

# Deep Learning

by DeepLearning.AI

Lucas Paiolla Forastiere

September 16, 2021

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Notation . . . . .	4
1.2	Logistic Regression as a Neural Network . . . . .	6
1.2.1	Gradient Descent . . . . .	7
<b>2</b>	<b>Basic Concepts</b>	<b>10</b>
2.1	Representation . . . . .	10
2.2	Activation Functions . . . . .	12
2.3	Backpropagation . . . . .	13
2.4	Random Initialization . . . . .	14
2.5	Notation for deep neural networks . . . . .	14
<b>3</b>	<b>Practical Aspects of Deep Learning</b>	<b>16</b>
3.1	Regularization . . . . .	16
3.2	Optimizing the Problem . . . . .	17
3.2.1	Normalizing input . . . . .	17
3.2.2	Vanishing / Exploding Gradients . . . . .	18
3.2.3	Gradient Checking . . . . .	19
3.3	Optimization Algorithms . . . . .	19
3.3.1	Mini-batch Gradient Descent . . . . .	19
3.3.2	Exponentially weighted averages . . . . .	20
3.3.3	Gradient Descent with Momentum . . . . .	22
3.3.4	RMSprop . . . . .	22
3.3.5	Adam . . . . .	23
3.3.6	Learning Rate Decay . . . . .	24
3.4	Hyperparameter tuning . . . . .	24
3.4.1	Selection the appropriate scale for hyperparameters . . . . .	25
3.5	Batch normalization . . . . .	25
3.6	Multiclass problems . . . . .	26
3.7	Deep Learning Frameworks . . . . .	26

# 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.

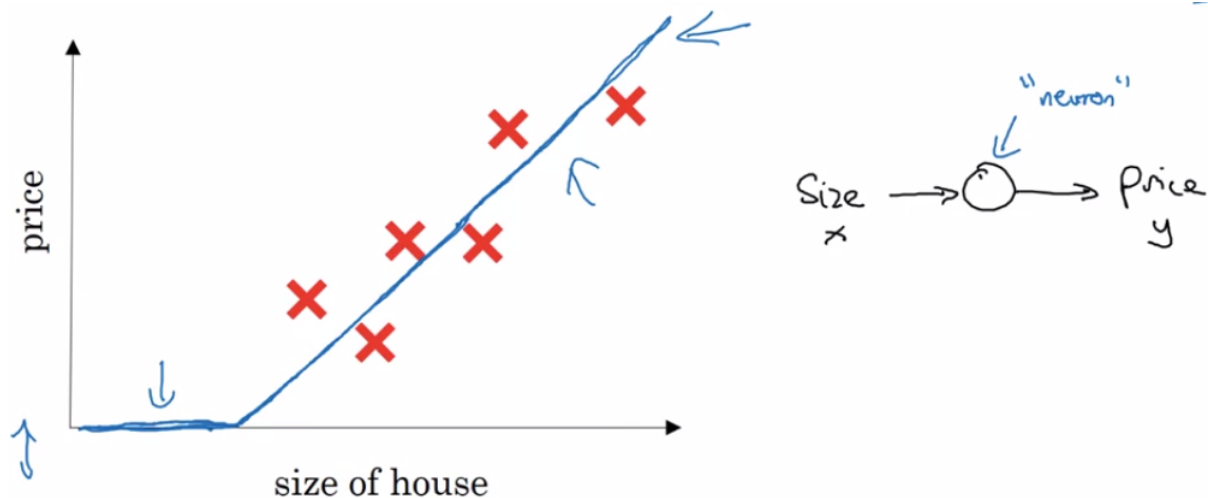


Figure 1.1: Here we see the graph of the problem we described.

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

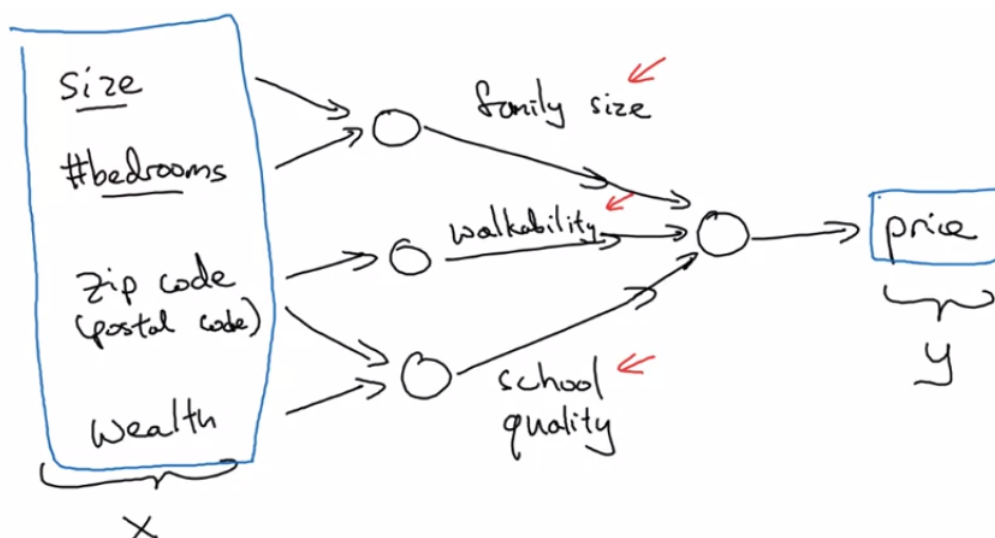


Figure 1.2: Now we have a more complex neural network, which is the stack of many ReLUs.

