

Identification and Classification in Food Images

Tom Barrett

March 27, 2018

Abstract

Health conditions in the modern age are progressively getting worse with a high percentage of the population being obese. In order to combat this problem, a dietary assessment smartphone application would be invaluable. The task of recording ones food intake could be quite tedious and therefore, utilizing computer vision to automatically classify and calculate nutritional value of a users food image could help to make the process more practical for use. One approach that could be attempted, is by using Convolutional Neural Networks (CNNs) for the classification of food images. This is due to the recent high success in using CNNs for image classification. Retraining of the Inception-V3 model architecture, using tensorflow, the Food-101 dataset and the ImageNet dataset was carried out, resulting in a classification model. A Top 1 accurcay of 66.6% and a Top 5 accuracy of 85.96% was achieved using a dataset of 108 classes. It was found that CNNs are very successful in classifying images with over a 100 classes and that network fine tuning could result in even better results.

Acknowledgements

Contents

1	Introduction	1
1.1	Overview	1
1.2	Objectives	3
1.3	Methodology	5
1.4	Overview of Report	6
1.5	Motivation	7
2	Background	9
2.1	Introduction to Machine Learning	9
2.2	Neural Computing	10
2.3	Convolutional Neural Networks Overview	15
2.4	Convolution Neural Networks Extended	18
2.5	Support Vector Machine	20
2.6	Evaluating the Output	21
2.7	Dietary Assessment using Computer Vision	23
2.8	APIs and Libraries	38
2.9	Conclusion	39

3	Introduction to Using Tensorflow	40
3.1	Template for Experiments	40
3.2	Udemy Tutorial	40
3.3	Udemy Tutorial 2	44
3.4	Using the Food 101 dataset	48
4	Training Using the Inception-V3 Model Architecture	52
4.1	Retrain ImageNet Inception V3 Model	52
4.2	Retrain with Extended Dataset	58
4.3	Retrain with Parameter Tuning	61
4.4	MobileNet	68
4.5	Food 101 subset	70
5	Analysing the Trained Model	72
5.1	Sliding Window	72
5.2	Recursive Refinement	78
5.3	Impact of Background	82
5.4	Alternative Test Image	84
5.5	Analysing Results of Recursive Refinement Further	85
5.6	Scaling Down Images	91
5.7	Effect of Colour	94
5.8	Visualising Images Through the Network	98
6	Prototype Application	99
6.1	Requirements	99
6.2	Design	99
6.3	Implementation	103

6.4	Testing	110
6.5	Backend	113
6.6	Resources Utilised	115
7	Discussion and Conclusion	116
A	Appendix	122
A.1	Image Segmentation	123

List of Figures

1.1	Pre-segmented Image	2
1.2	Segmented Image	2
2.1	Perceptron	11
2.2	Linearly Separable, adapted from Mitchell, (1997)	12
2.3	Multi Layer Perceptron	13
2.4	Gradient Descent	15
2.5	CNN Architecture	16
2.6	CNN Overview	17
2.7	Support Vector Machine Applicability	20
4.1	Inception module	55
4.2	Inception V3 Architecture	55
4.3	Pizza	58
4.4	Pizza not classified correctly by the model	59
4.5	Banana	61
4.6	Graph of accuracy of the test dataset during training	65
4.7	Graph of accuracy of the validation dataset during training	65
4.8	Comparison of accuracy	66
4.9	Comparison of accuracy	66

5.1	Bowl of fruit	73
5.2	Grid based window	74
5.3	Row based window	74
5.4	Column Based Window	74
5.5	Fruit with Color Overlay	75
5.6	Recursive refinement 1	79
5.7	Recursive refinement 2	80
5.8	Recursive refinement 3	80
5.9	Bowl of fruit with background removed	81
5.10	Alternative Bowl of fruit	84
5.11	Fruit image taken on a mobile phone	87
5.12	Image after sliding window - 13 classes	88
5.13	Image after sliding window - 108 classes	89
5.14	Banana Image Pre-Resolution	91
5.15	Apple Pie Image Pre-Resolution	92
5.16	Pizza Image Pre-Resolution	92
5.17	Banana Image Greyscale	94
5.18	Apple Pie Image Greyscale	95
5.19	Pizza Image Greyscale	95
6.1	Landing Activity	101
6.2	Image Submission Activity	101
6.3	Classification Activity	101
6.4	Food Logs Activity	101
6.5	Package Diagram	102
6.6	Automated Unit Tests	110

A.1	Detection Average Precision Donahue, (2017)	125
A.2	Mean Average Precision Donahue, (2017)	126
A.3	Distribution of top-ranked false positives Donahue, (2017) . .	126
A.4	Segmentation Mean Accuracy Donahue, (2017)	127
A.5	Per-category segmentation accuracy Donahue, (2017)	127
A.6	Per-class segmentation accuracy Donahue, (2017)	127

List of Tables

2.1	DeepFood Results	25
2.2	Summary of results in CNN based methods	26
2.3	Summary of accuracy in dietary assessment methods	33
2.4	Results from Region Based CNN Research	38
3.1	Experiment Template	41
4.1	Comparison of parameters	66
4.2	Accuracy of other studies	71
5.1	Grid Based Sliding Window Results	77
5.2	Row Based Sliding Window Results	77
5.3	Column Based Sliding Window Results	77
5.4	Comparison of fruit image sliding window results with and without background	82
5.5	Comparison of fruit bowl images	83
5.6	Results of recursive refinement segment classifications using two models	90
5.7	Results of Down Scaled Images	93
5.8	Effect of Colour	96
A.1	FCN Resulys Shelhamer, Long, and Darrell, (2017)	124

A.2	Results	124
A.3	Project Plan	127

Chapter 1

Introduction

The chapter outlines an introduction to the subject area of this Final year Project (FYP), the objectives of this project, the methodology used to complete this project, an overview of the report and an explanation of the motivation behind this FYP.

1.1 Overview

This project explores identification and classification of food images for use in a calorie measurement android application. Food calorie consumption is a huge problem in the modern world. Over 25% of the population in Ireland is obese and this figure is likely to rise over the coming years. A mobile application that could help keep track of a user's calorie intake by taking pictures of their meals would be a great help. The area of computer vision is a very difficult topic to address as it is a very hard task for computers to undertake. We, as humans, take vision for granted as we can soon see, from the study of computer vision, that there are many difficult steps that have to be made for full identification and classification of an image.

When looking into calorie measurement using an image, there are three questions that have to be answered:

- Where are the Regions of Interest (ROI) in this food image?



Figure 1.1: Pre-segmented Image



Figure 1.2: Segmented Image

- What food types are in these ROI's?
- What is the portion size of each food type?

In this project, the main focus will be on the first two questions, 'Where are the Regions of Interest (ROI) in this food image?' and 'What food types are in these ROI's'. The first step is normally achieved through image segmentation. Image segmentation is the process in which you divide an image into multiple segments as per Figure 1.2 and 1.1.

Many researchers in various computer vision labs have attempted to solve this problem using different methodologies. There has been promising results from some papers but these are mostly under highly constrained circumstances. When mixed foods are introduced to the problem, many of the methods fail.

Convolutional Neural Networks (CNN) have had very promising results in the field of image classification in the recent years and have been shown to work best in the area of face detection.

Extensive research on many different methods of image segmentation and classification has been carried out but it seems that CNN's have had the best results for multiple objects in one image and therefore, this method will be applied to many foods in an image. This is because we will rarely want to classify an image that only has one food item in it. Therefore, the one shot approach may not be the most successful by which we build a classifier that takes in an image and gives back the most likely food in that image. In contrast to this, there could be an element of sub-sampling of the image to rectify this.

The system proposed to solve the problem statement would be able to integrate with an Android mobile phone application. The idea is, that when a user is about to eat their meal, they can simply take a picture of their meal for computation. From here, the application would take an image, and send it to a server to be analysed. On the server side the image would be classified using a Tensorflow model. Once this is done, the size of each food type would be measured and through this an overall calorie count would be displayed for the user. This could be logged for user metrics. The full system may not be possible to implement due to time constraints so therefore, classification will be the furthest step looked in to.

1.2 Objectives

Primary Objectives

Use Convolutional Neural Networks for Food Image Classification

There has been a large paradigm shift in food image recognition in recent years to using convolutional neural networks. This paradigm will be used to answer the problem statement.

Tune the CNN, replicating previous work

There are many different approaches to food identification and classification, many of which we will see in Chapter 2. An approach will be selected that has shown promising results in the past and replication of these results will be attempted. In addition to this, there are many different network architectures and parameters that can be adjusted when building CNNs and these will be tuned in order to get the best outcome possible.

Develop a mobile application to leverage the CNN

Another objective for this project is to develop an application that can be used for dietary assessment. This application would be able to take an image and then send the image to a server to identify and classify the foods within the image. Size estimation will not be explored for this project. Alternatively, a previously developed application will be used to demonstrate the CNN developed.

Secondary Objectives

Understanding of Convolutional Neural Networks

In the project, Convolutional Neural Networks (CNN's) will be used for object identification in Food Images. A machine learning library will be used for this due to time constraints but it is a key objective to develop a deep understanding of CNN's as they are quite pivotal in the current computer vision industry and bio-inspired systems are very interesting.

Learn about different image identification and classification techniques

Although, CNN's will be used for implementation, other methods of identification and classification will not be ignored. It is very important to learn about other methods as different methods are better suited for some situations and

it would be best to know about these methods due to the inevitability of their use.

1.3 Methodology

The following methodology were adapter for this project:

Define the research question

The first step to this project was to define the research question. The general area of a computer vision supported nutritional assessment application was known at conception but the scope of this is too broad for an FYP. Therefore it was decided that the research question would be to look at the food identification and classification aspect.

Literature Review

Once the research question was defined, finding related work was the next milestone. There are many attempts at Food Image Classification and these were not difficult to find, using Google Scholar, but many of these papers glossed over the segmentation aspect and relied on third parties for this step. Because of this, quite a few references to different papers had to be followed that focused solely on image segmentation and classification.

