

AMATH 482: Neural Networks: My Personal Black box

Anh-Minh Nguyen

March 10, 2020

Abstract In this report, we will be tinkering with neural networks to help us classify a huge dataset concerning fashion items, or just clothes. What we are about to witness is, in my personal opinion, a defiance of human perseverance and reasoning for the acceptance of a far superior computational force, that is my laptop.

1 Introduction and Overview

Computers have given us the gift of being a trillion times faster than we could ever be in terms of calculation. We, the humans, have done the theoretical work, and now it's up to these small metal chips to do the brunt of the work. One thing we have been trying to speed up and refine is the decision-making process of a computer. In comes my new dear friend and enemy, neural networks. We will be using the fashion_mnist data set, which we will training our neural network on. Our goal is to achieve at least a 90% accuracy.

2 Theoretical Background

I want to start by speedrunning through two statistical staples, Linear Regression and Logistic Regression (the algorithm is regression while the decision process is classification).

Linear Regression is fitting a line of best fit, where the data looks like $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$, where N is the number of observations. We see that \mathbf{x} can be a vector of inputs while each y is our output based on those inputs. We assume that the relationship between our predictors and response is linear. We often don't know the true value of the predictors, so we have to predict our parameters. To measure how well we have done, we use the mean squared error (MSE). Our prediction is:

$$\hat{y}_j = m\mathbf{x}_j + \mathbf{b} \quad (1)$$

and our MSE is:

$$\frac{1}{n} \sum_{j=1}^N (y_j - \hat{y}_j)^2 \quad (2)$$

Let's look at a multivariate linear regression.

$$\begin{aligned} y_1 &= a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n + b_1 \\ y_2 &= a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n + b_2 \\ &\vdots \\ y_m &= a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n + b_m \end{aligned}$$

Which we can relabel as

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b} \quad (3)$$

\mathbf{A} is the matrix of our parameters weight.

Now, let's consider Logistic Regression. The data is exactly the same as our Linear Regression. However, we also have the logistic function, which looks like a curved Z being pulled apart sideways.

$$\sigma = \frac{1}{1 + e^{-x}} \quad (4)$$

To plug in our parameters, we can input a linear function into where the x is. Because the values of our function, will fall between 0 and 1, we can write them as probabilities:

$$p = \frac{1}{1 + e^{-mx+b}} \quad (5)$$

One big difference between this and Linear Regression is our loss function. For the sake of brevity, we will briefly explain it:

$$-\frac{1}{N} \sum_{j=1}^N [y_j \ln(p_j) + (1 - y_j) \ln(1 - p_j)] \quad (6)$$

This is called the cross-entropy loss. The reason we take the log is so we can make it convex, so we can find the local minimum. It helps work around the fact that we don't have a close form solution for the most basic loss function of logistic regression. We want our result to be a vector of probabilities that add up to 1, so we have the softmax function:

$$\mathbf{p} = \sigma(\mathbf{y}) = \frac{1}{\sum_{j=1}^m e^{y_j}} (e^{y_1} \ e^{y_2} \ \dots \ e^{y_m})^T \quad (7)$$

Now, that we have covered these topics. Let's start with a perceptron, a network made up of linear threshold units. The perceptron has input and output layers, will possible hidden layers in between. Each layer, from left to right, has a bias term and weights that connect it to each component of the next layer. We see that 1 will represents an active neuron while 0 represents an inactive one. For the output neuron to activate, the weighted sum of the inputs should be greater than some threshold determined by our activation function. The reason we mentioned Linear and Logistic Regression is that we see that the functions we use to represent these calculation resemble the calculations we perform within those statistical methods. For example, we can write the equation of the multilayer perceptron below as:

$$\mathbf{x}_2 = \sigma(\mathbf{A}_1 \mathbf{x}_1 + \mathbf{b}_1) \quad \mathbf{y} = \sigma(\mathbf{A}_2 \mathbf{x}_2 + \mathbf{b}_2) \quad (8)$$

We have that \mathbf{b}_1 are the biases between input and hidden layers, \mathbf{b}_2 are the biases between the hidden and output layers, \mathbf{A}_1 and \mathbf{A}_2 are the respective weight matrices between the layers, and \mathbf{y} is our output layer.