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.

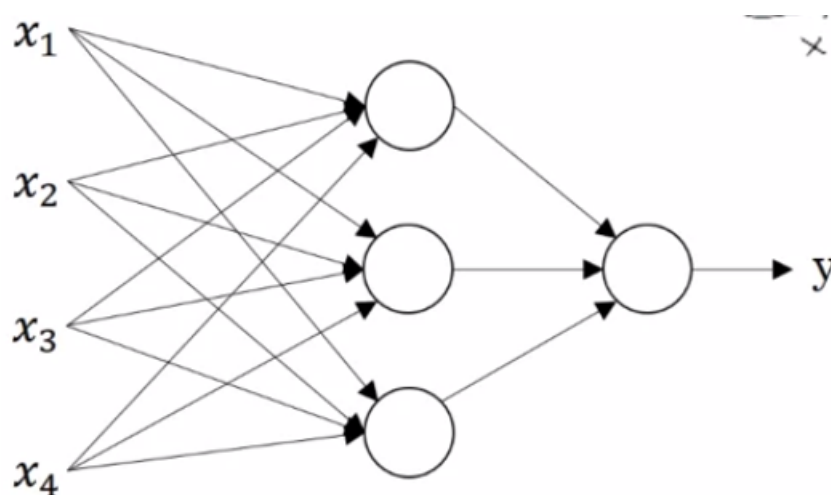


Figure 1.3: The generic form of a neural network.

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 kinds 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

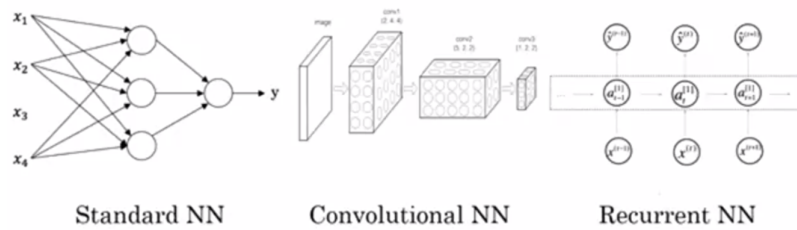


Figure 1.4: Examples of neural networks.

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.

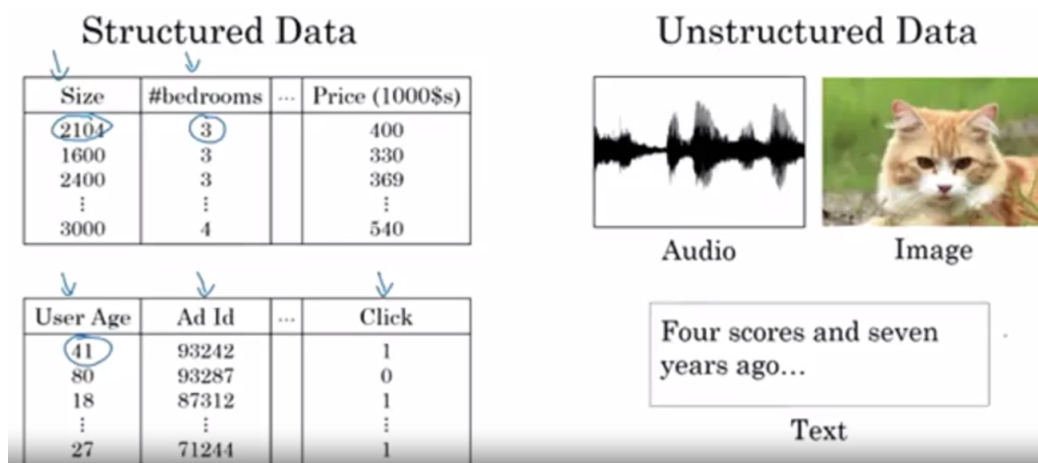


Figure 1.5: The two kinds of data.

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.

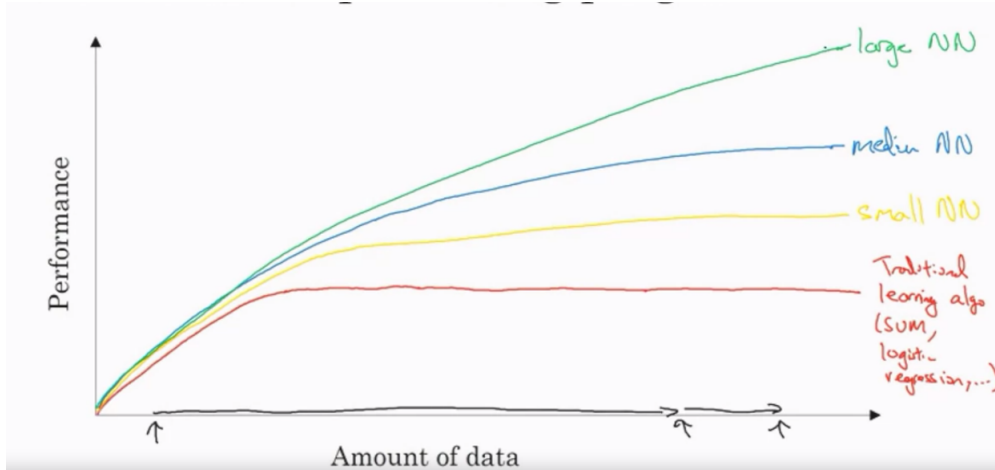


Figure 1.6: The performance of machine learning algorithms in respect to the data we provide to them.

