

Deep Learning

by DeepLearning.AI

Lucas Paiolla Forastiere

August 19, 2021

Contents

1	Introduction	2
1.1	Notation	4
1.2	Logistic Regression as a Neural Network	6
1.2.1	Gradient Descent	7
2	Basic Concepts	10
2.1	Representation	10
2.2	Activation Functions	12
2.3	Backpropagation	13

Chapter 1

Introduction

The term deep learning refers to training *neural networks*, sometimes very big neural networks. But what are neural networks?

So let's suppose we want to predict housing prices based on the size of the house. And let's say we'll use Logistic Regression to do that. But as we know, house prices can't be negative, so we simply say the value of the house is 0 if the Logistic Regression would predict something negative.

That's indeed the simplest neural network we can have, we have a single input **size** and a single output **price** and in the middle we have a single neuron: the logistic regression.

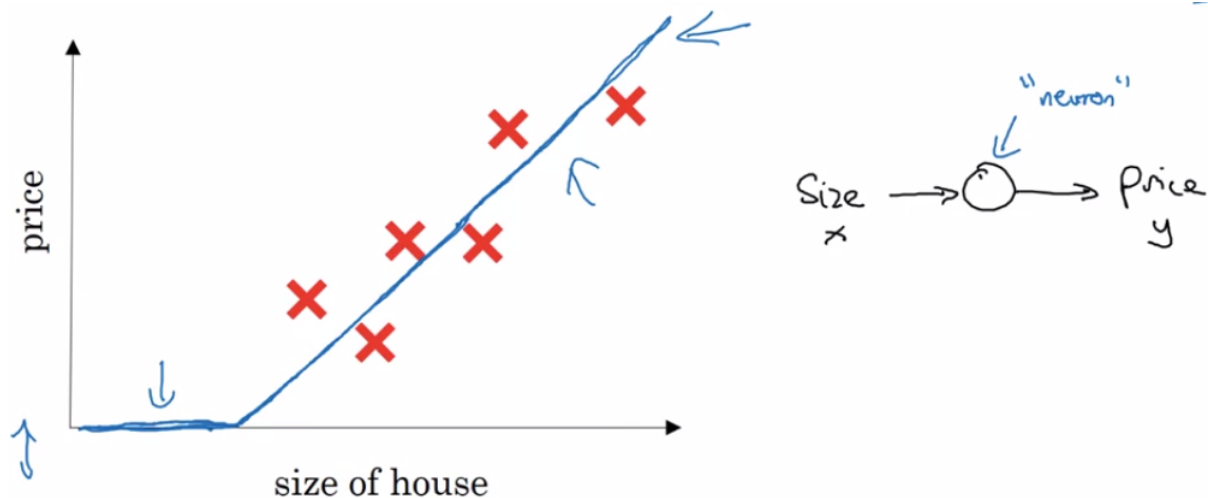


Figure 1.1: Here we see the graph of the problem we described.

That function which is zero and then linear is called *ReLU* and it's used a lot in neural networks. It stands for *Rectified Linear Unit*.

So to get a bigger neural network, we stack these neurons. Instead of predicting using only the size of house, we could use the number of bedrooms, zip code and wealth. We could use the size and number of bedrooms to predict the family size; use the zip code to predict the walkability; and use the zip code and wealth to predict the school quality. And then, we could use the family size, walkability and school quality to predict the price. See in the picture:

However, in general what we have is something a little more complex than that. We would have something like figure 1.3. Here we see that the internal nodes (which are called **hidden nodes** or **hidden neurons** or **hidden units**) receive the output of all the

