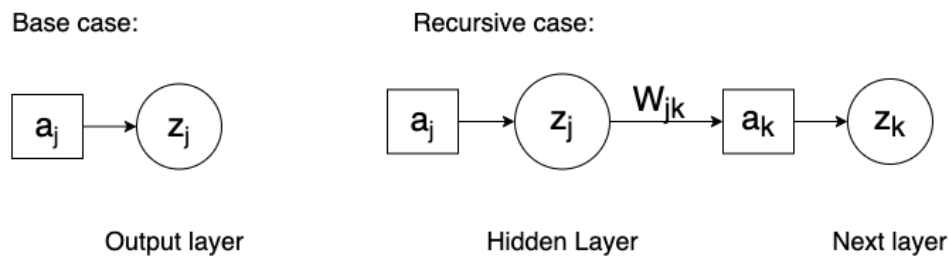
Note that the gradient for the weights in each layer has a similar expression. Each expression is a product of two terms: the delta value and the input going into the weight.

The remaining question is, how do we calculate the delta values? I am showing you the expressions for  $\delta_2$  for the output layer and  $\delta_1$  for the hidden layer. Note that we need the  $\delta_2$  values to calculate  $\delta_1$ . This gives you a hint of what the recursive relationship looks like.

### 3.5 The recursive relationship in the delta values

Let's take a closer look at the recursive relationship of the delta values. The recursive relationship allows us to propagate the error backward through the network and calculate the gradients efficiently. This recursive relationship is also how the back-propagation algorithm got its name.

Consider a unit  $j$ . This unit may be in the output layer or any hidden layer.



In general,  $\delta_j$  is the partial derivative of  $E$  w.r.t. to  $a_j$ , where  $a_j$  is the weighted sum of the values from the previous layer.

$$\delta_j = \frac{\partial E}{\partial a_j}.$$

To calculate  $\delta_j$ , we need to consider two cases.

- When  $j$  is an output unit, we have the base case and we will calculate  $\delta_j$  directly.
- When  $j$  is a hidden unit, we have the recursive case. In this case, calculating  $\delta_j$  requires  $\delta_k$ , which is the delta value for the next layer — the layer to the right, on the side closer to the output layer.



$$\delta_j = \begin{cases} \frac{\partial E}{\partial z_j} \times g'(a_j), & \text{base case, } j \text{ is an output unit} \\ \left( \sum_k \delta_k W_{jk} \right) \times g'(a_j), & \text{recursive case, } j \text{ is a hidden unit} \end{cases} \quad (5)$$

Let's think about the recursive case intuitively. The hidden unit  $j$  is connected to multiple units  $k$  in the next layer. Therefore, unit  $j$  is responsible for some fraction of the error  $\delta_k$  in each unit  $k$  that  $j$  connects to. For each unit  $k$ , we will weigh each error  $\delta_k$  by the strength of the connection between  $j$  and  $k$  — this is the weight  $W_{jk}$ . This weighted sum allows us to take the errors  $\delta_k$  from the next layer and propagate them back to calculate the error  $\delta_j$  in the current layer.

Also note that, the expressions for the two cases are quite similar. In particular, the second terms are identical — they are both the derivative of the activation function  $g$ .

Since our network only has one hidden layer, we only need to use the recursive case once. If a network has multiple hidden layers, we would need to apply the recursive case multiple times.