Explore different image identification methods

Various image identification methods were collected from the literature review that was carried out, so there were many options to evaluate. Convolutional Neural Networks were the clear choice due to recent popularity and successful results, so more traditional methods of identification using colour and texture was not so strenuously explored.

Select an image identification and classification method

After the various methods to identify and classify food images had been collated in the literature review, one of these had to be selected.

Research technologies and develop skills in these technologies

As Convolutional Neural Networks (CNN) are used in this project, many resources have been leveraged to enrich understanding of the process. Tensorflow is the main resource utilised in creating a CNN so on-line tutorials for this technology were greatly beneficial. A Deep Learning Course on Udacity was also used to enhance understanding and skills.

Build a prototype of the application

A prototype application had to built for demonstration for this project.

Compare and analyse results to other implementations

Once a model had been successfully trained, comparison to previous work was carried out.

1.4 Overview of Report

This report is broken down into various chapters. The introduction chapter is to give an overview of what this project is about, how the project will be approached and why it is being carried out. Following this, some information on the background of the subject of this Final Year Project (FYP) will be outlined. Once the chapters of introduction and background have been completed, a chapter outlining the learning of Tensorflow will be explored. Tensorflow is a machine learning library that is used to develop CNNs. In this chapter, the

goal is to demonstrate the activities undertaken to learn about the technologies used and show tutorials that have been completed to aid with this learning. In chapter titled 'Training Using the Inception-V3 Model Architecture' , the purpose is to outline experiments that have been carried out in relation to training Tensorflow models. A Tensorflow model is the term used to describe the neural network produced by Tensorflow. The Inception-V3 architecture is proven to be a very successful architecture for classifying images and will be explained in detail in said chapter. Once the model has been trained, the following chapter 'Analysing the Trained Model' will consist of experiments that analyse and test the models trained. A prototype smart phone application was developed to demonstrate the efficacy of the system and an overview of the process will be explained in the chapter titled 'Prototype Application.' The final chapter's purpose is to provide a discussion on the results obtained and offer a conclusion to the findings.

1.5 Motivation

I find the topic of Computer Vision a very interesting one. It excites me, to be able to 'teach' a machine how to see as we do. For this reason, I really wanted to learn about Neural Networks and this was a large motivator for this project.

Once I had a topic that I wanted to research, I needed a focus or problem statement for this research. I find that it is much more rewarding to work on something that positively impacts both myself and other people so I decided that I wanted to research something that fit this requirement.

Food calorie consumption is a very big problem in the modern world. Over 25 % of the population in Ireland are obese. A mobile application that could help keep track of a user's calorie intake by taking a picture of their meals would be a big help to combat this problem. This problem statement works very well for me because of it's application use and because of its complexity. Identifying and recognising food is much more difficult than say recognising

faces as it has no uniform shape. Therefore, this problem would also be very beneficial to developing skills in the computer vision area.

I would like to develop these real world skills so that I can partake in Computer Vision projects in industry or to do further research in academia. This is because machine learning has really taken off in the last few years and is used by many in industry. While machine learning as a whole has become very popular, computer vision is probably the most prominent that has come out of it. Face detection, for one, is being researched extensively for use in personalised advertising and also secure access to devices and systems.

Chapter 2

Background

This chapter outlines background information on the subject of nutritional assessment using computer vision. An overview of neural networks is explored, followed by a literature review of the subject area. Analysis into evaluating the efficacy of CNNs is covered and finally a conclusion to the information gained.

2.1 Introduction to Machine Learning

In Mitchell, (1997), machine learning is defined as "the question of how to construct computer programs that automatically improve with experience". Machine learning has blossomed in recent years with applications across multiple domains using vastly different paradigms and technologies. Some of the different approaches used in machine learning are: Artificial Neural Networks, Genetic Algorithms, Decision Tree Learning and Bayesian Learning Mitchell, (*ibid.*).

There are many ways in which machine learning can be used in the modern world, many of which are being utilised to great affect. Some of these applications are image recognition, natural language processing, medical diagnosis and many more. There may be fear that machine learning will start to take away many jobs from humans but this may not be the case. There are many

practical uses such as security and safety that could be leveraged such as face detection for security in airports or autonomous cars.

One of the most exciting avenues in machine learning is computer vision. Computer vision is the process of extracting high-dimensional data from an image to produce useful information, which in terms of classification usually results in labeling. It can be used in many areas to improve our lives. As mentioned earlier, autonomous cars are only possible when a machine can determine what objects are around it. Computer Vision can allow a machine to recognise skin diseases in an image. The applications are nearly limitless.

2.2 Neural Computing

The main area of my focus for this project is in Artificial Neural Networks (ANN). This is because extensive research has been carried out into Convolutional Neural Networks (CNNs) which are based on ANNs.

Artificial Neural Networks

An Artificial Neural Network is a bio-inspired system that is used to model the human brain in how it learns from experience. The ANN uses this model to build a very complex web of connected units called artificial neurons. These neurons are connected by certain weights which determines the processing capacity of the network and these weights are created by learning a dataset.(reference Malachy) An ANN has a set of inputs that take in a value, sometimes from network outputs and produce a single result or classification. While an ANN is bio-inspired from the human brain, there are many elements of the brain that are not present in ANN and many new elements in ANN that are not modeled from the human brain.

Before Convolution Neural Networks can be explored, which are vital to image processing, the perceptron learning algorithm, the multi layer perceptron, backpropogation must be analysed.

Perceptron Learning - Artificial Neuron

In our Artificial Neural Network a Perceptron is an Artificial Neuron. It is called an Artificial Neuron because it is a bio-inspired neuron which models a neuron in the human brain in terms of inputs and output.

In Perceptron learning, we can take two inputs which are put towards an activation function with a bias attached as seen in Figure 2.1. These inputs are multiplied by the weights that connect the input to the activation function and depending on the result, the activation function may fire an output. These inputs are either 1 or -1.

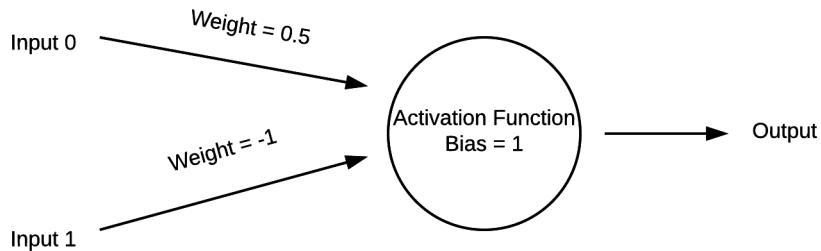


Figure 2.1: Perceptron

The Perceptron Training Rule is the means by which weights are selected to produce the correct output during training. As in Mitchell, (ibid.), a common way to train a perceptron is to start with random weights and change them during training as per the training rule. This rule follows the formula below:

$$w_i \leftarrow w_i + \Delta w_i$$

where x_i is the input and the equation below is valid:

$$\Delta w_i = n(t - o)x_i$$

In the above formula, "t" is the target output for the current training example, o is the output generated by the perceptron, and n is the positive constant called

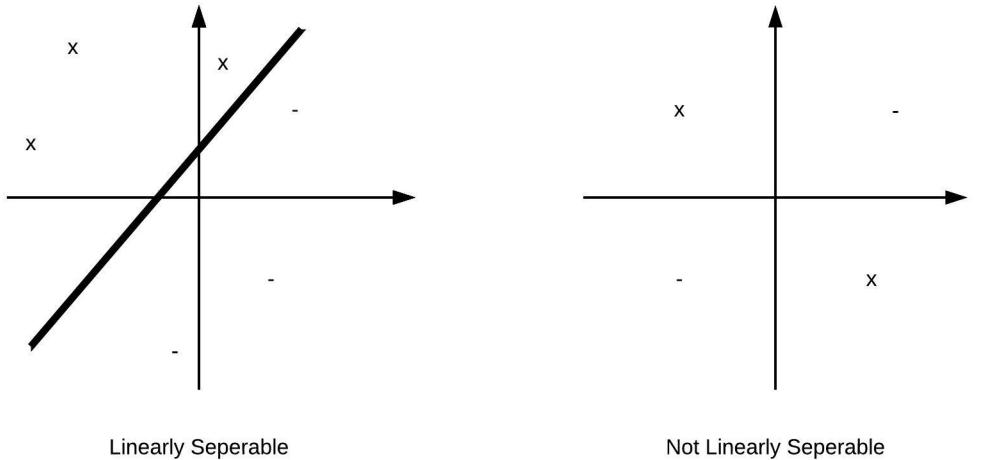


Figure 2.2: Linearly Separable, adapted from Mitchell, (1997)

the learning rate” Mitchell, (1997). This Perceptron Training Rule assumes that there are two sets of instances, a positive and negative set, and that they are linearly separable, as in demonstrated in Figure 2.2.

A perception is trained using supervised learning. When the perceptron classifies a results, it is told if it is correct or not. If the result is incorrect, weights are changed in value so that this error can be reduced Luger, (2005).

There is one major problem with perceptron learning and that is, it can't solve a problem if there is not a clear linear separation between the classes. There is a way in which we can attempt to solve this, through the delta rule. The delta rule utilizes gradient descent to find the best weight for the training samples Mitchell, (1997). We will discuss gradient descent in the next section.

Multi Layered Perceptron

Multi Layer Perceptrons (MLP) are made up of multiple layers of perceptrons connected together and are used to combat non-linearly separable classes. Firstly, we have an input layer, followed by one or more hidden layers and then finally an output layer. Any Neural Network with more than three hid-

den layers is categorised as a deep neural network.

