

Understanding & Developing Artificial Neural Network with Objective-C and Python

Deerfield Academy Advance Computer Science Research

Yongyang (Neil) Nie

Feb, 2016

Contents

1 Abstract:	2
2 Introduction:	2
3 Forward Feed:	2
4 Quantifying and Minimizing Cost:	2
5 Gradient Descent:	3
6 Back propagation:	3
6.1 Back propagation overview	3
6.2 Mathematics behind backpropagation	4
6.2.1 An equation for the error in the output layer δ^L	4
6.2.2 An equation for the error δ^l in terms of the error in the next layer, δ^{l+1} in particular:	4
6.2.3 An equation for the rate of change of the cost with respect to any bias in the network	4
6.2.4 An equation for the rate of change of the cost with respect to any weights in the network	4
6.2.5 Backpropagation with simple neural network example	4
7 Improve neural network training result	5
7.1 Making good decision	5
7.1.1 Hyperparameters	5
7.1.2 Hidden nodes	6
7.2 Methods to improve training result	6
8 Examples and Applications:	6
8.1 Create an application with Xcode to recognize handwritten digits	7
8.1.1 Introduction	7
8.1.2 Building the application	7

1 Abstract:

Machine learning is a subset of artificial intelligence that provides computers with the ability to learn without being explicitly programmed. Machine learning focuses on the development of computer programs that can change when exposed to new data.

There are many methods and algorithms for learning algorithms, from Support Vector Machine¹ to Artificial Neural Networks². The purpose of them are similar, they can all be used to classify complex data, such as images and DNA samples.

2 Introduction:

There are two types of machine learning, one is supervised machine learning³. In this research, we will mainly focus on artificial neural networks and supervised machine learning. Supervised learning is the machine learning task of inferring a function from labeled training data. The training data consist of a set of training examples. In supervised learning, each example is a pair consisting of an input object (typically a vector) and a desired output value (also called the supervisory signal).

You will discover that machine learning is beautiful and can be very simple. We will be able to implement a simple algorithm that can recognize hand written digits in less than a hundred lines of code. Only a small amount of mathematical proves and equations will be covered. The paper will discuss the principles behind designing, building, and debugging an artificial neural network. The projects will be written mainly in Objective-C or Python. You can find the resources on [Github](#).

3 Forward Feed:

The neural network will begin be taking in some inputs and making some predictions based on the inputs and the weights between nodes. We can think of it like a one by two matrix. There will be a weight connecting every input layer node to every output layer node, therefore, there are six weights between the input layer and the second layer. The matrix calculation should yield us some result

$$\begin{bmatrix} 0 \\ 1 \end{bmatrix} [w_0^0 \ w_0^1 \ w_0^2 \ w_1^0 \ w_1^1 \ w_1^2] = \begin{bmatrix} z \\ z \\ z \end{bmatrix}$$

If we give the matrices some names, call inputs x , weights $w_{(l)}$ and output $z_{(l)}$. In our case, l indicate the layer, for example, the first layer weights are $w_{(1)}$. $z_{(l)}$ represent the output matrix with layer l , in this case, the output z is the hidden layer output.

$$xw_{(n)} = z_{(n)} \quad (1)$$

Afterwards, we have to apply an activation function⁴. The activation function that I used in this research and this paper will be the sigmoid function. This is a complete cycle, there is one more layer to go in order to yield a result. Note, in a multilayer neural network, this process will be repeated until we yield some output. In our case, we only have to do this twice. The activation function is shown as below.

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$
$$a_{(2)} = f(z_{(2)}) \quad (2)$$

By applying the same process as before, this time, the inputs will be $a_{(2)}$ then, $z_{(3)}$ will be our final result. This process is known as forward feed. It's relatively straightforward.

$$z_{(3)} = a_{(2)}w_{(2)}$$

After we yield some result, we can more on and look at how far off we are from the expected result. We need some methods to quantify this and so adjustments to the network to minimize error.

$$\sigma'(z) = \sigma(x)(1 - \sigma(z))$$

4 Quantifying and Minimizing Cost:

Now the neural network can make calculation/predictions, however, the result is far from desired. In almost all learning algorithms, the input data cannot be altered, therefore the x term is constant in equation one. In order to

¹https://en.wikipedia.org/wiki/Support_vector_machine

²https://en.wikipedia.org/wiki/Artificial_neural_network

³Mehryar Mohri, Afshin Rostamizadeh, Ameet Talwalkar (2012) *Foundations of Machine Learning*, The MIT Press ISBN 9780262018258.