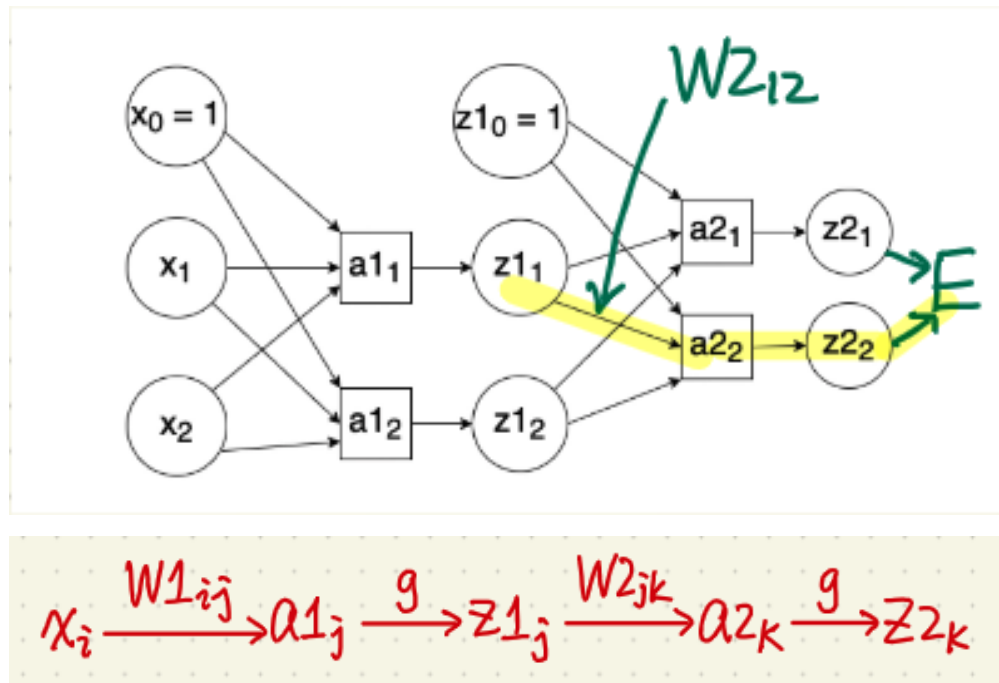
Here is a practice question for you. Construct a 3-layer neural network and write down the delta values for all three layers.

So far, I've shown you all the formulas you need to execute the back-propagation algorithm. If your goal is to implement the algorithm and run it, this is all you need. However, I hope that you are also interested in how the expressions were derived. In the next video, I will derive the partial derivatives by using the chain rule many many times. Once you understand the derivations, you will realize that the recursion with the delta values is nothing mysterious. The back-propagation process emerges directly from the derivation of the gradients.

### 3.6 Deriving the gradients for $W_2$

First, let's derive the partial derivatives for the gradients for  $W_2$ .

To remember what the network looks like, you can look at the picture or the chain. The chain is a simpler illustration of how the values flow through the network from inputs to outputs. In the picture, I have highlighted one path — the one for the weight  $W_{212}$ .



Before I start the derivations, let's do a super quick review of the chain rule. Deriving the gradient boils down to repeatedly applying the chain rule. Suppose that we have an expression in which we have applied several functions in sequence. In this case, we applied  $f$  first, and then  $g$  second. Taking derivative of the expression is equivalent to taking derivatives of the functions in reverse order, from the last one applied to the first one applied. I often think of this process as peeling off the functions one by one. For this example, the derivative is equal to the derivative of  $g$  w.r.t. to  $f$  multiplied by the derivative of  $f$  w.r.t to  $x$ .

$$\text{chain rule } \frac{d}{dx} g(f(x)) = \frac{dg}{df} \frac{df}{dx}$$

We are ready to derive our first gradient expression — the partial derivative of  $E$  w.r.t  $W_{2jk}$ . Our expression is basically tracing through the highlighted path from right to left. First,  $E$  is a function of  $z_2$  and  $y$ , so we have  $\frac{\partial E}{\partial z_{2k}}$ . Next,  $z_2$  is  $g$  of  $a_2$ , so we have  $\frac{\partial z_{2k}}{\partial a_{2k}}$ . Next,  $a_2$  is a function of the weights  $W_2$  and the hidden unit values  $z_1$ . Since we want to take derivative w.r.t. to the weights, we have  $\frac{\partial a_{2k}}{\partial W_{2jk}}$ .

**Example:**

$$\frac{\partial E}{\partial w_{2jk}} = \frac{\partial E}{\partial z_{2k}} \frac{\partial z_{2k}}{\partial a_{2k}} \frac{\partial a_{2k}}{\partial w_{2jk}} \quad (6)$$

If you have got here and understood how the expression was derived, great job! You got

through the most difficult part — deriving the most general and abstract expression for this derivative. It only gets easier from here. If you still have questions, no worries. I will show you some more concrete examples right now and they should help clear up some confusion.

Let's make some of the terms more concrete. What is the middle term?  $z_2$  is  $g(a_2)$ , so the middle term is  $g'(a_2)$ . How about the last term?  $a_2$  is the sum of  $W_2$  multiplied by  $z_1$ . So the last term should be  $z_{1j}$ , the input corresponding to the weight  $W_{2jk}$ . This expression is equivalent to the previous one, but it's easier on the eyes.