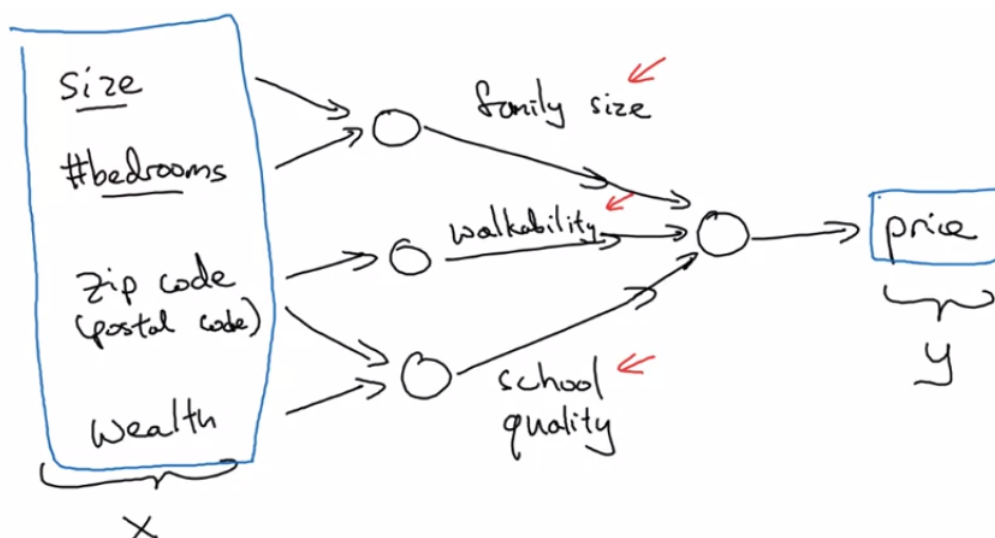


Figure 1.2: Now we have a more complex neural network, which is the stack of many ReLUs.

previous nodes to make its predictions. These hidden nodes don't really have a meaning like the example we gave. We don't try to predict family size or walkability or whatever, we simply let the neural network decide what that neuron will output in order to predict the final output *price* in the better way it can.

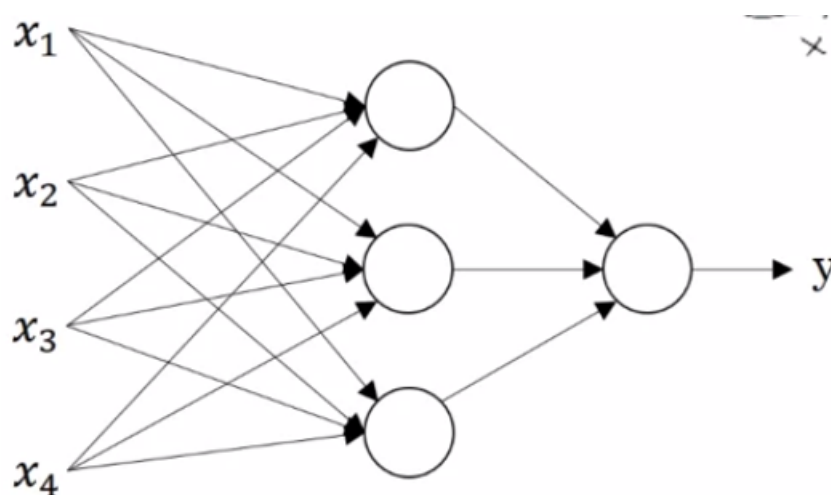


Figure 1.3: The generic form of a neural network.

We can use neural networks in many applications, here we're going to focus on **supervised learning**, which are problems that you have a set of variables called input (represented by x) and an output (y) related to that input. In order to solve these kind of problems, there are many kinds of neural networks. The one we saw is the most common one, but there are others, like convolutional nn or recurrent nn.

Another thing that's important to decide what kind of nn we'll use is knowing if the data we're dealing with is *structured* or *unstructured*.

Structured data is data in the form of a table. We have a very clear set of input variable X and a set of output variables y . Each line of our table represents one instance

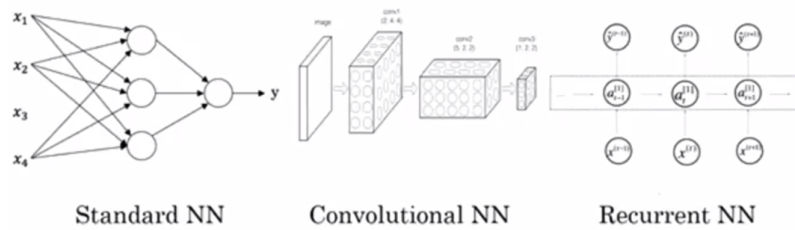


Figure 1.4: Examples of neural networks.

of data with many inputs and one or more outputs related to those inputs.

Unstructured data is all the other kinds of data: audio, video, texts, images, etc.

