

Machine Learning Engineer Nanodegree

Capstone Project

Spandan Roy

July 14th, 2019

I. Definition

We recognize images and patterns at every moment in our life. When individuals detect an object, it seems like our brain doesn't take a lot of efforts to tell what we are looking at. However, it did take a long time for each person to learn what is what since we were born, not only because there are so many different categories in the world for us to learn, but also due to the fact that our intelligence was growing gradually to support our learning process.

Object recognition [1] is a technology in the field of computer vision for finding and identifying objects in an image or video sequence. There are many well-researched domains of object detection include face detection and pedestrian detection. Many of them have been successfully achieved using artificial neural network [2].



Figure1: Face detection. Retrieved from <https://towardsdatascience.com/face-detection-for-beginners-e58e8f21aad9>

Artificial neural networks are computing systems vaguely inspired by the biological neural networks that constitute animal brains. Such systems mimic the learning

process of living creatures, which do not need to be programmed with any task-specific rules. They "learn" to perform tasks from examples, e.g. learn to classify dog breeds with pictures of different dog breeds and their names (labels). Neural networks were deployed on a large scale, particularly in image and visual recognition problems, which are known as "deep learning [3]".



Figure 2: Traffic sign detection. Retrieved from <http://www.mdpi.com/1424-8220/17/6/1207>

The problem that this project is focusing on is one of the key tasks in autonomous vehicles areas -- Traffic Sign Recognition [4]. Traffic-sign recognition (TSR) is a technology by which a vehicle is able to recognize the traffic signs put on the road, such as "speed limit" or "turn ahead". Modern traffic-sign recognition systems are being developed using convolutional neural networks [5], which is a class of deep, feedforward artificial neural networks. The trained neural net can then be used in real life driving to detect traffic signs in real time. The predictions outputted by such systems will immediately be processed by other systems in self-driving car to make further decisions, such as whether it should perform a full-stop or drop its speed below 50 km/h.

Project Overview

[Traffic sign recognition](#) is one of the most important parts in self-driving car field. It enables vehicles to recognize modern traffic signs on the road during driving such as speed limit and stop sign. On November 8th, 1968. The Vienna Convention on Road Signs and Signals was held to increase road safety and aid international road traffic by standardising the signing system for road traffic, such as road signs, traffic lights and road markings (Wikipedia, 2018). The first TSR systems which recognized speed limits were developed appeared at the end of 2008, on the redesigned BMW 7 Series. Deep learning neural networks method is one of the major algorithms for pattern recognition problems currently. It is also widely used on Traffic signal recognition on self-driving cars.

This project is designed to find a solution to the traffic sign recognition problem in real life scenario. In specific, I will be mainly focusing on training artificial neural networks to classify traffic signs within the range of the given dataset.

Moreover, implemented solutions will be evaluated according to its efficiency, accuracy and utility. For example, how long does it take for the trained neural network to make one prediction, what is the accuracy score evaluated by defined metrics, how useful is the implementation in real life, etc. Adjustments to the methods or parameters may also be made to maximize the model strength.

Reference:

Academic paper - [Radu Timofte*, Markus Mathias*, Rodrigo Benenson, and Luc Van Gool, Traffic Sign Recognition - How far are we from the solution?, International Joint Conference on Neural Networks \(IJCNN 2013\), August 2013, Dallas, USA.](#)

Problem Statement

The problem to be solved is how to teach computers recognize traffic sign when it receives an image of a traffic sign. For instance, we want the machine to output “40 km/h speed limit” when receiving a given photo of such a traffic sign. Furthermore, we expect our solution still does a good job when the given traffic sign image isn’t in good quality, such as an image with a relatively small traffic sign, or an unclear picture.

The main method we are going to use is training a Convolutional Neural Network using traffic sign datasets. In general, a CNN takes an image, which may be seen as a bunch of pixels in different colors, each color is made of blue, green and red, extract features from these input images, update the weights of this feature to the target through training, to minimize its prediction error.

Let’s use stop sign as a simple example. A CNN can recognize the features, such as the “S”, “T”, “O” and “P”. Then, it serves as a classifier on top of these extracted features, by assigning a probability for the object on the image.

The main steps of the project can be divided into the following:

1. **Data exploration:** Discover and download traffic sign images in decent quality.
2. **Data preparation:** Reshape them into input data that can be feed into our neural network. Split them into training dataset, validation dataset and test dataset.
3. **Pre-trained model exploration:** We are going to use a pre-trained model as the “base” of our model. Thus, we need to look for such a model that fits our problem.
4. **Fine-tune [\[6\]](#) the chosen model:** Fine-tune the chosen pre-trained model.
5. **Training:** Train our model using training dataset from step 2
6. **Testing:** Test our trained CNN using testing dataset from step 2
7. **Improve:** Adjust parameters to improve test score
8. **Evaluate:** Evaluate its comprehensive performance and explore potential improvements

By the end of this project, we are expected to see a well-trained traffic sign classifier that is capable of recognizing the traffic sign in any given related pictures.

Metrics

Due to the fact that the dataset I chose is not evenly distributed, I chose weighted F1 score as the metric to this problem. F1 score is given by the following formula:

$$F1_score = 2 * \frac{precision * recall}{precision + recall}$$

To understand how to calculate precision and recall, we need this chart to understand what are True Positive, False Positive, False Negative and True Negative:

Actual Class	Predicted class	
	Class = Yes	Class = No
	Class = Yes	Class = No
	Class = Yes	Class = No
	True Positive	False Negative
	False Positive	True Negative

Figure 3: metric chart. Retrieved from <http://blog.exsilio.com/all/accuracy-precision-recall-f1-score-interpretation-of-performance-measures/>

The calculation of precision and recall are given by the following formula:

$$Precision = \frac{True\ Positive}{True\ Positive + False\ Positive}$$

$$Recall = \frac{True\ Positive}{True\ Positive + False\ Negative}$$

The final F1 score is given by $F1 =$

$$\frac{F1score_1 * weight_1 + F1score_2 * weight_2 + ... + F1score_n * weight_n}{\sum_{i=1}^n weight_i}$$

II. Analysis

Data Exploration

The dataset I used is a subset of the [Belgium Traffic Signs database](#), which is a multi-class, single-image classification challenge. There are 62 classes in total, and is a large, lifelike database. The images are taken in different angles and lighting conditions. There is one directory for each class. Each directory contains the corresponding images of the sign in .ppm format (RGB color), as well as a csv file that has annotations, but the number of images in each class is not same ,i.e., the classes are not balanced. There are [training dataset](#) and [testing dataset](#) already separated in advance. Each of the two directories contain 62 subdirectories, named sequentially from **00000** to **00061**. The directory names represent the labels, and the images inside each directory are samples of each label.