**Example:**

$$\frac{\partial E}{w_{2jk}} = \frac{\partial E}{z_2} \frac{\partial z_2}{a_2} \quad (7)$$

$$\frac{\partial a_2}{w_{2jk}} = \frac{\partial E}{z_2} g'(a_2) z_{1j} \quad (8)$$

$$z_2 = g(a_2) \Rightarrow \frac{\partial z_2}{a_2} = g'(a_2) \quad (9)$$

$$a_2 = \sum_j w_{2jk} z_{1j} \Rightarrow \frac{\partial a_2}{w_{2jk}} = z_{1j} \quad (10)$$

Let's do an example for a particular weight. Consider the weight going from  $z_1$  to  $a_2$  — the highlighted one in the picture. All we need to do is take the general expression and plug in  $j = 1$  and  $k = 2$ . Here is the result.

**Example:** An example with  $j = 1$  and  $k = 2$

$$\frac{\partial E}{w_{212}} = \frac{\partial E}{z_2} \frac{\partial z_2}{a_2} \frac{\partial a_2}{w_{212}} = \frac{\partial E}{z_2} g'(a_2) z_{11} \quad (11)$$

If you want to make this expression more concrete, you'd have to pick the error function  $E$  and the activation function  $g$ . Here are two examples. If the error function is the sum of the squared difference, then the partial derivative looks like this. If the activation function is the sigmoid function, then the partial derivative looks like this. Interestingly, the partial derivative of the sigmoid function can be written as an expression involving two copies of the sigmoid function. If you don't believe it, please verify it yourself.

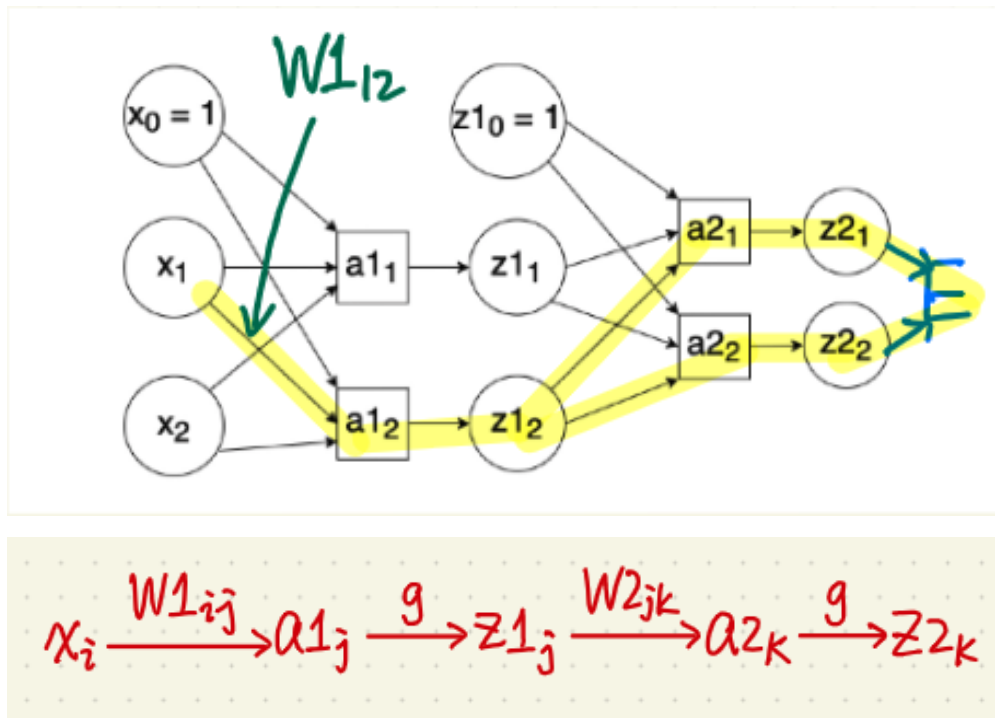
**Example:**

$$E = \frac{1}{2} \sum_k (z_{2k} - y_k)^2 \Rightarrow \frac{\partial E}{\partial z_{2k}} = (z_{2k} - y_k) \quad (12)$$

$$g(x) = \frac{1}{1 + e^{-x}} \Rightarrow g'(x) = g(x)(1 - g(x)) \quad (13)$$

### 3.7 Deriving the gradients for W1

Next, let's derive the expression for the partial derivative of E w.r.t to W1. These expressions are more complex but also more interesting.