- $(x, y)$  will denote a single training input;
- $m$  or  $m_{\text{train}}$  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);
- $m_{\text{test}}$  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  $\sigma$  is the **sigmoid function**:

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

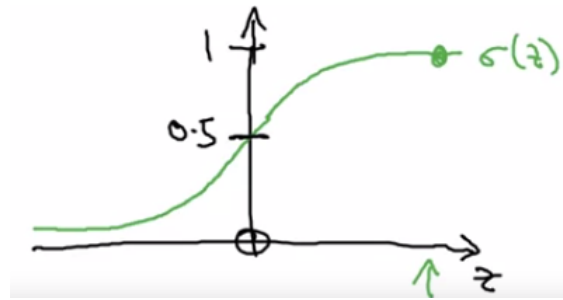


Figure 1.7: A sigmoid graph.

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

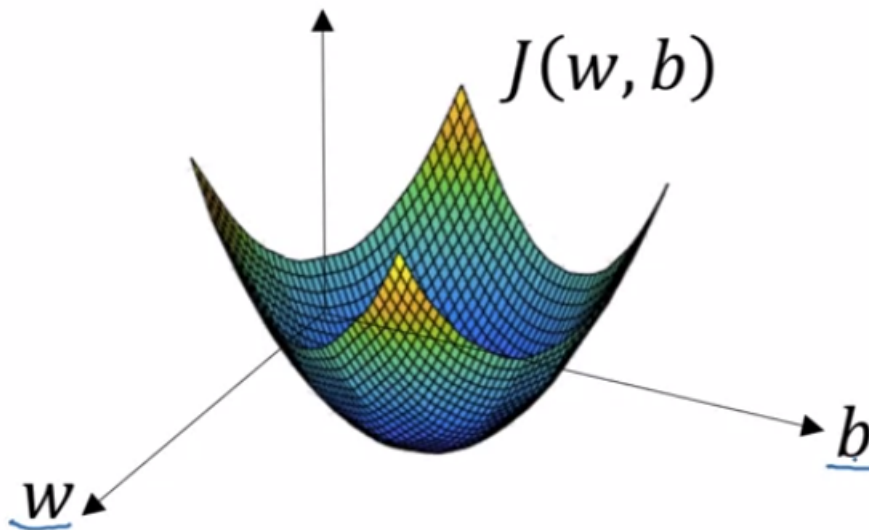


Figure 1.8: A generic graph of the cost function.

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,  $\alpha$  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.

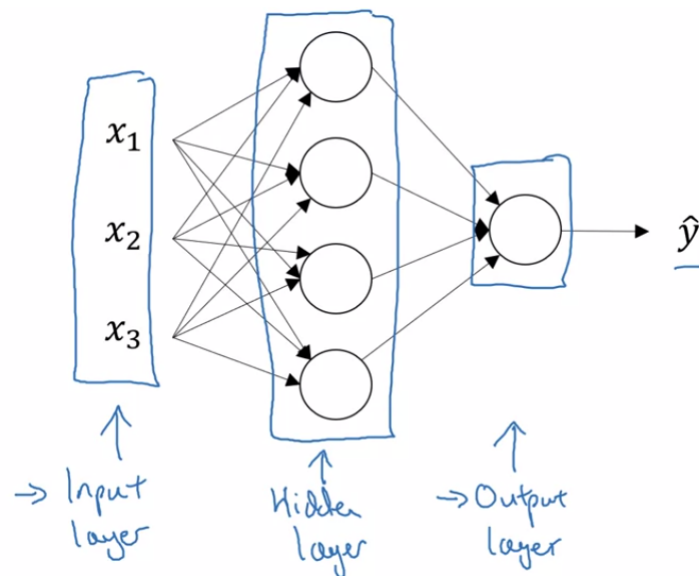
Thus, we have the following important equation:

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

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

### 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

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

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  $j$ ).

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

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*  $\sigma$ , 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 choice. 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 they 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  $\sigma$ , big numbers will have a very flat derivative and, thus, the algorithm will learn slower.

## 2.5 Notation for deep neural networks

We'll use  $L$  to describe the number of layers the neural network have (remember, the input layer is not a real layer).

$n^{[l]}$  denotes the number of nodes on layer  $l$ .

$a^{[l]}$  denotes the activation values on layer  $l$ .

$g^{[l]}$  denotes the activation function used on layer  $l$ .

Also, the dimensions of  $W^{[l]}$  and  $b^{[l]}$  are  $(n^{[l]}, n^{[l-1]})$  and  $(n^{[l]}, 1)$ .



# Chapter 3

## Practical Aspects of Deep Learning

In this chapter we're going to focus on the practical aspects of deep learning, like how to choose the hyperparameters, how to make the code faster, how to apply regularization and others.