Note: the delimiter of the csv file is semicolon (;).

The csv file contains the following details,

Filename – Image file the following information applies to

Width, Height – Dimensions of the image

Roi.x1, Roi.y1, Roi.x2, Roi.y2 – Location of the sign within the image

ClassId – The class of the traffic sign. **Note: ClassID starts from 0.**

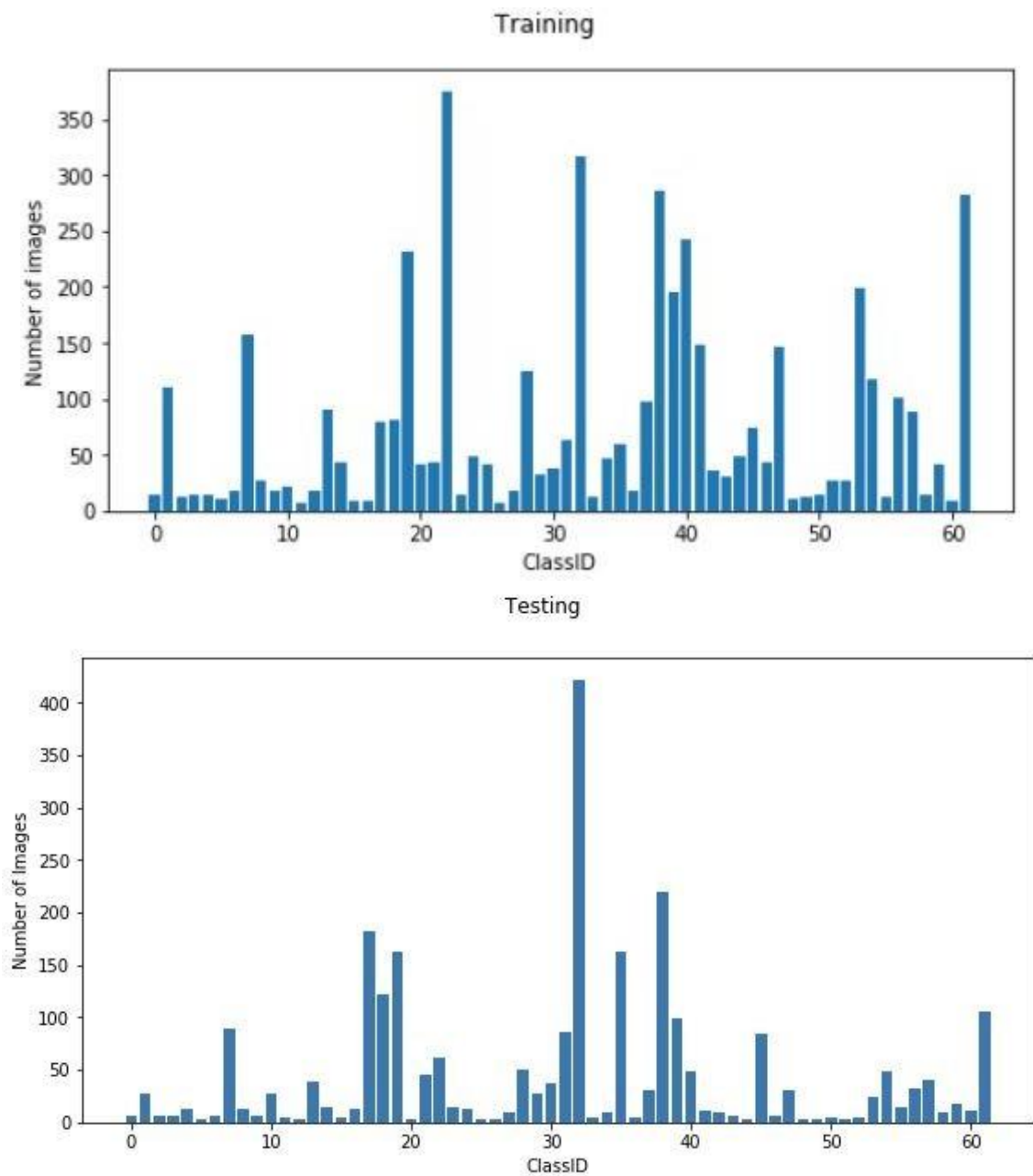
The input images are the keys to train our machine to learn how to recognize traffic signs. The csv file also helps us focus on the right spot on the image and classify different traffic signs. It should be clear how the dataset(s) or input(s) will be used in the project and whether their use is appropriate given the context of the problem.

Reference:

[Radu Timofte*, Markus Mathias*, Rodrigo Benenson, and Luc Van Gool, Traffic Sign Recognition - How far are we from the solution?, International Joint Conference on Neural Networks \(IJCNN 2013\), August 2013, Dallas, USA.](#)

Exploratory Visualization

Below are plots that shows the numbers of images of each class.



From these two plots, we realize that the data we got is not very evenly distributed. For instance, Class 22 has about 350 training images whereas class 0 only has about 5 images.

