

Info.6205 Project Report

A Neural Network Implementation For Handwritten Digits Recognize

Yuxuan Yang(001389098),

Haoqi Huang(001835259)

ABSTRACT

Neural networks are used as a method of deep learning, one of the many subfields of artificial intelligence. In this project, a small subsection of object recognition—digit recognition will be implemented, based on the MNIST database of handwritten digits. It includes building a neural network structure, training the model and accuracy evaluating.

STRUCTURE

1.Network

There are three layers in the neural network: Input, hidden and output layers, which is built by a one-dimensional array. Each element(an integer in this case) in the array represents one single node in the layer. Because the resulting images contain grey levels as a result of the anti-aliasing technique used by the normalization algorithm. the images were centered in a 28x28 image by computing the center of mass of the pixels, and translating the image so as to position this point at the center of the 28x28 field[1], there are 28x28=784 input Neural and output layers contains 10 nodes that show the similarity between result and actual number(0-9).

2.Activation and cost functions

In this project, the activation interface is set up to normalize the process, which includes a guidance method for predicting and reverse guidance for weight fixing. The different activators could be switched by changing the parameter and sigmoid is the default. Also, the quadratic cost function to monitor how close we get to the expected result.

MODELING

1.Feedforward

In a word, we feed in the data and calculate the result. And that's where activation function comes in, we use the activation function to consider a neuron "fired" or not. The reason why we chose sigmoid function is that we want small changes in the weights and bias cause only a small change in the output. That is an important fact that will allow our sigmoid neurons to learn. Training and learning are crucial for an artificial neural network.

```
Start feed the hidden layer
feed value: 0.9998246002161049
origin value:8.648267294928193
feed value: 5.170704782225562E-5
origin value:-9.869864755842498
feed value: 0.07240761889380155
origin value:-2.5502808645976156
feed value: 0.9995783961215072
origin value:7.77102266854018
feed value: 0.9998546604219798
```

We use input and its weights plus bias to calculate the weighted sum, and take the weighted sum into our sigmoid function to compute the activation. Using the activation as the input to calculate the next layer's activation, we do this Layer by layer, at last, we will get our output.

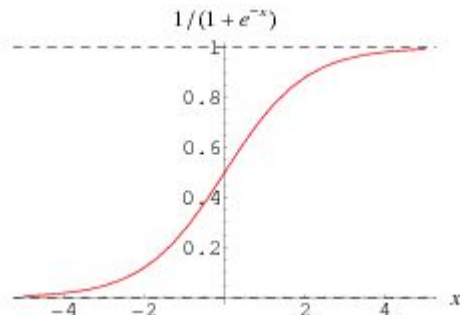
2.Back-propagation

The backpropagation approach is to model a given function by modifying internal weightings of input signals to produce an expected output signal[2]. The

most important part of the training process, in the feed-forward process we get our output, but we want to know how far we are from the target. Error is calculated between the expected outputs and the outputs forward propagated from the network.

Our aim is to decrease the distance, in other words, the minimum the value of cost function.

In the previous feed part, We choose the sigmoid function, and its formula shows like[3]:



In order to calculate how much we need to fix between prediction and reality, we count the error into the derivatives of our activation method- in this case(sigmoid) it is

$$\frac{d}{dx} S(x) = S(x)(1 - S(x))$$

We can use the property of derivative, therefore, we can know which direction is the correct way to get to the minimum. So back-propagation is surrounded by one key factor that is partial derivative. The partial derivative of the cost function with respect to any weight and bias in the network. it will tell us how quickly the cost changes when we change the weights and bias.

Since we don't know the expected value until arriving at the very last layer, the direction would be the back layer to the front. In other words, we calculate the hidden layer error using the result we got in the output layer.