The first thing to understand is that, in practical, we divide our data set into three:

- *Training set*, in which we're going to train the model;
- *Hold-out cross validation set (or dev set)*, in which we're going to evaluate our model and choose hyperparameters;
- *Test set*, in which we're going to evaluate our final model, with the hyperparameters setted.

In the previous era of Machine Learning, we would use 70% for the training/dev sets and 30% for testing, or 60/20/20%. But now in the big data era, we can use much small fractions for the testing and hold-out sets (like 1% or even less, since this can be equivalent to about 10000 examples).

Another important thing to keep in mind is that the training set and test set must have the *same distribution* of that. For instance, if we train on images of cats coming from the web, than we should not test on images of cats coming from mobile cameras.

### 3.1 Regularization

We've seen that the cost function of our neural network is a function  $J$  such that:

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

**Regularizing** a neural net model means adding a penalty when the model makes  $W^{[l]}$  bigger. This helps the model to prevent overfitting.

Thus, we have:

$$J(W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|W^{[l]}\|^2$$

Where

$$\|W^{[l]}\|^2 = \sum_{i=1}^{n^{[l]}} \sum_{j=1}^{n^{[l-1]}} (W_{i,j}^{[l]})^2$$

is called the *Frobenius norm* (sometimes we write  $\|W^{[l]}\|_F^2$ ).

Here  $\lambda$  is called the **regularization factor** and the bigger it is, the more we penalise the model for having bigger  $W^{[i]}$ .

Another kind of regularization is called **dropout regularization**. In this kind of regularization, we set a small probability of removing a node at all and have a new smaller network.

One way to implement it is called *inverted dropout*.

```

1 l = 3
2 keep_prob = 0.8
3 d3 = np.random.rand(a3.shape[0], a3.shape[1]) < keep_prob
4 a3 = np.multiply(a3, d3)
5 a3 /= keep_prob

```

The last line is used to keep the expected value of  $a$  the same, so we keep the same scale of  $z$  when doing  $z = wa + b$ .

When using dropout, though, we *don't apply it when using the neural net to predict values*.

To explain intuitively why dropout works, we need to think about the nodes connected to a particular node. If the NN just give importance to one of the those nodes, then it might be dropped and the model will perform poorly. Therefore, to work well, the model needs to spread the weights into all the neurons, shirking them.

There are other ways to perform regularization.

**Data augmentation** is when we change the data to get more data. For instance, if we have a image dataset, we can flip the images, reduce saturation, blur the image, and do many other things to get more images to train the make our mode better.

If we have a sound dataset, we can make the sound loud, or add noise, supression and so on.

**Early stopping** is when we stop iterating our model using the error in the dev set. While the error in the train error one decreases, the error in the dev set tends to decrease for a while and than starts increasing. What early stopping does is stop iterating the model when the error in the dev set starts increasing.

## 3.2 Optimizing the Problem

### 3.2.1 Normalizing input

Normalizing the input is very important because it garantees to us that the input data will have zero mean and variance one. To perform normalization, we calculate:

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

and

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m x^{(i)2}$$

where  $x^{(i)2}$  mean elementwise power.

Please note that these values are calculated using the *training set*. And then we'll apply the below formula to *every* input that we'll feed into the network (it can come from the training, dev or test set).

$$x := \frac{x - \mu}{\sigma}$$

Normalizing the input is important because it guarantees that Gradient Descent will converge much faster, since  $J$  will have a more circular bow shape, since of a elliptical one.

### 3.2.2 Vanishing / Exploding Gradients

In this section we're going to cover a problem that we can face when traning neural networks. Sometimes, the gradients become too small or too big and we can't training anything at all.

To understand that, suppose we have a very deep neural network. It's not so hard to imagine that as computations are performed upon computations, we can get very large or very small numbers for values like the activation function.

For instance, suppose all  $W^{[l]}$  have values of 1.5 then if we don't use an activation function, our final output will be something proportional to  $W^{[1]}W^{[2]} \dots W^{[L]}$ , which will me proportional to  $1.5^L$ . If  $L$  is very big, this can be huge.

If we changed 1.5 by 0.5, then we would have amoust zero.

To solve (parcially) this problem, we need to have a careful initialization of our weights.

Let's just consider a single neuron. If we ignore  $b$ , then we have

$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n$$

The later  $n$  is, the smaller we want  $w_i$  to be, because we want  $z$  to be approximately in the range  $[-1, 1]$ . One thing that we would wish is to have

$$\text{Var}(w_i) = \frac{1}{n}$$

So we can initialize it like:

$$W^{[l]} = \text{np.random.randn(shape)} \cdot \text{np.sqrt}\left(\frac{1}{n^{[l-1]}}\right)$$

It turns out that if we are using ReLU, it's better to have a variance of  $\frac{2}{n}$ , so we can replace the one by 2 in the sqrt.

The variatian with 1 is called the *Xavier initialization*, but we can use another variatian in which we multiply by:

$$\sqrt{\frac{2}{n^{[l-1]} + n^{[l]}}}$$

These factors could all be used as a base model, but we can really tune this as a hyperparameter if we want, although it's not a very important hyperparameter.

