



MASTER CORO M1  
MASTER IN CONTROL AND ROBOTICS

PROJECT REPORT

June 2018

**New labs for the Computer Vision  
module**

*Author:*

Ashok Muralidharan

Thuc-Long Ha

*Supervisor:*

Bogdan Khomutenko

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Objective of the project . . . . .	2
<b>2</b>	<b>Background</b>	<b>3</b>
2.1	Image classification . . . . .	3
2.2	Convolutional Neural Networks . . . . .	3
<b>3</b>	<b>Dataset</b>	<b>7</b>
3.1	Structure . . . . .	7
3.2	Train Set . . . . .	7
3.3	Test Set . . . . .	8
<b>4</b>	<b>Implementation</b>	<b>9</b>
4.1	Software and Technologies used . . . . .	9
4.2	Architecture of CNN . . . . .	10
4.3	Major steps in implementation . . . . .	10
<b>5</b>	<b>Results</b>	<b>13</b>
5.1	Result from our developed architecture . . . . .	13
5.2	Report by class . . . . .	16
5.3	Learning curve & Training Cost . . . . .	19
5.4	Comparing with other models . . . . .	20
<b>6</b>	<b>Conclusion &amp; Futher Development</b>	<b>21</b>
<b>7</b>	<b>Appendix</b>	<b>24</b>
7.1	References . . . . .	24

# 1 Introduction

With the advent of high computing power and large amount of images being generated every day, deep learning techniques are starting to play a crucial role in computer vision applications. A number of challenging applications such as image classification, object detection and localization, face recognition, etc.. are now being solved efficiently, thanks to deep learning. The key essence of deep learning are two things - large amount of training data and good computing power. Deep learning algorithms are proven to perform well with more training data, thus has the ability to achieve accuracy in certain problems at par or some times better than humans.

## 1.1 Objective of the project

The primary objective of our project is to design a new lab material for the computer vision module. The main aim of the lab material was decided to illustrate the importance of deep learning in computer vision applications, different techniques and technologies used for developing such algorithms quickly and efficiently. The task that was chosen was image classification. The task was to build an image classifier that should predict the road signs using a deep learning technique called Convolutional Neural Networks(CNNs). Few example images of the task are shown in the figure 1

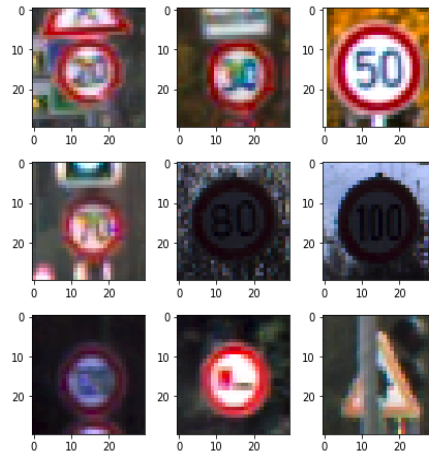


Figure 1: Traffic Signals for classification

## 2 Background

### 2.1 Image classification

### 2.2 Convolutional Neural Networks

Convolutional Neural Networks (CNNs) are the specialized Neural Networks which have been very effective in artificial intelligence tasks such as image classification and recognition. CNNs have been successful in identifying objects, human faces and traffic signs.

LeNet was one of the CNNs which started the researches of Deep Learning. LeNet5[6], which is developed by Yann LeCun, was the result after many previous successful development since the year 1988. During the past, LeNet architecture was used mainly for character recognition tasks such as reading zip codes, digits, etc

Below is the explanation how LeNet architecture is used to learn to recognize images.

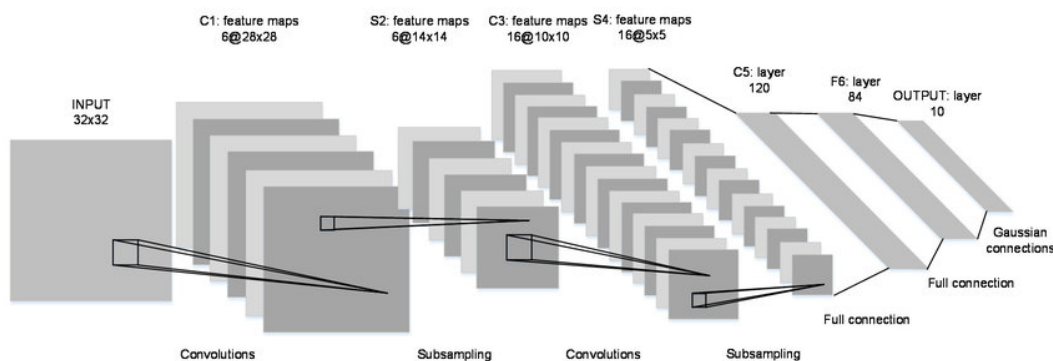


Figure 2: Lenet5 Architecture[6]

The Convolutional Neural Network above is identical to the original LeNet used for character recognition works.