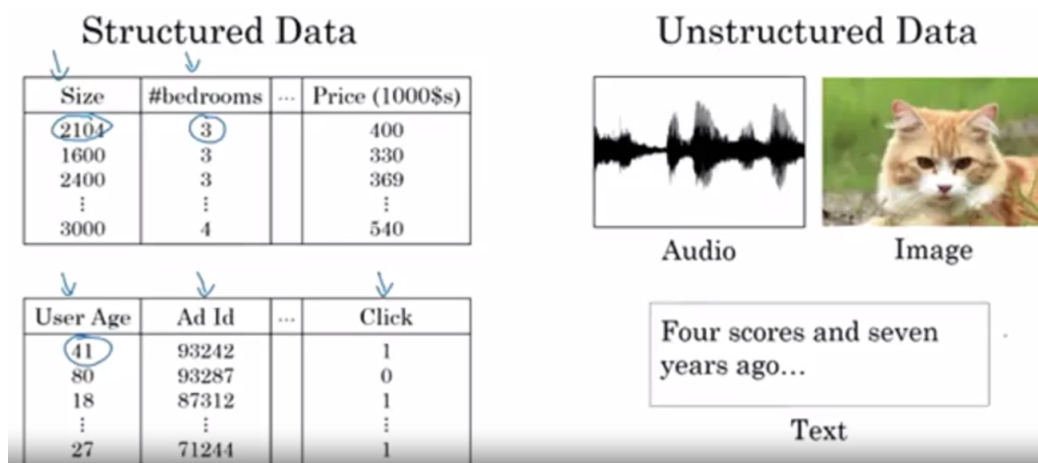


Figure 1.5: The two kinds of data.

It turns out that machine learning algorithms performed better on structured data over the years and more recently neural networks are performing better also on unstructured data.

Why is Deep Learning taking off? This is one of the questions we must ask ourselves when beginning to learn deep learning. Let's see the graph of the performance of the machine learning algorithms versus the amount of data that we provide to them. We see that traditional learning algorithms have a plateau where they can't improve anymore, which neural networks can lead with that data as we make them bigger and bigger.

We also see in the graph that when we don't have a large amount of data, NNs all algorithms perform pretty much the same.

So in order to answer our question, we have to understand the evolution of three things: *data*, *computation* and *algorithms*.

Through the years, the amount of data available was increased a lot, so NNs can take advantage from that. Also the computation power was increased with the use of GPUs to make a large amount of computations. And finally new algorithms have been developed to make NNs faster. That's the main reason why deep learning is taking off.

1.1 Notation

Before continuing, we need to define the notation we're going to use.

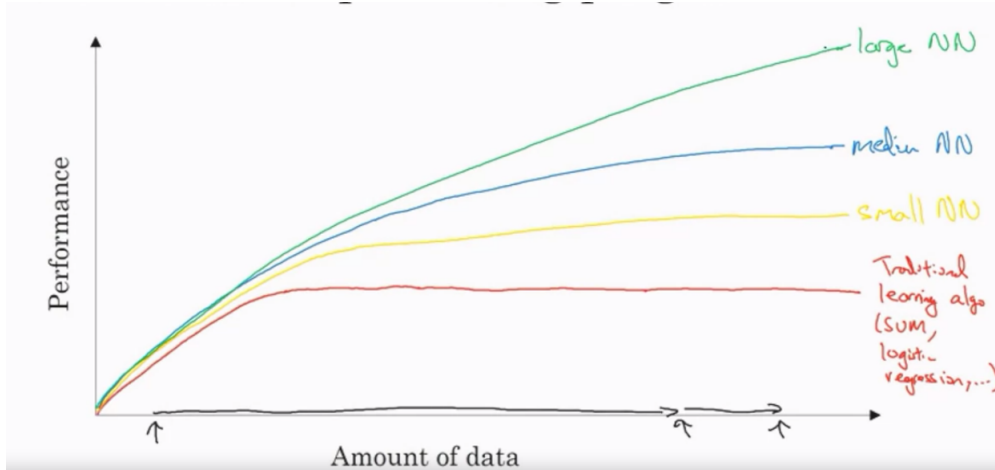


Figure 1.6: The performance of machine learning algorithms in respect to the data we provide to them.

- (x, y) will denote a single training input;
- m or m_{train} denotes the number of training examples;
- $x^{(i)}$ denotes the i -th training input;
- $y^{(i)}$ denotes the i -th training output;
- n_x or n denotes the number of dimensions x has (or the number of features);
- m_{test} denotes the number of testing examples;
- X is the matrix of all training examples. It's defined as:

$$X = \begin{bmatrix} | & | & \dots & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & \dots & | \end{bmatrix}$$

X is an $m \times n$ matrix;

- Y is the matrix of all outputs. It's defined as:

$$Y = [y^{(1)} \quad y^{(2)} \quad \dots \quad y^{(m)}]$$

Y is a $1 \times m$ matrix.