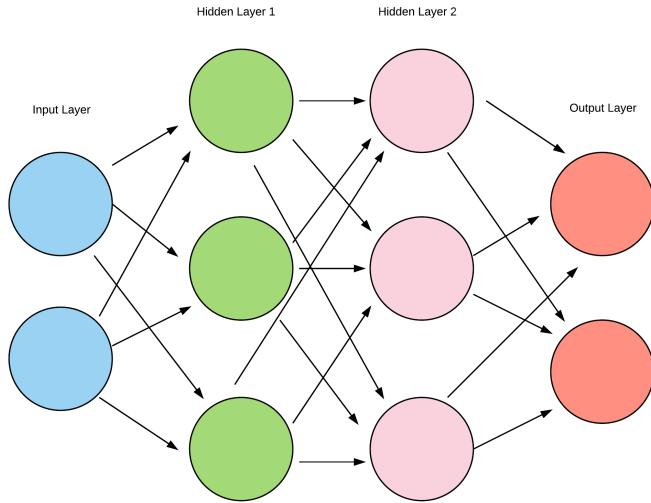


Figure 2.3: Multi Layer Perceptron

The input layer of your network consists of the data you feed into the network in order to classify it. The input layer passes this data to a hidden layer whose purpose is to transform this data into something that the output layer can understand. This transformation results in a linearly separable space that can be classified. The output layer normally consists of a class prediction.

Multi Layer Perceptrons are a class of feed forward Artificial Neural Networks. These means that the output of each perceptron feeds into an input in the next layer of the network, example seen in Figure 2.3.

During training, using backpropagation, for each step, the output of every neuron is calculated in each layer and then passed to the next layer. Then the error of the output is calculated and the network calculates how each neuron

in the last hidden layer effected the error. This is continued back through all the layers until the input layer Aurélien, (2017). The weights are then altered to try and reduce this error.

There is one large problem with MLP's and this is why Convolutional Neural Networks (CNN) were created. If you attempting to classify images with a MLP then each pixel in that image would have to be a separate input. This creates a massive amount of neurons through all the layers and this isn't feasible. CNN's solve this problem which we will discuss later.

Gradient Descent

Gradient Descent is an algorithm used to find the optimal weights to produce the smallest prediction error. It is used to overcome problems of non linearly separable classes. Gradient descent search selects a random weight value and then modifies it gradually to minimize the error. "At each step, the weight vector is altered in the direction that produces the steepest descent along the surface" Mitchell, (1997). This step is iterated until the lowest value is met.

There is an error function used for the perecptron which finds the lowest error for that neuron, but it can't be used here because, since we have many neurons, there could be an error in multiple neurons. Gradient Descent is mathematically based on the derivitive of a function. The gradient of a function can be calculate by differentiating it. As the weights are what is being controlled, "they are what we differentiate in respect to" Marsland, (2015). The negative gradient of this function is followed to find the lowest possible point, hence the name gradient descent Marsland, (ibid.).

One problem with Gradient Descent is that if we look at Figure 2.4, we may never get to the optimal point, point B. This is because we will find point A without too many problems but when the weights change we will get too high a slope of error and therefore will never reach point B.

Another variation of Gradient Descent is Stochastic Gradient Descent (SGD). SGD is different because it updates "weights incrementally, following the calculation of the error of each individual example" Mitchell, (1997).

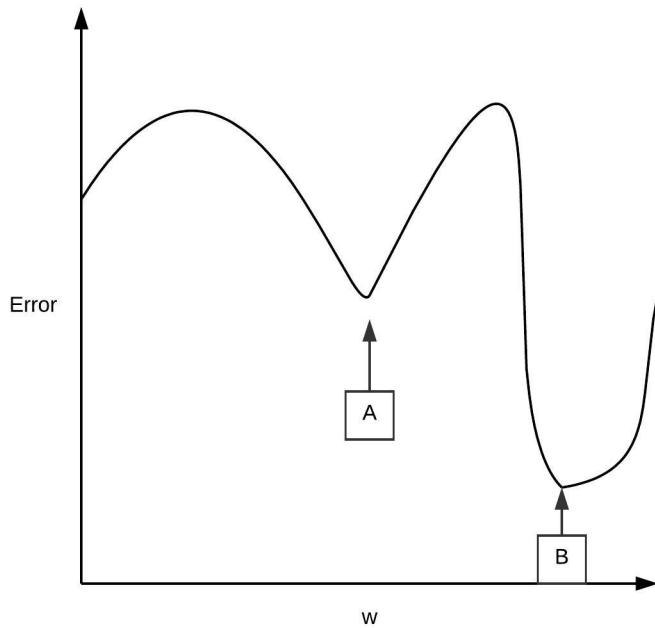


Figure 2.4: Gradient Descent

Backpropagation

”The Backpropagation algorithm learns the weights of a multilayer network, given a network with a fixed set of units and interconnections” Mitchell, (ibid.). Backpropagation attempts to minimise the mean squared error between the target output and the output of a network.

Backpropagation works by starting at the output layer of the network and going back through previous hidden layers, updating weights as it goes ie. it propagates back through the network, updating the weights to try and reduce the error.

2.3 Convolutional Neural Networks Overview

Convolutional Neural Networks (CNN’s) are essentially a Multi Layered Perceptron with a special structure. CNN’s have one major difference from a

MLP, they have extra layers of convolution and pooling. The architecture of a convolution network can be seen in Figure 2.5.

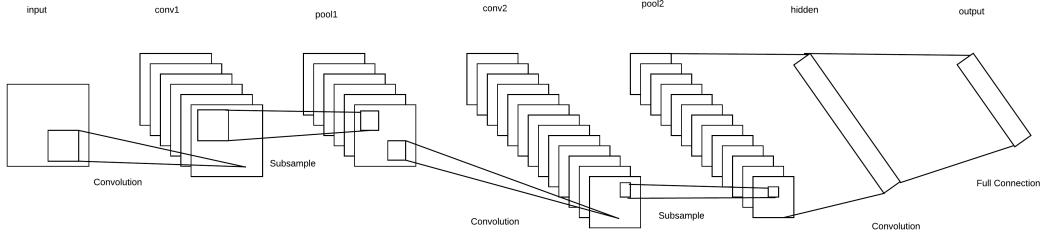


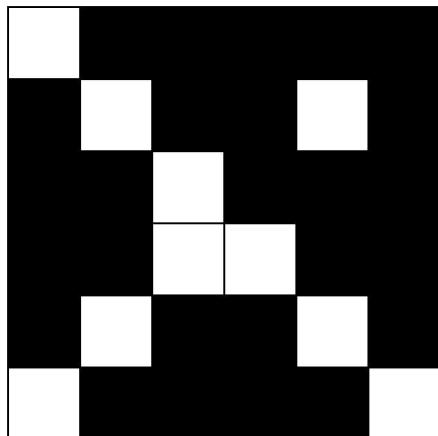
Figure 2.5: CNN Architecture

Figure 2.6 (a) shows an image that we want to compare against Figure 2.6 (b). For humans, it is quite easy to determine that these images are very similar but for a computer this task is surprisingly difficult.

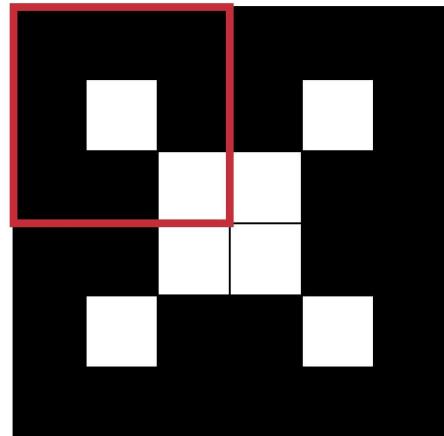
So what a CNN does, to combat this problem, is to take a small feature from Figure 2.6 (a) and compare it to a subsection of Figure 2.6 (b). The CNN multiplies the feature and a section of Figure 2.6 (c), adds up the results and divides by 9 (The number of pixels in the feature). This then gives a decimal value of how likely it is that the feature is in the part of the image, as seen in Figure 2.6 (d). This is called filtering. The Convolutional layer is composed of carrying out this filtering for every single possible location in Figure 2.6 (a).

Next is the Pooling Layer, this layer takes the convoluted layer output (you can use Figure 2.6 (d) as reference) and from a user defined size ie. 2x2, gets either the highest decimal value (Max pooling), the average value (Mean pooling) or the lowest decimal value (Min pooling) and records that as the new value for the section. This is then applied to the entire image. As we can see in Figure 2.6 (e) we now have a much smaller image stack in which to classify, thus making the computation easier. Pooling also makes features of the image clearer which helps to result in more accurate classification.

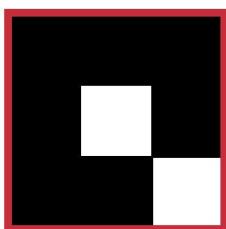
In between the Convolution and Pooling layer, there is sometimes a Normalization layer. This Normalization layer creates Rectified Linear Units (ReLU's). In other words, if we take Figure 2.6 (d), it changes all minus values to zero.



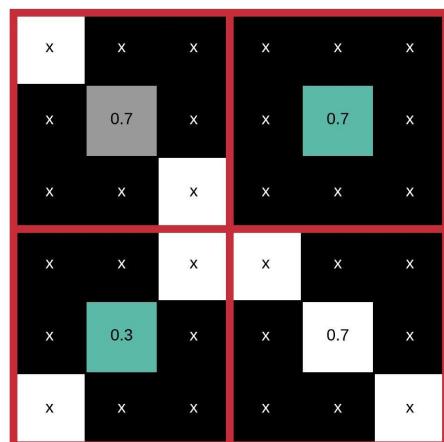
(a) Image to Classify



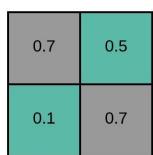
(b) Image to Compare



(c) Image Feature to Search



(d) Convolved Image



(e) Pooled Image

Figure 2.6: CNN Overview

There are some problems with CNN's however. One of the main problems is that you need a very large dataset in order to produce an accurate model and the training can be very time consuming without a GPU.

2.4 Convolution Neural Networks Extended

A basic overview of the operations of Convolutional Neural Networks has been explained in the previous section. Some further points are explored below.