The main operations in the CNN shown in Figure 1 are:

- Convolution
- Non Linearity (ReLU)
- Pooling or Sub Sampling
- Fully Connected Layer

**Convolution:**

We can imagine convolution in deep learning as a dot product of two vectors. The pixels of the image is convolved with corresponding values in the filter (i.e) the pixels are multiplied with the corresponding filter values are added to form the output. Thus the convolved output will have reduced dimensions than that of the original input. It is just like applying a filter in image processing. A typical example is shown in the figure 3

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	3	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

\*

1	0	-1
1	0	-1
1	0	-1

=

-5	-4	0	8
-10	-2	2	3
0	-2	-4	-7
-3	-2	-3	-16

Figure 3: Convolution - Vertical edge detection

**ReLU - Nonlinear Activation Function**

As we can see that neural networks consists of bunch of neurons connected in linear fashion performing multiplication and addition ultimately. But many functions are non-linear and linear model might just not be good enough to represent the data. To add non-linearity to the CNNs we use the ReLU(Rectified Linear Unit) as activation function. ReLU is a function that thresholds the value to zero. The function is can be represented as  $f(x) = \max(0, x)$ . Graphically representation of the function is shown in the figure 4. Apart from ReLU, there are many other activation functions such as sigmoid, tanh, etc.. that are used in deep learning.

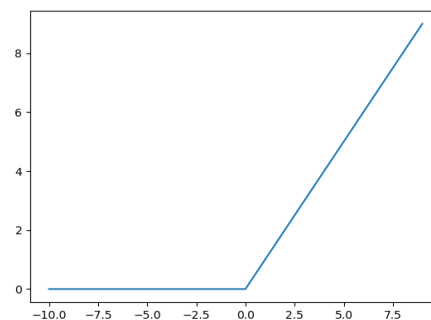


Figure 4: ReLU

**Pooling or Sub Sampling:**

The next major part of CNNs are pooling layers. The commonly used pooling functions are max-pooling and average-pooling. But max-pooling is the most adopted one. We can imagine max-pooling as something like taking the max value among a region that we see in the input. This way we try to reduce the dimension of the input, by generalizing based on the max information considering the adjacent values. The figure 5 should give a clear idea of max pooling.

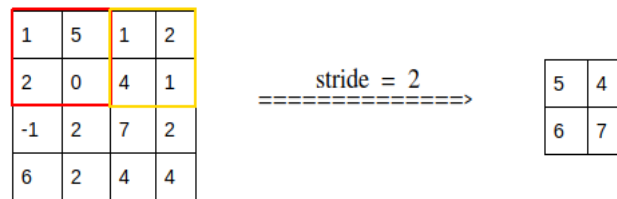


Figure 5: Max Pooling - 2X2 with stride 2

**Fully connected layers:**

The last major component in CNNs are fully connected networks. Fully connected networks are just the feed forward deep neural networks. Generally, in CNNs, one or two layers of fully connected neural network layers are added. The feed forward networks consists of number of neurons which are connected fully with each other, which performs multiplication(with weights) and addition(with biases) operation. Again a non-linearity is necessary here, which is called activation function. The figure 6 should give a brief idea of deep neural networks.

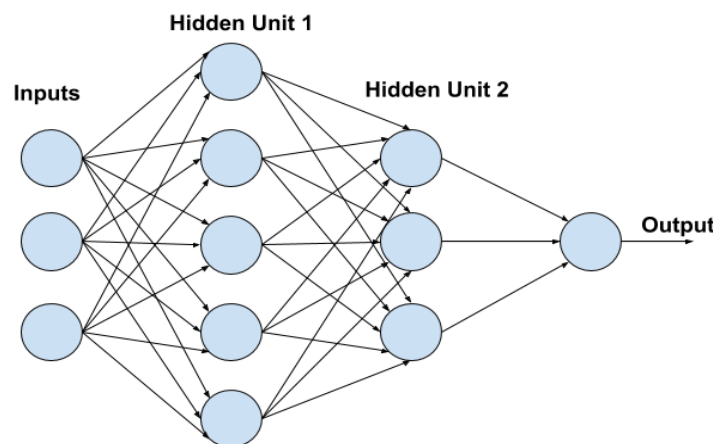


Figure 6: Fully Connected Networks

Apart from full connections, we also have biases which is basically a number that is

added. Biases are used to shift our activation function to left or right. For brevity, we are not getting into the details of deep neural networks working, but here is what a single neuron(say in hidden layer 1) does

$$\begin{aligned} Inputs &= X1, X2, X3 \\ ConnectionWeights &= W1, W2, W3 \\ Bias &= b1 \\ Output &= ActivationFunction(W1 * X1 + W2 * X2 + w3 * X3 + b1) \end{aligned}$$

Similarly, every neuron will multiply the inputs and weights that it receives, add a bias and apply activation function such as ReLU it. This forward computation in neural networks is called as the forward propagation.

### Back Propagation:

The learning algorithm that is generally used for the neural networks is back propagation. As the name suggests, back propagation involves propagation of error backwards in the network by computing the gradients at each step and updating the weights and biases which are the learning parameters. The main intuition here is to consider the CNN as one giant optimization problem, have all weights and biases as the system variables and error to decrease or optimize.