Observation. In other courses we might see X defined as the transpose of the matrix we've just defined. But it turns out that when using this definition, it's much easier to implement algorithms, so remember the definition we're going to use throughout the course.

The same thing for Y . We see that here Y is the transpose of that it's tends to be in other courses.

1.2 Logistic Regression as a Neural Network

To end this introduction, we'll see the basics of neural network programming using the simplest NN we can: a logistic regression.

So let's recall what's logistic regression and why it's useful. Logistic Regression is used in binary classification, the kind of problem where we have an input and want to predict between 0 or 1. An example could be an image and we want to say it what's a cat (1) or not (0).

Basically we want an algorithm to estimate the probability of $y = 1$ given x . In math we write:

$$\hat{y} = P(y = 1 | x), \quad x \in \mathbb{R}^n$$

Logistic Regression estimates this quantity using the formula:

$$\hat{y} = \sigma(w^T x + b),$$

where w and b are parameters to be discovered and σ is the **sigmoid function**:

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

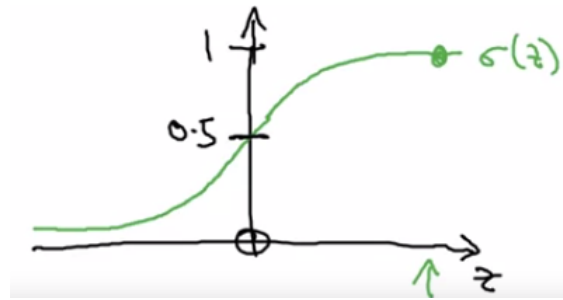


Figure 1.7: A sigmoid graph.

It's also common to create a new input $x_0 = 1$ and use the x vector as $x \in \mathbb{R}^{n+1}$ and use the formula $\hat{y} = \sigma(\theta^T x)$, where

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \quad \theta_0 = b \quad w = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

To find the parameters b and w , we need to define a **cost function**, which is a function that says how badly our algorithm is performing. This is a function that we want to minimize and when we minimize, we find the best values of b and w .

The **cost** function is a function of all training examples, while a **loss function** or **error function** is a function of a single training example that measures how well our algorithm is performing.

For logistic regression, we use the loss function:

$$\mathcal{L}(\hat{y}, y) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

Notice that this is the same as:

$$\mathcal{L}(\hat{y}, y) = \begin{cases} -\log(\hat{y} - 1), & \text{if } y = 0 \\ -\log \hat{y}, & \text{if } y = 1 \end{cases}$$

That give us the cost function:

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

1.2.1 Gradient Descent

We know have:

- A way of predicting the classes 0 or 1 using the sigmoid function;
- A way of measuring the error of our predictions.

What we need now is a way of changing our parameters b and w in order to minimize the error. That's what the **gradient descent** algorithm does.

Let's first see a general graph of the cost function. In general, it looks like figure 1.8

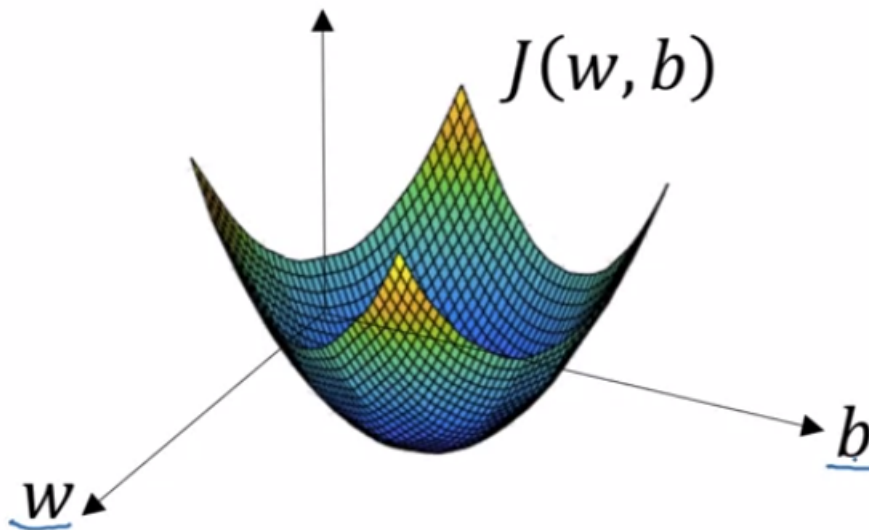


Figure 1.8: A generic graph of the cost function.