### 3.2.3 Gradient Checking

A very common technique when using Gradient Descent is to use **Gradient Checking** to make sure our gradient computations are been performed right. It will not go to the final model, but we use it to make sure we've coded everything right.

To do it, basically we approximate the gradients numerically and compare with the gradient we're computing using the formulas. Of course we don't do that when we want to predict any value or use our model in the real world because we're computing the gradient two times (and computing it numerically is not efficient). We just use this technique when coding the backpropagation to make sure we're doing it right.

To approximate the gradient, it's very simple, we just use:

$$g(\theta) \approx \frac{f(\theta + \varepsilon) - f(\theta - \varepsilon)}{2\varepsilon}$$

which is called *two-sided approximation*.

So the first thing in order to use gradient checking is to reshape our parameters  $W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}$  into just one big vector  $\theta$ . And its derivative will be  $d\theta$

Now we're going to compute  $d\theta_{\approx}$  which can be computed doing:

$$d\theta_{\approx i} = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \varepsilon, \dots, \theta_k) - J(\theta_1, \theta_2, \dots, \theta_i - \varepsilon, \dots, \theta_k)}{2\varepsilon}$$

For each value of  $i$ .

And what we want is

$$d\theta_{\approx} \approx d\theta$$

To now that, we can calculate:

$$\frac{\|d\theta_{\approx} - d\theta\|}{\|d\theta_{\approx}\| + \|d\theta\|} < \delta$$

for a small value of  $\delta$ . It could be  $10^{-7}$  to make sure it's really correct or  $10^{-5}$ . If the difference is greater than something like  $10^{-3}$  then probability we have a bug somewhere.

## 3.3 Optimization Algorithms

In this section, we're going to view many algorithms that allow us to train much faster.

### 3.3.1 Mini-batch Gradient Descent

What we've been doing in terms of gradient descent so far is to process the whole training set to take a step in the opposite direction of the gradient.

If  $m$ , our number of training examples, is very large (which is normal in Big Data applications), then a single gradient step is very expensive.

One way to optimize that is to use what's called **mini-batches**. Basically we split our training set into multiple training sets (if say 1000 examples) and we run each step of gradient descent using one of these sets.

That's what we call an **epoch** of gradient descent (a pass of forward-prop, cost computation, back-prop and parameters update).

Of course, by doing this, we're not optimizing the original  $J$  function (which is computed for every training example). However, we'll get a low cost model at the end.

If we plot the cost of  $J$  at each iteration (epoch), we'll not have that classical curve in which  $J$  decreases at each iteration. The curve will be a little noisy, increasing and decreasing. But in the long run, it'll converge to some low value. And not just that, it does it much faster than the original gradient descent (which is also called **batch gradient descent**).

One important concept to mention is the size of the mini-batches. If the size is one, then we call the algorithm **stochastic gradient descent**. This will make our epochs really fast, but we'll lead to an algorithm with a much higher cost than we could get using greater batches. Stochastic gradient descent will be *around* the minimum optima of  $J$ . The greater our mini-batches are, the more closer the that minimum optima we'll lead. That's the trade-off between speed and precision.

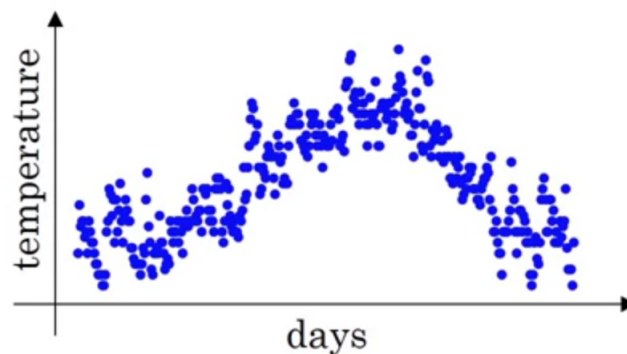
We can also train with mini-batch until the cost is very close to the minimum optima and then change to batch gradient descent and run some more epochs to make sure that our algorithm is really precise.

In practical, because how the computer memory and vectorization works, people tend to use powers of two for the size of the mini-batches. So that's a thing to keep in mind. Common values are  $2^6$  up to  $2^9$ .

### 3.3.2 Exponentially weighted averages

Now we're going to cover other optimization algorithms that are better than gradient descent. But to understand them, we need to talk about a concept of statistics which is called **exponentially weighted (moving) averages**.

To understand that, let's use an example. We'll plot the temperatures over the days.



Here we see that the values are very noisy, going up and down on each day. However, if we want to visualize the *trend* of the values, we can apply the moving average technique.

Basically, instead of plotting these points (let's call them  $\theta_i$ ), we'll plot  $V_i$ .

First, we initialize  $V_0$  as zero.

$$V_0 = 0$$

then, we get the next value by averaging the previous value of  $V_i$  and the current value of  $\theta_i$ .

$$V_1 = 0.9V_0 + 0.1\theta_1$$

a general formula would be:

$$V_i = 0.9V_{i-1} + 0.1\theta_i$$

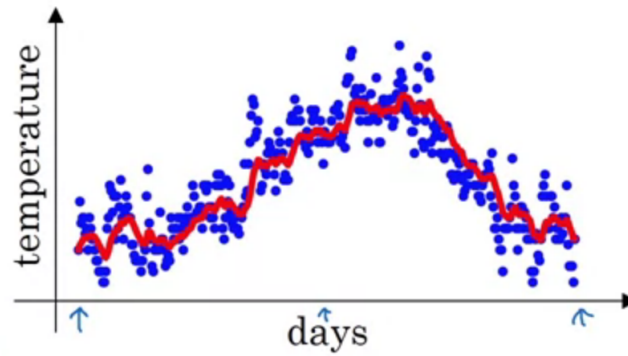


Figure 3.1: The moving average of the temperatures.