This process happens during the training phase. Initially we randomly initialize weights and biases. We pass the input through the network(forward propagation) and compute our prediction. We then find the error(prediction value - actual value). We then compute the gradient with respect to weights in the previous layer that affects the error.

$$\frac{\partial error}{\partial W1}$$

We can use the chain rule for updating the weights deep in the layer. An example is shown below.

$$\frac{\partial error}{\partial W2} = \frac{\partial W3}{\partial W2} * \frac{\partial error}{\partial W3}$$

Thus by updating the weights towards the direction of gradient, and doing this iteratively, we will minimize our error, thus move close to the actual output.

## 3 Dataset

### 3.1 Structure

The data set that we used in the lab were obtained from German Traffic sign dataset [9]. The images are in ppm format. Training and test images are divided into two different folders for facilitating easier retrieval of images during the process.

**Training Images Location:** /GTRSB/Final\_Training/Images

**Testing Images Location:** /GTRSB/Final\_Test/Images

Each Image has size that varies between 15x15 to 250x250 pixels. So we can observe that preprocessing of images is a necessary step here to make the dimensions of the images similar. Also the actual traffic sign is not always centered within the image.

### 3.2 Train Set

The training folder contains 43 folder named from 0 to 42 indicating the classes to which the images belong. There is a 'csv' file in every folder which can be used to read the image and the corresponding class of the image. There are a total of 39209 images as a part of the training dataset. Some of the sample images are shown in figure 7. As we can see, the dataset has lots of variance in lighting, orientation, scale, etc.

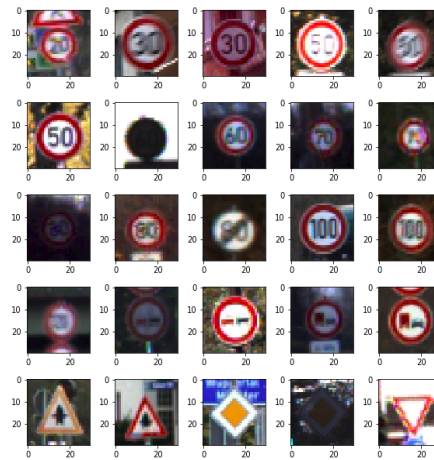


Figure 7: Training Images



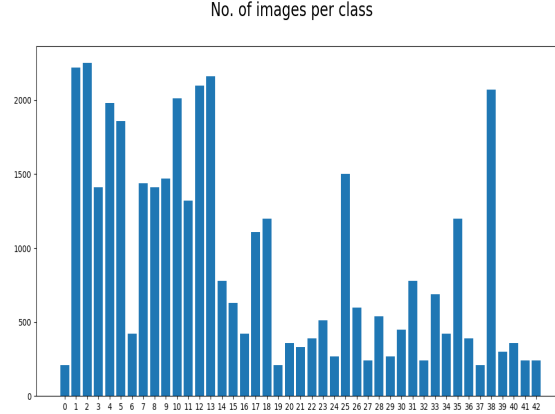


Figure 8: No. of images per class

As we can see from 8, the number of images per class is definitely not uniform. But the distribution is not very bad either. One way to overcome this uneven distribution is to find more data to add to the training samples. This is always possible when we take our training data on our own. But here since we are working on a provided dataset, we are not adding additional images. We can also perform image augmentation from the original image to get duplicate augmented images and add them to the dataset. But adding a lot of augmented images might affect the network generality in real world. We have not augmented images here, but have used augmentation during the run time which will be discussed later.

### Validation Set:

In deep learning, there is always a validation set, which is generally an extract from the training set. This set will be used during the training process for validating our model for every epoch (one complete pass through the training set). This becomes essential for early stopping of the training if our training process does not improve after few iterations. This validation set will not get involved in the training procedure.

## 3.3 Test Set

The test folder contains 12630 images. The file GT-final.test.csv labels the class and properties of each image. Some sample images are shown in figure 9. These images will not get involved in the training procedure and will be purely used for testing purpose.

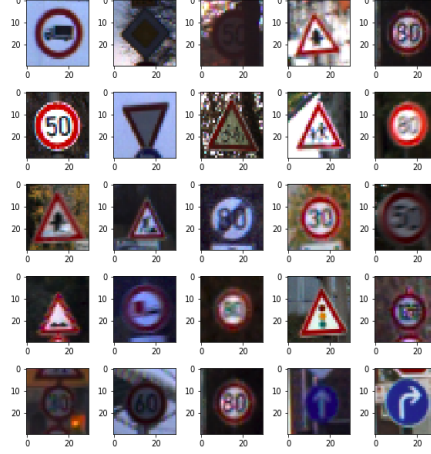


Figure 9: Test Images

## 4 Implementation

### 4.1 Software and Technologies used

The following are the programming languages and major technologies used for our purpose

- Python - programming language
- Numpy - Scientific computation library
- TensorFlow - open source machine learning framework from Google
- OpenCV - Open Source Computer Vision Library
- Git - Version Control System(VCS)

Python is the preferred language for deep learning for computer vision applications since it facilitates fast prototyping and has excellent community support. TensorFlow is a popular framework for machine learning applications which was developed and has support from Google. TensorFlow plays the major part of CNN implementation. This software stack is very popular and prevalent in the computer vision community. We also use GIT for the purpose of version control of the code.