⁴https://en.wikipedia.org/wiki/Activation_function

change the output z the only option is to change the weights w .

First of all, we have to come up with ways to quantify the cost.

$$(3) \quad C = \sum_j \frac{1}{2} (\hat{y} - y)^2 \quad (3)$$

C is the cost, which equals to the sum of all the differences between calculated result and actual result squared and times one half. We can take advantage of the equation derived above and substitute for some of the variable.

$$(4) \quad C = \sum_j \frac{1}{2} (y - f(f(xw_{(1)})w_{(2)}))^2$$

Here we have it, a way to quantify the cost of the neural network. This function will be referred to as the cost function. Now, we have to solve the problem, how do we minimize C , will brute force work? It turns out, no, because in a three-node neural network we have to compute more than a million possible weights, which will be gruesome.

We can think of the equation above as a function of cost in terms of all possible weights. There will be one set of weights that will bring the cost to the lowest. Then, this becomes a minimization problem.

5 Gradient Descent:

The best way to minimize the cost is to use gradient descent, a very fast and classic way to solve problems like this. In fact, gradient descent is widely used in math, image process and machine learning.

Gradient descent is a first-order iterative optimization algorithm. To find a local minimum of a function using gradient descent, one takes steps proportional to the *negative* of the [gradient](#) (or of the approximate gradient) of the function at the current point.

Gradient descent is also known as steepest descent, or the method of steepest descent. The process can be seen as a ball rolling down a hill⁵ and trying to find the lowest point. Note that actual physics doesn't apply here and we will define our own movement of the ball.

⁵<https://iamtrask.github.io/2015/07/27/python-network-part2/> by Andrew

⁶Proof and definition of convex functions: <http://mathworld.wolfram.com/ConvexFunction.html>

⁷http://www.mit.edu/~dimitrib/Incremental_Survey_LIDS.pdf Dimitri P. Bertsekas Report LIDS - 2848

There are limitations to this method. First of all, what if we are stuck in a local minimum, our goal is to find the global minimum for the cost function. In another word, this method will not work for a non-convex function. In fact, this problem is solved in equation (3), by squaring the difference in $(\hat{y} - y)^2$, we transformed this function to a convex function⁶. Even though there are tens of thousands of variables (dimension), the function will still be convex, thus, we can apply gradient descent to find the global minimum.

There are other types and variation of gradient descent as well. One of the most commonly used one is Stochastic gradient descent (SGD), also known as incremental gradient descent, is a stochastic approximation of the gradient descent optimization method for minimizing an objective function that is written as a sum of differentiable functions. In other words, SGD tries to find minima or maxima by iteration.⁷

6 Back propagation:

With gradient descent, we can create a set of routines that can help us to change the weights of the network to minimize the cost.

6.1 Back propagation overview

Phase 1: Propagation: Each propagation involves the following steps:

1. Forward propagation of a training pattern's input through the neural network in order to generate the network's output value(s).
2. Backward propagation of the propagation's output activations through the neural network using the training pattern target in order to generate the deltas (the difference between the targeted and actual output values) of all output and hidden neurons.

Phase 2: Weight update: For each weight, the following steps must be followed:

1. The weight's output delta and input activation are multiplied to find the gradient of the weight.
2. A ratio (percentage) of the weight's gradient is subtracted from the weight.

6.2 Mathematics behind backpropagation

Backpropagation is based around four fundamental equations. Together, those equations give us a way of computing both the error δ^L and the gradient cost of the function. In fact, the backpropagation equations are so rich that understanding them well requires considerable times and patience.⁸

6.2.1 An equation for the error in the output layer δ^L

$$\delta_J^L = \frac{\partial C}{\partial a_J^L} \sigma'(z_J^L) \quad (\text{BP1})$$

The first term on the right, $\frac{\partial J}{\partial a_J^L}$ measures how fast the cost is changing as a function of j^{th} output activation. Everything in (BP1) is easily calculated. We computed z_J^L while computing the behavior of the network. Depending on the cost function, in our case, the quadratic cost function is relatively easy to compute.

Equation (BP1) is a perfectly good expression, however, it's not matrix based, form that backpropagation desires. The fully matrix form becomes.

$$\delta^L = (a^L - y)(\sigma'(z^L)) \quad (4)$$

6.2.2 An equation for the error δ^l in terms of the error in the next layer, δ^{l+1} in particular:

$$\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l) \quad (\text{BP2})$$

where $(w^{l+1})^T$ is the transpose of the weight matrix (w^{l+1}) for the $(l+1)^{\text{th}}$ layer. Suppose we know the error δ^{l+1} at the $(l+1)^{\text{th}}$ layer. When the transpose weight matrix is applied $(w^{l+1})^T$, we can think of this as moving the error backward through the network, giving us some sort of measure of the error at the output of the l^{th} layer. Finally, we take the Hadamard product $\odot \sigma'(z^l)$.