History

Convolutional Neural Networks are a bio-inspired approach to solving intensive tasks. They have been applied to image recognition, along with natural language processing and voice recognition, with image recognition being explored since the 1980's Aurélien, (2017). CNN's have become increasingly successful in recent years due to increased computational power in the hardware and availability in datasets. CNN's have been recently applied to many complex tasks with exceptional results such as self-driving cars, facial recognition and many more.

Receptive Field and Kernel

In the previous section where an overview of CNN's was addressed, there was a discussion focused on filtering where features were searched for across the images. This filtering resulted in windows sliding across the image, looking for common features. This window is called the receptive field and this is the section of the image that is being looked at.

As the image is convoluted and processed through the network, the changed image can be called the kernel. The kernel of the CNN represents what features are extracted from the image throughout the network.

Activation Function

The activation function used in a neural network calculates the output for a neuron. The most commonly used activation function in CNN's is the ReLu activation function or the rectified linear unit Aurélien, (2017). The softmax activation function is sometimes used for the output layer Aurélien, (2017).

Dropout

Dropout is a type of regularisation technique used to prevent overfitting in a network Aurélien, (ibid.). Overfitting in a network is when the network cannot generalise to new data after over learning the training data. In dropout, each neuron is given a probability (normally set to 0.5) and this is the likelihood that the neuron will be used for that training step. It has proven quite successful for increasing a networks accuracy Aurélien, (2017).

Different Sizes of Convolutions

Zero padding can be added to an output so that it is same size as the previous layer by adding zeros around the feature Aurélien, (ibid.).. Through training, it is "possible to connect a large input layer to a much smaller layer by spacing out the receptive fields" Aurélien, (2017). The stride between two different receptive fields is the spacing between them Aurélien, (2017)..

Fully Convolutional Networks

A Fully Convolutional Network is one that does not have a fully connected layer and in a fully connected layers place is another convolution layer. This can be achieved by making the size of the receptive field equal to the size of the input *Image Segmentation Using DIGITS 5* (n.d.). They are proven to be very successful in object detection.

2.5 Support Vector Machine

A Support Vector Machine (SVM) is a machine learning algorithm that has been very popular before the use of a CNN was mainstream. Many of the texts that will be analysed later use SVMs for their classification.

A SVM works by creating an n-dimensional space, with n as the number of inputs you have *Support Vector Machines for Machine Learning* (n.d.). The SVM algorithm finds the hyperplane that splits this space. This hyperplane can then be used for classification. An SVM casts the problem to a higher dimensional space and this can solve a problem when the classes are not linearly separable. This can be seen in Figure 2.7, where on the left there is no clear way to separate the classes but once the problem is cast to another dimensional, the separation is clear.

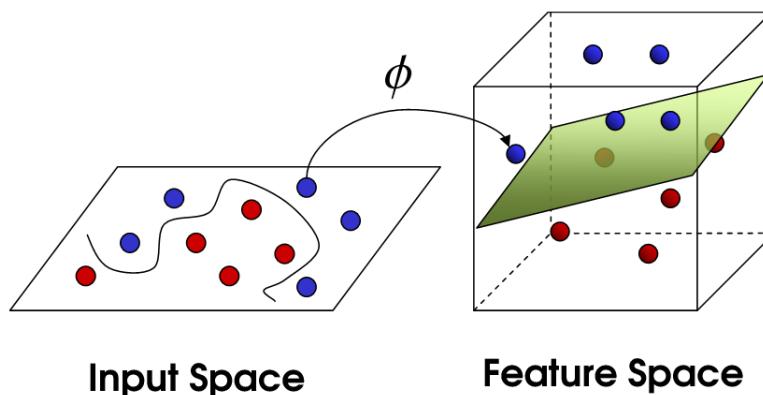


Figure 2.7: Support Vector Machine Applicability

There are benefits and liabilities to using a SVM. They can be very accurate and can work efficiently with small datasets but unfortunately it can take a large amount of time to train.

2.6 Evaluating the Output

There are various different metrics that can be used for evaluating the output of a classifier or segmentation algorithm, many of which we have seen in previous sections. These metrics will be examined below.

Research into Diagnosing Errors in Object Detectors

There has been some research into the question of how to evaluate object detectors, one of which be discussed in detail Hoiem, Chodpathumwan, and Dai, (2012). This paper in question ”analyzes the influences of object characteristics on detection performance and the frequency and impact of different types of false positives” Hoiem, Chodpathumwan, and Dai, (ibid.). They found that there were many effects that had influence on detectors as follows:

- occlusion
- size
- aspect ratio
- visibility of parts
- viewpoint
- localization error
- confusion with semantically similar objects
- confusion with other labeled objects
- confusion with background

The research team goes on to analyse false positives in object detectors. Localization errors were a large factor. This is where bounding boxes overlap to other objects in the image. Confusion with similar objects had a large influence on false positives also by which, for example, a dog detector had a high score for

a cat Hoiem, Chodpathumwan, and Dai, (2012). Confusion with dissimilar objects and confusion with background are the categories of the rest of the false positives they measured.

In conclusion the team would that "Most false positives are due to misaligned detection windows or confusion with similar objects" Hoiem, Chodpathumwan, and Dai, (ibid.). They had some recommendations towards improves detectors as follows:

- Smaller objects are less likely to be detected
- Localization could be improved
- Reduce confusion with similar categories
- Robustness to object variation
- More detailed analysis

Detection Average Precision

The average detection accuracy of a system. See A.1.

Mean Average Precision

The mean accuracy of a system across all results. See A.2.

Distribution of top-ranked false positives

This metric is used to tell where most errors occur when false positives are evident. These errors are broken down into four categories of localization, similar objects, background confusion and others. See A.3.

Segmentation Mean Accuracy

This is the mean accuracy of segmentation. See A.4.

Per-category segmentation accuracy

This metric measures the accuracy of segmentation at a category level. See A.5.

Per-class segmentation accuracy

The accuracy of segmentation at a class level is analysed. See A.6.

Top-1 and Top-5 Accuracy

When a classifier is given an image it normally responds with a list of predictions along with a decimal representation of the likelihood that it is of that class. The Top-1 accuracy is the top valued prediction and Top-5 accuracy is the top 5 predictions.

2.7 Dietary Assessment using Computer Vision

Convolution Neural Networks

Many researchers have used convolutional neural networks for image classification with various network architectures and many have used a food image dataset. Some of these papers will be looked at below. A summary of their results can be seen in Table 2.2.

Mao, Yu, and Wang, (n.d.) focused on a deep learning approach to food image recognition based their neural network architecture on Inception-ResNet and Inception V3. Deep learning is a term usually given to algorithms based on neural networks. They also used the Food-101 dataset. For this system, Google's Tensorflow was used for image preprocessing. Preprocessing was needed as the environmental background is different in many food images.

Because of these "Grey World method and Histogram equalization" Mao, Yu, and Wang, (*ibid.*) were used. AWS GPU instances were used for training and the results on completion were quite impressive with a Top-1 Accuracy of 72.55% and a Top-5 Accuracy of 91.31%.

Another research team in Japan, Yanai and Kawano, (2015) researched this topic. This was built off previous research they had carried out in the field Kawano and Yanai, (2014). They were aware of how difficult the problem was and therefore employed many techniques to solve the problem such as "pre-training with the large-scale ImageNet data, fine-tuning and activation features extracted from the pre-trained DCNN". In conclusion, they found that the "fine-tuned DCNN which was pre-trained with 2000 categories" from ImageNet was the best method. A DCNN is a Deep Convolution Neural Network. A network can become a DCNN when the number of hidden layers is larger than three. While many of the CNNs discussed have been DCNN, most are just labeled as CNNs. The achieved results of 78.77% for Top-1 Accuracy in the UECFOOD100 dataset.

Kagaya, Aizawa, and Ogawa, (2014) also employed the use of convolutional neural networks for image detection. They used a CNN for the "tasks of food detection and recognition through parameter optimization". They found that a CNN is much better suited to the task than a Support Vector Machine (SVM). They achieved an overall classification accuracy of 93.8% against their baseline accuracy of 89.7%. This accuracy was calculated using a dataset that they created specifically for this task. When they had completed the task they analysed the trained convolutional kernels and came to an interesting conclusion. They found that "color features are essential to food image recognition".

The last paper that will be analysed, Liu et al., (2016), oriented around using a Convolutional Neural Network for food image recognition, focuses on developing a dietary assessment application for use on a smart phone. They used the UEC-256 and Food-101 dataset for their experiments and achieved impressive results. They used a Convolutional Neural Network but "with a few major optimizations, such as optimized model and an optimized convolution technique". They used the Inception module for their CNN. After the

inception module was complete, they made the GoogleNet architecture by combining modules. In total, the network had 22 layers. They achieved the results shown in Table 2.1.

Table 2.1: DeepFood Results

	Top-1	Top-5
UEC-256	54.7%	81.5%
UEC-100	76.3%	94.6%
Food-101	77.4%	93.7%
UEC-256 With Bounding Box	63.8%	87.2%
UEC-100 With Bounding Box	77.2%	94.8%

Mezgec and Koroušić Seljak, (2017) developed a new neural architecture specifically for detecting food and drink images using deep convolutional neural networks called NutriNet. The trained network was to be used to aid patients with Parkinson’s disease in monitoring their diet. NutriNet was created based off of the AlexNet architecture. The dataset used for this study consisted of approximately 500 images for each of over 500 classes. Through this dataset a Top-1 accuracy of 86.72% and a Top-5 accuracy of 94.47% was recorded. A smart phone application was used for real world testing which brought a Top-5 accuracy of 55%. The application also saved these real world images from the smart phone to increase their dataset size. In conclusion, the team found that there were modifications that could be made to the NutriNet architecture as real world images didn’t perform incredibly well due to occlusion and background noise in the images. The detection and recognition steps were separated in this architecture. The team also acknowledges that joining these steps into a single DCNN may be successful and should be explored.