Let's look at a specific example. Consider  $W1_{12}$ , the weight between  $x_1$  and  $z1_2$ . I've highlighted the weight in the picture. How does this weight influence the error/loss? This weight affects  $a1_2$ , which affects  $z1_2$ . From there, it affects the error/loss E through 2 paths, one through each of the two output units. Because of this, we will have a summation with 2 terms in our expression.

Let's derive our expression now. To apply the chain rule, we'll look at the highlighted path and go from right to the left. We'll start with a summation over k since we need to consider two paths, one for each output unit  $z2_k$ .

E is a function of both  $z2$  values. So the first term is  $\frac{\partial E}{\partial z2_k}$ . Next,  $z2$  is  $g(a2)$ , so we have  $\frac{\partial z2_k}{\partial a2_k}$ . After that,  $a2$  is a weighted sum of  $z1$ , so the next term is  $\frac{\partial a2_k}{\partial z1_j}$ . At this point, the two paths merge into one, so we can end our summation.

Let's keep following the single path from  $z1$ .  $z1$  is  $g$  of  $a1$ , so the next term is partial  $z1$  over partial  $a1$ .  $a1$  is a weighted sum of  $x$  using the weights W1, so the final term is partial  $a1$  over partial  $W1_{ij}$ . That's the entire expression.

**Example:**

$$\frac{\partial E}{\partial w1_{ij}} = \left( \sum_k \frac{\partial E}{\partial z2_k} * \frac{\partial z2_k}{\partial a2_k} * \frac{\partial a2_k}{\partial z1_j} \right) \frac{\partial z1_j}{\partial a1_j} \frac{\partial a1_j}{\partial w1_{ij}} \quad (14)$$

Similar to the previous derivation, you might find it helpful to make it more concrete. If we simplify some terms, we'll get the second expression.

**Example:**

$$\frac{\partial E}{\partial w_{1ij}} = \left( \sum_k \frac{\partial E}{\partial z_{2k}} * \frac{\partial z_{2k}}{\partial a_{2k}} * \frac{\partial a_{2k}}{\partial z_{1j}} \right) \frac{\partial z_{1j}}{\partial a_{1j}} \frac{\partial a_{1j}}{\partial w_{1ij}} \quad (15)$$

$$= \left( \sum_k \frac{\partial E}{\partial z_{2k}} g'(a_{2k}) * w_{2jk} \right) g'(a_{1j}) x_i \quad (16)$$

Whenever we have a derivative of  $z$  with respect to  $a$ , the result is the derivative of the activation function  $g$ . We have two other similar terms: the derivative of  $a$  with respect to  $z$  or  $W$ .  $a_2$  is a sum of  $z_1$  weighted by  $W_2$ . So partial  $a_2$  partial  $z_2$  is equal to  $W_2$ . Similarly,  $a_1$  is a sum of  $x$  weighted by  $W_1$ . So partial  $a_1$  partial  $W_1$  is equal to the input value  $x$ .

This is all the simplification we can do without knowing the expressions for the error function  $E$  and the activation function  $g$ . You can also derive the specific expression for  $W_{112}$ . Here it is.

**Example:** An example with  $i = 1$  and  $j = 2$

$$\frac{\partial E}{\partial w_{112}} = \left( \sum_k \frac{\partial E}{\partial z_{2k}} g'(a_{2k}) * w_{22k} \right) g'(a_{12}) x_1 \quad (17)$$

### 3.8 Defining the delta values

We have derived some complicated derivative expression so far. Finally, let's connect these to the delta values.

Here are the simplified gradient expressions again.

**Example:**

$$\frac{\partial E}{\partial w_{2jk}} = \frac{\partial E}{\partial z_{2k}} g'(a_{2k}) z_{1j} \quad (18)$$

$$\frac{\partial E}{\partial w_{1ij}} = \left( \sum_k \frac{\partial E}{\partial z_{2k}} g'(a_{2k}) w_{2jk} \right) g'(a_{1j}) x_i \quad (19)$$

Note that, the last term in each expression is the input value for that layer. Let's disregard the last term in both expressions. Take the rest of the expression and define it to be delta. I've denoted them as  $\delta_{2k}$  and  $\delta_{1j}$ .