Statistics for datasets:

```
print('There are %d total traffic sign categories.' %  
len(set(all_training_targets)))  
  
print('There are %d total traffic sign images.' %  
len(traffic_sign_files_total))  
  
print('There are %d training traffic sign images.' % len(all_training_targets))  
print('There are %d validation traffic sign images.' %  
len(all_validation_targets))  
  
print('There are %d test traffic sign images.' % len(all_test_targets))
```

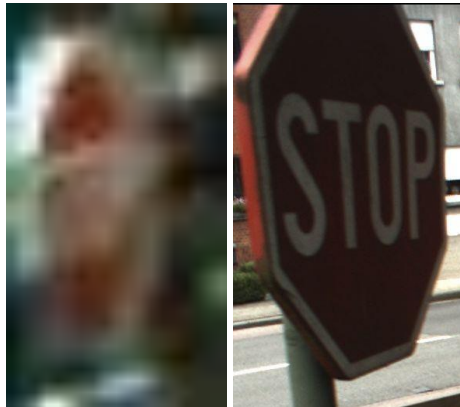
```
There are 62 total traffic sign categories.  
There are 7125 total traffic sign images.  
There are 4591 training traffic sign images.  
There are 1267 validation traffic sign images.  
There are 1267 test traffic sign images.
```

For display and future use, I convert all .ppm images into .png images and store them in another folder under the same structure.

The screenshot below shows several sample traffic sign images for class 00021.



An example of image in RGB channel:



Our dataset indeed has abnormal images such as dark and unclear images.

In my opinions, I don't think they need to be eliminated or fixed from the dataset. Our machine can still extract features (lines, curves...) from such dark or unclear images. It is expected to tell the traffic sign it's looking at under such conditions.

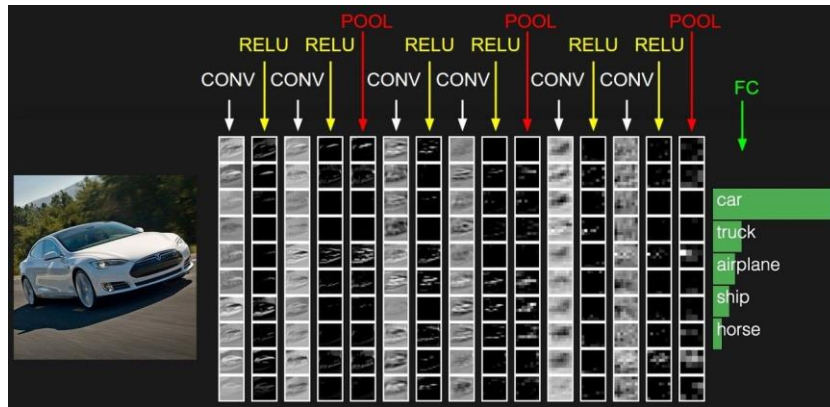
Therefore, pictures like these won't hurt if the proportion is relatively small.

During training, we need all pictures to be converted into the same ratio (224 x 224). This means some images will be stretched vertically or horizontally. I don't think this is an issue because the features in the picture won't get distorted or ruined as the whole picture doesn't change significantly.

Algorithms and Techniques

Convolutional Neural Network:

Convolutional Neural Networks is a class of deep, feedforward artificial neural networks. It is most commonly applied to analyzing visual imagery. In this section, I will be talking about how CNN works on image analysis in our problem, in a step-by-step way. First of all, I would like to show a picture of what an CNN looks like.



Before we dive into how CNN works, let's see what is CNN made of. CNN is made of layers of neurons that have learnable weights and biases to the extracted features.

Again, Artificial Neural Networks are like our brain. Weights and biases here, are the most important parts. Imagine when you see zebra, what is the most important factor that you consider to classify a zebra?

Black stripes, isn't it? Without those black stripes, I would definitely think I am looking at a horse. In this example, black stripe is a feature. How important you think it helps you to recognize it is a zebra, is the weight to this feature. CNN works exactly like that. By learning from examples, it extracts features and adjust weights of each feature to minimize its prediction error. Next, we are diving deeper to see how specific layers work.

Let's start from the very start -- the input file. We know our input file is a raw image of traffic sign. We see these as images, but computer sees them as a group of numbers. For example, for a 224×224 image, computer sees $224 \times 224 = 50,176$ pixels, each pixel is made of 3 color -- Red, Blue and Green.



Figure 5. Zebra. Retrieved from <https://en.wikipedia.org/wiki/Zebra>

The input layer is the first layer in CNN. This layer contains exactly what we discussed in the last paragraph. Since the neurons in a layer will be connected to a small region of the layer before it, the image data from this layer will be passed to the next layer through the connection.

CONV layer is usually the next layer to receive input data from input layer. It will compute the output of neurons that are connected to local regions in the input, each computing a dot product between their weights and the small region they are connected to in the input volume.

$$Y = \Sigma (\text{weight} * \text{input}) + \text{bias}$$

Basically, each neuron acts like a flashlight, which is for example, can light up 16x16 area of the 224x224 image. Note that the depth of neuron is the same as the depth of input, i.e. 3. Each neuron also has its own weights for the input. The output is simply one single number. This gif perfectly shows how convolutional layers work.