We see that J is what we call a **convex function**, which means that it is a function that has only one **local minimum** (or local maxima). This property is very important if we want to apply the gradient descent algorithm.

In Gradient Descent, we initialize b and w randomly and take steps into the direction that leads us to the lowest possible value of J . In order to do that, we calculate the **gradient** (the derivatives in each direction) of the function J and take a step in the opposite direction of the gradient.

Proposition 1.1

The gradient gives us the direction of the maximum increase of the a function.

Algorithm 1.1: Gradient Descent

```

Repeat {
     $w := w - \alpha \frac{\partial w}{\partial J(w, b)}$ 
     $b := b - \alpha \frac{\partial b}{\partial J(w, b)}$ 
}

```

In the algorithm, α is what we call the **learning rate**. It's how large we should step in the direction of the maximum decrease. If we take big steps, we can go much faster to the global minimum, but we might not be so accurate. On the other hand, we take small steps, we can find the global minimum accurately, but achieve it much slower.

Gradient Descent is a general optimization algorithm and can be applied to any convex function. So now we need to understand how to use it with logistic regression.

After calculating the derivatives, we'll have:

Algorithm 1.2: Gradient Descent for Logistic Regression

```

Repeat {
     $J = 0; dw = 0; db = 0$ 
    For  $i = 1 \dots m$ :
         $z^{(i)} = w^T x^{(i)} + b$ 
         $a^{(i)} = \sigma(z^{(i)})$ 
         $J += -[y^{(i)} \log a^{(i)} + (1 - y^{(i)}) \log(1 - a^{(i)})]$ 
         $dz^{(i)} = a^{(i)} - y^{(i)}$ 
         $dw += x^{(i)} dz^{(i)}$ 
         $db += dz^{(i)}$ 
     $J / = m$ 
     $dw / = m; db / = m$ 
     $w := w - \alpha dw$ 
     $b := b - \alpha db$ 
}

```

This version of the algorithm uses a for loop to compute J , dw and db . But when implementing the code into Python or other language, we always try to **vectorize** the code to make it faster.

Algorithm 1.3: Gradient Descent for Logistic Regression Vectorized

```

Repeat {
     $Z = w^T X + b$ 
     $A = \sigma(Z)$ 
     $dZ = A - Y$ 
     $dw = \frac{1}{m} X dZ^T$ 
     $db = \frac{1}{m} \sum dZ$ 
     $w := w - \alpha dw$ 
     $b := b - \alpha db$ 
}

```

}

Chapter 2

Basical Concepts

So let's start introducing some notation.

Whenever we use:

$$z^{[i]}$$

we're talking about the z values of the i -th *layer* of the neural network.

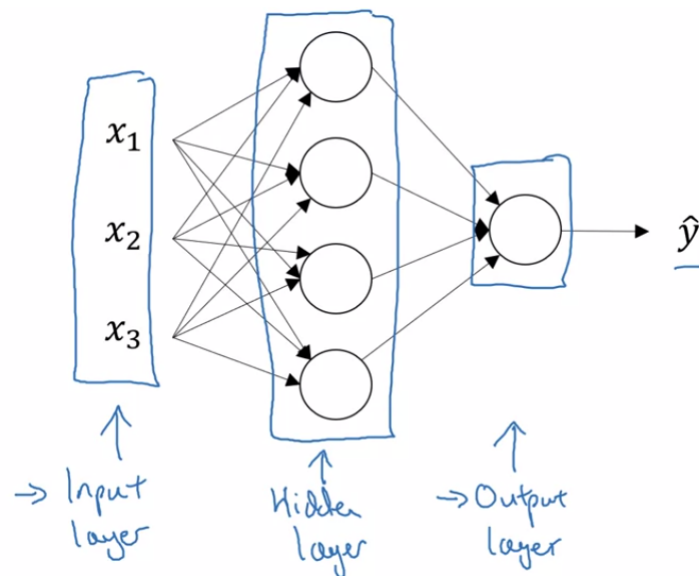
Thus, we have the following important equation:

$$z^{[i]} = W^{[i]}a^{[i-1]} + b^{[i]} \quad (2.1)$$

$$a^{[i]} = \sigma(z^{[i]}) \quad (2.2)$$

2.1 Representation

Let's illustrate a simple neural net:



This nn is divided into three parts:

- The **input layer** is the layer where we enter the data, the features. It's represented by a single vector x , which is a column of our inputs matrix X . Another way to denote the input is by $a^{[0]}$;

- The **output layer** is the layer where we get our answer. It can be one or more nodes and it's represented by a single vector y . It can be a binary vector, continuous vector, it could be even a codification of an image, sound, word or any codification just like the input;
- Finally, the **hidden layer** are all the nodes which are between the input and output, they are the intermediate computations that our neural net does. We don't need to stick to just one hidden layer, there could be many of them.