In this project, we only use sigmoid and tanh. The generic formula[4] is:

$$\delta^L = \nabla_a C \odot \sigma'(z^L).$$

BP1: Output error

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

BP2: Hidden layer error

$$\frac{\partial C}{\partial b_j^l} = \delta_j^l.$$

BP3: Bias Change Equation

$$\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l.$$

BP4: Weight Change Equation

BP1 is to get the error from the last layer (output), BP2 is used for calculating the error for the hidden layer, BP3 is to get the new bias and the last one BP4 is for weights.

3.Update weight and bias

To find out how we change the weights and bias, we have to know the error between the output and the expected result, each output neuron more or less have some difference with the expected result. We can know whether cost function depends much on a particular output neuron, by calculating the partial derivative of activation of specific output neuron. If the error is small, then the cost function does not depend on that neuron. if the error is big, it means when we change the weights connected to this neuron and the bias on this neuron will effectively reduce the cost function.

Since we got the error value for every node in the hidden layer and output layer, we could update the weight and bias now with the learning degree.

Because of the equation: result=input*weight+bias, which is obvious, we could simply regard the total training progress as fix the weight and bias.

The implementations are like below. Note that the learning degree (n) for bias is 1/10 since it turns out that bias is actually more sensitive than weights.

The update result shows like the below screenshots:

DELIVERABLES

Part1. the training process

Here you can see the change of cost function(Mean Square Error), as we mentioned before, to understand the variation of the cost function is a very important way to know how well our program has learned. We also can find our program bugs through MSE. So here is some screenshot.

Here is part of the output result divided by epoch(0, 1, 2), in each epoch, we have 5 MSE value, each MSE is the average MSE calculated from 100 random handwriting images, since each image through the program will generate an output, and through the expected result we can compute the MSE. We can clearly see that the MSE is reducing, that's mean our program is learning well and making a good self adjust of weights and bias.

[illegible]

Up there, I change the number of training images for each MSE, instead of 100, I decided to add to 1000, and then get the above result. As you can see the MSE change more quickly than the previous one. That's because we feed in much more training data than the first one. So the program will become more and more accurate when you train enough data. Also, we may notice that the MSE always fluctuates, sometimes it gets

```
time consume: 17208.0 millisecs
System gets 921 correct prediction from 1000 hand-writing images
```

```
time consume: 37041.0 millisecs
System gets 954 correct prediction from 1000 hand-writing images
```

If we increase hidden neurons to 150:

If we increase hidden neurons to 200:

[illegible]

It is easy to find out that the total time for the same amount of training data increase when you got more and more hidden neural. As for the correct rate predicted by the program, it increases when we have more hidden neurons but it seems to have a peak, after that the rate will decrease drastically. So we guess, that's maybe related to our learning rate. As a result, it is important to find out the right amount of hidden neurons, which will get the shortest time and highest correct rate.

Take a look at this training result of average MSE. In the first epoch of first 1000 “random” images, the MSE is too small to make sense, so there must be some problem with the random data we select, since we already pass the unit test of all the process in the feedforward and back-propagation part. After debugging, there is a problem in the random selection algorithm. We always chose the first data in the MNIST database, so that explains why we have such a small MSN and never fluctuate because we always training the same data! Therefore, the MSE can also help us detect the problem our program has.

The reason why the learning rate is important is that it affects the training process directly if we give the value to the learning rate too high, we will not get the result we want, even worse, the result is confusing. If it is too low, we will not see the remarkable changing in our correct rate. But, it is always recommended to set the learning rate lower.

We keep all the variable unchanged, except the amount of the hidden neurons. We feed in $10 * 5 * 2000$ random data, 10 epoch, each epoch loop 5 times, each loop train 2000 random images. And we got:

This is the result we got, for the learning rate is 0.3 at the first time we ran our program.

Reference

- [1]<http://yann.lecun.com/exdb/mnist/>
- [2]<https://machinelearningmastery.com/implement-backpropagation-algorithm-scratch-python/>
- [3]https://en.wikipedia.org/wiki/Sigmoid_function
- [4]http://neuralnetworksanddeeplearning.com/chap2.html#the_four_fundamental_equations_behind_backpropagation