

Object Classification in Road Images



UNIVERSITY *of* LIMERICK

OLLSCOIL LUIMNIGH

Final Year Project 2019

Author: Rory Egan

ID: 15178668

Bachelor of Science in Computer Systems (LM051)

Supervised by: J.J. Collins

Abstract

This project will be a research-oriented project centred around the use of Convolutional Neural Networks for object detection in images from the perspective of a self-driving car. The concept of self-driving cars is a rapidly expanding field, with significant investments having taken place over the past several years in this area. Many companies have chosen to use cameras as one of the primary sensors deployed on their prototype cars in conjunction with other sensors such as LIDAR, radar etc. to allow a vehicle to understand its environment and orient itself appropriately. This project will focus on the way these cameras can allow the vehicle to classify different objects in its environment. An important challenge within the field of machine vision in self-driving cars is reliably classifying multiple objects at once within an image. The current standard Machine Learning class of algorithms for object classification are CNNs (Convolutional Neural Networks). This project will use CNN's as well as some Computer Vision techniques to analyse colourised 2D images that depict typical scenes presented to a self-driving car and classify the different objects within them. The initial goal is to research and outline the basics of how CNN's operate, then choose a particular CNN paradigm and train a classifier using this paradigm. The classifier should be able to reliably classify major obstacles that are important to the function of a self-driving car such as other vehicles, pedestrians etc. The classifier should aim to be robust to adverse conditions creating noise within the images such as rain, snow etc. This project should then investigate the effects of parameter tuning and the consequences this reflects in the accuracy of the model.

Contents

| | | |
|----------|---|-----------|
| 1 | Overview | 3 |
| 1.1 | Introduction to Problem Area | 3 |
| 1.2 | Objectives | 4 |
| 1.2.1 | Primary Objectives | 4 |
| 1.2.2 | Secondary Objectives | 5 |
| 1.3 | Contribution | 5 |
| 1.4 | Methodologies | 5 |
| 1.5 | Project Plan | 8 |
| 1.5.1 | Table | 8 |
| 1.6 | Motivations | 9 |
| 2 | Background Research | 10 |
| 2.1 | Introduction to Machine Learning | 10 |
| 2.1.1 | Supervised Learning | 10 |
| 2.1.2 | Unsupervised Learning | 11 |
| 2.1.3 | History of Neural Computing | 11 |
| 2.2 | Introduction to Neural Networks and Deep Learning | 12 |
| 2.2.1 | Input Nodes | 14 |
| 2.2.2 | Hidden Nodes | 14 |
| 2.2.3 | Output Nodes | 15 |
| 2.2.4 | XOR Problem | 15 |
| 2.2.5 | Multilayer Perceptron | 15 |
| 2.2.6 | Gradient Descent and Backpropagation | 16 |
| 2.2.7 | Overfitting | 17 |
| 2.3 | Introduction to Machine Vision | 18 |
| 2.3.1 | Edge Detection | 19 |
| 2.4 | Introduction to Convolutional Neural Networks | 22 |
| 2.4.1 | Convolution Layer | 23 |

| | | |
|----------|---|-----------|
| 2.4.2 | Pooling Layer | 25 |
| 2.4.3 | Relevant Architectures | 26 |
| 2.4.4 | Evaluating the Classifier | 27 |
| 2.4.5 | Transfer Learning | 29 |
| 2.4.6 | Berkeley Deep Drive | 30 |
| 2.5 | Tensorflow Introduction | 32 |
| 2.5.1 | MNIST Experiment | 32 |
| 2.5.2 | CIFAR-10 Experiment | 34 |
| 2.5.3 | Initial Findings | 35 |
| 2.5.4 | Experiment Environments | 35 |
| 3 | Empirical Studies | 45 |
| 3.1 | SSD MobileNet V1 Experiments | 45 |
| 3.1.1 | Full Dataset Retraining Experiment | 46 |
| 3.1.2 | Partial Dataset Retraining Experiment 1 | 47 |
| 3.1.3 | Partial Dataset Retraining Experiment 2 | 50 |
| 3.2 | Faster RCNN Inception V2 Experiments | 54 |
| 3.2.1 | Initial Experiment | 55 |
| 3.2.2 | No Pretraining Experiment | 56 |
| 3.3 | Empirical Studies Results | 59 |
| 4 | Application Implementation | 61 |
| 4.1 | AWS Instance | 61 |
| 4.2 | Implementation | 61 |
| 4.3 | Code Snippets | 62 |
| 5 | Final Conclusion and Discussion | 64 |
| 5.1 | Summary | 64 |
| 5.2 | Reflections | 65 |
| 5.3 | Future Work | 65 |

Chapter 1

Overview

1.1 Introduction to Problem Area

This project deals with the area of classifying multiple objects within an image, specifically from the perspective of images presented to a self-driving car. A typical use case for a Machine Learning algorithm within the field of Computer Vision is image classification, whereby an image is presented to an algorithm and the algorithm decides which class that the image belongs to from a set of classes. For example, if an algorithm that is trained to tell the difference between dogs and cats is provided an image of a dog or a cat it should return a label for the image of either dog or cat.

A more complex problem however is the concept of object detection. If in the previous example an image that contains both a cat and a dog is supplied to the algorithm, it will typically label the image with the label that corresponds to the most prominent object within the image. A Machine Learning algorithm for object detection must be able to take in this image and correctly label both the dog and the cat within the image.

This class of Machine Learning algorithm is particularly useful in the field of self-driving cars. A self-driving car will typically be provided images which contain multiple objects, all of which must be correctly classified in order for the vehicle to function safely.

The standard type of Machine Learning algorithm that is utilised for image based problems is the Neural Network, a class of Machine Learning

algorithm modelled off the neural pathways of the brain. Within this class of algorithm, a class known as Convolutional Neural Networks (CNN's) has become increasingly popular. Both Neural Networks and CNN's will be explained further in this paper.

The overall aim of this project is to leverage CNN's to train an object classifier on a dataset extracted from cameras mounted on cars that consists of images containing multiple different objects such as pedestrians, vehicles etc.

1.2 Objectives

1.2.1 Primary Objectives

Utilise a Convolutional Neural Network for Object Classification

The core concept of this project is the use of CNN's for the detection of certain objects present within driving-based images. As such, the primary objective for this project is to research and implement a CNN that will return an acceptable level of accuracy on the testing data. The Berkeley Deep Drive dataset utilised by this project has a related paper in which object detection experiments are run on the dataset and performance metrics from these experiments are recorded (Yu et al., 2018). A major objective for this project therefore is to achieve similar results to those achieved in this paper. As an understanding of how CNN's are evaluated is required in order to properly expand upon the results obtained in the paper, the results will be discussed further in a later section.

Understanding of Convolutional Neural Networks

As this is a primarily research based project, one of the main objectives is to gain a deep understanding of how CNN's operate. The following sections of this report should demonstrate a depth of knowledge around the theory of how CNN's work, as well as demonstrate an ability to put these theories into practice. The section on Empirical Studies in particular should demonstrate an ability to forensically examine the results produced from CNN's and provide insight into how these results are achieved.

1.2.2 Secondary Objectives

Provide a Visual Interaction With The Trained Convolutional Neural Network

Once CNN's have been trained to acceptable accuracy levels during the Empirical Studies section, a basic application should be implemented in order for users to observe the CNN in action upon testing data. To this end, a Flask application is to be developed that will allow users to upload images upon which object detection will take place. As the Berkeley Deep Drive dataset contains videos as well as images, object detection should take place on a small number of test videos in order to provide users an example of how a CNN operates in real-time.

1.3 Contribution

Although results of applying a CNN to the Berkeley Deep Drive dataset have been published in (Yu et al., 2018), the authors of this study only used one particular architecture of CNN to obtain their results. This report outlines results obtained using two different types of CNN architecture across a range of different types of experiment. As the Berkeley Deep Drive dataset is a newly released dataset from 2018, a large body of research has yet to be carried out on the dataset. As such, this report may be of some small use to aid in avoiding some of the problems encountered during this project.

1.4 Methodologies

1. Define the Research Area: The first step in this project was to define which problem domain to base this FYP around. As an avid SCUBA diver, the initial area I investigated was the field of fish classification in images and videos collected from Irish waters. This area appealed to me as I wanted to investigate the possibility of collecting the dataset myself using underwater camera equipment. However I quickly abandoned this area for several reasons:
 - (a) Equipment: Underwater camera equipment that is capable of taking high resolution images and videos is extremely expensive.

Due to the strenuous nature of diving in rough Irish waters, equipment with excellent stabilisation software would be required in order to reliably obtain non-blurred images.

- (b) Domain Knowledge: As I am not an expert in the field of fish identification, I would be unable to reliably annotate species present in the training images. This would require enlisting the help of a third-party expert in the field of fish species.
- (c) Time Constraints: There are no datasets currently available that contain annotated fish species in underwater images or videos in Northern Atlantic waters. As such, the full dataset would need to be collated and annotated myself. With time constraints present for this project, it would not be feasible to carry out all this work in time to reach project deadlines.

With these issues present for my initial choice of problem domain, I decided to instead focus on the area of road images as there have been several large-scale annotated datasets released in recent years that I could leverage. Moreover, being due to start at Jaguar Land Rover as a graduate I decided that this project could serve as an excellent introduction into the problems faced by automotive companies in the field of Machine Vision.

2. Background Research: The next step was to carry out a literature review around the topic of Neural Computing, specifically CNN's. Papers dealing specifically with the challenges presented during object detection were the main focus of the literature review. In an attempt to demonstrate an understanding of how CNN's operate, the literature review has been summarised in the following sections of this report.
3. Gain a hands-on knowledge of tensorflow: As tensorflow was the tool utilised to implement CNN's during this project, a Udemy tutorial was first followed in an attempt to gain an understanding of using tensorflow in a hands-on manner (*Complete Guide to Tensorflow for Deep Learning with Python* 2018).
4. Carry out Empirical Studies: Empirical Studies were a major point of focus for this project. Experimentation was carried out on the Berkeley Deep Drive dataset using a range of different types of CNN.

Results from these experiments were documented and investigated in order to gain an understanding of what techniques work best for the dataset. The results from the experiments were then compared to results carried out in (Yu et al., 2018).

5. Build an application: Once Empirical Studies have been carried out and CNN's have been trained on the Berkeley Deep Drive dataset to a satisfactory degree, a prototype application should be built to allow users to observe the Object Detection taking place.

1.5 Project Plan

1.5.1 Table

| Chapter | Section | Est. Completion Date |
|---------------------------|--------------------------|----------------------|
| Introduction | Intro to Problem Area | 12/10/2018 - wk5 |
| | Motivations | 12/10/2018 - wk5 |
| | Objectives | 05/10/2018 - wk4 |
| | Contribution | End? |
| | Methodologies | 12/10/2018 - wk5 |
| | Project Plan | 05/10/2018 - wk4 |
| Background Research | Intro to ML | 19/10/2018 - wk6 |
| | Intro to Neural Networks | 26/10/2018 - wk7 |
| | Intro to CNN's | 26/10/2018 - wk7 |
| | Intro to Computer Vision | 02/11/2018 - wk8 |
| | Technology Investigation | 19/10/2018 - wk6 |
| Empirical Studies | CNN Architectures | 16/11/2018 - wk10 |
| | Image Processing | 30/11/2018 - wk11 |
| Classifier Implementation | Requirements | 14/12/2018 - wk13 |
| | System Design | |
| | Implementation | |
| | Analysis | |
| Final Conclusion | Review of Research | |
| | Reflections | |

1.6 Motivations

The main motivation behind this project for me is working within the field of self-driving cars. I find this particular area fascinating due to the broad range of technologies present within these cars. I have previously worked a summer internship at Jaguar Land Rover in Shannon working within the ADAS (driver assistance) team. I was exposed to a broad range of different technologies, from Computer Vision and Machine Learning based teams to Big Data pipelines concerned with offboarding data from test cars. This internship and the potential to work with the company further down the line has really motivated me to further my study within this field. I have found that these technologies are much more interesting to me than many more conventional potential areas of work. With ever increasing amounts of automotive manufacturers investigating this field, I feel like a final year project focusing on self-driving cars could be very beneficial to my career going forward.

I am also very interested in self-driving cars due to the far-reaching safety implications of the technology. Hundreds of people die on Irish roads every year, with driver error being the primary cause of fatalities. Self-driving cars have the potential to save countless lives in the future, and I find it very motivating to work within a field that has the potential to produce very tangible positive changes in the daily lives of many people.

Chapter 2

Background Research

2.1 Introduction to Machine Learning

A general definition of Machine Learning: “[Machine Learning is the] field of study that gives computers the ability to learn without being explicitly programmed“ - Arthur Samuel, 1959.

One of the goals for Machine Learning algorithms is automatically observing structures in data and fitting these structures to a model in order to allow people to interact with the data in a way that is humanly intuitive. Development within the field has progressed to the point where people interact with Machine Learning algorithms multiple times in their daily lives without noticing. Areas such as email spam detection, Facebook image tagging suggestions and voice-to-text are all examples of a broad range of Machine Learning algorithms that are commonly used.

Machine Learning can be roughly divided into two categories - Supervised and Unsupervised.

2.1.1 Supervised Learning

In Supervised Learning, algorithms are provided with some labelled input data which they attempt to learn patterns from. The algorithms will then attempt apply this learned experience to new unseen data and attempt to create their own labels for the data, with varying degrees of success.

Prominent examples of the Supervised category of Machine Learning include Support Vector Machines, Linear Regression and the focus of this project, Convolutional Neural Networks (O’Shea and Nash, 2015).

2.1.2 Unsupervised Learning

Unsupervised Learning algorithms differ from Supervised Learning algorithms in that they are given data with no labels. They must then attempt to find some structure in this input data themselves with no given direction or explicit programming. Popular Unsupervised Learning examples include clustering algorithms.

2.1.3 History of Neural Computing

Neural Networks are not a new concept - they were first proposed in 1943 by neurophysicists in the form of a primitive electrical circuit. The concept was studied up until the 1960's until it fell out of favour with researchers. Bold claims had been made by many researchers about the vast potential of the field, however a failure to back up these claims led to widespread skepticism about the true potential of Neural Computing. The area was in part hampered by the technology of the time - processing power available to researchers was very low. Additionally, in 1969 a paper entitled "Perceptrons: An Introduction to Computational Geometry" was published introducing the XOR problem, which will be explained later. The paper stated that the research being carried out at the time on Neural Networks was fundamentally flawed and that the field would not experience any major successes. A revival was seen in the 1980's, when the concept of using multiple layers of neurons to create a network began to emerge, and the issues presented by the XOR problem were solved. From 1989 to 1994 Yann LeCun developed the LeNet architecture, one of the first examples of a CNN, which was used to recognise handwritten postal addresses. In around the 2010 the field of Deep Learning (explained below) experienced a surge in popularity, primarily due to the increase in processing power made available through GPU's. In 2012, Alex Krizhevsky won the ImageNet competition, a popular image recognition competition. The architecture he used, known as AlexNet, achieved an error rate of 15%, which was approximately 10% better than the closest competition at the time. His architecture essentially scaled up the LeNet architecture into a larger, more complex network (Krizhevsky, Sutskever, and Hinton, 2012). This architecture caused a widespread adoption of CNN's within image based Machine Learning, and CNN's are now the de facto standard for many image processing tasks(O'Shea and Nash, 2015).