Support Vector Machines

While Convolutional Neural Networks have proven very successful in recent years, there are many other methods of food image identification and classification that have been employed by food image recognition researchers, such

Table 2.2: Summary of results in CNN based methods

Title	Dataset	Top-1 Accuracy
Mao, Yu, and Wang, (n.d.)	Food 101	72.55%
Yanai and Kawano, (2015)	UECFood101	78.77%
Kagaya, Aizawa, and Ogawa, (2014)	Own dataset	93.80%
Liu et al., (2016)	Food 101	77.40%
Mezgec and Koroušić Seljak, (2017)	Own dataset	86.72%

as SVMs. A summary of the results of these methods can be seen in Table 2.3.

A Food Image Recognition System with Multiple Kernel Learning

In this paper, Joutou and Yanai, (2009), a practical use for food image recognition in the form of a mobile phone application was proposed. In order to classify the images, multiple kernel learning(MKL) was used. MKL is similar to an SVM except that instead of a single kernel during training, MKL "treats with a combined kernel which is a weighted linear combination of several single kernels" Joutou and Yanai, (2009). The idea behind this is that different food types are distinguishable by different factors and using this method, the best of these factors can be used for classification of that food type. In the experiment carried out by Joutou and Yanai, (ibid.), three different factors were used for learning:

- Color Histograms
- Gabor Texture Features
- Bag-of-Features using Scale Invariant Feature Transformation (SIFT)

50 different classifiers were created in a SVM using MKL with "one category as a positive set and other 49 categories as a negative set" Joutou and Yanai, (ibid.). For each of these categories, a web scrape was carried out and then the best 100 images for each scrape was manually selected. Five-fold cross

validation was utilised in the paper. MKL proceeded to yield results of 61.34% on the 50 food types and a Top-3 accuracy of 80.05% Joutou and Yanai, (ibid.). The prototype mobile phone application resulted in a 37.55% user accuracy.

A Novel SVM Based Food Recognition Method for Calorie Measurement Applications

Another quite successful study was carried out using a SVM. Pouladzadeh, Villalobos, et al., (2012) had established that both colour and texture are very important but they also decided that shape and size are vital features to analyse. The proposed system has two main parts, segmentation followed by classification. In order to create a 'robust' system, a 'Robust Handling of Different Lighting Conditions' module is added to the system Pouladzadeh, Villalobos, et al., (2012). This is so that various lighting conditions don't cause color data to be distorted.

Since this paper calls for calorie estimation, the first step of the system calculates the size of the food portion. In order to do this, a coin or the users thumb is included in the image taken so that the pixel count of the thumb and the food can be compared to estimate the size. Following this the image is segmented into various portions. The following step classifies each segment of the image by extracting color, texture and shape features and inputting these into a SVM Pouladzadeh, Villalobos, et al., (ibid.).

12 different food types were trained for this SVM with an average classification accuracy of 92.6%.

Measuring Calorie and Nutrition From Food Image

Another study that employed both a SVM and an emphasis on colour, texture and shape features, was carried out by Pouladzadeh, Shirmohammadi, and Al-Maghribi, (2014). Size was also a factor in the calorie measurement module of the system. It was found that using all four of these features increases the overall accuracy.

In order to segment the image successfully, Gabor filters were applied to separate texture features while color was also utilised. For each segment established, size, shape, color and texture features were extracted and using a SVM, a classification was made. The SVM used the radial basis function kernel Pouladzadeh, Shirmohammadi, and Al-Maghribi, (ibid.).

Calorie estimation was also a large part of this paper, and the users thumb was taken with the food in order to calculate food size.

In the prototype application, once the classification had been made, the user can confirm or change the prediction. Another feature of the application was in regards to "Partially Eaten Food" Pouladzadeh, Shirmohammadi, and Al-Maghribi, (ibid.). This was accomplished by taking a picture before and after consumption and as a result only calculated the size of the food eaten and therefore more accurate calorie counts can be produced.

15 food types were trained using the SVM with 3000 images. The accuracy for the classifier averaged at 90.41% using 10 fold cross-validation Pouladzadeh, Shirmohammadi, and Al-Maghribi, (ibid.). There was also a calorie count accuracy of 86%. The best classification results were on single foods followed by non mixed and finally mixed foods produced the worst results.

Segmentation Assisted Food Classification for Dietary Assessment

Zhu et al., (2011) had a strong focus on the segmentation aspect of a dietary assessment system. The segmentation of the food images was achieved "using Normalized Cuts based on intensity and colour" Zhu et al., (2011). Normalized Cuts is a graph based segmentation method. To aid the segmentation aspect of this study, a common background colour was introduced to the images. Segmentation refinement was also an important module in the experiment. This is the process by which neighbouring segments with the same classification label are merged together. This also helps calculate a more accurate size estimation.

The classification of the segmented image was processed by using a SVM calculating colour and texture features. Gabor filters were used for the texture

feature extraction Zhu et al., (ibid.).

In the experimental results for this study, it was found that segmentation was not always successful "when the region of interest is camouflaged by making its boundary faint" Zhu et al., (ibid.). In their case, it was a can of coke that wasn't segmented correctly. The classification accuracy was of 56.2% and 95.5% with ground truth segmentation data Zhu et al., (ibid.).

Large Scale Learning for Food Image Classification

Abbirami.R.S et al., (2015) proposed a food image recognition system using a Bag of Features model. This study used over 5000 images separated into 11 classes.

A clustering algorithm was employed on this study before classification. For the classification step, experiments were carried out using different methods:

- SVM
- ANN
- Random Forests

The final accuracy of the system was 78% Abbirami.R.S et al., (ibid.).

A Personal Assistive System for Nutrient Intake Monitoring

Similar to other approaches seen thus far, Villalobos, Almaghrabi, Hariri, et al., (2011) employs the use of the users thumb in the image for size estimation.

Once a photo has been taken by the user with their thumb present, the system segments the food on the plate using shape, colour and texture detectors. The system then classifies the food type based on these features.

In this paper, it was decided to allow the users to change the prediction by the system. The thumb of each user is calibrated upon first use of the application so that size estimation can be as accurate as possible Villalobos, Almaghrabi, Hariri, et al., (ibid.).

Toward Dietary Assessment via Mobile Phone Video Camera

Another study into using computer vision for dietary assessment was carried out by N. Chen et al., (2010). They had a unique medium for the topic by using a video of the dishes in question and extracting frames from these videos to get the food from different angles.

N. Chen et al., (*ibid.*) then formed a region of interest in the image, where there were the most food items and extracted colour and image features. These image features were extracted using Maximally Stable Extremal Regions (MSER), Speeded Up Robust Features (SURF) and Star detector.

This research team also uses k-means clustering to build a bag-of-words model N. Chen et al., (*ibid.*).

The system had results as seen below across 20 categories using five images out of each video taken of the food:

- MSER - 95%
- SURF - 90%
- STAR - 90%

Automatic Chinese Food Identification and Quantity Estimation

There was a study carried out on food identification through a smart phone application M.-Y. Chen et al., (2012). This study resulted in an application that allows a user to send an image of their food to a server which can give them an automatic response in 12 seconds M.-Y. Chen et al., (*ibid.*). This back end service can have 34 threads working concurrently as stated at the time the paper was published M.-Y. Chen et al., (*ibid.*).

A SVM is used to classify the image across 50 categories trained on around 100 images each. The SVM uses SIFT and Local Binary pattern feature extractors M.-Y. Chen et al., (*ibid.*). A separate SVM was trained for each of these extractors and was merged together using a "Multi-class AdaBoost algorithm" M.-Y. Chen et al., (*ibid.*).

The study produced a top-1 accuracy of 68.3%. Accuracy of 80.6%, 84.8% and 90.9% were recorded using top-2, top-3 and top-5 accuracy respectively M.-Y. Chen et al., (2012).

An Image Processing Approach for Calorie Intake Measurement

Villalobos, Almaghrabi, Pouladzadeh, et al., (2012) researched the question of using a computer vision approach to this topic. They focus mostly on the segmentation and region of interest calculation of the system in their study.

The system in question requires two images of the food, one from above and one from the side. This helps with size estimation. The users thumb is required to be in the image for accurate size estimation. The application also requires an image after consumption as to not calculate calories for uneaten food.

The system segments the image and then extracts colour, size and shape information from each segment. This data is then used by a SVM for classification along with a nutritional database for calorie information.

Multiple segmentation methods were tested such as:

- Semi automatic contour definition
- Watershed transformation
- Colour rasterization
- Edge accentuation

The first two were dismissed due to poor results but the second two were used in conjunction for the segmentation aspect of the system.

Food Recognition and Nutrition Estimation on a Smartphone

An application called "Snap-n-Eat" was proposed by Zhang et al., (2015).

When an image is taken using this application, the system finds saliency regions to remove the background of the image. If the image has multiple food

types present, hierarchical segmentation takes place before proceeding to a SVM. Similar segments are merged together. These are found by using colour, texture and size.

SIFT and Histogram of Oriented Gradients (HOG) feature extractors are used on the image and these features are used by the SVM for classification. The SVM is trained using Scholastic Gradient Descent. Zhang et al., (2015) also uses a Bag of Visual Words model along with k-means clustering.

An accuracy of 85% on 15 classes was recorded using this method Zhang et al., (ibid.).

Merging dietary assessment with the adolescent lifestyle

Schap et al., (2014) proposed a system used by smart phones which sends an image of a users food to a back end system for computation.

Once this has been completed, the image is segmented, features are extracted from each segment and these segments are classified. Colour and texture features are used for classification. The user has the ability to confirm or amend predictions of the food type.

Size estimation is also an important aspect of this system. In contrast to previous studies, Zhang et al., (2015) uses food type shape and then those shape's geometric properties to estimate size.

This study produced results of 94% out of 32 test cases.

Other Methods

Methods outside of CNNs and SVMs are outlined below. A summary of the results can be seen in Table 2.3 along with SVM method results.