If we compute and plot that, we'll get figure 3.1

We see that we got a much more smooth curve because we're always considering the previous value to create the next one. We could make the formula even more general:

$$V_i = \beta V_{i-1} + (1 - \beta)\theta_i$$

where  $\beta$  is the percentage of how much we would like to consider the previous value.

Indeed the value of  $\beta$  will dictate approximately how many previous iterations we're considering on each point. To have an idea of how many previous points we're considering, we use the formula

$$\frac{1}{1 - \beta}$$

Therefore, for  $\beta = 0.9$ , we're considering approximately the 10 previous temperatures in order to get the current value. That's why it's called a *moving* average or a *local* average.

The lower  $\beta$  is, the less temperatures we'll be considering on each average and, therefore, the faster those averages will change, leading to a noisier curve, but that adapts faster. And if  $\beta$  is big, we'll be considering more days and, therefore, the average will change more slowly, leading to a smoother curve that has some latency to adapt to new trends.

This value will be a hyperparameter in our models.

To understand a little better this algorithm, let's try to get a non-iterative formula. We'll get something like:

$$V_i = 0.1\theta_i + 0.1 \times 0.9\theta_{i-1} + 0.1 \times 0.9^2\theta_{i-2} + \dots + 0.1 \times 0.9^k\theta_{i-k} + \dots$$

So, indeed, we're calculating  $V_i$  as an average of *all* previous values. But the farther the values are from the current we're calculating, the less impact it has on the current one. Indeed, we see that the influence decays *exponentially* (that's why this term appears in the name).

If we sum all coefficients, we would expect to get exactly 1, since this is an average. But unfortunately, that's not what happens. In fact, we say that this is approximately an average, because we're missing an **bias term**.

However, using this kind of average is very good computationally, because we don't have to keep track of (say) the last  $n$  terms to compute a local average. We can do that in  $O(1)$  for each element.

If we would implement this algorithm, it would be like:

**Algorithm 3.1: Exponentially weighted averages**

```

 $V_\theta = 0$ 
Repeat {
   $V_\theta := \beta V_\theta + (1 - \beta)\theta_t$ 
}

```

As one can see, this algorithm is very memory and time efficient.

Finally, we can make a correction to the bias term by dividing  $V_i$  by  $1 - \beta^i$ .

So we can use the formula:

$$V_i = \frac{\beta V_{i-1} + (1 - \beta)\theta_i}{1 - \beta^i}$$

This will help a lot in the first iterations where we still don't have many values in the window. Think of it as  $1 - \beta^i$  being a good approximation for the sum of the coefficient of the average.

### 3.3.3 Gradient Descent with Momentum

The first algorithm we'll see is **gradient descent with momentum**. Basically, instead of using the derivatives in each iteration, we're going to average the derivatives using exponentially weighted averages to get a more smooth descent to the minimum. This algorithm works pretty much always better than normal gradient descent.

**Algorithm 3.2: Gradient Descent with Momentum**

```

On iteration  $t$ :
  Compute  $dW, db$  on current batch (or mini-batch)
   $V_{dW} = \beta V_{dW} + (1 - \beta)dW$ 
   $V_{db} = \beta V_{db} + (1 - \beta)db$ 
   $W = W - \alpha V_{dW}$ 
   $b = b - \alpha V_{db}$ 

```

This averaging over the last  $k$  derivatives makes our algorithm take this *momentum*, which is going to create speed into the minimum direction, avoiding oscillations in other directions like standart gradient descent does.

The downside of this algorithm is that now we have two hyperparameters instead of just one:  $\alpha$  and  $\beta$ . In practice, values close to  $\beta = 0.9$  work good.

One more note is that in practice people tend to not use *bias correction* because after just 10 or so iterations the algorithm will have a very small bias value.

### 3.3.4 RMSprop

**RMSprop** which means for Root-Mean-Square prop. It also speed gradient descent by “normalizing” the different directions. To do that, it takes the element-wise square of the gradient, average it over the last  $k$  squares and than divide the gradient by the squareroot of that.

**Algorithm 3.3: RMSprop**

On iteration  $t$ :

Compute  $dW, db$  on current batch (or mini-batch)

$$S_{dW} = \beta V_{dW} + (1 - \beta) dW^2$$

$$S_{db} = \beta V_{db} + (1 - \beta) db^2$$

$$W = W - \alpha \frac{dW}{\sqrt{S_{dW} + \varepsilon}}$$

$$b = b - \alpha \frac{db}{\sqrt{S_{db} + \varepsilon}}$$

The  $\varepsilon$  value here is just to guarantee we're not dividing by zero.

**3.3.5 Adam**

**Adaptative moment estimation** optimization algorithm will basically take the two previous ideas (momentum and rmsprop) and put them together. It's one of the best algorithms known so far along RMSprop.