2.2 Introduction to Neural Networks and Deep Learning

Artificial Neural Networks (ANN's) are a particular Machine Learning paradigm modelled after the neural pathways present in the brain. The building blocks for every Neural Network is the perceptron or node, an artificial neuron. It receives inputs with weights associated with them that show the importance of each input relative to the others. The node applies a particular activation function to the weighted sums of the input, and an output is generated, as shown in Figure 2.1 (Géron, 2017). The Perceptron Traini

A single Perceptron is only able to classify linearly separable data (Kotsiantis, Zaharakis, and Pintelas, 2007). In order to classify data that is not linearly separable, techniques such as Multilayer Perceptrons must be used. These will be explained in a following section. Linear separability refers to the ability for a single line to separate all members of a given class A from all members of a given class B. This is illustrated in Figure 2.3.

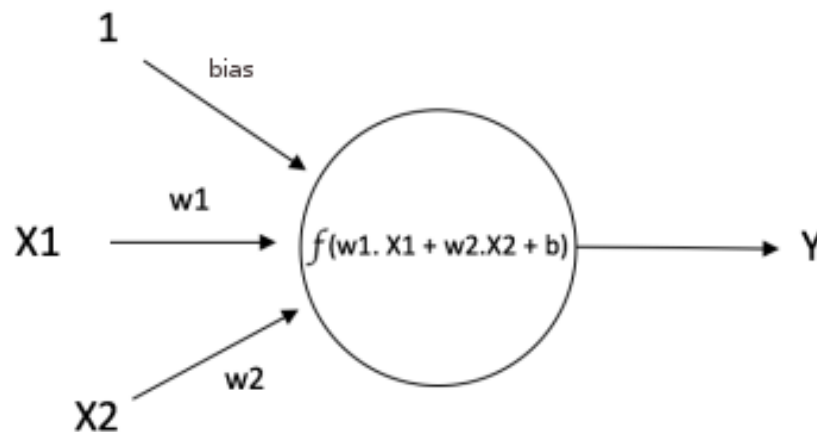


Figure 2.1: An Artificial Neuron

There are several different types of activation function, however every activation function takes in a single number as input and performs a certain mathematical operation on the number. The three most common activation functions generally encountered are ReLU, tanh and sigmoid. Rectified Linear Unit, or ReLU, takes in an input and replaces negative numbers with zero. The tanh function squashes the input to between the range between -1 and 1. Finally the sigmoid activation function takes the input and squashes it to the range between 0 and 1. In the field of CNN's ReLU is commonly used, as training times are significantly better when using this activation function. (Krizhevsky, Sutskever, and Hinton, 2012). These different types of activation function are visualised in Figure 2.2.

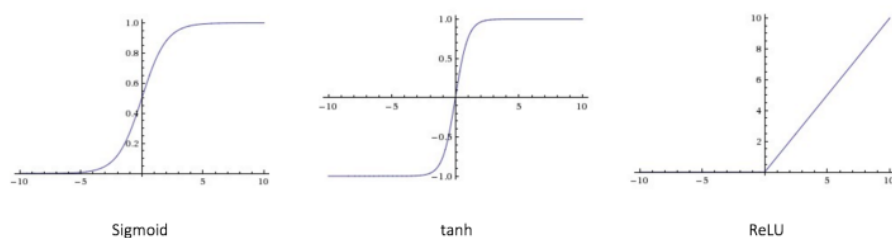


Figure 2.2: Activation Functions Visualised

A Neural Network consists of a series of interconnected layers of these artificial neurons, with the output of each layer of neurons serving as input for the next layer of neurons. The simplest and most common type of Neural Network is the feedforward neural network. It consists of multiple layers of neurons, with connections to all of the neurons in the preceding layer. Each connection or edge has a weight associated with it. There are three types of nodes - input nodes, hidden nodes and output nodes. As the name suggests, data in a feedforward neural network only moves forward through the network - into the input layer, through the hidden layers and then on to the output layer.

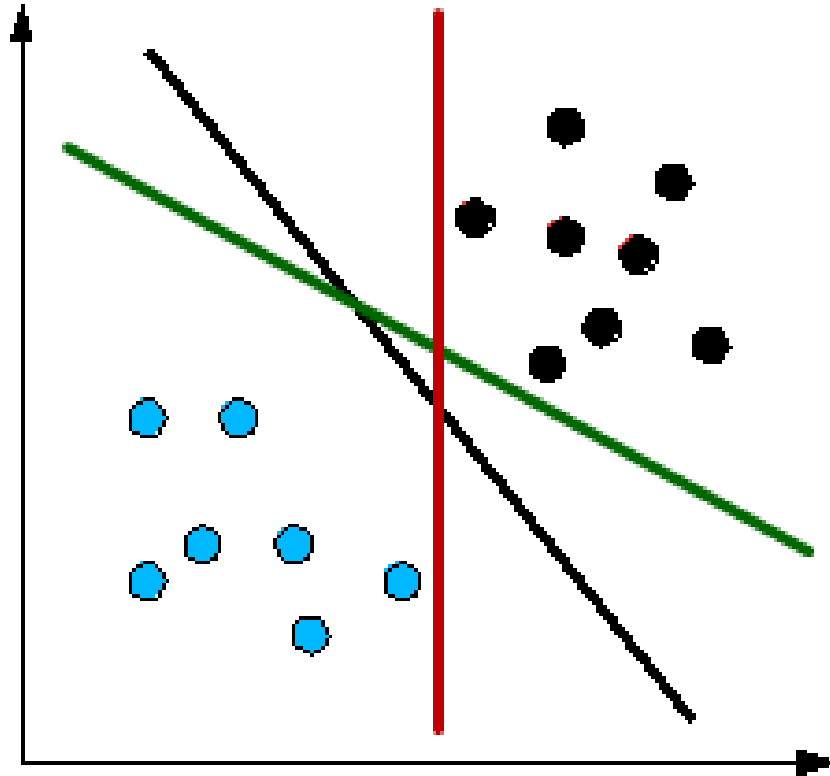


Figure 2.3: Linearly Separable Data

2.2.1 Input Nodes

Input nodes take in the input data fed into the network without performing any computation and pass the data on to the hidden nodes. These nodes make up the input layer of a neural network.

2.2.2 Hidden Nodes

The hidden nodes perform the actual computation, taking data from the input nodes and providing information to the output nodes. There can be multiple or zero hidden layers in a network, whereas there will only ever be one input and output layer.

2.2.3 Output Nodes

Like the input nodes, the output nodes do not perform any computation on the data - they simply take information from the hidden layers and expose this to the outside. This will generally consist of a prediction.

2.2.4 XOR Problem

In the paper mentioned above, "Perceptrons: An Introduction to Computational Geometry", the Xor problem was first introduced. Xor is a function that given two binary inputs returns 0 if the inputs are equal and 1 if they are not. Xor is a classification problem with known expected results, therefore it is appropriate to utilise a Supervised Learning approach to solving it. However, the Xor problem is an example of a problem that is not linearly separable. As mentioned previously, a single perceptron is not capable of predicting data that is not linearly separable, therefore a single layered architecture of perceptrons is simply not capable of solving this problem. The only way for a Neural Network to solve this problem is to expand the number of layers in the architecture, adding a hidden layer. This type of architecture is known as the Multilayer Perceptron.

2.2.5 Multilayer Perceptron

A Multilayer Perceptron (MLP) consists of an input layer, an output layer and one or more hidden layers, as seen in 2.4. MLP's are feedforward neural networks. The MLP architecture solves the Xor problem by introducing linear separability (Singh and Pandey, 2016).

Deep Learning is when there is more than one hidden layer present in a network (O'Shea and Nash, 2015). MLP's encounter issues when attempting to carry out classification tasks on large input images. As will be explained later, each pixel of an input image must correspond to an input node in the network. As every node is connected to every other node in its preceding layer, large images quickly require a large amount of nodes in the network. A small image of size 28x28 such as the images found in the MNIST dataset will require a network with 784 input nodes. Images sourced from the Berkeley Deep Drive dataset are of size 1280x720, which would require an input layer consisting of 928080 nodes. This will quickly cause issues due to overfitting and network training time.

Learning takes place in a MLP through changing the connection weights for each perceptron in the network based on the error between the expected and true output of each perceptron. This is carried out through a concept known as backpropagation.

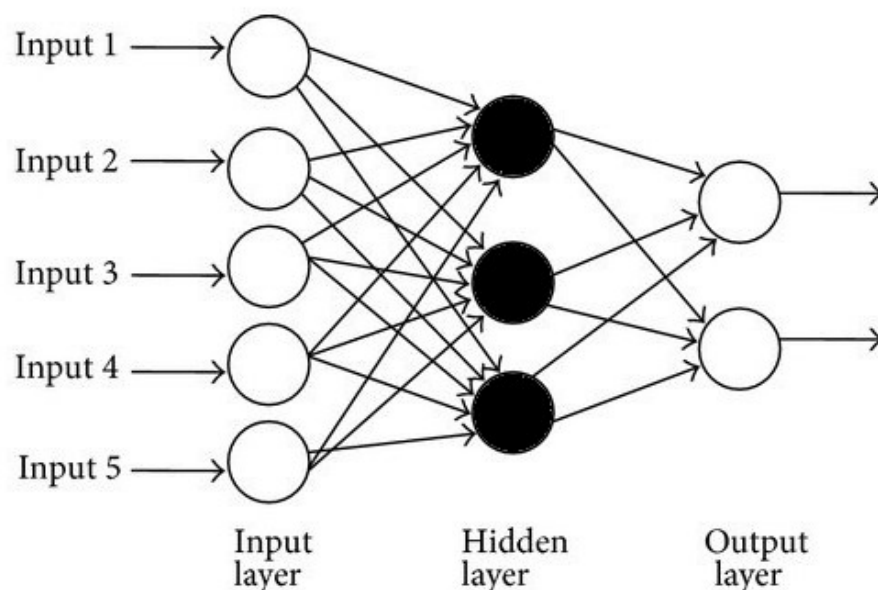


Figure 2.4: Multilayer Perceptron Architecture

2.2.6 Gradient Descent and Backpropagation

Gradient descent is an algorithm used to optimise the weights between neurons in such a manner that creates the least possible amount of error. When the network is training, a cost function is used to keep track of how the network is performing. The cost function looks at the discrepancies between the training output and the true values to determine an error. When the network trains, the goal is to therefore get this cost function as low as possible to ensure the lowest possible error. The way the cost function is minimised in an MLP is through gradient descent and backpropagation.

Every time the weights must be updated, the derivative of the cost function with regards to the weight itself scaled to a learning rate is subtracted. As the network trains the derivative should decrease with each training iteration. This is known as gradient descent. When the cost function cannot be reduced any more, it has converged.

The backpropagation algorithm works in conjunction with this by starting at the output layer of the network and stepping back through all the layers, updating the weights for each neuron as it goes (Rumelhart, Hinton, and Williams, 1985). This is in an attempt to minimise the overall error. These steps should allow the network to reduce its error and converge on an optimal solution over time as the network trains.

2.2.7 Overfitting

An important point to consider when training any Machine Learning model is how well the model generalises to new data (Domingos, 2012). An algorithm should be able to apply concepts learned from the training data to any previously unseen data in the problem domain in order to make accurate predictions. Overfitting is the concept of a model learning the details and noise in the training domain too well, and thus failing to provide accurate predictions on new data. Overfitting tends to occur when a model learns noise and randomness in the training data as concepts that it attempts to apply to data in the problem domain. These concepts will not apply to the new data however, and poor accuracy generally results. Although a noisy dataset will increase the severity of overfitting, it is not simply a result of a noisy dataset and can occur in any dataset (Domingos, 2012). Within the field of Neural Networks, a common approach to avoid severe overfitting is to keep the network architecture relatively simple. The less parameters required to train the model, the lower the overall chance that the network will overfit and fail to generalise (O'Shea and Nash, 2015). However, this approach is not always feasible when dealing with large and complex datasets, and other approaches must also be taken. Another popular method to reduce overfitting is the introduction of dropout in the network. Dropout is when nodes are randomly dropped from the network along with their connections (Srivastava et al., 2014). This serves to constrain the adaptation of the network while it trains in order to prevent the model from over-learning the training data.

2.3 Introduction to Machine Vision

Images are represented as matrices of values, with each value corresponding to a pixel in the image. Images can be broadly broken down into two groups, grayscale and colour. A grayscale image will be represented as a 2D array of values ranging from 0 to 255, with each value representing the intensity of that pixel, as seen in Figure 2.5. A value of 0 represents black and a value of 255 represents white. In grayscale images, each pixel therefore represents one colour channel (gray). In colour images however we have multiple different channels for red, green and blue. In order to represent this, a colour image will be represented as a 3D array of pixel values, with each pixel having 3 values between 0 and 255 associated with it. These three values represent the intensity of the three colour channels at that pixel. The idea is that a CNN should be able to receive these image representations as input data and return values for the probabilities of different classes.

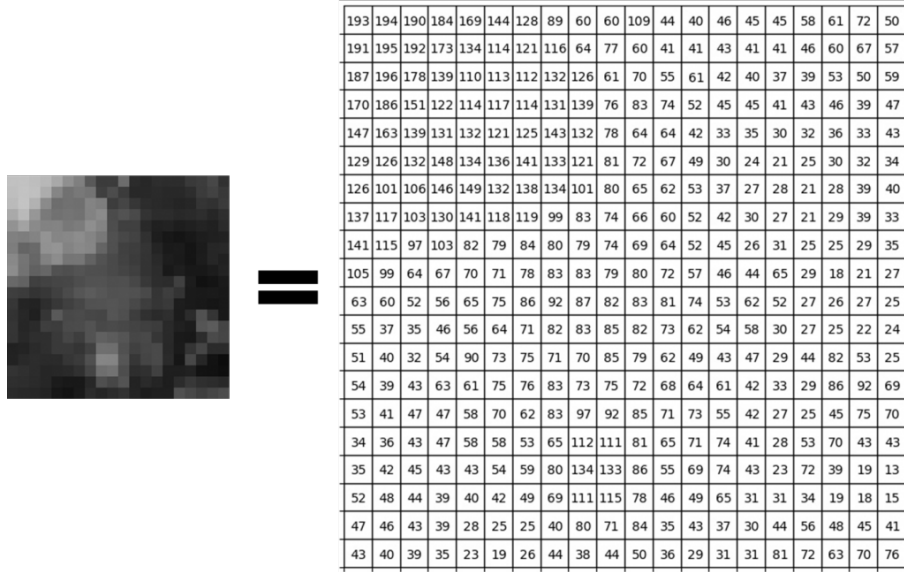


Figure 2.5: Pixel Representation of Grayscale Image

When humans look at an image of say, a dog, we automatically extract the features that make a dog unique in order to allow us to recognise it. For

example, we may see a tail, fur and a snout and recognise that we are looking at a dog. At an extremely high level this what a CNN will do - it will look for certain low-level features such as curves and edges and build these up into more abstract concepts such as a leg or a tail in order to classify the object (LeCun and Bengio, 1995).