One more thing to keep in mind is that the neural network drawn in the picture is called a *two layer* neural network, because we don't count the input layer as a "real" layer (it's the layer zero).

Each layer (real layer) will have **parameters** associated with them, which we'll denote $W^{[i]}$ and $b^{[i]}$ for the i -th layer.

To understand what a node computes given an input, let's focus in the first node of the hidden layer in the figure. We can see that the inputs x_1, x_2, x_3 are all passed to that node, which uses them to output something to the node in the output layer (or it could be outputted to the next hidden layer).

So basically a node is a logistic regression, it computes:

$$a = \sigma(Wx + b)$$

But as we have many computations like this, we have to use indices to denote the first hidden layer and the first node of that layer. So the computations would be described as following:

$$\begin{aligned} z_1^{[1]} &= W_1^{[1]T} x + b_1^{[1]} \\ a_1^{[1]} &= \sigma(z_1^{[1]}) \end{aligned}$$

The superscript means we're talking about the first (real) layer and the subscript means we're talking about the first node of that layer. Therefore $a_i^{[l]}$ is the **activation value** (output value) of the i -th node in the l -th layer.

Of course, when coding a neural network, we don't compute each one of these equations using a for loop, we vectorize in order to compute the whole vector $z^{[l]}$ at once using equation 2.1.

Vectorizing the input Now, one more thing that we want is to be able to predict *multiple* inputs at once. Indeed we can do that with some vectorization. Let's see how.

So let's recall that

$$X = \begin{bmatrix} | & | & \dots & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix}$$

and let's define a matrix $Z^{[l]}$ in which the j -th column is the vector $z^{[l](j)}$ (i.e., the value before the activation function of the layer l when applied to the input j).

$$Z^{[l]} = \begin{bmatrix} | & | & \dots & | \\ z^{[l](1)} & z^{[l](2)} & \dots & z^{[l](m)} \\ | & | & & | \end{bmatrix}$$

and also define the matrix $A^{[l]}$, which is $\sigma(Z^{[l]})$:

$$A^{[l]} = \begin{bmatrix} | & | & & | \\ a^{[l](1)} & a^{[l](2)} & \dots & a^{[l](m)} \\ | & | & & | \end{bmatrix}$$

We can see that each row of matrix A tell us the activation value for a particular node of that layer for each input. For instance, the value $A_{i,j}^{[l]}$ is the value of the i -th neuron in the l -th layer of the neural network when we input the j -th input.

Given these matrices, we can train over all inputs using the vectorized formula:

$$\boxed{Z^{[l]} = W^{[l]}A^{[l-1]} + b^{[l]}} \quad (2.3)$$

$$\boxed{A^{[l]} = \sigma(Z^{[l]})} \quad (2.4)$$

and recall that $X = A^{[0]}$.

2.2 Activation Functions

The **activation function** is the function g applied to z . We've been currently using the *sigmoid function* σ , but there are some other functions that can be used and can significantly change our results and performances.

Another function one can use is the tanh function, given by the formula:

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

which is a shifted version of the sigmoid function. We have $\tanh(0) = 0$ instead of $\sigma(0) = 0.5$.

It turns out that the tanh function frequently works better than the sigmoid, because the values are between -1 and 1 instead of 0 and 1 and thus the average values tend to be closer to zero.

The exception is for the output layer, where it's good to have a number between 0 and 1 in many cases, so we could keep using the sigmoid function there and use the tanh on the hidden layers.

One of the down sides of these two functions is that when z is very large, the gradient is very small, because they're very flat for large (positive or negative) values of z .

That's why many times we see the **ReLU** (rectified linear unit) function being used:

$$ReLU(z) = \max(0, z)$$

Indeed the ReLU is the most used activation in practice and should be the first choise. One of the only down sides of it is that the derivative for values lower than zero is zero, but it works fine in practice.