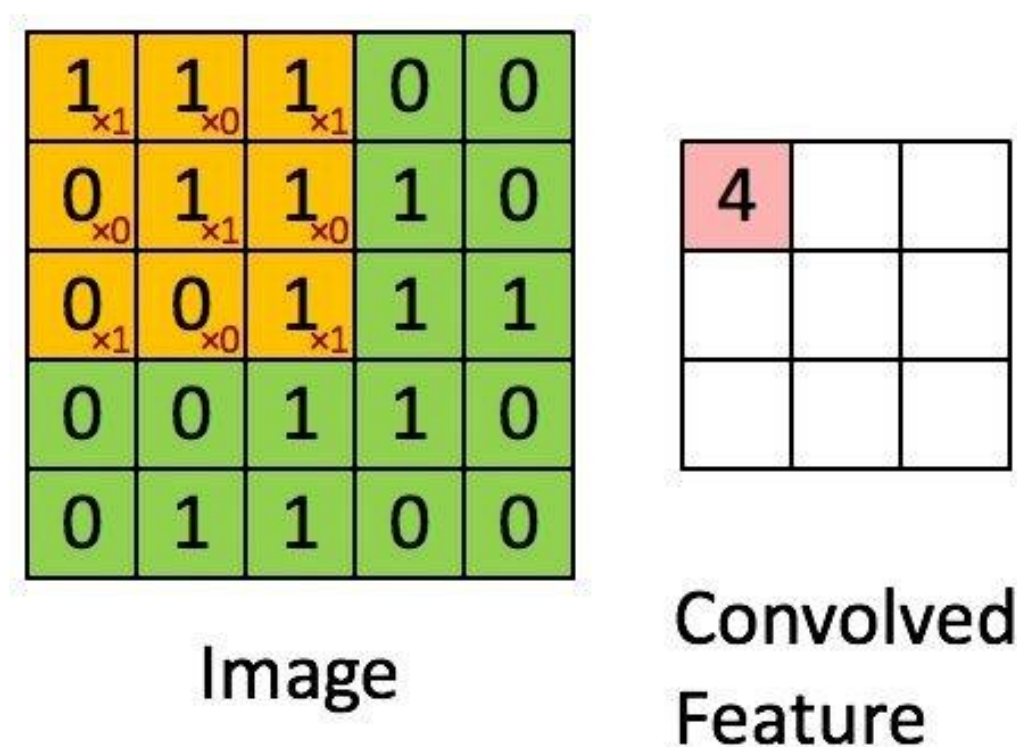


Figure 5. Convolutional Layer. Retrieved from <https://hackernoon.com/visualizing-parts-of-convolutional-neural-networks-using-keras-and-cats-5cc01b214e59>

As the flashlights (neurons) move 1 by 1 (or 2 by 2. This number is defined by the stride parameter), they are effectively checking for patterns in that section of the image. Thus, the task of convolutional layer is actually **filtering**.

You may be wondering, for all these values calculated by the neurons here, do they have a range? What if some number is extremely negative? The activation function in **activation layer** is responsible for squashing these values and give them a range.

There are a bunch of common activation functions. The most popular one is **ReLU (Rectified Linear Unit)**.

The reason why it is widely used in CNN is that it is easy to perform: If the number is negative: zero, else: the number. The simple operation makes it much faster to train networks.

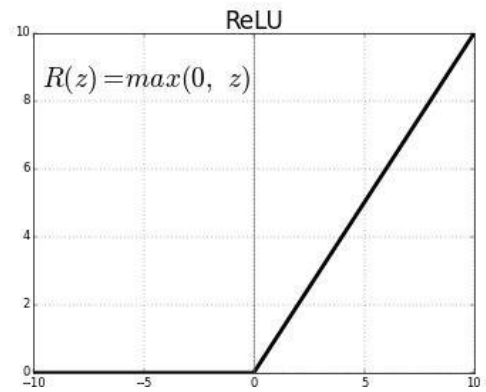


Figure 5. ReLU. Retrieved from <https://medium.com/@kanchansarkar/relu-not-a-differentiable-function-why-used-in-gradient-based-optimization-7fef3a4cecec>

Our next question is, how are we going to effectively reduce the spatial size of the visual representation, so that we can reduce the number of parameters and computation in the network. **Pooling layer** is the one who does the job.

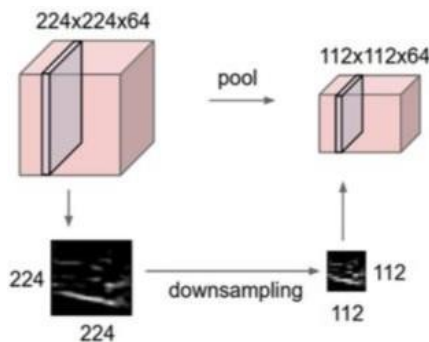


Figure 6. Pooling. Retrieved from <http://cs231n.github.io/convolutional-networks/>

Max pooling and **Average pooling** are the most common pooling functions. Max pooling takes the largest value from the window of the image, while average pooling takes the average of all values in the window. Both of them can effectively extract the important information from a image with a reduced space.

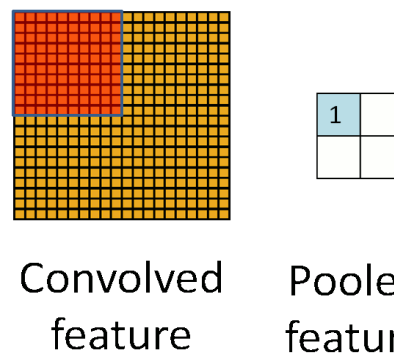


Figure 7. Pooling layer. Retrieved from <https://hackernoon.com/visualizing-parts-of-convolutional-neural-networks-using-keras-and-cats-5cc01b214e59>

Lastly, fully-connected layer. Neurons in a FC layer have full connections to all activations in the previous layer. It hence can compute the class scores, resulting in volume of size $[1 \times 1 \times \text{num_class}]$, where each number correspond to a class score.

So, this is basically how CNN works. A CNN can be any combination of the layers, from what we mentioned above and beyond. Design a CNN can be quite tricky: you need to add certain pooling layers to avoid overfitting, add proper activation layer to make sure values are in appropriate range, etc. Next, I will be explaining what is fine-tuning, and why it is a good option for this project.

Fine-tuning:

As I mentioned in the earlier section, I am going to fine-tune [6] a pre-trained deep learning model in Keras [7]. Fine tuning is a process to take a network model that has already been trained for a given task, train it on another similar task. The principle of fine-tuning is taking advantage of the feature extraction that happens in the front layers of the network without developing that feature extraction network from scratch.

The first thing we need to do, is to find a pre-trained model that fits our problem. We want to find such a model with pre-trained weights on a common dataset like the ImageNet. I have finally chosen [ResNet-50 Pre-trained Model for Keras](#), trained on ImageNet dataset.

In general, the technique of fine-tuning is replacing the output layer originally trained to recognize (in the case of ImageNet models) 1,000 classes, with a layer that recognizes the number of classes we require (in our case 62 classes). After attaching the new output layer, the model is then trained to take the lower level features from the front of the network and map them to the desired output classes, using SGD (Stochastic Gradient Descent [8]).