Interestingly enough, CNN's do take some inspiration from how the brain processes images. The visual cortex of the brain contains many fields that are sensitive to different specific elements of the input. Some groups of neurons will only respond when certain elements of the input are present, for example certain vertical edges, and other groups of neurons will respond to different elements, such as horizontal edges (Hubel and Wiesel, 2011). The concept of distinct groups of neurons looking for certain features is one that has translated very well into Machine Vision through the use of CNN's.

2.3.1 Edge Detection

Edge detection is the ability to recognise object boundaries within an image, with most edge detection techniques using changes in pixel intensities to identify potential edges (Arbelaez et al., 2011). Although edge detection may appear straightforward in concept, in practice it can be a difficult task. Input images typically suffer from focal blur due to a finite point-of-field, as well as blur caused by shadows. This can cause smoothing in the variations in pixel intensity at edge points, and make it difficult to define what actually is an edge (Ziou and Tabbone, 1998). The magnitude of a gradient will determine if a point in the image is an edge or not - a high gradient implies that an edge is likely present. The direction of a gradient shows how the edge is oriented within the image.

Approaches

There is a multitude of different techniques used for edge detection that can be split into two categories, search based and zero-crossing based. Search based techniques work by computing gradient edge strength, then searching for the direction of the edge by computing the gradient orientation. The Sobel operator, mentioned below in the explanation of the Canny edge detector, is an example of this type of technique. Zero-crossing based techniques differ in that they search for zero-crossing points on a

second-order derivative function calculated from the image. Generally smoothing is applied to the image prior to any of these techniques (Ziou and Tabbone, 1998).

Canny Edge Detection

Although the Canny edge detector is just one technique used for edge detection, it is one of the more prolific techniques used, and serves as a good example to illustrate the concepts used in edge detection. It consists of multiple stages. The first stage, preprocessing, involves smoothing the image to help reduce noise. Next, the gradient magnitudes and directions are calculated throughout the image via the use of the standard sobel edge detector. The magnitude of the gradient is calculated as $m = \sqrt{G_x^2 + G_y^2}$ and the direction of the gradient is calculated as $\theta = \arctan\left(\frac{G_y}{G_x}\right)$, where G_x and G_y are the derivatives of the x and y coordinates of the point under investigation. Once these are calculated the edge detection begins through the process of nonmaximum suppression - if a pixel is not at maximum value it is suppressed. The orientation of the pixel is placed into one of four bins representing the four possible directions the edge could be present within as illustrated in Figure 2.6. The possible edge directions at the grey pixel point could be either north to south (represented by the green pixels), east to west (represented by the blue pixels) or one of the diagonal directions (represented by the yellow and red pixels). Nonmaximum suppression is then applied to these four different possible areas. If the gradient orientation is at a value such as the one represented in Figure 2.7, then the edge is moving diagonally perpendicular to the gradient direction. To check if a pixel belongs to an edge, its value is checked against its neighbours along the same orientation to see if the gradient is maximum and if it is above a defined threshold. If so, the pixel is marked as an edge. After this process of nonmaximum suppression, thin edges that may be broken at points are left. An edge thresholding process called hysteresis is used to fill in these broken points (Canny, 1986). Hysteresis counters the "streaking" effect of broken lines by setting both an upper and lower threshold for pixel gradient magnitudes. Values above the upper limit are accepted and values below the lower limit are rejected. Values which lie between the two limits are accepted if they lie between two points above the upper limit, reducing the overall gaps exhibited in the line. An example of Canny edge detection being used in a road environment is shown in Figure 2.8.

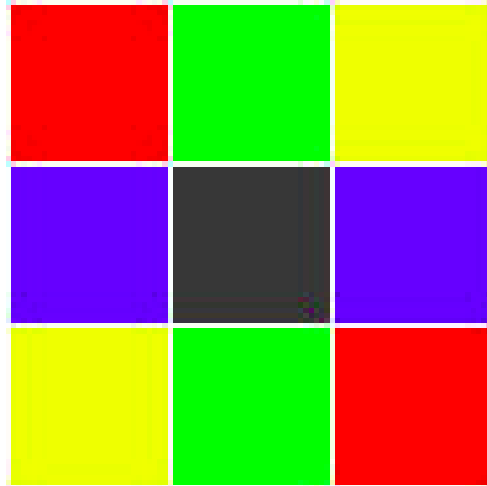


Figure 2.6: Possible Edge Directions

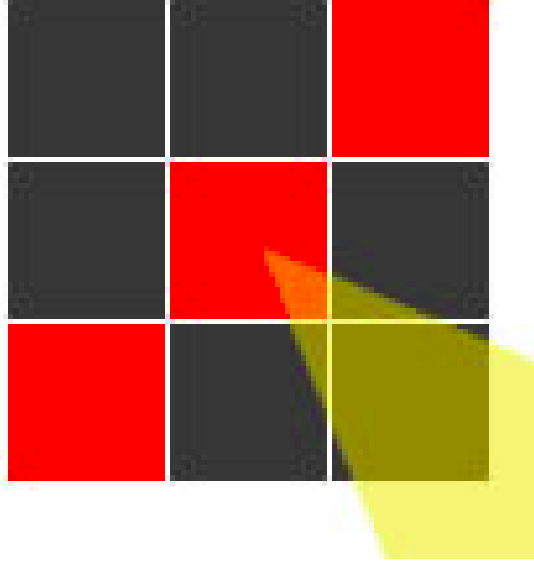


Figure 2.7: Orientation Value

2.4 Introduction to Convolutional Neural Networks

Convolutional Neural Networks have recently become the standard Neural Network used within the field of image-based Machine Learning tasks (O'Shea and Nash, 2015). They follow essentially the same architecture as a standard Multi Layer Perceptron, except that they include extra layers of Convolution and Pooling. When these layers of Convolution and Pooling are stacked with the Fully-Connected layers of the MLP, a CNN architecture has been formed. Once Convolution and Pooling has been performed the fully-connected layer will compute the overall class scores, resulting in a final output of $1 \times 1 \times n$, where n is the number of possible classes.



Figure 2.8: Canny Edge Detection

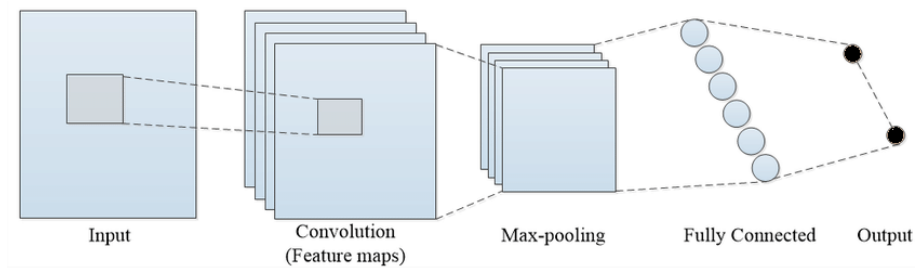


Figure 2.9: CNN Architecture

2.4.1 Convolution Layer

There are three important points of note regarding the Convolution layer: the input image, the feature detector and the feature map. The input image is the image which the CNN is given. The feature detector is a

matrix (usually 3x3 or 7x7), also called the kernel or filter, which is multiplied with the matrix values of the input image to create the feature map. The aim of this is to capture important features of the image, maintaining these important features while losing some unimportant elements and reducing overall image size. The way this works is the feature detector slides over the whole image. The particular point on the image that the feature detector is over is known as the receptive field. The parameter known as the stride is the value of the number of pixels by which the feature detector slides over the input image. For example, a stride of 2 will slide the feature detector over the input image jumping 2 pixels at a time. As the feature detector slides around the input image, the pixel values of the feature detector are multiplied against the pixel values of the input image. These multiplications are all summed to give a single value, as shown in Figure 2.10. The higher the value, the higher the likelihood that the particular feature represented in the feature detector is present in that section of the input image. The feature detector then slides on to carry out this process for every section of the input image.

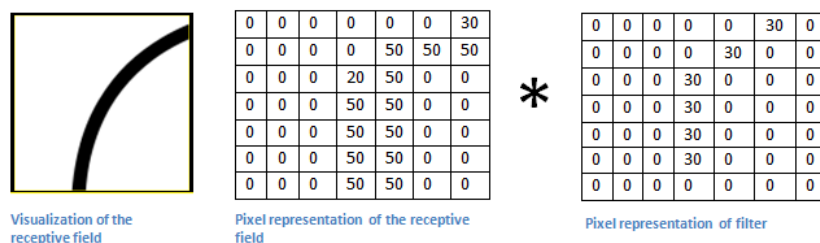


Figure 2.10: Convolution Layer

2.4.2 Pooling Layer

The pooling layer performs downsampling to extract key important features and reducing the number of parameters. The pooling layer operates over every feature map and by placing a matrix over the feature map and selecting the min value, max value or mean value, depending on the type of pooling being utilised, as shown in Figure 2.11. These extracted values form the pooled feature map. Pooling helps to reduce image size while also extracting important features.

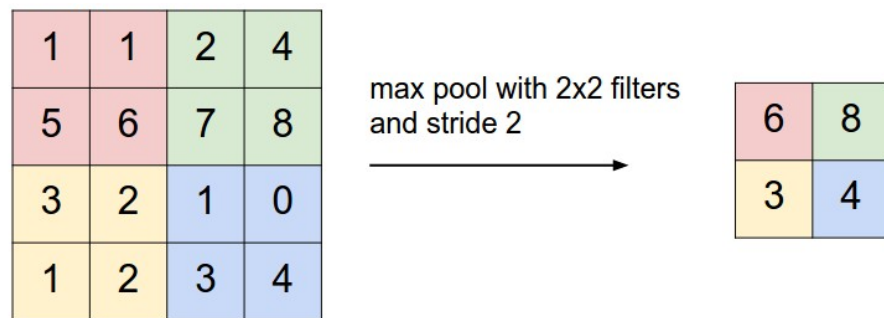


Figure 2.11: Pooling Layer

2.4.3 Relevant Architectures

This section will outline some important CNN architectures that are commonly used. These architectures will later be implemented and investigated during the Empirical Studies section.

MobileNet

Following the success of AlexNet in 2012, CNNs followed a general trend towards deeper and more computationally expensive architectures. However, increased network complexity does not guarantee increased network performance (Szegedy et al., 2016), and in computationally limited platforms such as smartphones and self-driving cars tasks need to be carried out quickly with a low computational overhead. The MobileNet architecture arose out of this demand for more lightweight and efficient architectures (Howard et al., 2017). MobileNet differs from conventional CNN architectures in that it makes use of depthwise separable convolutions in the place of regular convolutions. A depthwise separable convolution consists of two operations - a depthwise convolution followed by a pointwise convolution.

A standard convolution works on the spatial dimension of the feature maps as well as the input and output channels.

Inception

In a similar vein to the MobileNet architecture, the Inception architecture arose from a desire to create more complex networks to increase performance, without leading to the issues presented by simply stacking more layers into a network. At the current time there are four main versions of the Inception architecture, with each improving slightly over the preceding version. The core concept behind the creation of this network was that important objects or regions within a given image can vary greatly depending on the distance between the object and the camera. Due to this variance, selection of the correct kernel size is difficult. A larger kernel size is suitable for the larger objects, however this will not suffice for the smaller objects. Rather than simply adding extra layers, a possible solution is to have multiple filters operate on the same level, creating a "wider" rather than "deeper" architecture. This was the driving concept behind the creation of the Inception v1 (also known as GoogLeNet) architecture.

2.4.4 Evaluating the Classifier

There are various techniques which can be used to evaluate the operation of a classifier. Some of the more prevalent ones, many of which will be seen throughout the rest of this report, are explained below.

Top-1 and Top-5 Accuracy

First, the CNN makes a classification. Then in the case of Top-1, the class of the highest probability is checked against the target label. In the case of Top-5, the target label is compared with the top five highest probability predictions. Following this, in both cases the classifier score is calculated as the number of times that the prediction matched in either Top-1 or Top-5, divided by the total number of data points.

Accuracy Issues

Generally speaking using accuracy alone (number of correctly labelled classes) is not a very reliable metric for an object detection classifier. For example, if there were 95 dogs and 5 cats to be classified a classifier may label all objects present as dogs. This would lead to an overall accuracy score of 95%, however this classifier is obviously quite flawed. This is known as class imbalance, and in real life scenarios tends to be the norm. The decisions reached by the classifier can be explained as follows:

1. True positive: a positive example classified as positive
2. True negative: a negative example classified as negative
3. False positive: a negative example classified as positive
4. False negative: a positive example classified as negative

False negatives are particularly an issue for the field of self-driving cars, where this result could have potentially fatal ramifications. An example of a false negative classification was seen during a fatal crash of a Tesla model S in May 2016, whereby the autopilot system failed to recognise a white truck against a clear, brightly lit sky. Thus, we need to do more than just measure accuracy in order to properly evaluate the classifier (Rogers and Girolami, 2016).

| | | Actual | |
|-----------|---|-----------------|-----------------|
| | | + | - |
| Predicted | Y | True Positives | False Positives |
| | N | False Negatives | True Negatives |

Confusion Matrix

A Confusion Matrix, also referred to as an Error Matrix, is a relatively simple technique for classifier evaluation. It is essentially a table that plots the number of correct and incorrect predictions for all of the different classes. The false positives, false negatives, true positives and true negatives are plotted then the average value for all classes combined is calculated, as demonstrated in the confusion matrix below.

Other Evaluation Metrics

There are a number of different metrics of evaluation that can be used other than accuracy alone, with the major ones being summarised in the section. An important point to note is that each of these metrics is derived from the confusion matrix explained above, and can be thought of as different ways of summarising the matrix.

Precision

The precision is the number of correct positive predictions compared to the actual number of positive examples. It can be represented by the formula $precision = \frac{\#TP}{\#TP + \#FP}$. This is a good metric to evaluate when there is a high cost associated with a false positive. For example, an email spam detector must have high precision in order to avoid mistakenly classifying legitimate emails as spam.

Recall

The recall is the ratio of correct positive examples to the number of positive predictions, and can be represented as $recall = \frac{\#TP}{\#TP + \#FN}$. Recall is an important metric when there is a high cost associated with false negatives and as mentioned above are an extremely important metric for this project.

F-score

Another metric can be derived from precision and recall, known as F-score. F-score is the harmonic mean of precision and recall, and is defined as $F = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$. The F-score can be thought of as representing the balance between the precision and the recall. Although this is an important metric overall, it does not take precedence in this project, where a good recall score is of utmost importance.

AUC-ROC Curve

The AUC (Area Under Curve) ROC (Receiver Operating Characteristics) curve is another method of classifier evaluation, enabling the visualisation of results from a classification problem. The ROC curve is used to plot the true positive rate against the false positive rate, resulting in a probability curve, while the AUC represents how well the model distinguishes between classes. The higher the AUC, the better the model is. This is visualised in Fig 2.12.