## 4.2 Architecture of CNN

The architecture which we used is shown in the figure 10. We did not go with industry standard, state of art architectures like inception [10], VGG-16 [8], etc. because the network in these architectures are quite big and training will take a lot of time which might not be feasible to do in the labs. So we followed a pattern in which the network has CONV, CONV, POOL, CONV, CONV, POOL layers stacked in order. This architecture was chosen after some trial experiments with layers and with some inspirations from the internet. This architecture is fairly simple, so that this can quickly run in the labs as well as produce decent results. You can see the architecture in the figure 10.

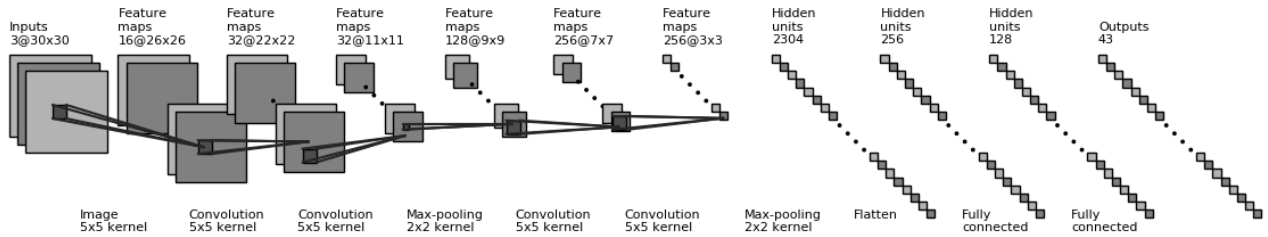


Figure 10: CNN Architecture [3]

As we can see the image that we use is of dimension (30X30X3). Thus it is an RGB image with width and height of 30 pixels. The image is passed through a series of CONV, CONV, POOL, CONV, CONV, POOL layers followed by two fully connected layers. The dimensions after each operation performed in each layer is shown in the figure 10. The operation performed is shown in the bottom of the layer and the resultant dimension is shown in the top of the layer.

## 4.3 Major steps in implementation

- Image Preprocessing and Augmentation
- Data Preparation
- Forward Propagation

- Back Propagation
- Model Evaluation, Logging and Saving

**Image Preprocessing and Augmentation:**

Since our model expects a (30X30X3) image, we need to resize the image to that dimension. We use OpenCV for this purpose. Initially, we tried using gray scale image for the purpose, but since traffic signs have colored content in them with specific meaning, we decided to go with RGB image later. Also with RGB image, the accuracy seem to slightly increase. We actually got a very good inspiration from an online article regarding the augmentation of images [11] where the images are augmented during the computation process. During initial phases of training, we add more augmentation but at later stages of training we reduce the augmentation so that the model fits to the training data set well. This had a small impact on our accuracy improvement.

**Data Preparation:**

Data preparation is an important step in the deep learning applications. Since deep learning works with data, proper structuring of data is necessary for getting desired results. Here we need to collect the image and split them to training set, validation set and testing set. These three sets are standards that are followed in deep learning community to enable proper evaluation and testing of the model. In order to predict the class, we design the outputs of the network as a onehot encoding. A onehot encoding is a typical structure used for classification problems. A onehot encoding has only 1s and 0s in it indicating the presence and the absence of the particular class(in our case vector of size 43).

**Forward Propagation:**

As explained earlier, the forward propagation of input through the network is done to get the predicted output. The forward propagation in general consists of series of multiplication and addition. The image is passed through a series of CONV, CONV, POOL, CONV, CONV, POOL layers followed by two fully connected layer in our architecture. We have already seen about the functions of each layer. We can imagine CONV layers as number of filters and these information are merged to get our final prediction. Our final output here is a vector of dimension 42(0-43) which indicates the presence or absence of particular class. TensorFlow lets us build these steps in a very easy and elegant fashion. The codes are available in the appendix section for reference. Thus the output of our forward propagation is one pass through the entire network producing output.

**Back Propagation:**

The learning algorithm in neural networks is back propagation. The algorithm was explained earlier where we compute the gradients of the error with respect to

every parameter of the model. It is during the back propagation where our model generalizes to the given training inputs and outputs.

**Model Evaluation, Logging and Saving:**

In deep learning applications, proper evaluation of the model at every steps during the training process becomes crucial as it provides us with a metric of knowing how well our model is performing or improving. Since CNNs involves lots of calculations with lots of images, the training process is painfully slow in nature. Not having a evaluation of model might result in spending enormous amount of time and resource on a bad model. With the evaluation in place, it helps us in early stopping of the training process, when we do not see any improvement in our learning. A model evaluation is something like determining the accuracy of the batch, loss function in the batch etc.. This is a very good practise that is commonly followed in the deep learning community. Also logging results such as accuracy, loss, etc.. onto the console during the training helps us get a glimpse of what is happening with each epoch. We can also write the results to a file, but it is not implemented as a part of the project. Also saving the model after every few iterations is a nice way to have different checkpoints which can be used to retrieve the model later. TensorFlow provides elegant ways for performing all these operations. We saved a copy of the model for every five epochs. This way we can use the best performing model from the list of models.