Just like training other Neural Networks, we need to define the following parameters before training:

```
# row number of pixels for input_shape
img_row = 224
# col number of pixels for input_shape
img_col = 224
# three for RGB channel = 3
# number of classes, in our case 62
```

```

num_class = len(set(all_training_targets))
# the number of training examples in one
forward/backward pass. Note: the higher the
batch size, the more memory needed
batch_size = 16
# One Epoch is when an entire dataset is passed forward
and backward through the neural network only once.
epochs = 10
# Learning rate:
Lr = 0.01
# Decay: Learning rate decay over each update.
Decay = 1e-6
# momentum: Parameter that accelerates SGD in the relevant
direction and dampens oscillations.

momentum=0.9

```

Benchmark

Initially, I am going to set a benchmark model that we are trying to beat. I have eventually chosen the ResNet50 pre-trained model on ImageNet weights. We only expect this model to recognize street sign, which is defined as class 919.

```

ResNet50_model = ResNet50(weights='imagenet')
def sign_detector(img_path):

    prediction = ResNet50_predict_labels(img_path)
    #print(prediction)

    return (prediction == 919) # 919 is street sign
#####

''' Test predictor performance '''
def test_sign_detector_performance(images_needed_check):
    # number of success of predicting there is traffic signs
    # in images that contain traffic sign

    positive_success = []
    all_success = [1]*len(images_needed_check)

```



```
for i in images_needed_check:
    if sign_detector(i):
        positive_success.append(1)
    else:
        positive_success.append(0)
print(f1_score(all_success, positive_success))
# training: 0.42355531095211885
# validation: 0.3816120906801007
# testing: 0.4027346177750155
```

It achieved F1 score of **42%** on the testing dataset, **38%** on the validation set and **40%** on testing dataset, which is good to serve as a benchmark.

I will also try to train the final model with testing dataset to see how it varies with the model trained with the training dataset.

III. Methodology

Data Preprocessing

In this section, I will show all the preparations I made before feeding training inputs into the model.

My first step is to convert every traffic sign picture as well as its information into a python object defined by this Traffic_Sign class.

```

class Traffic_Sign:
    """ A Single Traffic sign model """
    def __init__(self, filePath, fileName, width, height, x1, y1,
x2, y2, classTxt, classID):
        self.filePath = filePath
        self.fileName = filename

        self.width = width
        self.height =
height self.x1 =
x1 self.y1 = y1
        self.x2 = x2
        self.y2 = y2

        self.classTxt = classTxt
        self.classID = classID

```

To load dataset into a list of such objects, I iterated through every annotation csv files. After shuffling the list of Traffic_Sign objects, I then extract filePath and classID out as “training files” and “targets”.

```

# a list of training traffic_sign_collection
all_training_objs = load_dataset(training_dataset_directory)

# a list of testing traffic_sign_collection
all_testing_and_validation_objs =
load_dataset(testing_dataset_directory)

all_training_files, all_training_targets =
extract_X_and_Y(all_training_objs)

```

Even though our dataset has already been divided into training dataset and testing dataset, we still need to manually take validation dataset. Compare to the number of training images, we seem to have an abundant number of test files. Therefore, I split the test dataset into validation dataset and testing dataset.

```
# extract testing files and targets and split into validation and
test

X_temp, Y_temp = extract_X_and_Y(all_testing_and_validation_objs)

# merge X_temp and all_training_files to get total files
traffic_sign_files_total = X_temp + all_training_files
traffic_sign_targets_total = Y_temp + all_training_targets
# get all validation field and all test files

all_validation_files, all_test_files,
all_validation_targets, all_test_targets =
train_test_split(X_temp, Y_temp, test_size=0.5,
random_state=42)
```

Next, we need to define a function to convert the path to the image (string) into tensor to fit the input shape -- a 3D array in (224, 224, 3). Any pictures in different ratio will be converted into 224 x 224.

```
def path_to_tensor(img_path):
    # loads RGB image as PIL.Image.Image type
    img = image.load_img(img_path, target_size=(224, 224))
    # convert PIL.Image.Image type to 3D tensor with shape (224,
    224, 3)

    x = image.img_to_array(img)
    # convert 3D tensor to 4D tensor with shape (1, 224, 224, 3)
    and return 4D tensor

    return np.expand_dims(x, axis=0)

def paths_to_tensor(img_paths):
    list_of_tensors = [path_to_tensor(img_path) for img_path in
    tqdm(img_paths)]

    return np.vstack(list_of_tensors)
```

Implementation

Fine-Tuning ResNet50:

The following source code is from [ResNet50 implementation](#). Finding the right parameters for fine-tuning the pre-trained model was one of the most time-consuming tasks in the whole project.

```

x = ZeroPadding2D((3, 3))(img_input)

x = Conv2D(64, (7, 7), name="conv1", strides=(2, 2))(x)
x = BatchNormalization(axis=bn_axis, name='bn_conv1')(x)
x = Activation('relu')(x)

x = MaxPooling2D((3, 3), strides=(2, 2))(x)
x = conv_block(x, 3, [64, 64, 256], stage=2, block='a',
               strides=(1, 1))
x = identity_block(x, 3, [64, 64, 256], stage=2, block='b')
x = identity_block(x, 3, [64, 64, 256], stage=2, block='c')

x = conv_block(x, 3, [128, 128, 512], stage=3, block='a')
x = identity_block(x, 3, [128, 128, 512], stage=3, block='b')
x = identity_block(x, 3, [128, 128, 512], stage=3, block='c')
x = identity_block(x, 3, [128, 128, 512], stage=3, block='d')

x = conv_block(x, 3, [256, 256, 1024], stage=4, block='a')

x = identity_block(x, 3, [256, 256, 1024], stage=4, block='b')
x = identity_block(x, 3, [256, 256, 1024], stage=4, block='c')
x = identity_block(x, 3, [256, 256, 1024], stage=4, block='d')
x = identity_block(x, 3, [256, 256, 1024], stage=4, block='e')
x = identity_block(x, 3, [256, 256, 1024], stage=4, block='f')

x = conv_block(x, 3, [512, 512, 2048], stage=5, block='a')
x = identity_block(x, 3, [512, 512, 2048], stage=5, block='b')
    x = identity_block(x, 3, [512, 512, 2048], stage=5,
                       block='c')

    # original fully connected layer
x_orig = AveragePooling2D((7, 7), name='avg_pool')(x)
x_orig = Flatten()(x_orig)

    x_orig = Dense(1000, activation='softmax',
name='fcImageNet')(x_orig)

# Create model