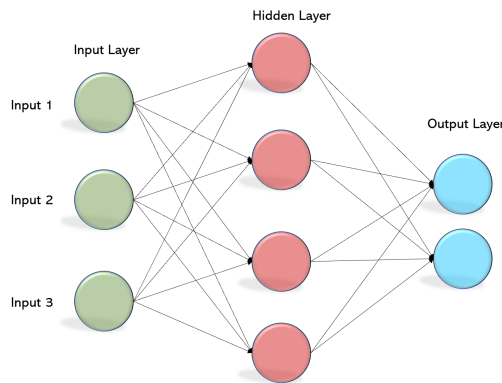


Figure 1: Example of Perceptron

(a) Convolutional Networks

Like the original perceptrons, convolutional neural networks also take inspiration from biology, specifically the eyes this time. But to greatly simplify it, specific neurons are responsible for specific regions, *receptive fields*, of our eyes. We can say that there are layers of neurons that take in the information, and process the images, from simple to complex.

This is replicated in machine learning, except we do it with numbers. Consider a matrix of numbers, or an image essentially, with the value of each pixel scaled between 0 and 1. We can look at different receptive fields of the matrix, usually $n \times n$ with a *filter kernel* to create a *feature map*. To center our filter at every value if desired, we can use zero padding, which creates a boundary around of 0 around the image. Sometimes, we can increase the *stride*, that lets us adjust the shift of our filter and size of our feature map. Next, we can do pooling, which is like our filter kernel, except we're fundamentally sub-sampling. Some of most popular options include *max pooling*, grabbing the max value of the pool, or *average pooling*, calculating the average of the pool.

3 Algorithm Implementation and Development

For our assignment, we will be building a fully-connected, feed-forward neural network, where every neuron is connected to every neuron in the previous layer and we move in one direction. We will be experimenting with the options of the neural network like the width of the layers, the step function, and the L2 Regularization parameters. It's going to be an adventure, tumultuous and emotional. Here, is where we go from the civilized, wealthy towns of the East Coast to the wild, wild West.

But, here's the situation. After we spend hours trying to build a neural network using fully connected layers, we'll experiment with convolutional layers and mess with some more options like stride, padding, and normalization layers. How many hours will we spend trying to hit our accuracy goal? Find out in the next section.

4 Computational Results

In this section, we go over the fruit of my laptop's labor, over a couple hours of nearly overheating. We always input our data as 28x28x1.

(a) Fully Connected Only

Confusion Matrix										
Output Class	0	1	2	3	4	5	6	7	8	9
0	5204 9.5%	0 0.0%	9 0.0%	2 0.0%	1 0.0%	0 0.0%	100 0.2%	0 0.0%	0 0.0%	97.9% 2.1%
1	1 0.0%	5430 9.9%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	100.0% 0.0%
2	14 0.0%	0 0.0%	5108 9.3%	0 0.0%	51 0.1%	0 0.0%	72 0.1%	0 0.0%	0 0.0%	97.4% 2.6%
3	134 0.2%	12 0.0%	34 0.1%	5474 10.0%	103 0.2%	0 0.0%	119 0.2%	0 0.0%	2 0.0%	93.1% 6.9%
4	6 0.0%	2 0.0%	214 0.4%	21 0.0%	5303 9.6%	0 0.0%	73 0.1%	0 0.0%	1 0.0%	94.4% 5.6%
5	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	5435 9.9%	0 0.0%	0 0.0%	1 0.0%	100.0% 0.0%
6	176 0.3%	0 0.0%	130 0.2%	2 0.0%	53 0.1%	0 0.0%	5141 9.3%	0 0.0%	1 0.0%	93.4% 6.6%
7	3 0.0%	0 0.0%	0 0.0%	0 0.0%	72 0.1%	0 0.0%	5485 10.0%	1 0.3%	155 0.3%	96.0% 4.0%
8	5 0.0%	0 0.0%	1 0.0%	1 0.0%	0 0.0%	1 0.0%	0 0.0%	5505 10.0%	0 0.0%	99.9% 0.1%
9	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	3 0.0%	5338 9.7%	99.9% 0.1%
93.9% 99.7% 92.9% 99.5% 96.2% 98.7% 93.4% 99.9% 99.9% 97.2% 97.1% 6.1% 0.3% 7.1% 0.5% 3.8% 1.3% 6.6% 0.1% 0.1% 2.8% 2.9%										
Target Class										