**Algorithm 3.4: Adam**

$$V_{dW} := S_{dW} := V_{db} := S_{db} := 0$$

On iteration  $t$ :

Compute  $dW, db$  on current batch (or mini-batch)

$$V_{dW} = \beta_1 V_{dW} + (1 - \beta_1) dW$$

$$V_{db} = \beta_1 V_{db} + (1 - \beta_1) db$$

$$S_{dW} = \beta_2 V_{dW} + (1 - \beta_2) dW^2$$

$$S_{db} = \beta_2 V_{db} + (1 - \beta_2) db^2$$

$$V_{dW} = \frac{V_{dW}}{1 - \beta_1^t}$$

$$V_{db} = \frac{V_{db}}{1 - \beta_1^t}$$

$$S_{dW} = \frac{S_{dW}}{1 - \beta_2^t}$$

$$S_{db} = \frac{S_{db}}{1 - \beta_2^t}$$

$$W = W - \alpha \frac{V_{dW}}{\sqrt{S_{dW} + \varepsilon}}$$

$$b = b - \alpha \frac{V_{db}}{\sqrt{S_{db} + \varepsilon}}$$

So as one can see, we're basically doing what we're doing in RMSprop, but now instead of using  $dW, db$ , we're using the average over the last  $k$  values, just like in momentum.

One more thing here is that we're using bias correction, which is very common to do in this algorithm.

Now we have three hyperparameters:  $\alpha, \beta_1, \beta_2$ . In practice,  $\alpha$  still needs to be tuned, and common values for the other two are  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . There's also the value of  $\varepsilon$ , but this is very rare to chance. We usually use something like  $\varepsilon = 10^{-8}$ .



### 3.3.6 Learning Rate Decay

One more technique that can be used through all the optimization algorithms is **learning rate decay**. Basically, we slowly reduce the  $\alpha$  parameter in our algorithm by the time. In the begining, we're far away from the minimum, so we can set our learning rate higher and as we go towards the minimum, we reduce this learning rate to make it smaller and more precise. One way to do that, is to use the formula:

$$\alpha = \frac{1}{\text{decay\_rate} \cdot \text{num\_epoch}} \cdot \alpha_0$$

Other formulas use exponentially decay:

$$\alpha = 0.95^{\text{num\_epoch}} \cdot \alpha_0$$

$$\alpha = \frac{k}{\sqrt{\text{num\_epoch}}} \cdot \alpha_0$$

## 3.4 Hyperparameter tuning

We've seen so far that the algorithms have lots of hyperparameters to turn, like the number of layers, number of nodes in each layer, mini-batch sizes,  $\alpha$ ,  $\beta_1$ ,  $\beta_2$ . Also the decision of using one optimization algorithm instead of other could be considered as a hyperparameter.

We need to understand how to select the values of the hyperparameters we're going to try. In the early machine learning times, we used a grid to try the hyperparameters. For instance, if we had  $h_1$  and  $h_2$ , then we would say in each values of them we wanted to try. Let's say we want to try  $h_1 = \{v_1, v_2, v_3\}$  and  $h_2 = \{u_1, u_2, u_3\}$ . Then we would try all the combinations of these values:

$$h_1 \times h_2 = \{ (v_1, u_1), (v_1, u_2), \dots, (v_3, u_3) \}$$

This is good when the number of parameters is small. But in deep learning, the number of parameters tend to be large and, thus, we use random points (i.e., for each hyperparameter, we pick a random value in some region).

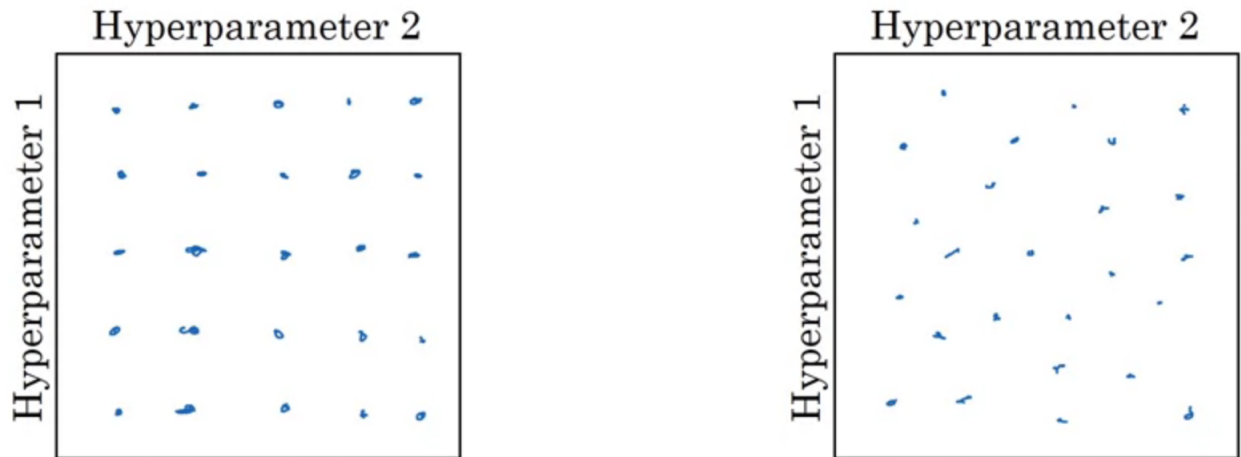


Figure 3.2: In the left we see the Grid Search and in the right, we see the random search

Another good thing to do is to perform many hyperparameter searches. In the first one, we can use larger scales for each hyperparameter and discover that the best value for, let's say,  $\alpha$  is between 0.9 and 1.1. Then we can perform a new search sampling values of  $\alpha$  just in that region and discover that the best value is approximately 0.96.