## 5 Results

### 5.1 Result from our developed architecture

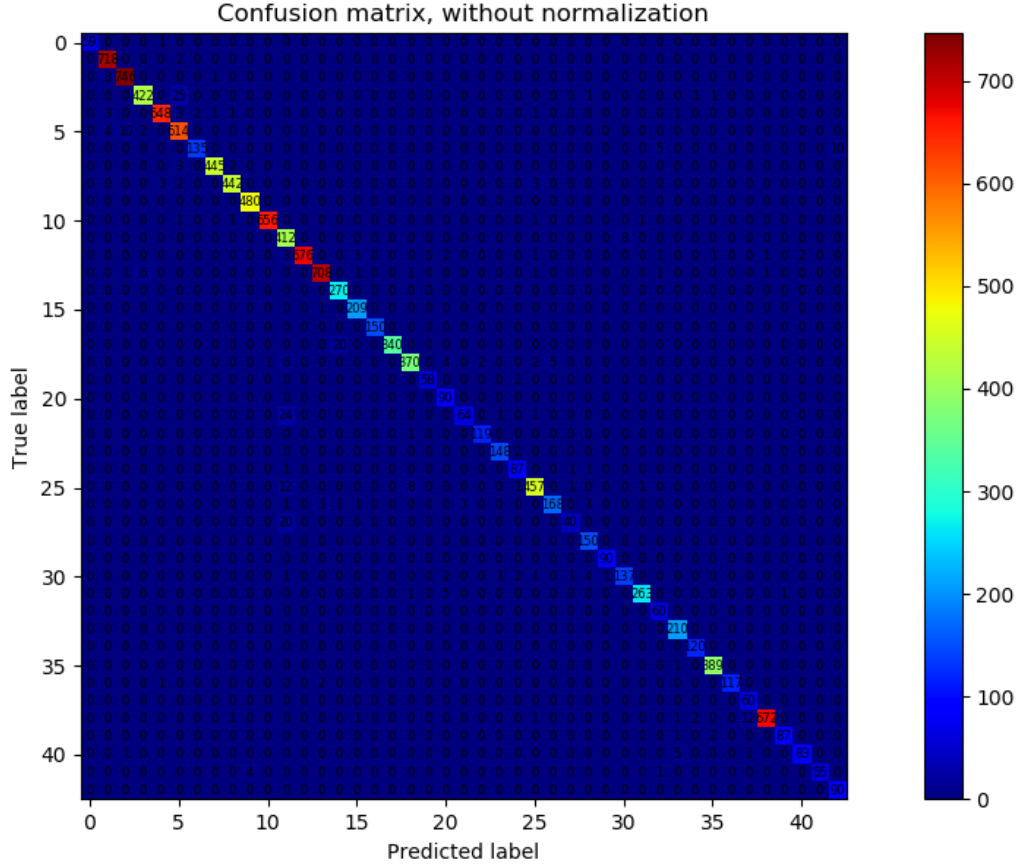


Figure 11: Confusion Matrix [4]

The confusion matrix in the figure 11 gives us a glimpse of overall classification performance. From the confusion matrix, we could see that the most prediction result concentrate in the diagonal line which is true positive. However, test data are uneven distributed in some classes. It is a matrix of actual values and the predicted values with count as the each cell of matrix. Thus we can see which class was wrongly predicted(confused) with which class. Hence the name confusion matrix.

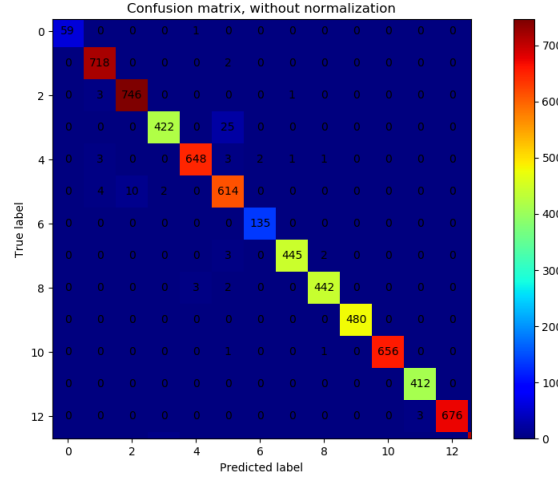


Figure 12: Zoomed - Confusion Matrix

The zoomed version of the confusion matrix is shown in the figure 12. As we can clearly see that the class 3 was wrongly predicted with class 5, 25 times.(i.e) 25 images that belong to class 3 was predicted class 5 by the network. Now let us examine what classes 3 and 5 are. The figure 13 shows the classes that were confused. It is evident that these two classes have similar features in them (like circle with red border, number 8 and 6 having similar curves, etc.). Thus we should do a some study on our analysis and not just settle on our accuracy. These classifications can be tricky.