Figure 2: Fully Connected Only - Train

Confusion Matrix										
Output Class	0	1	2	3	4	5	6	7	8	9
0	378 7.6%	1 0.0%	6 0.1%	3 0.1%	0 0.0%	0 0.0%	55 1.1%	0 0.0%	3 0.1%	84.8% 15.2%
1	2 0.0%	544 10.9%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	99.5% 0.5%
2	4 0.1%	2 0.0%	398 8.0%	1 0.0%	22 0.4%	0 0.0%	28 0.6%	0 0.0%	1 0.0%	87.3% 12.7%
3	27 0.5%	7 0.1%	13 0.3%	483 9.7%	28 0.6%	0 0.0%	24 0.5%	0 0.0%	3 0.1%	82.6% 17.4%
4	2 0.0%	1 0.0%	50 1.0%	6 0.1%	415 8.3%	0 0.0%	28 0.6%	0 0.0%	0 0.0%	82.7% 17.3%
5	0 0.0%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	472 9.4%	0 0.0%	3 0.1%	1 0.0%	98.3% 1.7%
6	41 0.8%	1 0.0%	33 0.7%	6 0.1%	20 0.4%	0 0.0%	352 7.0%	0 0.0%	1 0.0%	77.5% 22.5%
7	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	17 0.3%	1 0.0%	502 10.0%	2 0.4%	92.3% 7.7%
8	3 0.1%	0 0.0%	3 0.1%	2 0.0%	3 0.1%	2 0.0%	0 0.0%	478 9.6%	3 0.1%	96.0% 4.0%
9	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	2 0.0%	7 0.1%	478 9.6%	98.0% 2.0%
82.7% 97.8% 79.0% 96.4% 85.0% 95.7% 71.4% 98.0% 97.6% 94.5% 90.0% 17.3% 2.2% 21.0% 3.6% 15.0% 4.3% 28.6% 2.0% 2.4% 5.5% 10.0%										
Target Class										

Figure 3: Fully Connected Only - Valid

In our final model of using only fully connected layers, we ended up using 3 fully connected layers activated by relu layers, with widths of 700, 500 and 10, then a softmax layer and classification layer. I discovered the *sgdm* optimizer. I did not use a bias and had an initial learning rate of 0.01, with 50 epochs. So, for the training, validation, and testing set accuracies, we respectively got: 97.1%, 90.0%, and 88.8%.

Confusion Matrix

0	812	1	21	9	0	0	86	0	5	0	86.9%
	8.1%	0.0%	0.2%	0.1%	0.0%	0.0%	0.9%	0.0%	0.1%	0.0%	13.1%
1	1	973	0	3	0	0	1	0	0	0	99.5%
	0.0%	9.7%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.5%
2	18	0	786	12	59	0	64	0	4	0	83.4%
	0.2%	0.0%	7.9%	0.1%	0.6%	0.0%	0.6%	0.0%	0.0%	0.0%	16.6%
3	45	20	24	941	46	1	50	0	7	0	83.0%
	0.4%	0.2%	0.2%	9.4%	0.5%	0.0%	0.5%	0.0%	0.1%	0.0%	17.0%
4	5	3	93	16	837	0	66	0	3	0	81.8%
	0.1%	0.0%	0.9%	0.2%	8.4%	0.0%	0.7%	0.0%	0.0%	0.0%	18.2%
5	2	0	1	0	0	929	0	5	3	7	98.1%
	0.0%	0.0%	0.0%	0.0%	0.0%	9.3%	0.0%	0.1%	0.0%	0.1%	1.9%
6	104	2	75	14	55	0	722	0	5	1	73.8%
	1.0%	0.0%	0.8%	0.1%	0.5%	0.0%	7.2%	0.0%	0.1%	0.0%	26.2%
7	0	0	0	1	0	48	0	983	8	56	89.7%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.5%	0.0%	9.8%	0.1%	0.6%	10.3%
8	13	1	0	4	3	1	11	1	964	0	96.6%
	0.1%	0.0%	0.0%	0.0%	0.0%	0.0%	0.1%	0.0%	9.6%	0.0%	3.4%
9	0	0	0	0	0	21	0	11	1	936	96.6%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.2%	0.0%	0.1%	0.0%	9.4%	3.4%
	81.2%	97.3%	78.6%	94.1%	83.7%	92.9%	72.2%	98.3%	96.4%	93.6%	88.8%
	18.8%	2.7%	21.4%	5.9%	16.3%	7.1%	27.8%	1.7%	3.6%	6.4%	11.2%
	0	1	2	3	4	5	6	7	8	9	
	Target Class										