⁸The equations were created by Michael A. Neilson "Neural Networks and Deep Learning", Determination Press, 2015

⁹The equation was reference from Michael A. Neilson "Neural Networks and Deep Learning", Determination Press, 2015

¹⁰The equation was referenced from Michael A. Neilson "Neural Networks and Deep Learning", Determination Press, 2015

¹¹The equation was referenced from Michael A. Neilson "Neural Networks and Deep Learning", Determination Press, 2015

By combining (BP1) with (BP2), the error δ^l for any layer in the network can be computed. We start by using (BP1) to compute δ^L , then apply Equation (BP2) to compute δ^{L-1} , then Equation (BP2) again to compute δ^{L-2} , and so on, all the way back through the network.

6.2.3 An equation for the rate of change of the cost with respect to any bias in the network

$$\frac{\partial C}{\partial b_j^l} = \delta_j^l \quad (\text{BP3})$$

This is the error δ_j^l is exactly equal to the rate of change $\frac{\partial C}{\partial b_j^l}$. This equation can be rewritten as

$$\frac{\partial C}{\partial b} = \delta$$

6.2.4 An equation for the rate of change of the cost with respect to any weights in the network

$$\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l \quad (\text{BP4})$$

This equation can help us to compute the partial derivative $\frac{\partial C}{\partial w_{jk}^l}$ in terms of the quantity δ^l and a^{l-1} . It can be rewritten in a less index-heavy notation:

$$\frac{\partial C}{\partial w} = a_{in} \delta_{out} \quad (5)$$

6.2.5 Backpropagation with simple neural network example

In a simple three-layer neural network, below are the four equations that we derive with the principle of gradient descent that will help us minimize the error.¹¹

$$\delta_{(3)} = -(y - \hat{y})(\sigma'(z_{(3)}))$$

Calculate the δ for the third layer of the network.

$$\frac{\partial J}{\partial w_{(2)}} = (a^{(2)})^T \delta_{(3)}$$

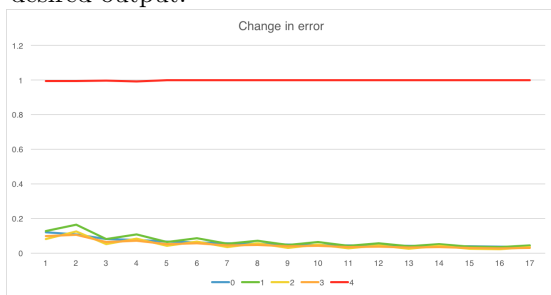
Use the calculated result to help us modify the second layer of weights in the network.

$$\delta_{(2)} = \delta_{(3)} (w_{(2)})^T \sigma'(z_{(2)})$$

Calculate the δ for the second layer of the network.

$$\frac{\partial J}{\partial w_{(1)}} = x^T \delta_{(2)}$$

Finally, we can modify the first layer of weights in the network. This process is often repeated until the accuracy of the network reaches a threshold. Here is a visualization of how individual output node is approaching the desired output.



Line 0-3 are lowering as we train the network. Line 4, is approaching 1 as we keep train the network. This figure displayed that the error of the network C is being minimized.

Here is a graph of the correction of the neural network over times of training. Each time the network will train with 1000 out of 60,000 training data. After, The program will shuffle the training data and repeat the training process for 197 times, or until the network reaches 95% accuracy.

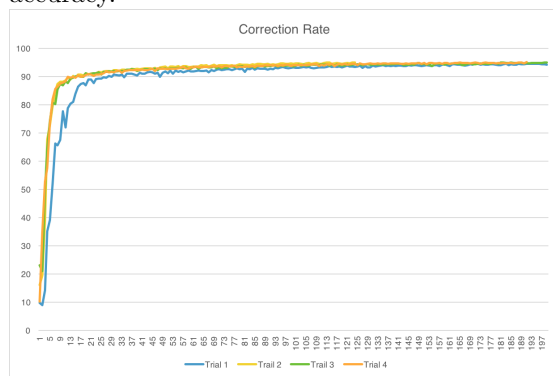


Figure 1

The neural network is trained four times with the same training and testing data.