Figure 13: Class 3 and 5

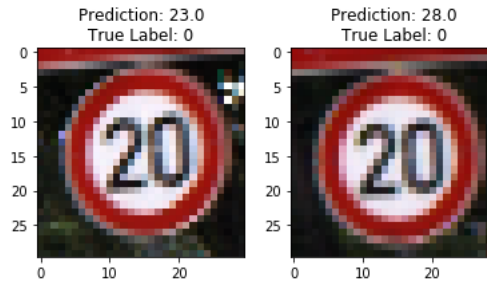


Figure 14: Failed prediction for class 0

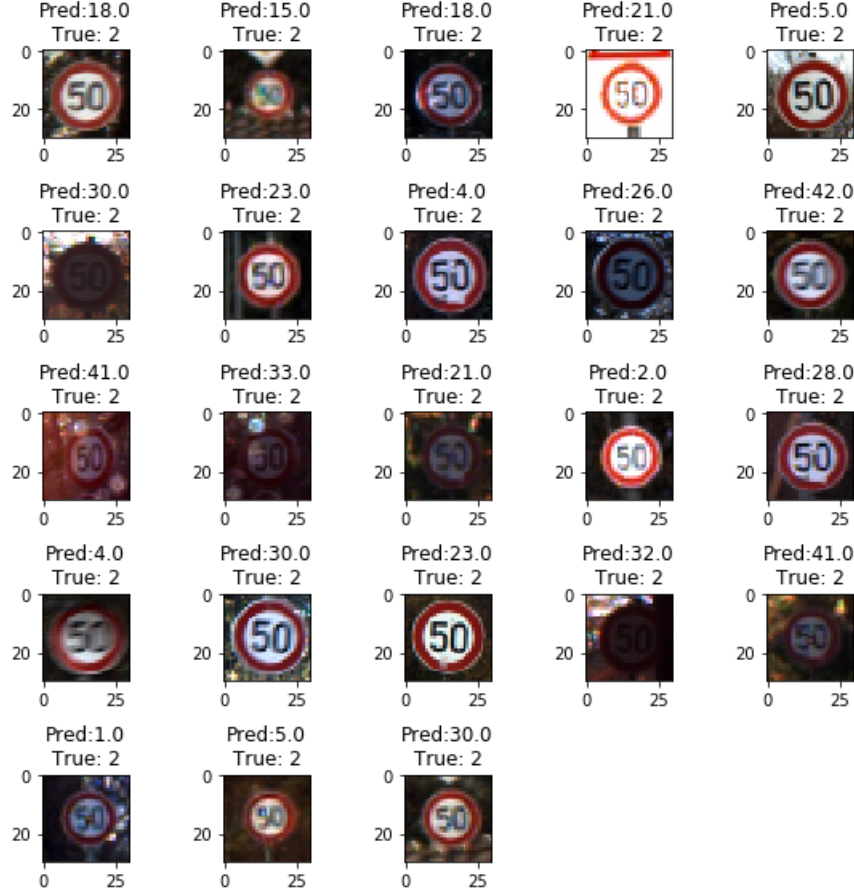


Figure 15: Failed prediction for class 2

From the failed prediction cases, we could see that there is still some wrong predictions among the classes with similar features such as shape, color, etc. This is definitely a drawback of deep learning applications as debugging the reason for failure in these cases is quite difficult task. Further tuning of hyper parameters such as learning rate, type of optimizer, etc., modifying the model architecture, might result in improved accuracy. We did not find time to play around with our architecture for improving accuracy but we managed to achieve model generalization and decent test accuracy of about 97.2%



## 5.2 Report by class

		Predicted Value	
		TRUE	FALSE
Actual Value	TRUE	True Positive	False Negative
	FALSE	False Positive	True Negative

Figure 16: True, False - Positives and Negatives

**True Positives (TP)**- These are the predicted positive values which are correct. So the value of real class is 'yes' and so is the value of predicted class. In our case, if actual class value indicates the sign and the sign you actual see the same thing.

**True Negatives (TN)** - These are the predicted negative values which are correct. It means the value of actual class is 'no' and so is the value of predicted class. In our case, if actual class value indicates the sign does not belong to one specific class and the sign you actual see is not that sign class also.

False positives and false negatives, these values happen when your actual class is in conflict with with the predicted class.

**False Positives (FP)** – If the actual class is 'no' and predicted class is 'yes'. For example, if actual class says that the sign is not 20km/h speed limit but the predicted class is 20km/h speed limit

**False Negatives (FN)** – If the actual class is 'yes' but predicted class is 'no'. For example, if actual class value indicates the sign is speed limit 20km/h but the predicted class classified otherwise.

**Accuracy** - Accuracy is the the ratio between the correctly predicted observation to the total of observations. It is intuitive indicator but is helpful only for the symmetric datasets where values of false positive and false negatives are almost same. Therefore, we have to observe other parameters to define the performance of the model. For our model, we have got 0.97 which means our model is approx. 97% accurate.