### 3.4.1 Selection the appropriate scale for hyperparameters

In many cases, we don't want to *uniformly* sample values from a particular range. Suppose we use the range  $[0.0001, 1]$  for the learning rate  $\alpha$ . Then it's clear that most (90%) of the time we'll select values between 0.1 and 1. That's not what we want. We want, to select values between  $(0.0001, 0.001)$ ,  $(0.001, 0.01)$ ,  $(0.01, 0.1)$  and  $(0.1, 1)$  in the same proportion.

Therefore, what we need to do is to sample using a *logarithmic* random scale.

To do that in python, one way is the following:

```
1 r = -4 * np.random.rand()
2 alpha = 10 ** r
```

**Scale for parameter  $\beta$**  Another important topic is how to pick a scale for the parameter  $\beta$ . We've seen that  $\beta = 0.9$  means we're averaging over the 10 previous gradients and if  $\beta = 0.999$ , we're averaging over the last 1000 gradients. So we want to have a linear sample *in the number of gradients*. To do that, we can think in  $1 - \beta$  and get the interval  $(0.1, 0.001)$ . And just like in the example from above, we can use a logarithmic scale, but now reversed. Therefore, we use:

$$\beta = 1 - 10^r$$

Instead of the code above.

## 3.5 Batch normalization

As we've seen before, input normalization can help a lot to optimize our algorithm and train faster. Another technique is called **batch normalization** in which normalize the values in between the hidden layers.

So basically for each intermediate value  $z^{[l]}$ , we normalize it using:

$$\begin{aligned}\mu &= \frac{1}{m} \sum_i z^{[l](i)} \\ \sigma^2 &= \frac{1}{m} \sum_i (z^{[l](i)} - \mu)^2 \\ z_{\text{norm}}^{[l](i)} &= \frac{z^{[l](i)} - \mu}{\sqrt{\sigma^2 + \varepsilon}}\end{aligned}$$

But indeed, we don't want all the intermediate values to be normalized, so we do:

$$\tilde{z}^{[l](i)} = \gamma z_{\text{norm}}^{[l](i)} + \beta$$

where  $\gamma$  and  $\beta$  are learnable parameters of the model (using their gradient as with  $W$  and  $b$ ).

This allows us to change the mean and variance of each layer to our benefit. Basically now, instead of doing this process:

$$X \xrightarrow{W^{[1],b^{[1]}}} z^{[1]} \xrightarrow{g^{[1]}} a^{[1]} \xrightarrow{W^{[2],b^{[2]}}} z^{[2]} \xrightarrow{g^{[2]}} a^{[2]} \rightarrow \dots$$

We're doing:

$$X \xrightarrow{W^{[1],b^{[1]}}} z^{[1]} \xrightarrow{\gamma^{[1],\beta^{[1]}}} \tilde{z}^{[1]} \xrightarrow{g^{[1]}} a^{[1]} \xrightarrow{W^{[2],b^{[2]}}} z^{[2]} \xrightarrow{\gamma^{[2],\beta^{[2]}}} \tilde{z}^{[2]} \xrightarrow{g^{[2]}} a^{[2]} \rightarrow \dots$$

And we can compute the gradient of all parameters to update each of them using Gradient Descent.

**Observation.** We also need to say that indeed, when using batch normalization, we don't need to use the  $b$  parameter at all, because  $\beta$  will always change it before computing  $a$ . Therefore, the only parameters we need are  $W, \gamma, \beta$ .

## 3.6 Multiclass problems

So far, we've seen only how to do binary classification with neural nets. Now, let's see a logistic regression variant that allow us to have multiple classes: **softmax regression**.

### Notation 3.1

- $C$  will denote the number of classes;
- The classes will be numbers from 0 to  $C - 1$ .

The first modification we'll do is in the output layer. Now  $\hat{y}$  will be a vector instead of a single number and  $\hat{y}_i = p(i | x)$ , i.e., the probability of class  $i$  given the input  $x$ .

We'll use an **softmax** layer in the output layer to achieve such goal. To do so, we'll compute  $z^{[L]}$  as usual and use the **softmax activation function**, which is:

$$t = e^{(z^{[L]})}$$

$$\hat{y} = a^{[L]} = \frac{e^{(z^{[L]})}}{\sum_i t_i}$$

As you can see, we basically compute this vector  $t$  and divide by the sum of its components (it's like we were normalizing  $t$ ).

Indeed, if we sum the coordinates of  $\hat{y}$ , they'll add up to 1.

## 3.7 Deep Learning Frameworks

We've learned so far how to program deep learning algorithms from scratch but there are several python frameworks that are used in practice when implementing these deep learning algorithms.

Some of the most used DL frameworks are:

- Keras
- TensorFlow
- Torch