Large Scale Learning for Food Image Classification

Abbirami.R.S et al., (2015) proposed a food image recognition system using a Bag of Features model. This study used over 5000 images separated into 11

Table 2.3: Summary of accuracy in dietary assessment methods

Title	Classes	Accuracy
Food Image Recognition with Multiple Kernel Learning	50	61.34%
A Novel SVM Based Food Recognition Method	12	92.6%
Measuring Calorie and Nutrition From Food Image	15	90.41%
Segmentation Assisted Food Classification	N/A	77.4%
Large Scale Learning for Food Image Classification	11	78%
Toward Dietary Assessment via Mobile Phone Video Camera	20	~92%
Automatic Chinese Food Identification and Quantity Estimation	50	68.3%
Food Recognition and Nutrition Estimation on a Smartphone	15	85%

classes.

A clustering algorithm was employed on this study before classification. For the classification step, experiments were carried out using different methods:

- SVM
- ANN
- Random Forests

The final accuracy of the system was 78% Abbirami.R.S et al., (2015).

A Personal Assistive System for Nutrient Intake Monitoring

Similar to other approaches seen thus far, Villalobos, Almaghrabi, Hariri, et al., (2011) employs the use of the users thumb in the image for size estimation.

Once a photo has been taken by the user with their thumb present, the system segments the food on the plate using shape, colour and texture detectors. The system then classifies the food type based on these features.

In this paper, it was decided to allow the users to change the prediction by the system. The thumb of each user is calibrated upon first use of the application so that size estimation can be as accurate as possible Villalobos, Almaghrabi, Hariri, et al., (ibid.).

Toward Dietary Assessment via Mobile Phone Video Camera

Another study into using computer vision for dietary assessment was carried out by N. Chen et al., (2010). They had a unique medium for the topic by using a video of the dishes in question and extracting frames from these videos to get the food from different angles.

N. Chen et al., (*ibid.*) then formed a region of interest in the image, where there were the most food items and extracted colour and image features. These image features were extracted using Maximally Stable Extremal Regions (MSER), Speeded Up Robust Features (SURF) and Star detector.

This research team also uses k-means clustering to build a bag-of-words model N. Chen et al., (2010).

The system had results as seen below across 20 categories using five images out of each video taken of the food:

- MSER - 95%
- SURF - 90%
- STAR - 90%

Merging dietary assessment with the adolescent lifestyle

Schap et al., (2014) proposed a system used by smart phones which sends an image of a users food to a back end system for computation.

Once this has been completed, the image is segmented, features are extracted from each segment and these segments are classified. Colour and texture features are used for classification. The user has the ability to confirm or amend predictions of the food type.

Size estimation is also an important aspect of this system. In contrast to previous studies, Zhang et al., (2015) uses food type shape and then those shape's geometric properties to estimate size.

This study produced results of 94% out of 32 test cases.

Possible Applications in Region Based CNNs

Ross Girshik and other contributers had some very positive results in the area of object detection using region based convolutional neural networks. There were four iterations of papers based on this work by Ross and groups in UC Berkley, Mircosoft and Facebook. A PHD student at the time of Ross's first paper also completed his dissertation on the subject. These papers, their results (Table 2.4) and the changes made through each iteration will be analysed throughly in the proceeding sections.

RCNN

In the first paper written by Ross Girshik, while researching at UC Berkeley, focused on two main insights. These were that "one can apply high-capacity convolutional neural networks (CNNs) to bottom-up region proposals in order to localize and segment objects" and that "when training data is scarce, supervised pre-training for n auxiliary task, followed by domain-specific fine-tuning, yields a significant performance boost" Girshick et al., (2014).

The system that they developed followed these steps:

- Take image as input
- Extract approximately 2000 region proposals from the image
- Compute fixed length vectors of features for the regions using a convolutional neural network
- Use a Support Vector Machine (SVM) to classify these regions
- Bounding box regression for final region proposals

This system utilised selective search to gather these region proposals but they mention that a sliding-window detector is also an option. Ross Girshik and his team used the open source Caffe CNN library for this system. The system is quite efficient and scalable. It is scalable because of the fixed length vector

of features which will remain constant regardless of inputs and additional outputs.

The team evaluated their results on a few metrics and test sets as seen in Table 2.4.

Fast RCNN

Ross Girshik's next iteration of work on region based convolution neural networks took place in Microsoft Research. This paper was titled "Fast R-CNN" as it's aim was to decrease training and testing time "while also increasing detection accuracy" Girshick, (2015).

This paper analyses why RCNN Girshick et al., (2014) was slow and therefore how it could be improved. RCNN was classified to be slow because of three main factors:

- There are multiple stages to training as both a CNN and a SVM need to be trained.
- In training of the SVM, each region proposal must be written to disk and is therefore expensive.
- Object detection takes 47 seconds per image Girshick, (2015).

Due to these problems with RCNN, a new algorithm, titled Fast RCNN was proposed. The architecture is as follows. An image is taken as input along with a proposals for regions. The image is pushed through convolutional and pooling layers (using max pooling). A fixed-length vector of features is then extracted from each region proposal. These vectors are inputted to fully connected layers for bounding box location prediction Girshick, (2015).

At detection time, a pass through of the net is all that is needed so this runtime is significantly less than RCNN.

Faster RCNN

Due to the success of RCNN and Fast RCNN, Faster RCNN was introduced to combat the problem of region proposal computation Ren et al., (2015).

The architecture for this system comprises of two modules. These consist of a convolutional neural network for region proposals (RPN) which feeds into a Fast RCNN detector. These combine to produce a single neural network for object detection.

Instead of training these networks separately, the team had to look at how to share layers between the two networks. There were three options available:

- Alternating training whereby RPN is trained, and then used to train Fast RCNN. The Fast RCNN network is then used to initialise RPN and the process is iterated Ren et al., (2015). This paper follows this approach.
- Approximate joint training.
- Non-approximate joint training.

Mask RCNN

The most recent paper on this topic was also written by Ross Girshik while working with Facebook AI Research K. He et al., (2017). Mask RCNN "extends Faster RCNN by adding a branch for predicting an object mask in parallel with the existing branch for bounding box regression" K. He et al., (2017).

Mask RCNN has two modules, similar to Faster RCNN, where the first module is the Region Proposal Network. In the second module, in parallel to classification, a binary mask is outputted for each region. Bounding box regression and classification are done in parallel.

Table 2.4: Results from Region Based CNN Research

	VOC07	VOC10	VOC11	VOC12	COCO15	COCO16
RCNN	58.5%	53.7%	47.9%	N/A	N/A	N/A
Fast RCNN	70.0%	68.8%	N/A	68.4%	N/A	N/A
Faster RCNN	78.8%	N/A	N/A	75.9%	42.7%	N/A
Mask RCNN	N/A	N/A	N/A	N/A	N/A	63.1%

2.8 APIs and Libraries

Tensorflow

Tensorflow is a deep learning software library for various machine learning paradigms. Tensorflow will be used to create neural networks. Tensorflow has two utilisations, through a Graphics Processing Unit (GPU) and also through a Central Processing Unit (CPU). GPU computation is recommended for CNN training.

Central Processing Unit Computation

It is quite easy to get tensorflow up and running if you are only using a CPU to train. Tensorflow CPU has been successfully installed on both Windows and Ubuntu, for the purpose of this project. For Windows you can download and install using the tensorflow website and on ubuntu you can use apt-get. Once installed, tensorflow can be imported into any python shell or script for use. Tensorflow can also be used in C++. There will be various python implementations of neural networks in Chapter 3.

Graphics Processing Unit Computation

For use with a GPU, the set up for tensorflow is a bit more complicated. Firstly you must check that the GPU in your machine is compatible for CUDA 8.0 using the NVIDIA website. If your GPU is compatible, you must install CUDA after signing up as an NVIDIA developer. CUDA 8.0 is compatible with tensorflow.

You also need to install cudnn6. The NVIDIA website contains tutorials to install these. Once these are installed, download and install tensorflow-gpu. This can be imported into python similar to CPU computation.

Caffe

OpenCV

OpenCV is an industry wide, open source library for computer vision and machine learning *About - OpenCV library* (n.d.). It has over 2500 algorithms that are available for use *About - OpenCV library* (ibid.). OpenCV is supported across multiple languages and platforms such as Python, C++, C, Java, Matlab, running on Windows, Android, Mac OS and Linux *About - OpenCV library* (ibid.).

There is not much of the library utilised in this project due to the nature of Tensorflow but some algorithms for image reading, writing and resizing were used due to the ease of use.

```
image = cv2.imread('image.jpg')

resized = cv2.resize(image, (299, 299))

cv2.imwrite('imageResized.jpg', resized)
```

Numpy

Jupyter

2.9 Conclusion

Chapter 3

Introduction to Using Tensorflow

The purpose of this chapter is to give an overview of what has been completed in order to gain understanding of how to create CNNs and also how to work with datasets. Three experiments are explored in this chapter, the first two of which are tutorials completed as an introduction to Tensorflow and CNNs. The final experiment explores how to feed datasets from file directories into CNNs.

3.1 Template for Experiments

The template followed for all experiments in chapters three, four and five is outlined in Figure 3.1.

3.2 Udemy Tutorial

Overview

There are many on line resources that are geared towards helping deep learning novices. One of these resources is a course on Udemy titled 'Complete Guide to

Table 3.1: Experiment Template

Experiment Section	Rationale
Overview	An explanation of the purpose of the experiment along with how it is carried out.
Network Architecture	An explanation of the network architecture used in the experiment. If the architecture has been explained in another experiment, the architecture will only be referenced by name.
Dataset	The dataset used for the experiment, with its source and an overview of its details. A brief mention will be made in following experiments.
API's	Reference to the technologies used as outlined in Chapter 2.
Script	Snippets of the script used for the experiment.
Results	The results acquired from the experiment, usually in the form of a percentage accuracy.
Empirical Analysis	States any information gained from the experiment and speculation for the reasoning of the results.

Tensorflow for Deep Learning with Python' *Complete Guide to Tensorflow for Deep Learning with Python* (n.d.). This course has a section on Convolutional Neural Networks which has been followed and completed.

Network Architecture

The architecture for this network is very simple as it is an introductory CNN. It consists of two convolutional layers, two max pooling layers and a fully connected layer.