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{FN} + \text{TN}}$$

**Precision** - Precision is the ratio between predicted positive observations which are correct to the total predicted positive observations. This number tells how many signs are actual belong to one specific class among those is classified into that class. We got 0.97 precision on average which is good.

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

**Recall** - Recall is the ratio between predicted positive observations which are correct to all of the observations in the actual class. In our case, among the actual 20km/h limit sign, how many we have predicted correctly.

$$\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$$

**F1 score** - F1 Score is the harmonic average of the precision and recall. This number takes both false positives and false negatives into account. Although it is not as easy to interpret as accuracy, but F1 is more useful than accuracy in some case, especially when the class distribution is uneven. Accuracy works best if false positives and false negatives have similar cost. If the cost of false positives and false negatives are very different, it's better to look at both Precision and Recall. In our model, F1 score is 0.97.

$$\text{F1 Score} = 2 * (\text{Recall} * \text{Precision}) / (\text{Recall} + \text{Precision})$$

Below are our statistics of the scores in each class.

	precision	recall	f1-score	support
0	1.00	0.98	0.99	61
1	1.00	0.99	0.99	727
2	1.00	0.97	0.98	774
3	0.97	0.98	0.97	446
4	0.99	0.99	0.99	659
5	0.97	0.98	0.98	623
6	0.81	0.99	0.89	123
7	1.00	0.99	0.99	453
8	0.98	0.96	0.97	459
9	1.00	1.00	1.00	478
10	0.98	0.99	0.99	653
11	0.95	0.98	0.96	404
12	0.96	0.99	0.98	667
13	1.00	0.99	0.99	726
14	1.00	0.99	0.99	273
15	1.00	0.93	0.96	224
16	0.99	0.99	0.99	150
17	0.98	0.99	0.99	355
18	0.93	0.94	0.94	387
19	1.00	0.97	0.98	62
20	1.00	0.86	0.92	105
21	0.96	0.71	0.82	121
22	0.88	0.99	0.93	107

23	0.99	0.86	0.92	174
24	0.89	0.99	0.94	81
25	0.99	0.98	0.98	485
26	0.99	0.94	0.96	191
27	0.50	0.88	0.64	34
28	0.97	0.94	0.95	155
29	1.00	0.95	0.97	95
30	0.74	0.87	0.80	128
31	0.99	0.99	0.99	271
32	1.00	0.91	0.95	66
33	1.00	0.96	0.98	219
34	0.99	0.98	0.99	121
35	0.97	0.99	0.98	382
36	0.93	1.00	0.97	112
37	1.00	0.98	0.99	61
38	0.97	0.99	0.98	672
39	0.96	1.00	0.98	86
40	0.97	0.90	0.93	97
41	1.00	0.94	0.97	64
42	1.00	0.91	0.95	99
avg / total	0.97	0.97	0.97	12630