2.4.5 Transfer Learning

Transfer learning is a popular technique in Machine Learning whereby an already trained model is applied to new data. As mentioned earlier, CNN's use their lower layers to extract primitive features such as edges, the middle layers detect shapes from these edges and the final layers learn the task specific abstractions of these shapes. In transfer learning the last layers of the network are retrained using data from the target domain, allowing the insights that have been learned from the original data to be applied to the new data. This helps to avoid training models from scratch on unseen yet related datasets. Transfer learning is a very popular field within Machine Learning, as it can help drastically cut down on training times, as well as enabling better results on smaller datasets.

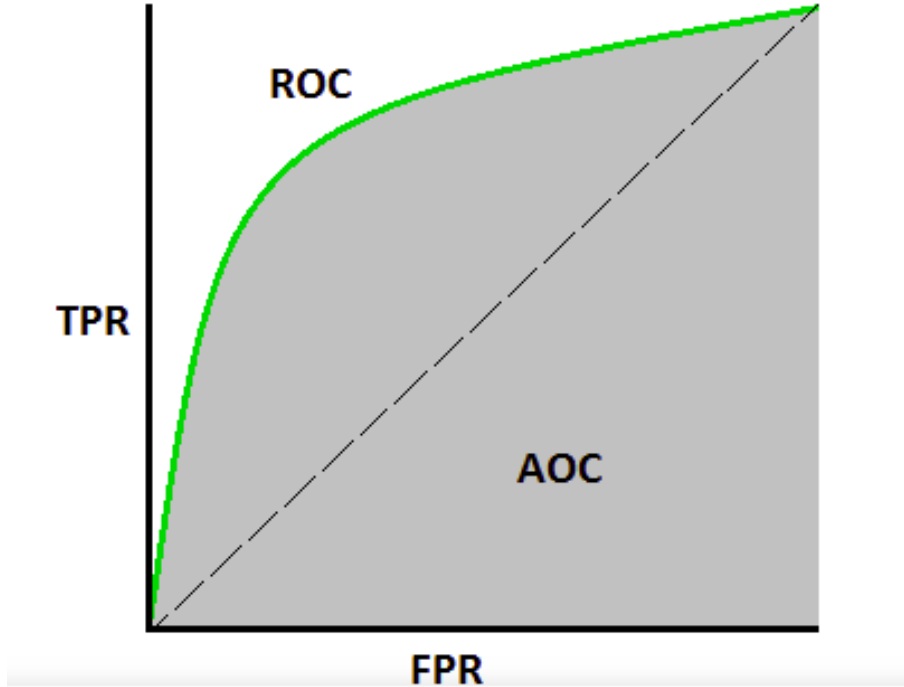


Figure 2.12: AUC-ROC Visualised

2.4.6 Berkeley Deep Drive

Dataset

As mentioned previously in this report, the dataset utilised for this project is the Berkeley Deep Drive dataset, a "large-scale diverse driving video dataset", also known as the BDD100K dataset. Released in 2018, the dataset consists of 100,000 1280x720 annotated images and 100,000 720p 30fps video sequences (Yu et al., 2018). The video was collected from a number of locations throughout the United States and consists of a range of weather conditions including sunny, overcast, and rainy. The data consists of a mixture of both daytime and nighttime driving conditions. The images section of the dataset has been created by selecting and annotating the frame at the 10th second of every video. Annotated objects range across 10 different classes - bus, light, sign, person, bike, truck, motor, car, train and rider. The number of instances of each of these annotations are displayed in

Fig 2.13.

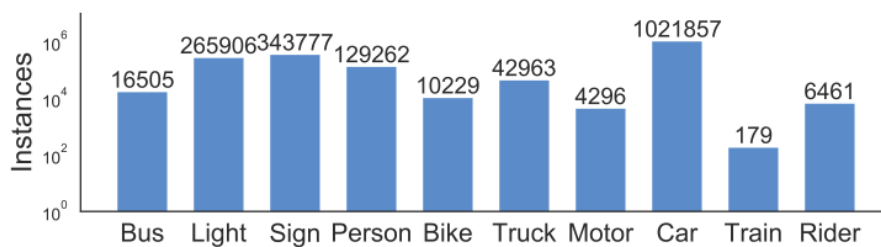


Figure 2.13: source: <https://bair.berkeley.edu/blog/2018/05/30/bdd/>

Related Results

The paper released with the BDD100K dataset contains a list of results from experiments carried out on the dataset. The experiments focused on the performance of a CNN on different domains within the dataset. For example, a Faster-RCNN was trained on daytime images and then its performance on nighttime images was tested and the discrepancies recorded. The performance metrics recorded from these experiments allow some basis upon which the CNN's trained in this project can be evaluated against.

2.5 Tensorflow Introduction

This section is dedicated to the background work done in order to gain an understanding of how to implement CNN's using tensorflow. The experiments illustrated in this section were carried out during the completion of the Complete Guide to TensorFlow for Deep Learning with Python course on Udemy (*Complete Guide to Tensorflow for Deep Learning with Python* 2018). This section consists of two experiments carried out on the MNIST and CIFAR-10 datasets. The first experiment is sample code provided as part of the course to illustrate the concepts learned during the course. The second experiment is an exercise for the course carried out to implement an understanding of these concepts.

2.5.1 MNIST Experiment

The first experiment carried out was the implementation of a very simple CNN. The architecture of the CNN consisted of two layers of Convolution, two pooling layers and a fully-connected layer. After carrying out 5000 training steps with a batch size of 50, a final Top-1 accuracy of 98.75% was achieved.

Code Breakdown

The functions displayed in Figure 2.14 are all helper functions. The `init_weights` function initialises the random weights for fully connected or convolution layers, with the shaper of the layer passed in as a parameter. The `init_bias` function performs the same operation for the bias. The `conv2d` function creates a convolution using a built in function from tensorflow. The `max_pool_2by2` function creates a max pooling layer, also using built in tensorflow functions. The `convolutional_layer` function uses the `conv2d` function to return an actual convolutional layer with a ReLu activation function. Lastly, the `normal_full_layer` returns a normal fully connected layer.

The functions displayed in Figure 2.15 are used to create the convolution and pooling layers of the network. In creating `convo_1`, we can see that a 6x6 filter is used from the parameters. An output value of 32 is used to

```

def init_weights(shape):
    init_random_dist = tf.truncated_normal(shape, stddev=0.1)
    return tf.Variable(init_random_dist)

def init_bias(shape):
    init_bias_vals = tf.constant(0.1, shape=shape)
    return tf.Variable(init_bias_vals)

def conv2d(x, W):
    return tf.nn.conv2d(x, W, strides=[1, 1, 1, 1], padding='SAME')

def max_pool_2by2(x):
    return tf.nn.max_pool(x, ksize=[1, 2, 2, 1],
                          strides=[1, 2, 2, 1], padding='SAME')

def convolutional_layer(input_x, shape):
    W = init_weights(shape)
    b = init_bias([shape[3]])
    return tf.nn.relu(conv2d(input_x, W) + b)

def normal_full_layer(input_layer, size):
    input_size = int(input_layer.get_shape()[1])
    W = init_weights([input_size, size])
    b = init_bias([size])
    return tf.matmul(input_layer, W) + b

```

Figure 2.14: Creating Helper Functions

represent the number of filters used. The parameter of 1 represents the original input of the image. This carries down to the other layers, with `convo_2` taking in an input image of 32 from `convo_1`.

The code displayed in Figure 2.16 shows the model being trained. A tensorflow session is created and data is read in from the dataset. The session is ran to begin model training. The current training step and current error calculated from a loss function is displayed as output as the model trains.

```

convo_1 = convolutional_layer(x_image,shape=[6,6,1,32])
convo_1_pooling = max_pool_2by2(convo_1)

convo_2 = convolutional_layer(convo_1_pooling,shape=[6,6,32,64])
convo_2_pooling = max_pool_2by2(convo_2)

convo_2_flat = tf.reshape(convo_2_pooling,[-1,7*7*64])
full_layer_one = tf.nn.relu(normal_full_layer(convo_2_flat,1024))

# NOTE THE PLACEHOLDER HERE!
hold_prob = tf.placeholder(tf.float32)
full_one_dropout = tf.nn.dropout(full_layer_one,keep_prob=hold_prob)

y_pred = normal_full_layer(full_one_dropout,10)

```

Figure 2.15: Creating Layers

```

steps = 5000
with tf.Session() as sess:
    sess.run(init)
    for i in range(steps):
        batch_x , batch_y = mnist.train.next_batch(50)
        sess.run(train,feed_dict={x:batch_x,y_true:batch_y,hold_prob:0.5})

        # PRINT OUT A MESSAGE EVERY 100 STEPS
        if i%100 == 0:
            print('Currently on step {}'.format(i))
            print('Accuracy is:')
            # Test the Train Model
            matches = tf.equal(tf.argmax(y_pred,1),tf.argmax(y_true,1))
            acc = tf.reduce_mean(tf.cast(matches,tf.float32))

            print(sess.run(acc,feed_dict={x:mnist.test.images,y_true:mnist.test.labels,hold_prob:1.0}))
            print('\n')

```

Figure 2.16: Training the Model

2.5.2 CIFAR-10 Experiment

The second experiment was carried out on the CIFAR-10 dataset, and the architecture used was consistent with the first experiment, consisting of two layers of Convolution and Pooling. After 5000 training steps with a batch size of 100, a final Top-1 accuracy of 72% was reached.

2.5.3 Initial Findings

The first two experiments were carried out using CNN's of the same architecture, however an accuracy discrepancy of 26.75% was observed. Both datasets contain images labelled into 10 classes. The differences arise in the dataset size and image complexity. Firstly, the CIFAR-10 dataset contains 50,000 training images, whereas the MNSIT dataset contains 60,000 training images. The images in the MNSIT dataset are grayscale, whereas images in the CIFAR-10 dataset consist of 3 colour channels. The images in the CIFAR-10 dataset are also of a higher complexity compared to the hand-written digits in the MNIST dataset. The CIFAR-10 images contain more complex real-world images with noisy backgrounds. These discrepancies between the datasets may explain the discrepancies in the accuracy achieved by the classifiers.

2.5.4 Experiment Environments

Initial environment setup to carry out the below experiments proved quite challenging. During the first attempts at carrying out the experiments, tensorflow was erroneously installed to run using the training machines CPU. This lead to untenable training times for the experiments. Although the training times were not recorded precisely, each experiment took over a day to finish training, during which time the training machine was rendered essentially unusable due to the high load on its CPU. There were two main take-aways from the initial iterations:

1. Tensorflow should be installed properly to use the GPU in order to manage training times.
2. Training times should be recorded for better understanding of the experiments.

Tensorflow Install

The initial attempts to setup tensorflow were carried out on a machine running Windows 7. Several versioning issues were encountered due to Bazel, the build tool used to compile tensorflow. These issues were solved via downgrading to an earlier version of Bazel, however new versioning issues were then encountered due to CUDNN, Nvidias GPU-accelerated

library for deep learning. Several different versions were installed without success. Due to these issues seemingly being exclusive to Windows, it was decided to attempt an install on a different machine running Ubuntu. The tensorflow install for the Ubuntu machine was carried out without issue. All necessary Nvidia drivers were installed on the system, and tensorflow was then installed and the experiments were rerun. The training times for running the two experiments with GPU acceleration were 12 and 16 minutes respectively, a vast improvement over the initial times. The GPU on the Ubuntu machine was an Nvidia Geforce GTX 650 with 2GB memory.

Memory Issues

Although the first two experiments ran without issues, problems were quickly encountered when modifications were made to the Cifar-10 experiment. In an attempt to improve accuracy results, two further layers of convolution and pooling were added. This led to out of memory (OOM) errors being thrown during training. This was caused by the graphical memory on the training machine being maxed out at the full 2GB. Several attempts at rectifying the issues were carried out. The first attempt tried was reduction of batch size from 100 images, causing less of the dataset being loaded into memory. At a batch size of 4 images, the three layer network successfully trained. The issue with this approach however is that a smaller batch size causes the gradient estimate of the network to be less accurate, and Top-1 accuracy of the network dropped to 60.46%. The second approach was to simply keep the number of layers low, resulting in a simpler network with less memory requirements. Both of these solutions were not viable for this project. Poor accuracy scores and an inability to train complex networks created an obvious demand for a more powerful training machine. Another point to note is image size disparities between the Cifar-10 dataset and the overall target dataset for this project, the Berkeley Deep Drive dataset. A machine struggling with the 32x32 low resolution images of the Cifar-10 dataset would surely run into a host of issues with the 1280x720 images of the target dataset. As such, it was decided that a much more advanced machine would be required. It was decided that utilising Amazon Web Services (AWS) products was the only way to quickly and easily access the required computational power demanded by this project.

AWS Setup

One of the many services offered by the AWS platform is Amazon Machine Images (AMI) for Deep Learning. These are virtual machine instances that launch with pre-installed pip deep learning frameworks such as tensorflow, keras etc. in different environments, as well as GPU acceleration. This allows development of models remotely from a client machine which can be then trained using the available GPU resources on the deep learning instance. These instances can be deployed and the desired environment (tensorflow in this case) activated. Jupyter notebooks can be run as normal on the AWS instance, and a rule is set up on the client machine to forward all requests on a certain port to the AWS instance. This allows the user to write code in a Jupyter notebook remotely on the AWS instance. The instance type chosen for this project was the p2.xlarge instance which contains 122GiB of memory, 16 virtual CPU cores and a Nvidia K80 GPU. The steps followed was from the official AWS documentation.

Cifar-10 Further Experiments

As mentioned above, further experiments were carried out on the Cifar-10 dataset in an attempt to improve the disappointing initial results. The base code provided as part of the tutorial (*Complete Guide to Tensorflow for Deep Learning with Python* 2018) was modified in a range of different ways in order to maximise results. In order to evaluate the network properly the following metrics were recorded:

1. Peak Top-1 accuracy
2. Precision
3. Recall
4. Training time

The baseline of these metrics achieved by the simple network provided from the tutorial are shown in the table below, with the accuracy, precision and recall being graphed in Fig 2.17. Interestingly enough, the accuracy and recall values for this experiment were identical throughout all of the training steps. This does not necessarily mean that there is an issue, however it is very unusual. A focus for the following experiments was to investigate what changes, if any, would cause these values to diverge.