¹²Michael A. Nielsen, "Neural Networks and Deep Learning", Determination Press, 2015. [Chapter 1 Using neural nets to recognize handwritten digits](#)

The results are similar and the network is steadily improving, thanks to gradient descent. There are many factors that contribute to a successful learning process. If these things are ignored or incorrect, the network will not be able to learn as well.

7 Improve neural network training result

Not all neural network behaves flawlessly. Generally, there are a few factors that will effect the training result of the network: hyperparameter, number of hidden nodes, human error, training data/testing data problem and more. Among them, hyperparameter, number of hidden nodes and be determined by estimation and trial and error. Huamn error can be eliminated by debugging. Training data and testing data problem is not easy to solve. That often contributes to the source of your data and the quality of your data, which are beyond the scope of this research. It's safe to assume that the training data and testing data that we are using are carefully chosen and considered.

7.1 Making good decision

In a pre-trained network, there are several critical hyperparameters and numbers that will help the network to improve. I mainly focused on learning rate and number of hidden nodes.

7.1.1 Hyperparameters

Learning rate is the hyperparameter that we will focus on. To make gradient descent work correctly, we need to choose the learning rate r to be small enough that Equation (9) is a good approximation. If we don't, we might end up with $\Delta C > 0$, which will take us to the opposite of minimizing. Meanwhile, r can't be too small, since that will make the changes Δv tiny, and thus the gradient descent algorithm will work very slowly. In practical implementations, r is often varied so that Equation (9) remains a good approximation, but the algorithm isn't too slow.¹²

This is a graph of the training result with a different learning rate. The batch size is 1000 and the training is repeated 100 times.

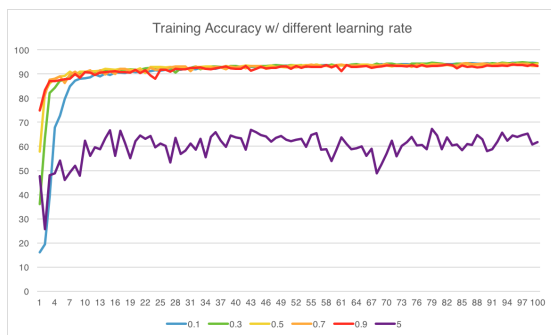


Figure 2

The range of the learning rate that I chose is from $[0.1, 5]$. When the learning rate is relatively low, the network began at a lower accuracy comparing to higher training rates. However, especially when the learning rate is 0.9, the training result became unstable, and the graph starts to oscillate.

On the other hand, when we set the training rate to 5.0, the network doesn't improve after 40%-50%. Then general rule of thumb for choosing the learning rate is somewhere between $[0.1-1]$.

With this observation in mind, we might be able to alter the learning rate as the network progress. The training rate should be high in the beginning to quickly bring us to the desired place. Then, we can lower the rate so the network learning result will not have oscillation.

7.1.2 Hidden nodes

In a multilayer neural network, you will likely encounter the problem of how many hidden layers and hidden nodes should the network have. Seemingly simple, but complex question will help you improve the learning rate of the network. In this research, we will only focus on one hidden layer. Multiple hidden layer is known as deep learning, which is out of our scope. So those few factors to keep in mind in order to finalize the number of layers and size (neurons/layer) for both the input and output layers.

It's difficult to form a good network topology just from the number of inputs and outputs. It depends critically on the number of training examples and the complexity of the classification you are trying to learn. There are problems with one input and one output that require millions of hidden units, and problems with a million inputs and a million outputs that require only one hidden unit, or none at all.

Some books and articles offer "rules of

thumb" for choosing a topology – inputs plus outputs divided by two, maybe with a square root in there somewhere – but such rules are not recommended. Other rules relate to the number of examples available: Use at most so many hidden units that the number of weights in the network times 10 is smaller than the number of examples. Such rules are only concerned with overfitting and are unreliable as well.

7.2 Methods to improve training result

There are many ways to improve the network's performance. In fact, gradient descent is a great optimization method to speed up the learning process. Here are some common ways to improve the accuracy.

- Cross-entropy
- Overfitting and regularization
- Altering the hyperparameters

In section 7.1.1, I stated that the learning rate will effect how the network learns in the beginning the fluctuation of the learning curve. This observation gives us the possibility to change the learning rate as the network learns. In the beginning, we will set the learning rate 0.9, then slowly lower the number as the network improves. This concept will seem intuitive in that you want the network slow down the learning rate when it's close to the desired result. Otherwise, a high learning rate will likely create some fluctuation like we have seen figure 2.

8 Examples and Applications:

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

8.1 Create an application with Xcode to recognize handwritten digits

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

8.1.1 Introduction

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst.

Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

8.1.2 Building the application

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.