```

```
model = Model(img_input, x_orig)
```

This fine-tune method doesn't use a "pop" method to remove the last layer. Instead, it connects to the final layer we want after loading the ImageNet weights.

```
# load image net weight
model.load_weights(imagenet_weights_path)

# replace softmax layer for transfer learning
x_new = AveragePooling2D((7, 7),
name='avg_pool')(x) x_new = Flatten()(x_new)

x_new = Dense(num_classes,
activation='softmax', name='fcTS')(x_new)

model = Model(img_input, x_new)
```

Train:

We pre-compiled the model having mentioned the optimizer and loss. Here, we proceed with the training of the model using fit() function and checkout the performance.

```
TS_model.fit(train_tensors, all_training_targets,
batch_size=batch_size, epochs=epochs, shuffle=True, verbose=1,
validation_data=(valid_tensors, all_validation_targets))
```

Metric:

```
TS_model.save("TS_resnet50.h5")
# Prediction
ts_prediction =
[np.argmax(TS_model.predict(np.expand_dims(tensor, axis=0)))
for tensor in test_tensors]

# Metric 1: F_1 score
test_accuracy_f1 = f1_score(np.argmax(all_test_targets,
axis=1), np.array(ts_prediction), average='weighted')
print("f_1 score: {}".format(str(test_accuracy_f1)))
# Output: f_1 score: 0.9940646687515156
```

We saved the model after training and predict the images of test dataset using predict() function and check the f1 score on the target dataset.

Refinement

Using the parameters, I defined in “Algorithms and Techniques” section, the model achieved a **90%** f1 score, which is pretty decent. The only refinement I made is changing number of epochs from **10** to **25**. The final f1 score on testing dataset achieved **99.4%**.

IV. Results

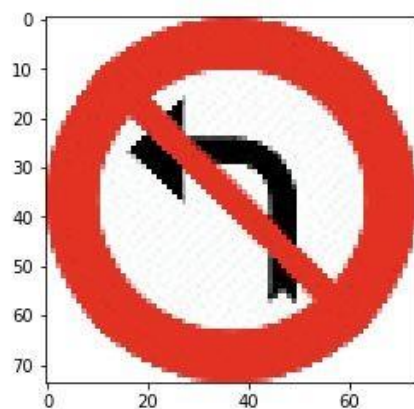
Model Evaluation and Validation

Below are some sample images downloaded from the Internet used for testing. I have classified them into different kinds of testing, so that we can see if the program is consistent:

Easy recognition test: traffic sign occupies most of the space of images

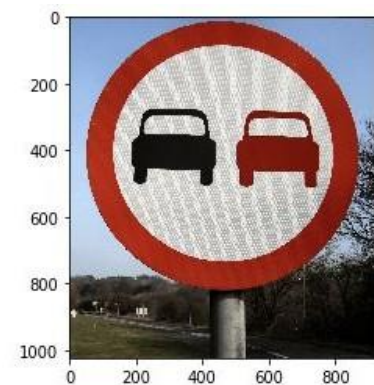
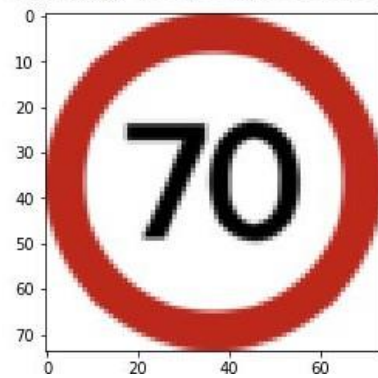


Confidence: 0.9993894. The classID for the traffic sign in this picture is: 29
Its class name in the database is: C31LEFT
Searching in the image database... It's refereneing to this one!

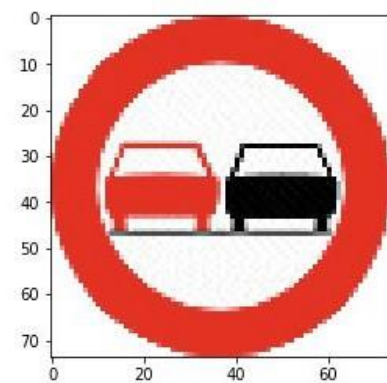


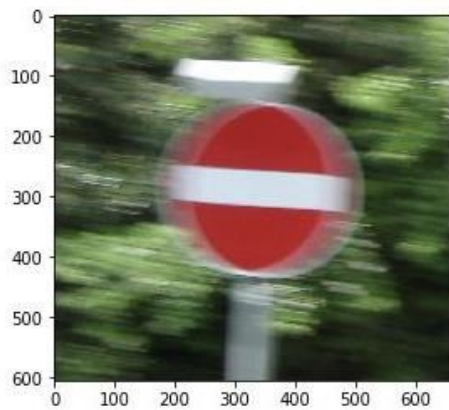


Confidence: 0.99999654. The classID for the traffic sign in this picture is: 32
Its class name in the database is: C43
Searching in the image database... It's referencing to this one!