Initial Experiment

| Peak Top-1 accuracy | Training Time | Precision | Recall |
|---------------------|---------------|-----------|--------|
| 70.16% | 220s | 71.67% | 70.16% |

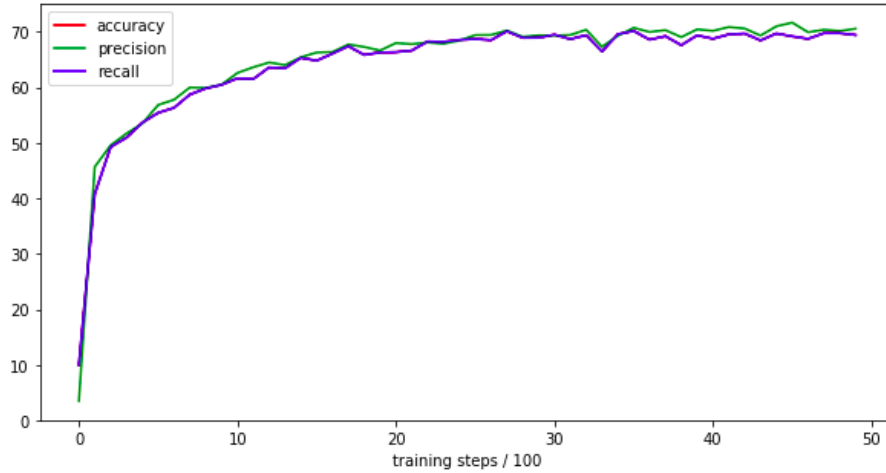


Figure 2.17: Baseline Plot

Experiment 1: Training Steps

Before investigating the effects of altering the network architecture, the optimal number of training steps was investigated. This parameter was investigated first as it has a large effect on training time. A network that does not have enough training steps will underlearn, whereby it does not learn all possible parameters. However having too many training steps will lead to increased training time with no increase in accuracy. The initial experiment carried out 5000 training steps, so this was doubled to 10,000 steps in an attempt to see if this would boost accuracy, precision and recall. The results are displayed in the table below and plotted in Fig 2.18. All three metrics levelled off at around the 5000 step mark, with no significant increases reported. It would therefore appear that 5000 training steps is the optimal value.

Experiment 1

| Peak Top-1 accuracy | Training Time | Precision | Recall |
|---------------------|---------------|-----------|--------|
| 70.46% | 466s | 70.71% | 70.46% |

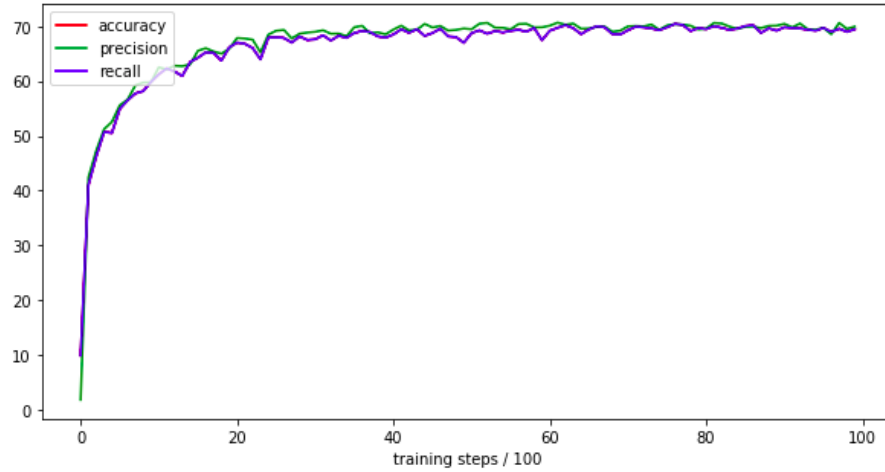


Figure 2.18: 10,000 Training Steps

Experiment 2: Convolution and Pooling Layers

The next experiment undertaken was to increase the number of convolution and pooling layers. For this experiment only one extra layer of convolution and pooling was added. As can be shown in the results displayed in the table below as well in Fig 2.19, this did not yield any significant results.

Experiment 2

| Peak Top-1 accuracy | Training Time | Precision | Recall |
|---------------------|---------------|-----------|--------|
| 70.38% | 192s | 70.82% | 70.38% |

Experiment 3: Fully Connected Layers

As the extra layers of convolution and pooling had failed to yield any benefits, the next step taken was to increase the number of fully connected layers in the network. One extra fully connected layer was added as well as a dropout layer. The reasoning behind this experiment was to investigate the possibility that the network was performing poorly with the class predictions. The images in the Cifar-10 dataset consist of non-occluded blurred objects with one class of object present per image. As such it was reasoned that only a small number of convolution and pooling layers were

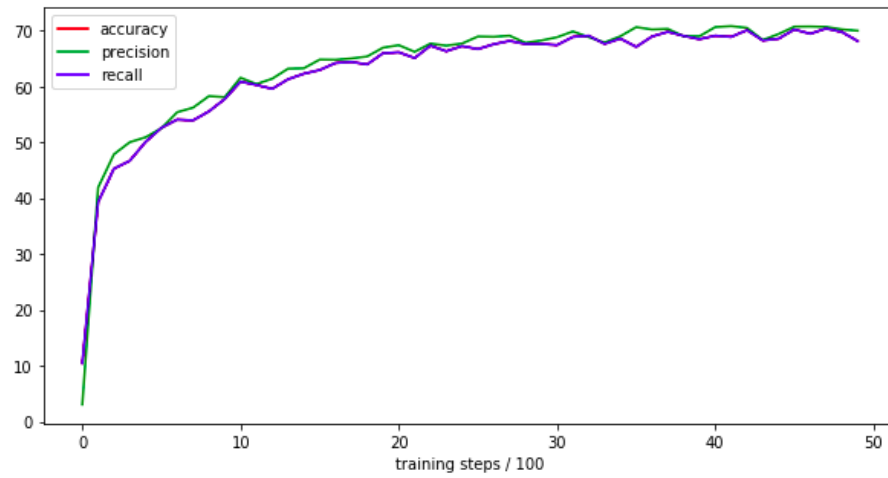


Figure 2.19: Experiment 2

required to extract feature maps from the images, and increased numbers of fully-connected layers would be more important. Dissapointingly, the results of this experiment were a marginal drop in accuracy, precision and recall.

Experiment 3

| Peak Top-1 accuracy | Training Time | Precision | Recall |
|---------------------|---------------|-----------|--------|
| 69.52% | 161s | 69.67% | 69.52% |

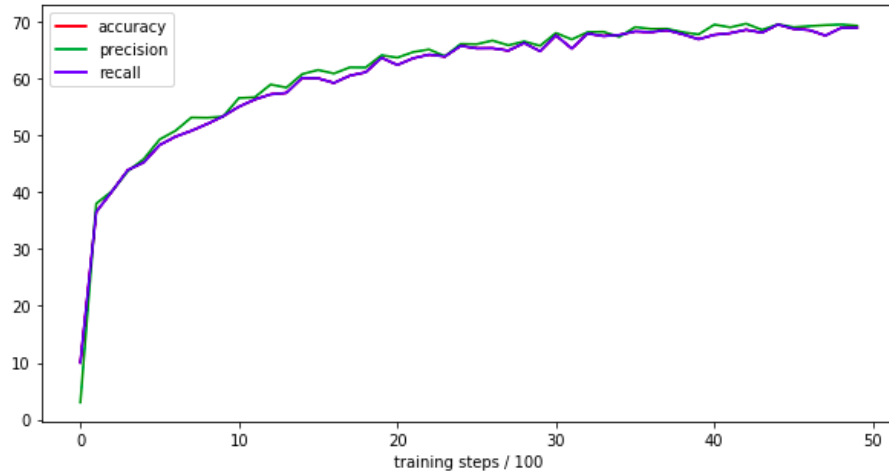


Figure 2.20: Experiment 3

Experiment 4: Convolution, Pooling and Fully Connected Layers

As Experiments 2 and 3 had failed to yield better results, 2 extra layers of convolution and pooling were added in conjunction with an extra fully-connected layer. As can be seen from the results table and Fig 2.21, results actually went down marginally.

Experiment 4

| Peak Top-1 accuracy | Training Time | Precision | Recall |
|---------------------|---------------|-----------|--------|
| 67.3% | 276s | 68.56% | 67.3% |

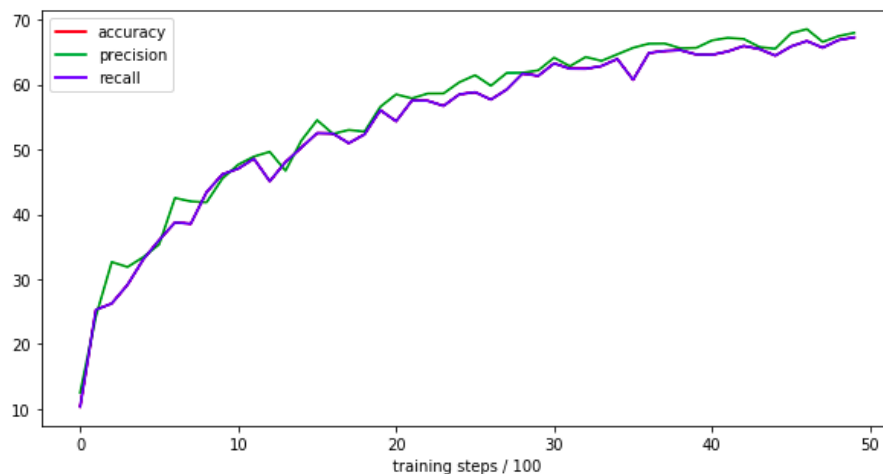


Figure 2.21: Experiment 4

Experiment 5: Extra Dropout Layers

Following the failure of experiment 4 to improve results, 2 extra layers of dropout with a probability of 0.25 were added to the experiment 4 architecture after the 2nd and 4th pooling layers in an attempt to tale any potential overfitting. Although the results summarised in the table below and Fig 2.22 do not show any meaningful increase in accuracy either, the accuracy curve from this experiment as well as experiment 4 suggest that the network is underlearning as the curves do not level off as much as in previous results. This implied that the network may simply need more time to train. Interestingly, during this experiment the accuracy and recall metrics diverged from one another, as can be observed in Fig 2.22. This would seem to imply that the network may have been overfitting slightly.

Experiment 5

| Peak Top-1 accuracy | Training Time | Precision | Recall |
|---------------------|---------------|-----------|--------|
| 60.47% | 425s | 60.60% | 60.57% |

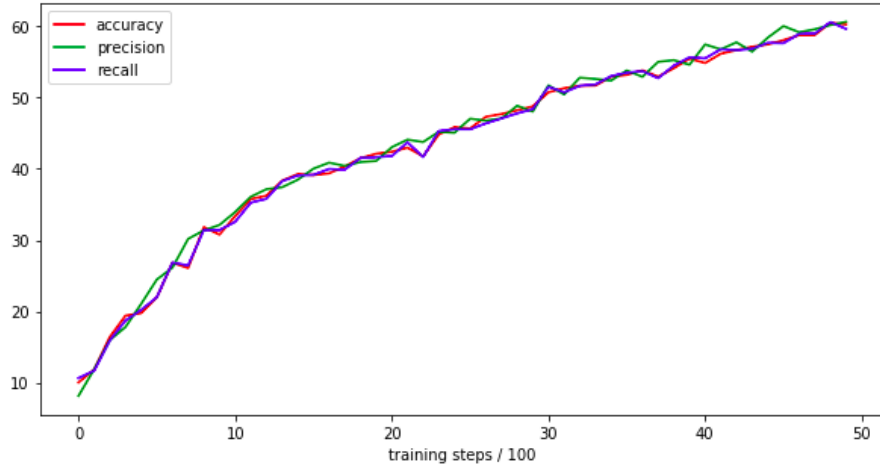


Figure 2.22: Experiment 5

Experiment 6: Increased Training Steps

The observations in experiments 4 and 5 suggested that the network might simply need more time to train. In experiment 1 it was concluded that 5000 training steps would be sufficient, however this experiment was carried out using the simple baseline network architecture. It was therefore a mistake to assume that 5000 training steps would be sufficient for all following experiments. The number of training steps was increased drastically to 25,000, with the architecture remaining the same as in experiment 5. Although training time was significantly increased, all metrics gained a slight increase with accuracy peaking at 73.25%, a decent improvement from the initial results.

Experiment 6

| Peak Top-1 accuracy | Training Time | Precision | Recall |
|---------------------|---------------|-----------|--------|
| 73.54% | 1778s | 74.03% | 73.48% |

Conclusions

These experiments conclude with a small yet significant increase in accuracy, precision and recall scores from the final experiment. These experiments were focused mainly upon the effects of network architecture rather than the effects of parameter tuning. Empirical studies in later

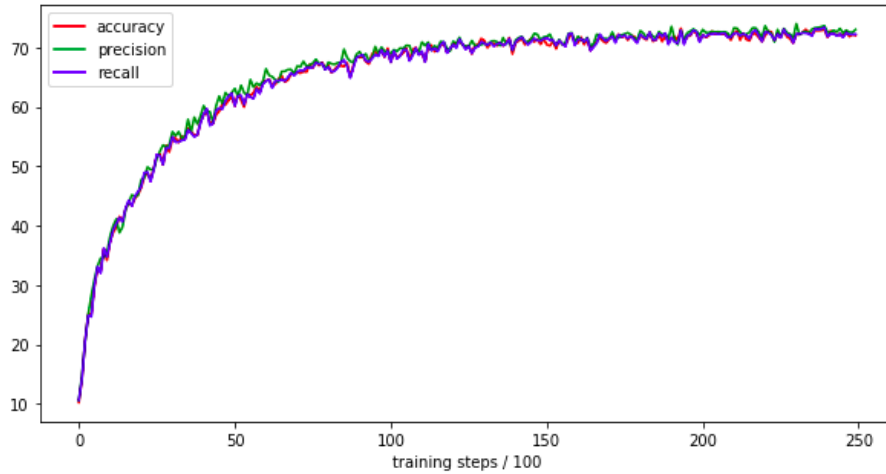


Figure 2.23: Experiment 6

sections will focus more on the effects of parameter tuning. Adding just extra layers of convolution and pooling without adding extra fully-connected layers did not result in improvements, and vice versa. Upon creation of a more complex network architecture with dropout layers results actually went down slightly. However, observation of the metric graphs showed that the network was possibly underlearning, and a drastic increase in training steps led to positive results. Accuracy, precision and recall all seemed to increase in proportion with one another for these experiments, however this is generally not guaranteed for all networks. Datasets containing more classes or high levels of class imbalance would typically not result in such regular accuracy, precision and recall metrics. If there is anything to take from these experiments moving forward with this project, it is that simply adding one or two extra layers is not sufficient to gain satisfactory results - adding multiple layers of convolution, pooling and fully-connected layers should be investigated in conjunction with each other in future experiments. As the unusual accuracy and recall results also demonstrate, it may also be beneficial to introduce dropout or some other form of regularisation in networks that are not obviously overfitting. Although the increases in performance shown are slight, I am confident that further increasing the complexity of the network and increasing training times can lead to further increases.

Chapter 3

Empirical Studies

The following sections will investigate the effects of retraining various different network architectures using transfer learning and investigate the effects of parameter tuning upon these networks. Models were retrained using the Tensorflow Object Detection API, with the steps followed being from a tutorial made available on GitHub (*TensorFlow Object Detection Model Training* 2018). All of the pretrained models utilised in the following sections have been trained on the COCO (Common Objects in Context) dataset, a dataset provided by Microsoft. Although the COCO dataset is updated every year, it currently consists of 164,000 images populated with 90 classes of objects. The dataset was chosen for this project as it contains 8 different classes of vehicle and displays objects in real-life contexts, with potential for occlusion and general noise in the images (Lin et al., 2014). Many other datasets instead provide an "iconic" view of objects, whereby the object appears centrally and unobstructed in the image. As the BDD100K dataset images are drawn from real-world scenarios from the perspective of a car, using the COCO dataset as the initial training dataset should allow the retrained models to generalise better to the BDD100K data than datasets with less variance in the images.