Figure 4: Fully Connected Only - Test

I think this is pretty good considering our final model didn't touch a single observation in the testing to train the model.

(b) Convolutional Neural Network

Confusion Matrix

0	5204	0	9	2	1	0	100	0	0	0	97.9%
	9.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.2%	0.0%	0.0%	0.0%	2.1%
1	1	5430	0	0	0	0	1	0	0	0	100.0%
	0.0%	9.9%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
2	14	0	5108	0	51	0	72	0	0	0	97.4%
	0.0%	0.0%	9.3%	0.0%	0.1%	0.0%	0.1%	0.0%	0.0%	0.0%	2.6%
3	134	12	34	5474	103	0	119	0	2	0	93.1%
	0.2%	0.0%	0.1%	10.0%	0.2%	0.0%	0.2%	0.0%	0.0%	0.0%	6.9%
4	6	2	214	21	5303	0	73	0	1	0	94.4%
	0.0%	0.0%	0.4%	0.0%	9.6%	0.0%	0.1%	0.0%	0.0%	0.0%	5.6%
5	0	0	0	0	0	5435	0	0	0	1	100.0%
	0.0%	0.0%	0.0%	0.0%	0.0%	9.9%	0.0%	0.0%	0.0%	0.0%	0.0%
6	176	0	130	2	53	0	5141	0	1	0	93.4%
	0.3%	0.0%	0.2%	0.0%	0.1%	0.0%	9.3%	0.0%	0.0%	0.0%	6.6%
7	3	0	0	0	0	72	0	5485	1	155	96.0%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.1%	0.0%	10.0%	0.0%	0.3%	4.0%
8	5	0	1	0	1	0	1	0	5505	0	99.9%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	10.0%	0.0%	0.1%
9	0	0	0	0	0	0	3	0	5338	99.9%	0.1%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	9.7%	0.0%	0.1%
	93.9%	99.7%	92.9%	99.5%	96.2%	98.7%	93.4%	99.9%	99.9%	97.2%	97.1%
	6.1%	0.3%	7.1%	0.5%	3.8%	1.3%	6.6%	0.1%	0.1%	2.8%	2.9%
	0	1	2	3	4	5	6	7	8	9	
	Target Class										