Dataset

This CNN is trained on the MNIST dataset. This dataset consists of 70000 handwritten digits LeCun and Cortes, (2010).

API's

This experiment was carried out in a jupyter notebook using tensorflow.

Script

```
1 #Helper Functions  
2  
3 #INIT WEIGHTS  
4 def init_weights(shape):  
5     init_random_dist = tf.truncated_normal(shape, stddev = 0.1)  
6     return tf.Variable(init_random_dist)  
7  
8 #INIT BIAS  
9 def init_bias (shape):  
10    init_bias_vals = tf.constant(0.1, shape = shape)  
11    return tf.Variable( init_bias_vals )
```

```

12
13 #CONV2D
14 def conv2d(x, W):
15     #x -> [batch, H, W, Channels]
16     #W -> [filterH, filterW, ChannelsIn, ChannelsOut]
17     return tf.nn.conv2d(x, W, strides = [1,1,1,1], padding = 'SAME')
18
19 #POOLING
20 def max_pool_2by2(x):
21     #x -> [batch, H, W, Channels]
22     return tf.nn.max_pool(x, ksize = [1,2,2,1], strides = [1,2,2,1],
23                           padding = 'SAME')
24
25     #NORMAL (FULLY CONNECTED)
26 def normal_full_layer (input_layer , size ):
27     input_size = int(input_layer .get_shape() [1])
28     W = init_weights([input_size , size ])
29     b = init_bias ([ size ])
30     return tf.matmul(input_layer, W) + b

1 #32 features for every 5 x 5 patch with 1(grayscale)
2 convo_1 = convolutional_layer(x_image, shape = [5,5,1,32])
3 convo_1_pooling = max_pool_2by2(convo_1)
4
5 convo_2 = convolutional_layer(convo_1_pooling, shape = [5,5,32,64])
6 convo_2_pooling = max_pool_2by2(convo_2)
7
8 convo_2_flat = tf.reshape(convo_2_pooling, [-1,7*7*64])
9 full_layer_one = tf.nn.relu( normal_full_layer ( convo_2_flat , 1024))

1 with tf.Session() as sess:
2     sess.run(init)
3

```

```

4   for i in range(steps):
5       batch_x, batch_y = mnist.train.next_batch(32)
6       batch_test = mnist.test.next_batch(32)
7       sess.run(train, feed_dict = {x:batch_x, y_true:batch_y,
8           hold_prob:0.5})
9       if i%500 == 0:
10          print("ON STEP: {}".format(i))
11          print("ACCURACY: ")
12          match = tf.equal(tf.argmax(y_pred, 1), tf.argmax(y_true, 1))
13          acc = tf.reduce_mean(tf.cast(match, tf.float32))
14          print(sess.run(acc, feed_dict = {x:batch_test,
15              y_true:mnist.test.labels, hold_prob:1.0}))
16          print('\n')
17
18 #Final Accuracy 0.973

```

Results

The Final Accuracy for this experiment was of 97.3%.

Empirical Analysis

3.3 Udemy Tutorial 2

Overview

Similar to the previous experiment 3.3, this experiment is a Udemy course exercise *Complete Guide to Tensorflow for Deep Learning with Python* (n.d.). In contrast, this does not use a dataset built into Tensorflow so therefore, there is extra configuration to be done on the dataset.

Network Architecture

Same architecture as in 3.2.

Dataset

The CIFAR-10 dataset is used here. This has 60000 images split into 10 classes and a test set of 1000 images Krizhevsky, Nair, and Hinton, (n.d.).

Script

```
1 def unpickle( file ):
2     import pickle
3     with open( file , 'rb' ) as fo:
4         cifar_dict = pickle.load(fo, encoding='bytes')
5     return cifar_dict
6
7 dirs =
8 ['batches.meta','data_batch_1','data_batch_2','data_batch_3','data_batch_4','data_batch_5','data_batch_6']
9 all_data = [0,1,2,3,4,5,6]
10
11 for i,direc in zip( all_data , dirs ):
12     all_data [i] = unpickle(CIFAR_DIR+direc)
13
14 batch_meta = all_data[0]
15 data_batch1 = all_data[1]
16 data_batch2 = all_data[2]
17 data_batch3 = all_data[3]
18 data_batch4 = all_data[4]
19 data_batch5 = all_data[5]
20 test_batch = all_data[6]
```

```

1 def set_up_images(self):
2
3     print("Setting Up Training Images and Labels")
4
5     # Vertically stacks the training images
6     self.training_images = np.vstack([d[b"data"] for d in
7         self.all_train_batches])
8     train_len = len(self.training_images)
9
10    # Reshapes and normalizes training images
11    self.training_images =
12        self.training_images.reshape(train_len ,3,32,32) .transpose (0,2,3,1) /255
13
14    # One hot Encodes the training labels (e.g. [0,0,0,1,0,0,0,0,0,0])
15
16    self.training_labels = one_hot_encode(np.hstack([d[b"labels"] for
17        d in self.all_train_batches]), 10)
18
19
20    # Vertically stacks the test images
21    self.test_images = np.vstack([d[b"data"] for d in self.test_batch])
22    test_len = len(self.test_images)
23
24
25
26    # Reshapes and normalizes test images
27    self.test_images =
28        self.test_images.reshape(test_len ,3,32,32) .transpose (0,2,3,1) /255
29
30    # One hot Encodes the test labels (e.g. [0,0,0,1,0,0,0,0,0,0])
31
32    self.test_labels = one_hot_encode(np.hstack([d[b"labels"] for d in
33        self.test_batch]), 10)
34
35
36    def next_batch(self, batch_size):
37        # Note that the 100 dimension in the reshape call is set by an

```

```

assumed batch size of 100

28     x =
29         self . training . images [ self . i : self . i +batch_size] . reshape(100,32,32,3)
30     y = self . training_labels [ self . i : self . i +batch_size]
31     self . i = ( self . i + batch_size) % len(self . training . images)

31 return x, y

1 cross_entropy =
2     tf . reduce_mean(tf . nn . softmax_cross_entropy_with_logits( labels =
3         y_true, logits = y_pred))
4 optimizer = tf . train . AdamOptimizer(learning_rate = 0.001)
5 train = optimizer.minimize(cross_entropy)
6 init = tf . global_variables_initializer ()
7 with tf . Session() as sess:
8     sess . run(tf . global_variables_initializer ())
9
10    for i in range(10000):
11        batch = ch.next_batch(100)
12        sess . run(train, feed_dict ={x: batch[0], y_true: batch[1],
13            hold_prob: 0.5})
14
15    # PRINT OUT A MESSAGE EVERY 100 STEPS
16    if i%1000 == 0:
17
18        # Test the Train Model
19        matches = tf.equal(tf.argmax(y_pred,1),tf.argmax(y_true,1))
20
21        acc = tf.reduce_mean(tf.cast(matches,tf.float32))
22        testSet = ch.next_test_batch(100)
23        print('Accuracy: ')
24        print(sess . run(acc, feed_dict ={x:testSet [0]
25            ,y_true: testSet [1], hold_prob:1.0}))
```

```
22 print('\'\n')
```

Results

A Final Accuracy of 71% was reached.

Empirical Analysis

3.4 Using the Food 101 dataset

Overview

For the next experiment, it was decided to use the Food-101 dataset. An on line tutorial was used to create a dataset in tensorflow using image directories on disk *Build an Image Dataset in Tensorflow* (n.d.). Tutorials used previously were also utilized for this experiment *Complete Guide to Tensorflow for Deep Learning with Python* (n.d.) and Krizhevsky, Nair, and Hinton, (n.d.).

Network Architecture

A similar network architecture to the previous two experiments (3.2 and 3.3) was used here.

Dataset

The Food-101 datset was used for this experiment Bossard, Guillaumin, and Van Gool, (2014).

Script

```
1 # Reading the dataset  
2 # 2 modes: 'file' or 'folder'
```

```

3  def read_images(dataset_path, mode, batch_size):
4      imagepaths, labels = list(), list()
5      if mode == 'file':
6          # Read dataset file
7          data = open(dataset_path, 'r').read().splitlines()
8          for d in data:
9              imagepaths.append(d.split(' ')[0])
10             labels.append(int(d.split(' ')[1]))
11     elif mode == 'folder':
12         # An ID will be affected to each sub-folders by alphabetical order
13         label = 0
14         # List the directory
15         try: # Python 2
16             classes = sorted(os.walk(dataset_path).next()[1])
17         except Exception: # Python 3
18             classes = sorted(os.walk(dataset_path).__next__()[1])
19         # List each sub-directory (the classes)
20         for c in classes:
21             c_dir = os.path.join(dataset_path, c)
22             try: # Python 2
23                 walk = os.walk(c_dir).next()
24             except Exception: # Python 3
25                 walk = os.walk(c_dir).__next__()
26             # Add each image to the training set
27             for sample in walk[2]:
28                 # Only keeps jpeg images
29                 if sample.endswith('.jpg') or sample.endswith('.jpeg'):
30                     imagepaths.append(os.path.join(c_dir, sample))
31                     labels.append(label)
32             label += 1
33     else:
34         raise Exception("Unknown mode.")