There's modified version of the ReLU called **Leaky ReLU**, which is given by:

$$LReLU(z) = \max(0.01z, z)$$

It's just like the ReLU function, but it's not zero for any negative value of z and, thus, the derivative is not null on those points. Actually, indeed there no reason for the 0.01 , it could be any small value and we can turn that into a hyperparameter of our algorithm.

Why activation functions Now that we've seen so many activation functions we might want to understand why do we even need them. And most importantly, why do we need them to be non-linear.

We could try using $a = g(z) = z$ (i.e., doing nothing). It turns out that if we have just the identity function, we can't express *complex* (non-linear) decision boundaries.

We'll not prove it here, but if we just use a linear function in all hidden nodes and a sigmoid, it is equivalent to a standard logistic regression in terms of what it can express.

$$a^{[1]} = w^{[1]}x + b^{[1]}$$

$$\begin{aligned} a^{[2]} &= w^{[2]}a^{[1]} + b^{[2]} \\ &= w^{[2]}(a^{[1]} = w^{[1]}x + b^{[1]}) + b^{[2]} \\ &= (w^{[2]}w^{[1]})x + (w^{[2]}b^{[1]} + b^{[2]}) = \\ &= w'x + b' \end{aligned}$$

Above we can see that from one layer to the next we still have a linear equation (a composition of two linear functions is a linear function).

2.3 Backpropagation

Now we're going to understand how to find the optimal values for W and b using backpropagation.

The first step in this direction is calculating the derivative of the activation functions. So let's do it.

Sigmoid

$$\begin{aligned} g(z) &= \frac{1}{1 + e^{-z}} \\ \frac{dg(z)}{dz} &= \frac{1}{1 + e^{-z}} \left(1 - \frac{1}{1 + e^{-z}} \right) \\ &= g(z)(1 - g(z)) \end{aligned}$$

Tanh

$$\begin{aligned} g(z) &= \tanh(z) \\ \frac{dg(z)}{dz} &= 1 - (\tanh(z))^2 \\ &= 1 - g(z)^2 \end{aligned}$$

ReLU

$$\begin{aligned} g(z) &= \max(0, z) \\ \frac{dg(z)}{dz} &= \begin{cases} 0, & \text{if } z < 0 \\ 1, & \text{if } z > 0 \\ \text{undef}, & \text{if } z = 0 \end{cases} \end{aligned}$$

Leaky ReLU

$$g(z) = \max(0.01z, z)$$

$$\frac{dg(z)}{dz} = \begin{cases} 0.01, & \text{if } z < 0 \\ 1, & \text{if } z > 0 \\ \text{undef}, & \text{if } z = 0 \end{cases}$$

Now let's understand how to do **gradient descent**.

Let's denote L the number of layers of our network. Then, the last value derivative (the first we're going to compute) is given by the formula:

$$dz^{[L]} = a^{[L]} - y \quad (2.5)$$

The values of dW and db , for any layer are given by the formulas:

$$dW^{[l]} = dz^{[l]} a^{[l-1]T} \quad (2.6)$$

$$db^{[l]} = dz^{[l]} \quad (2.7)$$

The other values of dz for any layer other than the output layer are given by the formulas:

$$dz^{[l]} = W^{[l+1]T} dz^{[l+1]} * g^{[l]'}(z^{[l]}) \quad (2.8)$$

where $*$ denotes the *element-wise* product.

There's also a vectorized way of computing this for all the training examples at once.

$$dZ^{[L]} = A^{[L]} - Y \quad (2.9)$$

$$dW^{[l]} = \frac{1}{m} dZ^{[l]} A^{[l-1]T} \quad (2.10)$$

$$db^{[l]} = \frac{1}{m} \text{sum}(dZ^{[l]}, \text{axis}=1) \quad (2.11)$$

$$dZ^{[l]} = W^{[l+1]T} dZ^{[l+1]} * g^{[l]'}(Z^{[l]}) \quad (2.12)$$

2.4 Random Initialization

Finally, we need to initialize the parameters of the neural network. It turns out that initializing everything as zero or everything as the same number, then all nodes of a layer will be exactly the same even with many iterations of gradient descent.

For the values of b , we can initialize with zero, but for W , we prefer to initialize with *small* random numbers. We like the numbers to be small though because if we're using \tanh or σ , big numbers will have a very flat derivative and, thus, the algorithm will learn slower.