For final result on test set, our architecture got 97.2% accuracy

### 5.3 Learning curve & Training Cost

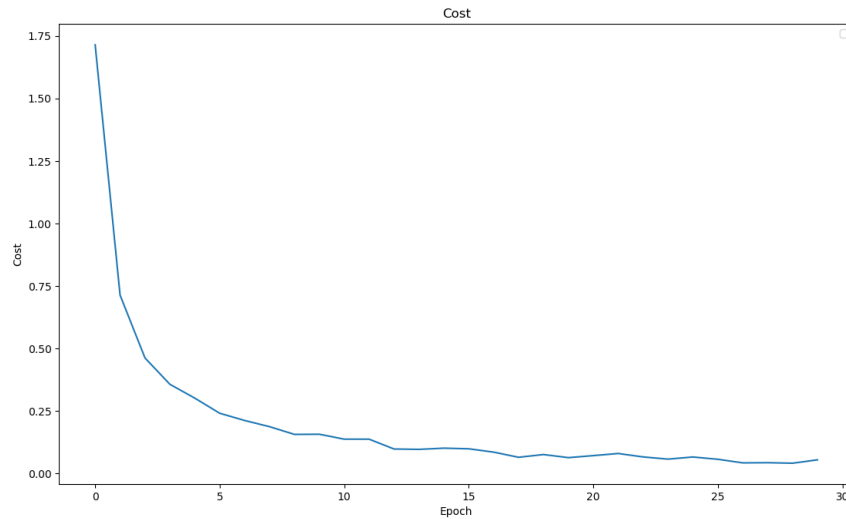


Figure 17: Training cost

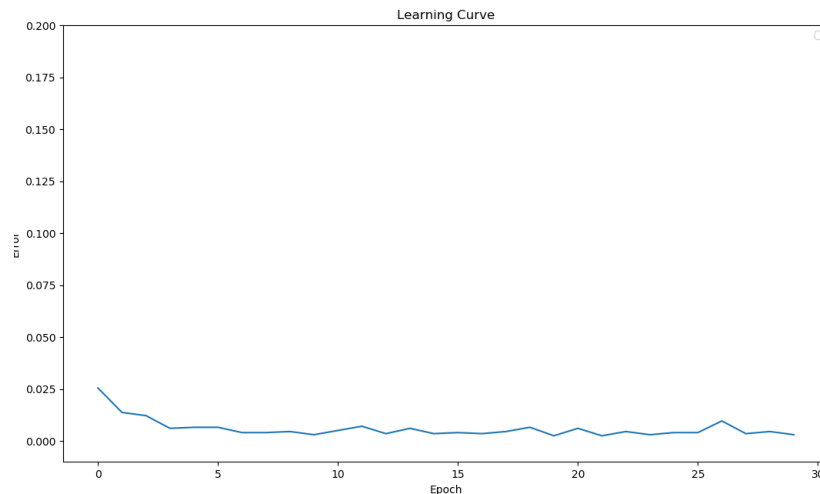


Figure 18: Training learning curve

Training cost curves can be seen above. The training cost decrease with high slop after each epoch. At the end of the training, cost converge and stabilized. While the error is high only at the first three epoch, stable at the following epochs.

## 5.4 Comparing with other models

TEAM	METHOD	TOTAL
DeepKnowledgeSeville[1]	CNN with 3 Spatial Transformers	<b>99.71%</b>
IDSIA[2]	Committee of CNNs	<b>99.46%</b>
COSFIRE[5]	Color-blob-based COSFIRE filters for object recogn	<b>99.87%</b>
INI-RTCV[9]	Human Performance	<b>99.84%</b>
Semanet[7]	Multi-Scale CNNs	<b>98.31%</b>
Our	Our Architecture	<b>97.23%</b>
CAOR[12]	Random Forests	<b>96.14%</b>
INI-RTCV [9]	LDA on HOG2	<b>95.68%</b>

Table 1: Our architecture result comparing of others using the same dataset

Comparing with other teams' techniques, our models surpass the feature based methods, while less perform than other CNN. The reason is that we lacked the processing steps like illumination.

## **6 Conclusion & Futher Development**

## References

- [1] Álvaro Arcos-García, Juan A. Álvarez-García, and Luis M. Soria-Morillo. “Deep neural network for traffic sign recognition systems: An analysis of spatial transformers and stochastic optimisation methods”. In: *Neural Networks* 99 (2018), pp. 158–165. ISSN: 0893-6080. DOI: <https://doi.org/10.1016/j.neunet.2018.01.005>. URL: <http://www.sciencedirect.com/science/article/pii/S0893608018300054>.
- [2] Dan Cireşan et al. “Multi-column deep neural network for traffic sign classification”. In: *Neural Networks* 32 (2012). Selected Papers from IJCNN 2011, pp. 333–338. ISSN: 0893-6080. DOI: <https://doi.org/10.1016/j.neunet.2012.02.023>. URL: <http://www.sciencedirect.com/science/article/pii/S0893608012000524>.
- [3] Weiguang (Gavin) Ding. *The figure is generated by adapting the code from*. URL: [https://github.com/gwding/draw\\_convnet](https://github.com/gwding/draw_convnet).
- [4] Eric. *Plotting confusion matrix*. URL: <https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix>.
- [5] Baris Gecer, George Azzopardi, and Nicolai Petkov. “Color-blob-based COS-FIRE filters for object recognition”. In: *Image and Vision Computing* 57 (2017), pp. 165–174. ISSN: 0262-8856. DOI: <https://doi.org/10.1016/j.imavis.2016.10.006>. URL: <http://www.sciencedirect.com/science/article/pii/S0262885616301895>.
- [6] Yann Lecun et al. “Gradient-based learning applied to document recognition”. In: *Proceedings of the IEEE*. 1998, pp. 2278–2324.
- [7] P. Sermanet and Y. LeCun. “Traffic sign recognition with multi-scale Convolutional Networks”. In: *The 2011 International Joint Conference on Neural Networks*. July 2011, pp. 2809–2813. DOI: 10.1109/IJCNN.2011.6033589.
- [8] K. Simonyan and A. Zisserman. “Very Deep Convolutional Networks for Large-Scale Image Recognition”. In: *CoRR* abs/1409.1556 (2014).
- [9] J. Stallkamp et al. “Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition”. In: *Neural Networks* 0 (2012), pp. -. ISSN: 0893-6080. DOI: 10.1016/j.neunet.2012.02.016. URL: <http://www.sciencedirect.com/science/article/pii/S0893608012000457>.
- [10] Christian Szegedy et al. “Rethinking the Inception Architecture for Computer Vision”. In: *CoRR* abs/1512.00567 (2015). arXiv: 1512.00567. URL: <http://arxiv.org/abs/1512.00567>.
- [11] Vivek Yadav. *Image Augmentation*. URL: <https://github.com/vxy10/ImageAugmentation>.
- [12] F. Zaklouta, B. Stanciulescu, and O. Hamdoun. “Traffic sign classification using K-d trees and Random Forests”. In: *The 2011 International Joint Con-*

*ference on Neural Networks*. July 2011, pp. 2151–2155. DOI: 10.1109/IJCNN.2011.6033494.



## 7 Appendix

### 7.1 References

1. Coursera; <https://www.coursera.org/specializations/deep-learning>
2. <https://www.tensorflow.org>
3. <http://www.numpy.org/>
4. <https://matplotlib.org/contents.html>