Figure 5: Convolutional - Train

Confusion Matrix

0	378	1	6	3	0	0	55	0	3	0	84.8%
	7.6%	0.0%	0.1%	0.1%	0.0%	0.0%	1.1%	0.0%	0.1%	0.0%	15.2%
1	2	544	0	0	0	0	1	0	0	0	99.5%
	0.0%	10.9%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.5%
2	4	2	398	1	22	0	28	0	1	0	87.3%
	0.1%	0.0%	8.0%	0.0%	0.4%	0.0%	0.6%	0.0%	0.0%	0.0%	12.7%
3	27	7	13	483	28	0	24	0	3	0	82.6%
	0.5%	0.1%	0.3%	9.7%	0.6%	0.0%	0.5%	0.0%	0.1%	0.0%	17.4%
4	2	1	50	6	415	0	28	0	0	0	82.7%
	0.0%	0.0%	1.0%	0.1%	8.3%	0.0%	0.6%	0.0%	0.0%	0.0%	17.3%
5	0	0	1	0	0	472	0	3	1	3	98.3%
	0.0%	0.0%	0.0%	0.0%	0.0%	9.4%	0.0%	0.1%	0.0%	0.1%	1.7%
6	41	1	33	6	20	0	352	0	1	0	77.5%
	0.8%	0.0%	0.7%	0.1%	0.4%	0.0%	7.0%	0.0%	0.0%	0.0%	22.5%
7	0	0	0	0	0	17	1	502	2	22	92.3%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.3%	0.0%	10.0%	0.0%	0.4%	7.7%
8	3	0	3	2	3	2	4	0	478	3	96.0%
	0.1%	0.0%	0.1%	0.0%	0.1%	0.0%	0.1%	0.0%	9.6%	0.1%	4.0%
9	0	0	0	0	0	2	0	7	1	478	98.0%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.1%	0.0%	9.6%	2.0%
	82.7%	97.8%	79.0%	96.4%	85.0%	95.7%	71.4%	98.0%	97.6%	94.5%	90.0%
	17.3%	2.2%	21.0%	3.6%	15.0%	4.3%	28.6%	2.0%	2.4%	5.5%	10.0%
	0	1	2	3	4	5	6	7	8	9	
	Target Class										

Figure 6: Convolutional - Valid

For our convolutional networks, we actually used starter code provided by MATLAB themselves. It features 3 sets of a convolution layer with padding, a batch normalization layer, and a relu activation layer. There is a max pooling layer with stride 2 after the first and second set. The convolution layers are size [3 3], and the number of filters go as 16, 32, and 64 as we proceed deeper. Finally, we have a fully connected layer with 10 neurons, a softmax layer, and finally, the classification layer.

We run it through 10 epochs, and an initial learning rate and L2 regularization of 0.0001. Also, we switched back to our default optimizer, adam. Our final model accuracies for train, valid, and test were: 92.8%, 90.5%, and 89.3%.

Output Class	0	1	2	3	4	5	6	7	8	9	
	868 8.7%	2 0.0%	22 0.2%	15 0.1%	1 0.0%	0 0.0%	152 1.5%	0 0.0%	3 0.0%	0 0.0%	81.7% 18.3%
	0 0.0%	972 9.7%	1 0.0%	4 0.0%	3 0.0%	0 0.0%	1 0.0%	0 0.0%	3 0.0%	0 0.0%	98.8% 1.2%
	26 0.3%	0 0.0%	873 8.7%	9 0.1%	123 1.2%	0 0.0%	92 0.9%	0 0.0%	6 0.1%	0 0.0%	77.3% 22.7%
	17 0.2%	15 0.1%	7 0.1%	916 9.2%	36 0.4%	0 0.0%	31 0.3%	0 0.0%	3 0.0%	0 0.0%	89.4% 10.6%
	3 0.0%	6 0.1%	48 0.5%	18 0.2%	781 7.8%	0 0.0%	63 0.6%	0 0.0%	1 0.0%	0 0.0%	84.9% 15.1%
	1 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	974 9.7%	0 0.0%	12 0.1%	3 0.0%	9 0.1%	97.5% 2.5%
	78 0.8%	1 0.0%	48 0.5%	37 0.4%	56 0.6%	0 0.0%	651 6.5%	0 0.0%	5 0.1%	0 0.0%	74.3% 25.7%
	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	20 0.2%	0 0.0%	965 9.7%	4 0.0%	32 0.3%	94.5% 5.5%
	7 0.1%	3 0.0%	1 0.0%	1 0.0%	0 0.0%	0 0.0%	10 0.1%	2 0.0%	972 9.7%	0 0.0%	97.6% 2.4%
0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	6 0.1%	0 0.0%	21 0.2%	0 0.0%	959 9.6%	97.3% 2.7%	
	86.8% 13.2%	97.2% 2.8%	87.3% 12.7%	91.6% 8.4%	78.1% 21.9%	97.4% 2.6%	65.1% 34.9%	96.5% 3.5%	97.2% 2.8%	95.9% 4.1%	89.3% 10.7%
	0	1	2	3	4	5	6	7	8	9	