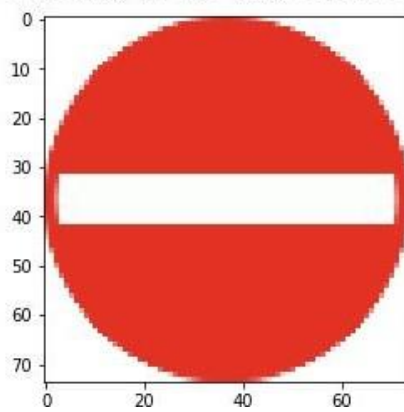


Confidence: 0.7758715. The classID for the traffic sign in this picture is: 31
Its class name in the database is: C35
Searching in the image database... It's referencing to this one!





Confidence: 0.9469988. The classID for the traffic sign in this picture is: 22
Its class name in the database is: C1
Searching in the image database... It's refereneing to this one!

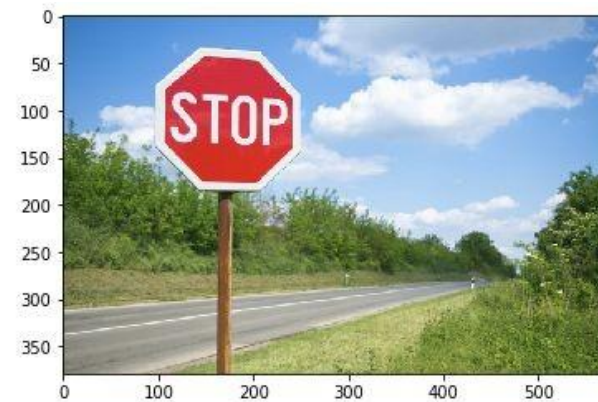
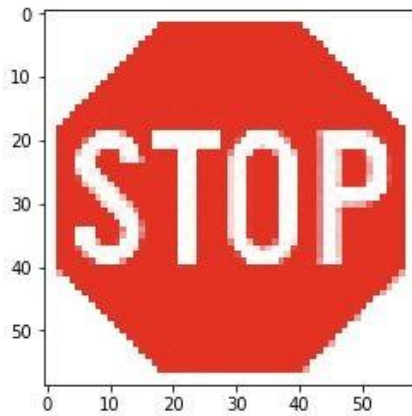


Our model is pretty solid when the target occupies more than half of a image. The confidence of prediction is generally greater than 90% for the samples I posted here as well as the tests I did locally. Some exceptions such as 3rd sample image have confidence lower than 80%. However, it's more because of the ambiguous color for the two cars in the picture. Overall, our model did a great job in this testing category.

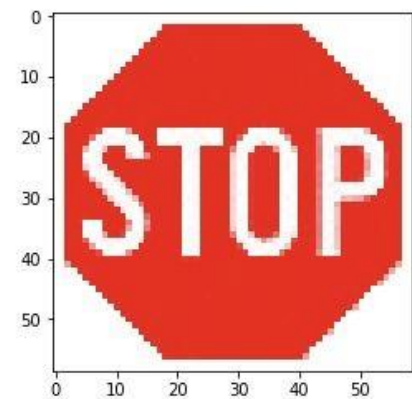
Small proportion test: traffic sign occupies less than half of the space of images



Confidence: 0.9919761. The classID for the traffic sign in this picture is: 21
Its class name in the database is: B5
Searching in the image database... It's referencing to this one!

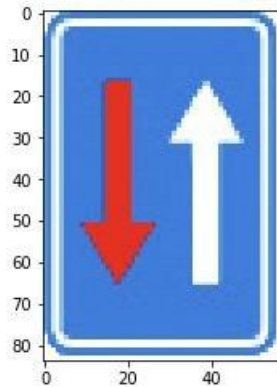


Confidence: 0.58831125. The classID for the traffic sign in this picture is: 21
Its class name in the database is: B5
Searching in the image database... It's referencing to this one!





Confidence: 0.16885376. The classID for the traffic sign in this picture is: 44
 Its class name in the database is: B21
 Searching in the image database... It's refereneing to this one!

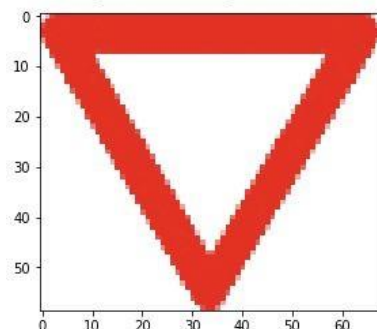


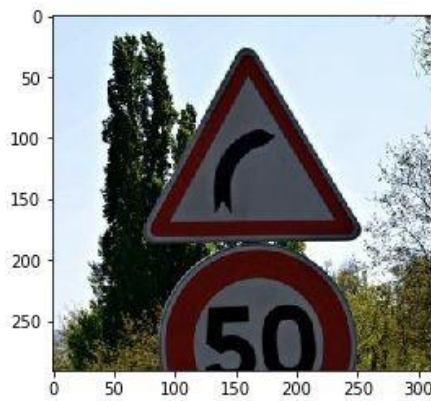
The model did a decent job most of the time, but with a lower confidence. Some wrong predictions can be found, such as the last sample. The main reason for this is that, the model takes the whole picture as input, and try to tell which traffic sign this picture looks like. For a model that have only “seen” cropped traffic sign images, tests from this section is a little bit hard for it, but still it performs not bad.

Multiple traffic signs test: more than one traffic signs in the image

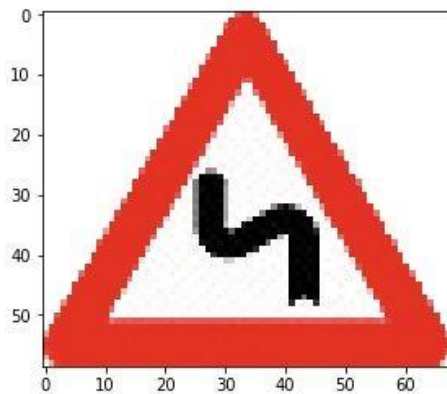


Confidence: 0.4284406. The classID for the traffic sign in this picture is: 19
 Its class name in the database is: B1
 Searching in the image database... It's refereneing to this one!





Confidence: 0.2042899. The classID for the traffic sign in this picture is: 5
 Its class name in the database is: A1C
 Searching in the image database... It's refereneing to this one!