3.1 SSD MobileNet V1 Experiments

The SSD MobileNet V1 architecture was selected for the first experiments due to its speed.

3.1.1 Full Dataset Retraining Experiment

Objectives

Initial attempts at retraining the SSD MobileNet V1 were carried out using the full 10k subset of the bdd100k dataset, containing 10,000 images annotated with 10 classes. The objective for this experiment was to retrain the MobileNet model with the new BDD100K data. With training times being a concern when dealing with large datasets, the MobileNet architecture was selected to be retrained first due to its quick training times. This main objective for this first experiment was to investigate how models would handle the large amount of data presented by the BDD100K dataset and to gain some understanding of how the object detection API works.

Setup

The tensorflow object detection API only accepts datasets in the tfrecord format, a binary file format for data storage. This requires parsing the images and their labels and writing them into tfrecord format. Luckily, an open source parser for the BDD100K dataset has been made available on GitHub (*Convert the Berkeley Deepdrive dataset to a TFRecord file* 2018). This tool was used to parse all images within the 10k section of the BDD100K dataset, resulting in a tfrecord file containing the data for the full 10,000 images and their annotations.

Results

After 5000 training steps, the results for this experiment were disappointing albeit not unexpected. As can be seen in Fig 3.3, loss values were extremely high, and did not follow any steady downward trend. Total training time to 10,000 steps took 2 hours and 48 minutes. These poor results were expected due to the very large and highly complex dataset. Fig 3.1 provides an example of this - the nighttime conditions combined with slight motion blur creates a very noisy image on which to make predictions. Notwithstanding dataset noise created by time of day or weather conditions, many of the images in the dataset are extremely noisy simply due to the scenes depicted within them. For example, in city scenes such as the example in Fig 3.2, rows of cars cause occlusion and people are hard to pick out between the cars. With the dataset consisting in large part of

images with similar levels of variance and noise, training times for the full dataset would likely be untenable if a CNN was to be trained to a point where it performed close to the results presented in (Yu et al., 2018).. As such, investigation into a different approach was required.



Figure 3.1: Example of an extremely noisy image

3.1.2 Partial Dataset Retraining Experiment 1

Objectives

Following the failure to produce adequate results on the full dataset, it was determined that an approach of drastic data reduction was required in order to investigate training on a problem domain of reduced complexity. The objective for this experiment was therefore to retrain the MobileNet V1 architecture again, this time with a much simpler dataset and with less classes required for detection. In order to achieve this, a subgroup of images would need to be extracted from the dataset and labelled manually.



Figure 3.2: Example of the image noise caused by city streets

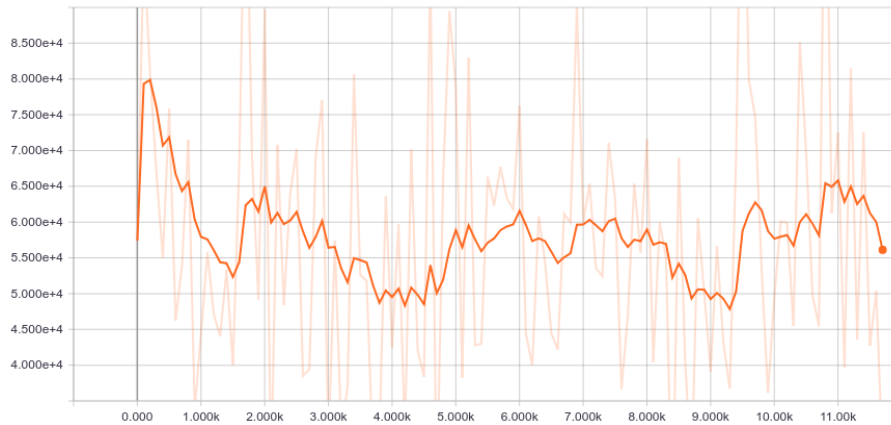


Figure 3.3: Loss values for full dataset

Setup

It was determined that 200 total images would be selected to be manually labelled from the bdd100k dataset, with 160 being utilised for training and 40 for testing. Classes labelled for detection were reduced in number to just 3 - cars, traffic lights and people. These classes were chosen as they are commonly present in most images in the dataset and are all unrelated to each other. The images selected were all images taken with good lighting,

no motion blur and no inclement weather present. The scenes depicted by the selected images were semi-urban - all three classes were present albeit in lesser numbers for traffic lights and people. Fully urban scenes were avoided due to the potential for occlusion and noise within the images. The images were annotated using the open source Labellmg tool, before being converted to tfrecord format. Training was then carried out exactly the same as the above experiment, this time with the reduced dataset.

Results

Although far from perfect, results from this experiment were much more encouraging. As evidenced in Fig 3.4, loss quickly declined before settling to values just below 2, with the minimum value being 1.788 and a total training time of just 1hr 25min. Mean average precision (3.5) and average recall (3.6), although slightly erratic, appeared to settle into an upward trend as training progressed. Maximum precision reached was 0.1995 and maximum recall reached was 0.2812. The model can be seen in action at 500 training steps in Fig 3.7 and then at 5000 training steps in Fig 3.8. The improvement in the model is clearly apparent in these two images. Please note however that these images appear blurred due to the fact that this MobileNet architecture was initially trained on images of size 300x300 pixels in order to reduce training times. This means that the BDD100K images that it has been retrained on have been scaled down to size 300x300. Scaling the images back up to full size results in the distortion present in the example images. These results are far from perfect however, and will need to be improved upon. Although the model at 5000 steps detects far more than it did at 500 steps, there are still several of the more occluded cars present in the image not being detected. As mentioned, mean average precision and recall did seem to be improving, however neither value increased to above 0.3, which is comparatively poor results. As such, further experimentation is required.

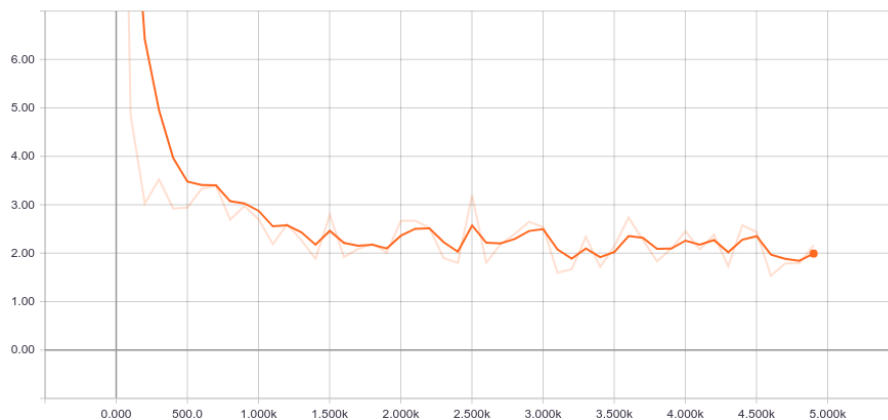


Figure 3.4: Loss values for reduced dataset

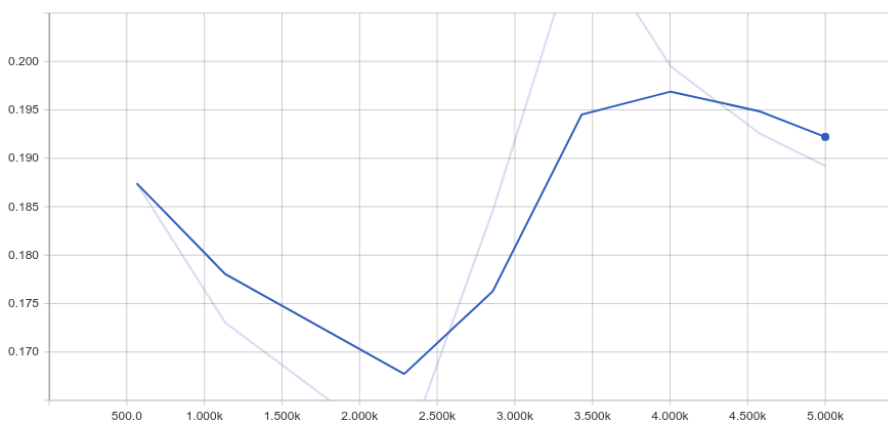


Figure 3.5: Mean average precision (mAP) values for reduced dataset

3.1.3 Partial Dataset Retraining Experiment 2

Objectives

Following the more encouraging results from the previous experiment, it was decided that the easiest way to boost performance was to simply increase the amount of training data. A further 200 images were hand annotated, with 120 being added to the training data and the final 80 being added to the validation data. This brought the size of the manually annotated dataset to 400 images. Although this is a comparatively small amount, the relatively close relationship between the COCO dataset and the BDD100K

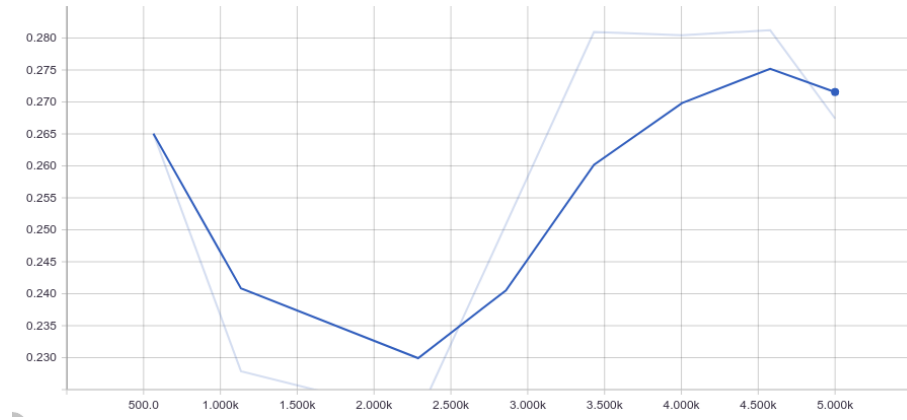


Figure 3.6: Average recall (ar) values for reduced dataset

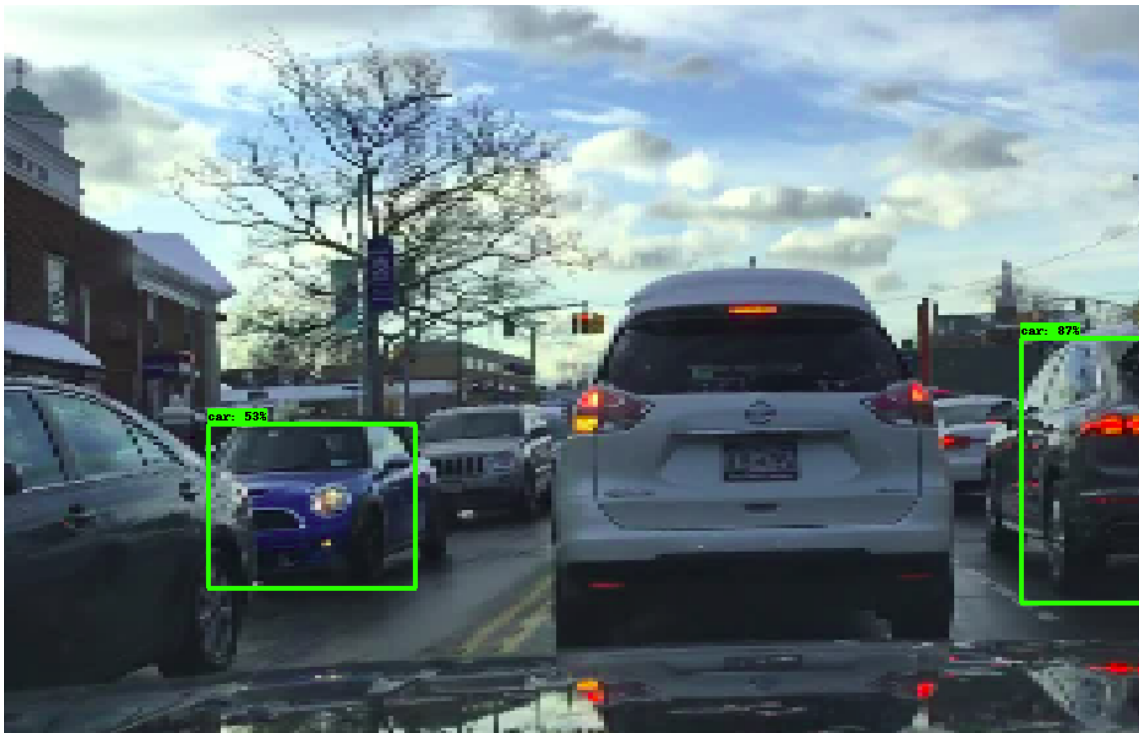


Figure 3.7: Object detector in action at 500 steps

dataset was a mitigating factor that allowed a smaller dataset.

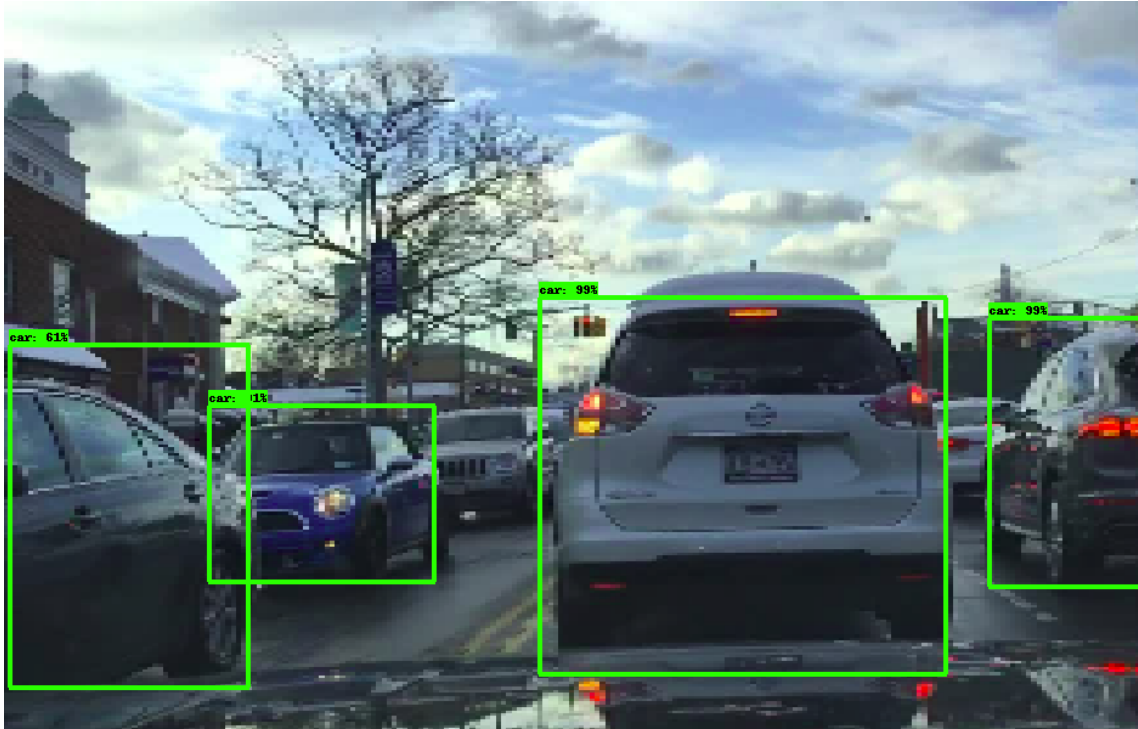


Figure 3.8: Object detector in action at 5000 steps

Setup

Setup was carried out almost exactly the same as the previous experiment. 200 total images were annotated manually using LabelImg across the three classes - car, traffic light and person. Images were selected from the same semi-urban scenes without inclement weather or poor lighting. As the dataset had now doubled in size, the number of training steps was increased. The CNN was trained for 14,000 steps, at which point the loss values appeared to not be dropping significantly.