**Example:**

$$\frac{\partial E}{\partial w_{2jk}} = \frac{\partial E}{\partial z_{2k}} g'(a_{2k}) z_{1j} = \delta_{2k} z_{1j} \quad (20)$$

$$\text{where } \delta_{2k} = \frac{\partial E}{\partial z_{2k}} g'(a_{2k}) \quad (21)$$

$$\frac{\partial E}{\partial w_{1ij}} = \left( \sum_k \frac{\partial E}{\partial z_{2k}} g'(a_{2k}) w_{2jk} \right) g'(a_{1j}) x_i = \delta_{1j} x_i \quad (22)$$

$$\text{where } \delta_{1j} = \left( \sum_k \frac{\partial E}{\partial z_{2k}} g'(a_{2k}) w_{2jk} \right) g'(a_{1j}) \quad (23)$$

$$\text{and } \delta_{2k} = \sum_k \frac{\partial E}{\partial z_{2k}} g'(a_{2k}) \quad (24)$$

$$\text{or } \delta_{1j} = \left( \sum_k \delta_{2k} w_{2jk} \right) g'(a_{1j}) \quad (25)$$

Note that,  $\delta_{2k}$  appears inside the expression for  $\delta_{1j}$ . If we re-write the delta expressions separately, we get the delta expressions that I showed you in the previous section.

**Example:**

$$\delta_{2k} = \frac{\partial E}{\partial z_{2k}} g'(a_{2k}) \quad (26)$$

$$\delta_{1j} = \left( \sum_k \delta_{2k} w_{2jk} \right) g'(a_{1j}) \quad (27)$$

Here is a final practice question for you. Construct a 3-layer network with 2 hidden layers. Each hidden layer has 2 real nodes and 1 dummy node. The output layer has 2 nodes. Calculate the gradients for all the weights in the 3 layers. Write down the general expressions and then write down the expressions using the delta values. When using the delta values, be sure to write down the expressions for the delta values for every layer. If you understood everything in this video, you should be able to complete this question.

## 4 Neural Networks vs. Decision Trees

In our short machine learning unit, we discussed two models: decision trees and neural networks. When should you use one versus the other?

### 4.1 When should we use a neural network?

- Neural networks are well-suited to high-dimensional or real-valued inputs, or noisy (sensor) data. You are likely familiar with the most recent stories of neural networks in applications with such data: images, videos, audio, and so on.
- A second situation is when the form of the target function is unknown (no model). There is a theorem which says that with enough levels and nodes, a neural network can approximate any arbitrary function. In the case that the function is unknown, you may want a very flexible model which allows you to discover the true function.
- The third case for using a neural network is when it is not important for humans to explain the learned function. Once we learn a neural network, it usually has high prediction accuracy but tends not to be interpretable.

### 4.2 Disadvantages of neural networks

- It is difficult to determine the network structure. There are many layer and many neurons.
- It is difficult to interpret weights, especially in multi-layer networks.
- Neural networks tend to over-fit in practice (predict poorly outside of the training data) since they have very high variance.

### 4.3 Choosing a decision tree or a neural network

There are many factors to consider when choosing between a decision tree and a neural network. Each model will work better under different circumstances. The following table provides a summary of the differences.

---

	<b>Decision trees</b>	<b>Neural networks</b>
Data types	Good for tabular data (rows and columns).	Good for images, audio, text, etc.
Size of data set	Work well with little data.	Require lots of data, but overfit easily.
Form of target function	Model nested if-then-else statements.	Model arbitrary functions.
Architecture	Only a few parameters.	Lots of parameters, all are critical.
Interpreting the learned function	Easy to understand.	Essentially a black box: difficult to interpret.
Time available for training and classification	Fast to train and classify.	Slow to train and classify.

---

You may be tempted to opt for a neural network because of how powerful it is, but bear in mind the drawbacks. It is often easier to begin with a simpler model, then increase the complexity as needed. Sometimes, a simpler model might do just as well as a complicated model, with less time required to implement and test.