Just like the previous test, the prediction comes with lower confidence and wrong predictions are made some time.

```
time_used = []
for eachfile in os.listdir(test_folder):
    start_time = time.time()

    if (eachfile[0]!='.'):
        predict_img(os.path.join(test_folder, eachfile))

    elapsed_time = time.time() -
    start_time
    time_used.append(elapsed_time)

print("Average time spend for 1 prediction:
{}".format(str(np.mean(time_used))))

# Output:
# Average time spend for 1 prediction: 0.5287921245281513
```

In terms of efficiency, the model takes an average of **0.53 second** to make one prediction, which is a reasonable time.

Justification

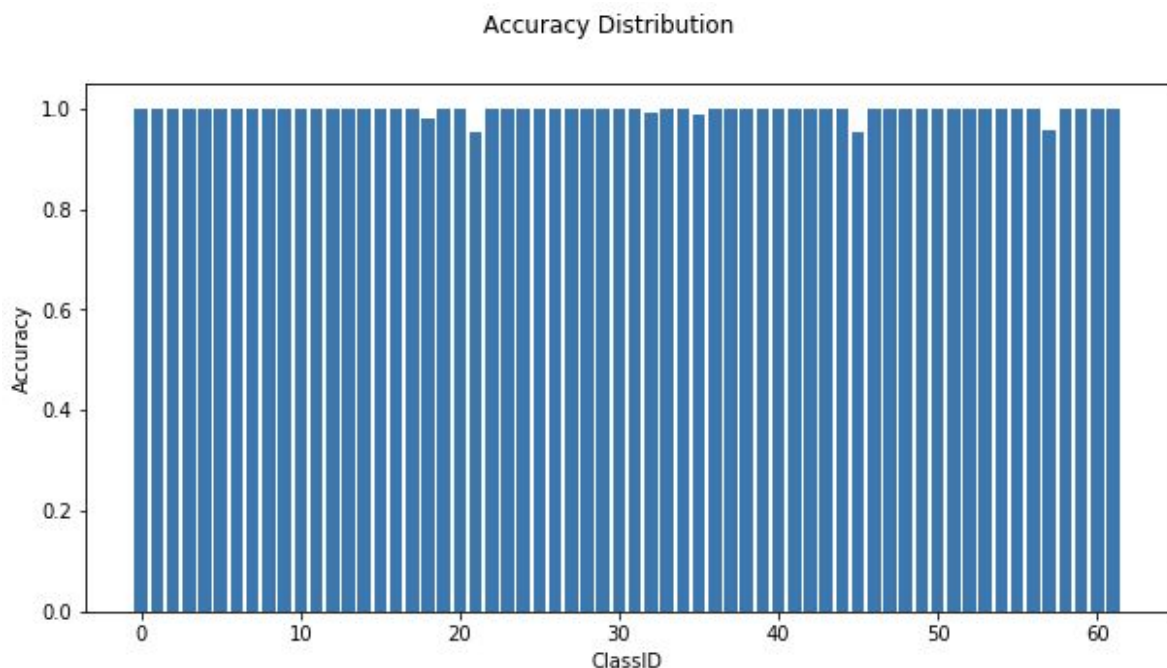
Compared to the benchmark we defined previously (40% f1 score on testing dataset), our final model performs way better (99.4%).

The final model made expected prediction most of the time, especially when I tested using images that have traffic sign part taking most of the spaces. However, low confidence and wrong predictions were made when traffic sign are relatively smaller in the image, or there are multiple traffic signs in one image (the reasons are actually the same). Thus, our model is robust to handle the job most of the cases when traffic sign occupies most of the image frame, but it is not reliable enough to be used in real time situation.

V. Conclusion

Free-Form Visualization

I have made a distribution chart for all 62 classes in the testing dataset, where each bar shows the percentage of the ratio of correctly predicted images out of the total images in each class. Here, the accuracy does not refer to any prediction accuracy score metric, it only serves the above purpose.



We can see that the model did a wonderful job on most of the classes as most of the images in each class were predicted correctly on the test dataset. Hence, the model is suitable for decent recognition of traffic sign, though it needs lot of improvements to be perfect.

Reflection

This capstone project not only helps me review my machine learning knowledges, but also offers me a great chance to gain hands-on experiences on a formal machine learning project.

Through this project, I found many interesting parts that I really enjoy. For example, manipulating dataset to fit in model, fine-tuning pre-trained model, analyzing and plotting data to see the statistics, etc.

The difficult part I found in this project, is searching the most correct and efficient way of doing what I want to do. For example, I took a long time to figure out what is the best way to fine-tune a ResNet50 model. It is not hard to make it work, but it takes time to figure out why it is implemented that way.

The final model and solution fit my expectations. However, it definitely needs improvements. I am confident to say I build a successful traffic sign classifier, but it needs more work to become a real-time traffic sign detector and classifier. This will be the next step I am working on.

Improvement

As I discussed in the previous section, the most important improvement is real-time detection support. Through my searches, it is difficult to make such an improvement on the model that we already trained. There are some mature object detection implementations such as object detection API of TensorFlow, YOLO etc. Therefore, this improvement may require changing the deep learning network we had.

Another improvement I could have made, is increasing the number of training dataset of classes that have less images. From the statistics listed in the report, I realized that the dataset is unbalanced. I could either replace our dataset with a better comprehensive dataset, or manually edit and add images to the dataset.

Reference

- [1]: https://en.wikipedia.org/wiki/Outline_of_object_recognition
- [2]: https://en.wikipedia.org/wiki/Artificial_neural_network
- [3]: https://en.wikipedia.org/wiki/Deep_learning
- [4]: https://en.wikipedia.org/wiki/Traffic-sign_recognition
- [5]: https://en.wikipedia.org/wiki/Convolutional_neural_network
- [6]: http://wiki.fast.ai/index.php/Fine_tuning
- [7]: <https://en.wikipedia.org/wiki/Keras>
- [8]: https://en.wikipedia.org/wiki/Stochastic_gradient_descent