Results

Results over the previous experiment were much improved. Over a total training time of 3hr and 50mins to reach 14,000 steps, the loss value fell to a lowest value of 1.356. Precision reached a maximum value of 0.2389 and recall reached a maximum value of 0.3402. Although these results are a

significant improvement over the initial experiments, they are still inadequate. Although the MobileNet architecture proved valuable to carry out some initial experiments and highlight the issues presented by the BDD100K dataset it is still unsuitable for this project for several reasons. The requirement to resize the images down to 300x300 pixels in order to provide training images of the same dimensions of the initial COCO training data means that the pretrained MobileNet models provided by tensorflow will never be suitable for achieving acceptable results on the larger BDD100K images. Moreover, the MobileNet architecture is simply not complex enough to perform well on a highly complex and variable dataset. As such, a more complex CNN architecture will need to be investigated in order to obtain a trained model that can reliably carry out object detection on the BDD100K dataset.

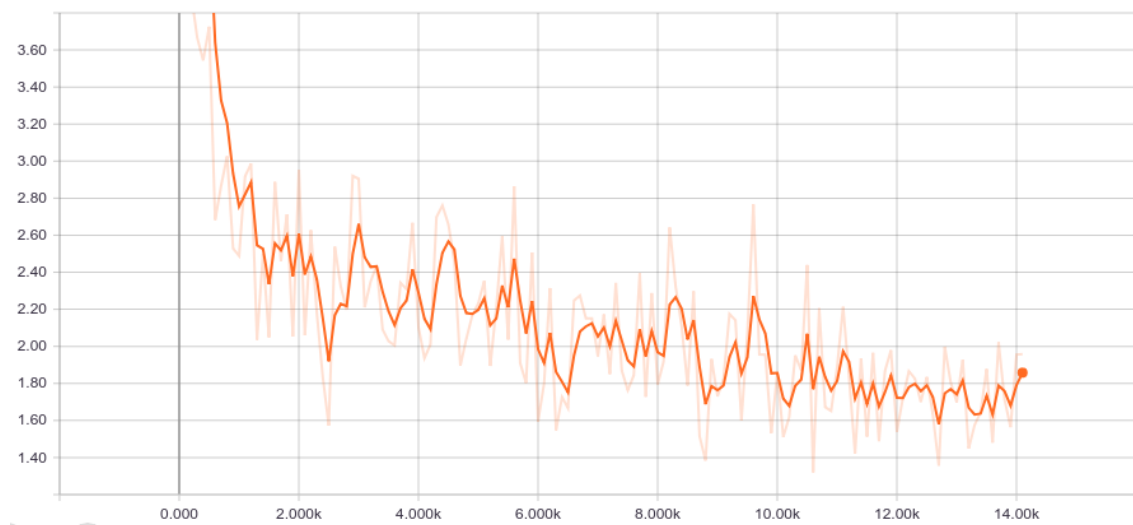


Figure 3.9: Loss values for MobileNet with 400 training images

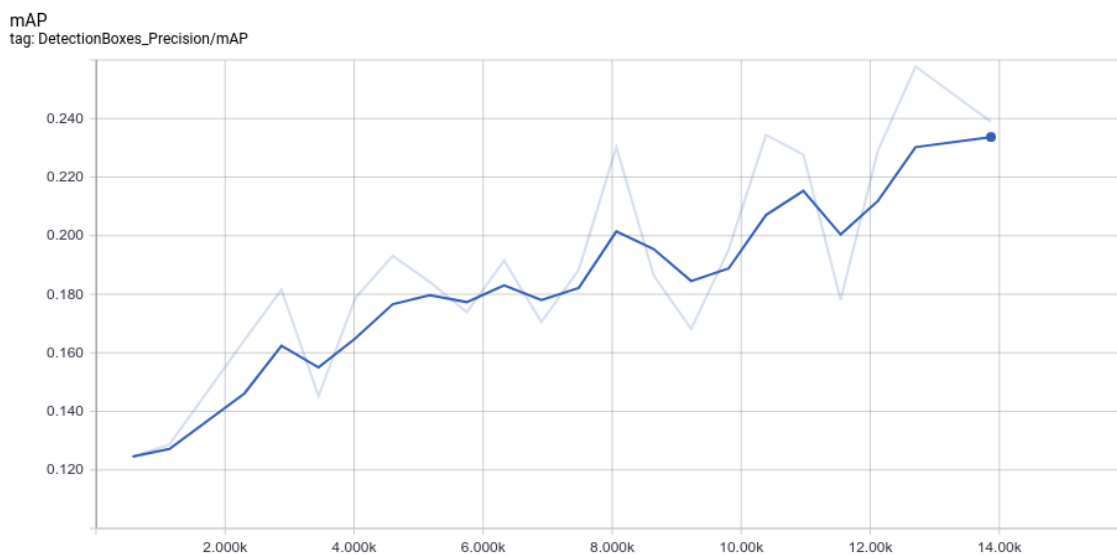


Figure 3.10: Precision values for MobileNet with 400 training images



Figure 3.11: Recall values for MobileNet with 400 training images

3.2 Faster RCNN Inception V2 Experiments

Following the relatively poor results of the MobileNet experiments, the Faster RCNN Inception V2 architecture was utilised.

3.2.1 Initial Experiment

Objectives

The objective for this experiment was to retrain the Inception V2 architecture using all of the subsampled data from the BDD100K dataset. As the Inception V2 architecture is a much more complex architecture than the MobileNet architecture, it was expected that the results obtained from this experiment would be much closer to those presented in (Yu et al., 2018).

Setup

Setup for this experiment was minimal as the training data for this experiment had already been annotated. The pretrained Inception V2 model was simply downloaded and the steps presented by (*TensorFlow Object Detection Model Training* 2018) were once again followed, this time for the new model.

Results

Results for this experiment were far better than the MobileNet experiments. Over 10,000 training steps taking a total time of 1d 3hr and 25min, loss fell to a minimum value of 0.1253. Precision reached a maximum value of 0.4392 and recall reached a maximum value of 0.2170. These metrics are graphed in Fig 3.12, Fig 3.13 and Fig 3.14. As these metrics did not seem to be improving significantly at the 10,000 steps mark, training was halted at this point.

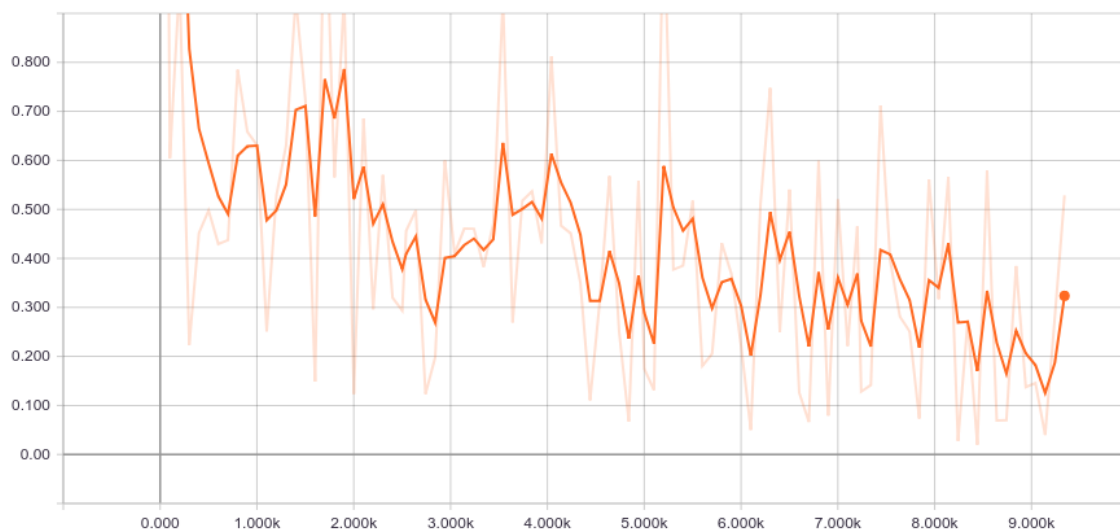


Figure 3.12: Loss values for Inception V2

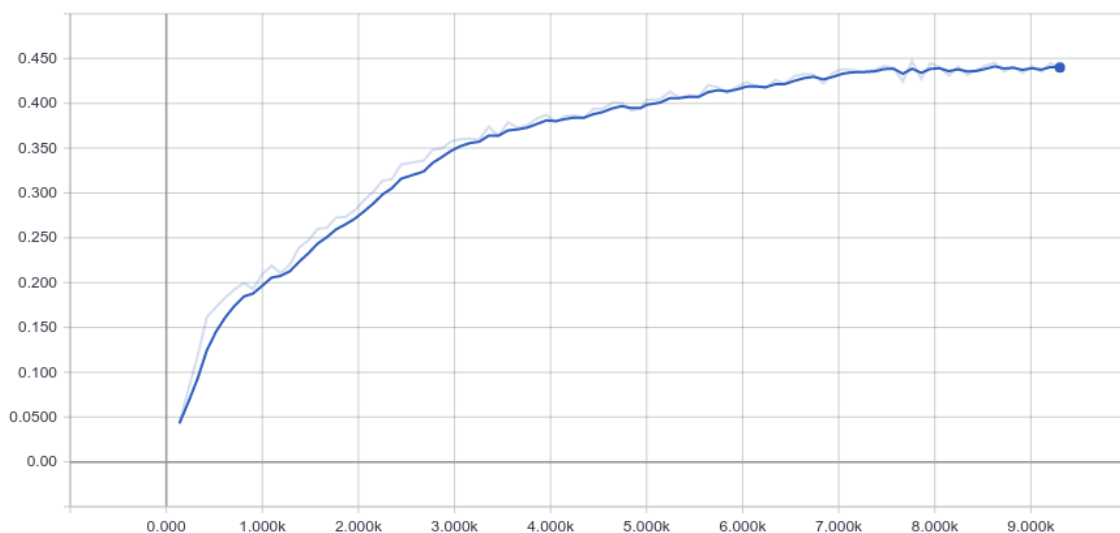


Figure 3.13: Precision values for Inception V2

3.2.2 No Pretraining Experiment

Objectives

As every experiment prior to this experiment was carried out on a pretrained model provided by tensorflow, it was decided to train a model

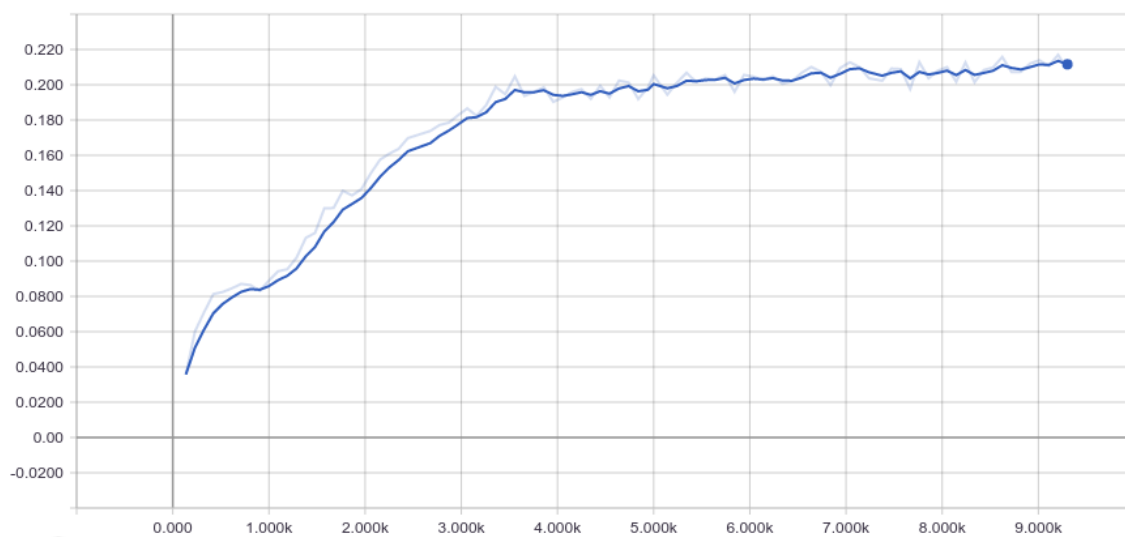


Figure 3.14: Recall values for Inception V2

using the Inception V2 architecture that had received no pretraining whatsoever. Results for this experiment were expected to be extremely poor - with a dataset consisting of a meager 400 images overfitting was almost guaranteed.

Setup

Results

As expected, results for this experiment were extremely poor over 3500 training steps totalling 5hr 35 minutes. Although Fig 3.15 may appear to be indicative of positive results, note the fact that this metric is simply the *training* loss. Training any CNN will lead to a lower training loss over time as the CNN adapts and learns the parameters of the training dataset. Observation of the validation loss in Fig 3.16, precision in Fig 3.17 and recall in Fig 3.18 clearly illustrate that this model is severely overfitting. The validation loss does not appear to be following any clear downward trend and although the precision and recall values appear to be following an upward trend, their values are extremely low. Training was halted at 3500 steps as it was clearly apparent that overfitting was taking place and further training would simply be a waste of resources.