6 Appendix

A MATLAB Functions

- *classify*: Compare fitted labels with the observed values
- *trainingOptions*: Specify options for our neural network
- *trainNetwork*: Find optimal weights and biases for our hyper parameters of the neural networks

B MATLAB Code

```
clear; clc; close all;

load('fashion_mnist.mat')
%%
X_train = im2double(X_train);
X_test = im2double(X_test);

X_train = reshape(X_train,[60000 28 28 1]);
X_train = permute(X_train,[2 3 4 1]);

X_test = reshape(X_test,[10000 28 28 1]);
X_test = permute(X_test,[2 3 4 1]);

X_valid = X_train(:,:,1:5000);
X_train = X_train(:,:,5001:end);

y_valid = categorical(y_train(1:5000))';
y_train = categorical(y_train(5001:end))';
y_test = categorical(y_test)';
%%
layers = [imageInputLayer([28 28 1])
          fullyConnectedLayer(784)
          reluLayer
          fullyConnectedLayer(200)
          reluLayer
          fullyConnectedLayer(100)
          reluLayer
          fullyConnectedLayer(10)
          softmaxLayer
          classificationLayer];
options = trainingOptions('sgdm', ...
    'MaxEpochs',15,...
    'InitialLearnRate',1.1e-2, ...
    'L2Regularization',1.7e-5, ...
    'ValidationData',{X_valid,y_valid}, ...
    'Verbose',true)
net = trainNetwork(X_train,y_train,layers,options);

%% Confusion for training
figure(2)
y_pred = classify(net,X_train);
plotconfusion(y_train,y_pred)
```

```

figure(3)
y_pred2 = classify(net,X_valid);
plotconfusion(y_valid,y_pred2)
%% Test classifier
figure(3)
y_pred3 = classify(net,X_test);
plotconfusion(y_test,y_pred3)

%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
PART 2
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
clear; clc; close all;

load('fashion_mnist.mat')
%%
X_train = im2double(X_train);
X_test = im2double(X_test);

X_train = reshape(X_train,[60000 28 28 1]);
X_train = permute(X_train,[2 3 4 1]);

X_test = reshape(X_test,[10000 28 28 1]);
X_test = permute(X_test,[2 3 4 1]);

X_valid = X_train(:,:,1:5000);
X_train = X_train(:,:,5001:end);

y_valid = categorical(y_train(1:5000))';
y_train = categorical(y_train(5001:end))';
y_test = categorical(y_test)';
%%
layers = [
imageInputLayer([28 28 1])

    convolution2dLayer(3,16,'Padding',1)
    batchNormalizationLayer
    reluLayer

    maxPooling2dLayer(2,'Stride',2)

    convolution2dLayer(3,32,'Padding',1)
    batchNormalizationLayer
    reluLayer

    maxPooling2dLayer(2,'Stride',2)

    convolution2dLayer(3,64,'Padding',1)
    batchNormalizationLayer
    reluLayer

    fullyConnectedLayer(10)
    softmaxLayer

```

```

        classificationLayer];
options = trainingOptions('adam', ...
    'MaxEpochs',10,...
    'InitialLearnRate',1e-4, ...
    'L2Regularization',1e-4, ...
    'ValidationData',{X_valid,y_valid}, ...
    'Verbose',false, ...
    'Plots','training-progress');

net = trainNetwork(X_train,y_train,layers,options);
%% Confusion for training
figure(1)
y_pred = classify(net,X_train);
plotconfusion(y_train,y_pred)

figure(3)
y_pred2 = classify(net,X_valid);
plotconfusion(y_valid,y_pred2)
%% Test classifier
figure(3)
y_pred3 = classify(net,X_test);
plotconfusion(y_test,y_pred3)

```

C References

<https://www.mathworks.com/solutions/deep-learning/examples/training-a-model-from-scratch.html>

Lecture Notes AMATH 482 Winter 2020

<https://stats.stackexchange.com/questions/376312/mnist-digit-recognition-what-is-the-best-we-can-get-with-a-fully-connected-nn-o>