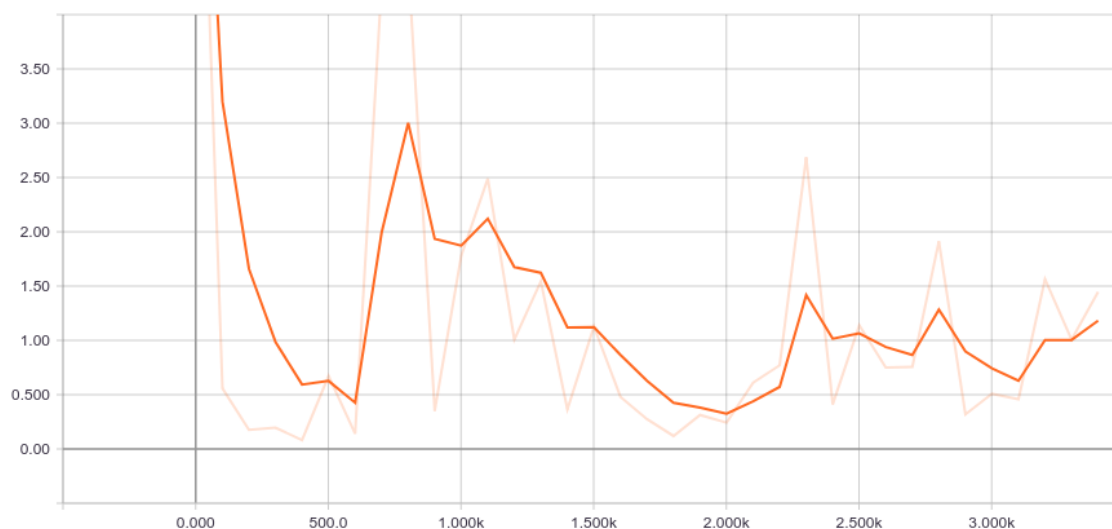


Figure 3.15: Training Loss with no Pretraining

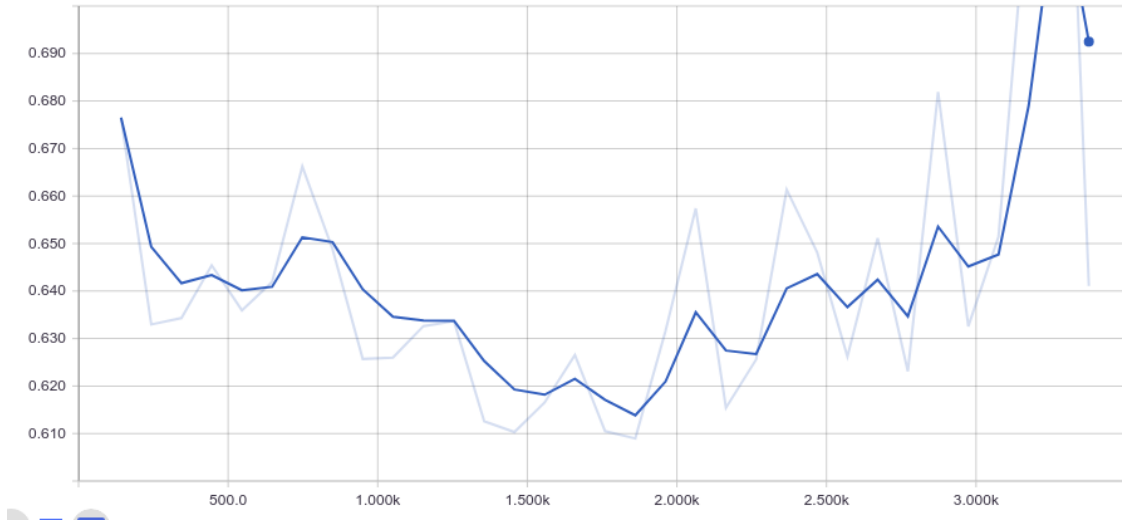


Figure 3.16: Validation Loss with no Pretraining



Figure 3.17: Precision with no Pretraining

3.3 Empirical Studies Results

The results from the Empirical Studies are shown in the below table along with the results provided in (Yu et al., 2018):

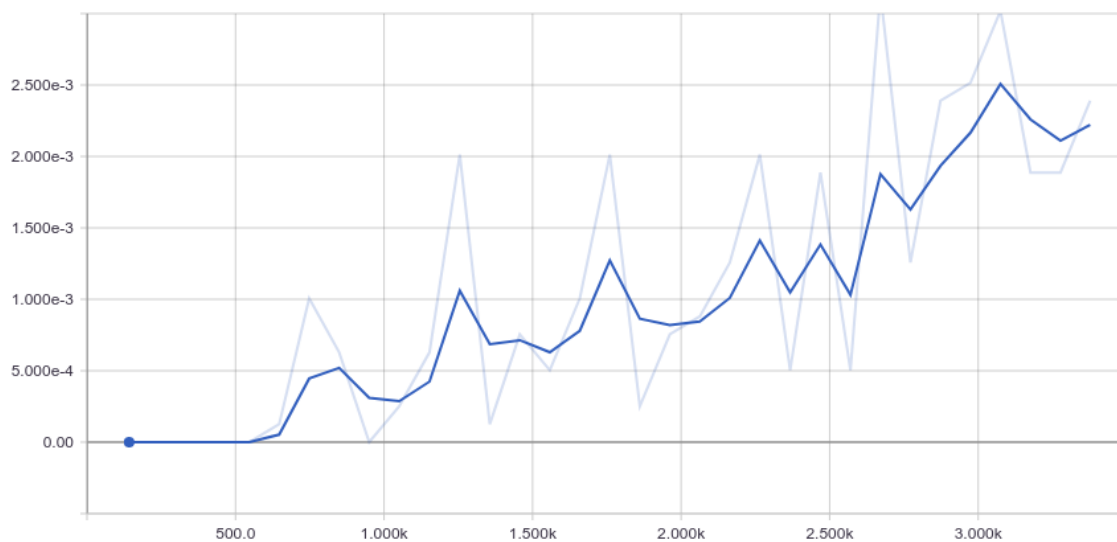


Figure 3.18: Recall with no Pretraining

| Experiment | Training Steps | Training Time | Precision | Recall |
|-----------------------------|----------------|-----------------|------------|-----------|
| Yu et al. | - | - | 0.53 | - |
| MobileNet V1 200 Images | 5000 | 1 hr 25 min | 0.1995 | tbd |
| MobileNet V1 400 Images | 14000 | 3 hr 50 min | 0.2389 | tbd |
| Inception V2 Pretrained | 10000 | 1 d 3 hr 25 min | 0.4392 | 0.2170 |
| Inception V2 Not Pretrained | 3500 | 5 hr 35 min | 1.7800 e-3 | 2.500 e-3 |

Chapter 4

Application Implementation

In order to provide an interaction with the models trained during the Empirical Studies section, a lightweight Flask application was created. This application was developed to allow users to upload images to a simple public-facing website and observe object detection taking place on the uploaded images.

4.1 AWS Instance

As the AWS Deep Learning AMI that was used to train the models was still in use, it was decided to simply use this virtual machine for the deployment of the application. All AWS EC2 instances are provided with a public-facing IP address by default, so the only work required on the instance was editing its security roles to allow HTTP requests to the instance for all IP addresses.

4.2 Implementation

The first step for the implementation of this prototype was providing some way for the trained model to annotate the provided images. Tutorial code provided by tensorflow to interact with their pre-trained models (tensorflow, 2019) was taken and adapted for use in the context of a Flask application. Modifications to the code were slight - the program was changed to accept files passed as a parameter instead of loading in images from a predefined location on the machine, and images were then converted

to bytes and return after annotation had taken place. Next a simple flat html page was created to be served by the Flask application that contained for image upload. Lastly a simple Flask script was written to accept images sent via POST request, send these images to the object detection code and finally display the newly annotated image. Fig 4.1 shows the design of the simple interface through which users can upload their images.

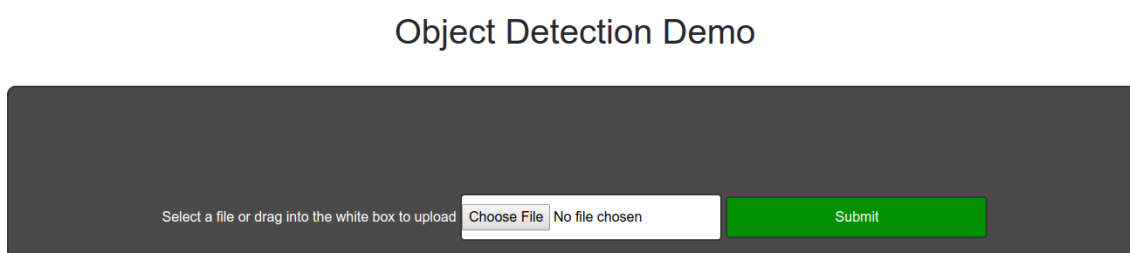


Figure 4.1: Simple Webpage to allow Image Upload

4.3 Code Snippets

Some interesting snippets of code from the application implementation have been provided below:


```

# Import method from the object detection script
from object_detection.ObjectDetector import detect_objects
import io

from flask import Flask, render_template, request

from PIL import Image
from flask import send_file

app = Flask(__name__)

# Calls the index.html file present in the "templates" folder
@app.route("/")
def index():
    return render_template('index.html')

# Accept jpg files via POST request, send to be annotated and then display result on-page
@app.route("/", methods=['POST'])
def upload():
    if request.method == 'POST':
        file = Image.open(request.files['file'].stream)
        img = detect_objects(file)
        return send_file(io.BytesIO(img), attachment_filename='image.jpg', mimetype='image/jpg')

# Accept requests through port 5000, this has been opened for HTTP requests for the AWS instance
if __name__ == "__main__":
    app.run(host="0.0.0.0", port=5000)

```

Figure 4.2: Flask Application Code, calls Annotation Code

```

def detect_objects(image):
    # the array based representation of the image will be used later in order to prepare the
    # result image with boxes and labels on it.
    image_np = load_image_into_numpy_array(image)
    # Expand dimensions since the model expects images to have shape: [1, None, None, 3]
    image_np_expanded = np.expand_dims(image_np, axis=0)
    # Actual detection.
    output_dict = run_inference_for_single_image(image_np, detection_graph)
    # Visualization of the results of a detection.
    vis_util.visualize_boxes_and_labels_on_image_array(
        image_np,
        output_dict['detection_boxes'],
        output_dict['detection_classes'],
        output_dict['detection_scores'],
        category_index,
        instance_masks=output_dict.get('detection_masks'),
        use_normalized_coordinates=True,
        line_thickness=8)

    img = cv2.cvtColor(image_np, cv2.COLOR_RGB2BGR)
    img = cv2.imencode('.jpg', img)[1]
    return(img.tobytes())

```

Figure 4.3: Method called by Flask Application, accepts and returns image

Chapter 5

Final Conclusion and Discussion

5.1 Summary

The primary purpose behind this project has been an investigation into how CNN's perform object detection and the issues encountered during this process. Once the research area for this project had been defined, a comprehensive literature review of the subject material was undertaken. One of the main objectives for the literature review was to develop and document an expansive knowledge of CNN's. A working knowledge of using tensorflow was also developed through a Udemy tutorial (*Complete Guide to Tensorflow for Deep Learning with Python* 2018). Once this knowledge had been gained, several experiments were undertaken in an attempt to replicate the results published in (Yu et al., 2018). Although the full BDD100K dataset quickly proved far too complex to use in its entirety, a subsampled version of the dataset was sufficient for training instead. Although the results provided by Yu et al. were not replicated perfectly, results that were within an acceptable margin of error were obtained. Unfortunately financial constraints related to using AWS services did not allow for further experimentation, however I am satisfied with the results obtained throughout the experimentation phase of this project. Finally, a lightweight Flask application was built and deployed in order to allow users to interact with the models trained during the experimentation phase.

5.2 Reflections

5.3 Future Work

Bibliography

- Arbelaez, Pablo et al. (2011). “Contour detection and hierarchical image segmentation”. In: *IEEE transactions on pattern analysis and machine intelligence* 33.5, pp. 898–916.
- Canny, John (1986). “A computational approach to edge detection”. In: *IEEE Transactions on pattern analysis and machine intelligence* 6, pp. 679–698.
- Complete Guide to Tensorflow for Deep Learning with Python* (2018). URL: <https://www.udemy.com/complete-guide-to-tensorflow-for-deep-learning-with-python/> (visited on 13/12/2018).
- Convert the Berkeley Deepdrive dataset to a TFRecord file* (2018). URL: https://github.com/meyerjo/deepdrive_dataset_tfrecord (visited on 01/03/2018).
- Domingos, Pedro (2012). “A few useful things to know about machine learning”. In: *Communications of the ACM* 55.10, pp. 78–87.
- Géron, Aurélien (2017). *Hands-on machine learning with Scikit-Learn and TensorFlow: concepts, tools, and techniques to build intelligent systems*. ” O’Reilly Media, Inc.”.
- Howard, Andrew G et al. (2017). “Mobilenets: Efficient convolutional neural networks for mobile vision applications”. In: *arXiv preprint arXiv:1704.04861*.
- Hubel, David and Torsten Wiesel (2011). *Hubel & Wiesel - Cortical Neuron - V1*. Youtube. URL: <https://www.youtube.com/watch?v=8VdFf3egwfg>.
- Kotsiantis, Sotiris B, I Zaharakis, and P Pintelas (2007). “Supervised machine learning: A review of classification techniques”. In: *Emerging artificial intelligence applications in computer engineering* 160, pp. 3–24.
- Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E Hinton (2012). “Imagenet classification with deep convolutional neural networks”. In: *Advances in neural information processing systems*, pp. 1097–1105.

- LeCun, Yann, Yoshua Bengio, et al. (1995). “Convolutional networks for images, speech, and time series”. In: *The handbook of brain theory and neural networks* 3361.10, p. 1995.
- Lin, Tsung-Yi et al. (2014). “Microsoft coco: Common objects in context”. In: *European conference on computer vision*. Springer, pp. 740–755.
- O’Shea, Keiron and Ryan Nash (2015). “An introduction to convolutional neural networks”. In: *arXiv preprint arXiv:1511.08458*.
- Rogers, Simon and Mark Girolami (2016). *A first course in machine learning*. CRC Press. Chap. 2.
- Rumelhart, David E, Geoffrey E Hinton, and Ronald J Williams (1985). *Learning internal representations by error propagation*. Tech. rep. California Univ San Diego La Jolla Inst for Cognitive Science.
- Singh, Vaibhav Kant and Shweta Pandey (2016). “Minimum configuration MLP for solving XOR problem”. In: *2016 3rd International Conference on Computing for Sustainable Global Development (INDIACom)*, pp. 174–179.
- Srivastava, Nitish et al. (2014). “Dropout: a simple way to prevent neural networks from overfitting”. In: *The Journal of Machine Learning Research* 15.1, pp. 1929–1958.
- Szegedy, Christian et al. (2016). “Rethinking the inception architecture for computer vision”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2818–2826.
- tensorflow (2019). *Object Detection Demo*. tensorflow. URL: https://github.com/tensorflow/models/blob/master/research/object_detection/object_detection_tutorial.ipynb.
- TensorFlow Object Detection Model Training* (2018). URL: <https://gist.github.com/douglasrizzo/c70e186678f126f1b9005ca83d8bd2ce> (visited on 01/03/2018).
- Yu, Fisher et al. (2018). “BDD100K: A diverse driving video database with scalable annotation tooling”. In: *arXiv preprint arXiv:1805.04687*.
- Ziou, Djemel, Salvatore Tabbone, et al. (1998). “Edge detection techniques-an overview”. In: *Pattern Recognition and Image Analysis C/C of Raspoznavaniye Obrazov I Analiz Izobrazhenii* 8, pp. 537–559.