```

```

1 # Convert to Tensor
2     imagepaths = tf.convert_to_tensor(imagepaths, dtype=tf.string)
3     labels = tf.convert_to_tensor(labels, dtype=tf.int32)
4 # Build a TF Queue, shuffle data
5     image, label = tf.train.slice_input_producer([imagepaths, labels],
6                                                 shuffle=True)
7
8 # Read images from disk
9     image = tf.read_file(image)
10    image = tf.image.decode_jpeg(image, channels=CHANNELS)
11
12 # Resize images to a common size
13    image = tf.image.resize_images(image, [IMG_HEIGHT, IMG_WIDTH])
14
15 # Normalize
16    image = image * 1.0/127.5 - 1.0
17
18 # Create batches
19    X, Y = tf.train.batch([image, label], batch_size=batch_size,
20                          capacity=batch_size * 8,
21                          num_threads=4)
22
23    return X, Y

```

```

1 # Start training
2 with tf.Session() as sess:
3
4 # Run the initializer
5 sess.run(init)
6
7 # Start the data queue
8 tf.train.start_queue_runners()

```

```

9
10     # Training cycle
11     for step in range(1, num_steps+1):
12
13         if step % display_step == 0:
14             # Run optimization and calculate batch loss and accuracy
15             _, loss, acc = sess.run([train_op, loss_op, accuracy])
16             print("Step " + str(step) + ", Minibatch Loss= " + \
17                 " {:.4f}".format(loss) + ", Training Accuracy= " + \
18                 " {:.3f}".format(acc))
19
20     else :
21         # Only run the optimization op (backprop)
22         sess.run(train_op)
23
24
25     print("Optimization Finished!")

```

Results

The final accuracy for this model was 18.8% after 5000 steps.

Empirical Analysis

The poor accuracy of this model is to be expected due to the simple architecture of the network and the fact that 101 classes is a lot to process.

Chapter 4

Training Using the Inception-V3 Model Architecture

This chapter consists of experiments which train various Tensorflow models. There are parameters that can be changed in the code that creates a Tensorflow model based on the Inception-V3 model architecture. These parameters were changed in order to see how the model accuracy would be affected. In addition to this, change in dataset size was explored in relation to model accuracy.

4.1 Retrain ImageNet Inception V3 Model

Overview

For the fourth experiment undertaken, it was decided to take inspiration from Yanai and Kawano, (2015), where pre-training was used training a model for food classification. In order to achieve this, the final layer of the Inception V3 model which was trained on the ImageNet dataset had to be retrained. This is called Transfer Learning. A tutorial, created by Google, on the tensorflow website was followed for direction on this process *How to Retrain Inception's Final Layer for New Categories* (n.d.).

Firstly, in order to retrain the final layer of a model, a dataset must be prepared in the correct way. The Food-101 dataset Bossard, Guillaumin, and Van Gool, (2014) was used for this experiment, which will be analysed below. This dataset is in contrast to Yanai and Kawano, (2015) as they used the UECFOOD100 dataset. The food-101 dataset was chosen due to larger number of images per class, a difference between 100 and 1000. The dataset must be structured so that there is a separated directory for each class with the directory name as the class name. These directories should contain all the images for this class.

Once this dataset has been set up correctly, a directory can be found on github which contains the necessary files for this tutorial. When the directory has been downloaded, the following command can be ran:

```
python tensorflow/examples/image_retraining/retrain.py \ --image_dir  
~/dataset_directory
```

The first thing that the script will do is create bottleneck files for the images. A bottleneck is a term used to define the final layer before the output layer. This is so that for each image, we do not have to push it through the entire network during training *How to Retrain Inception’s Final Layer for New Categories* (n.d.).

After, the bottlenecks are created, the training can be completed. The images are split into three sub directories of training, testing and validation. By default, these images are split into percentages of 80%, 10% and 10% respectively. The model is trained at a default of 4000 steps.

At the final stage of the script, the model is run on a batch of test images not yet seen and a final test accuracy is displayed. This can be seen in the Script section below.

The command used for using this model once it is trained is:

```
python tensorflow/examples/label_image.py --graph=/tmp/output_graph  
--labels=/tmp/output_labels.txt --input_layer=Mul --output_layer=fin  
--input_mean=128 --input_std=128 --image=~/image_directory
```

Network Architecture

The Inception V3 model network architecture was used for this experiment. The Inception V3 architecture was created by building on the existing Inception model aimed at efficient image classification Szegedy et al., (2016). This research was carried out due to the popularity of convolutional neural networks. The main aim of this study was to produce a model that would "scale up networks in ways that aim at utilizing the added computation as efficiently as possible by suitable factorized convolutions and aggressive regularization" Szegedy et al., (2016). The team did not want to simply rely on larger models for better results.

There are 4 main design principles that the research team followed in the paper Szegedy et al., (2016):

- Avoid representational bottlenecks
- It is easier to process higher dimensional representations locally
- Without much loss in power, spatial aggregation can be carried out over lower dimensional embeddings
- The width and depth of the network should be balanced and when one is increased it can be beneficial to increase the other.

The main idea behind the network is that different strands of the network are followed and the best result is used as seen in Figure 4.1.

In order to combat the problem of low resolution input, three separate experiments were carried out based on the fact that "One simple way to ensure constant effort effort is to reduce the strides of the first two layer in the case of lower resolution input, or by simply removing the first pooling layer of the network" Szegedy et al., (ibid.).

These experiments are as follows, with the first resulting in the most accurate:

- "299 x 299 receptive field with stride 2 and maximum pooling after the first layer." Szegedy et al., (2016)

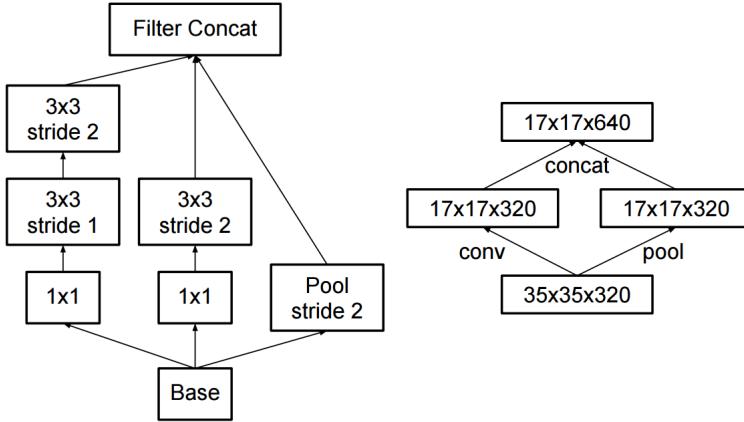


Figure 4.1: Inception module

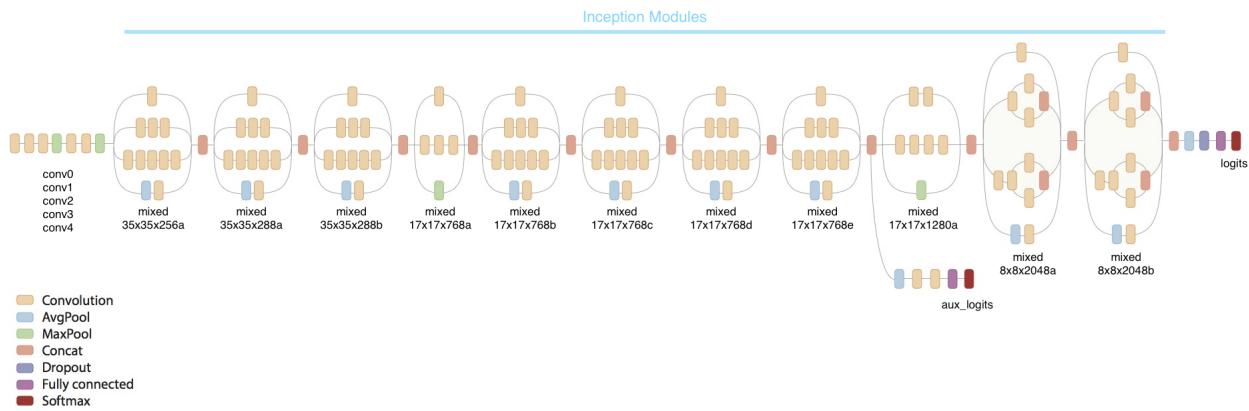


Figure 4.2: Inception V3 Architecture

- "151 x 151 receptive field with stride 1 and maximum pooling after the first layer." Szegedy et al., (ibid.)
- "79 x 79 receptive field with stride 1 and without pooling after the first layer." Szegedy et al., (ibid.)

The best results of the Inception-V3 model achieved 78.8% top-1 accuracy and 94.4% top accuracy of the ILSVR 2012 dataset Szegedy et al., (ibid.). This model also was the least computationally expensive of other published and successful models. The full model can be viewed in Figure 4.2.

Dataset

The dataset used for this experiment is the Food-101 dataset **food 101**. This dataset has 101 classes with 1000 images for each class.

Libraries

Tensorflow and Numpy were used to run this script.

Script

The following snippets of code are from the retrain.py script.

```
1 # Add the new layer that we'll be training.  
2 ( train_step , cross_entropy , bottleneck_input , ground_truth_input ,  
3   final_tensor ) = add_final_training_ops(   
4     len( image_lists.keys() ), FLAGS.final_tensor_name ,  
5     bottleneck_tensor ,  
6     model_info[ 'bottleneck_tensor_size' ] ,  
7     model_info[ 'quantize_layer' ] )  
8  
9 # Create the operations we need to evaluate the accuracy of our new layer.  
10 evaluation_step , prediction = add_evaluation_step(  
11   final_tensor , ground_truth_input )  
12  
13 # Merge all the summaries and write them out to the summaries_dir  
14 merged = tf.summary.merge_all()  
15 train_writer = tf.summary.FileWriter(FLAGS.summaries_dir + '/train' ,  
16                                         sess.graph)  
17  
18 validation_writer = tf.summary.FileWriter(  
19   FLAGS.summaries_dir + '/validation' )  
20
```

```

21 # Set up all our weights to their initial default values.
22 init = tf.global_variables_initializer()
23 sess.run(init)

1 # We've completed all our training, so run a final test evaluation on
2 # some new images we haven't used before.
3 test_bottlenecks, test_ground_truth, test_filenames =
4     get_random_cached_bottlenecks(
5         sess, image_lists, FLAGS.test_batch_size, 'testing',
6         FLAGS.bottleneck_dir, FLAGS.image_dir, jpeg_data_tensor,
7         decoded_image_tensor, resized_image_tensor, bottleneck_tensor,
8         FLAGS.architecture))
9 test_accuracy, predictions = sess.run(
10    [evaluation_step, prediction],
11    feed_dict={bottleneck_input: test_bottlenecks,
12               ground_truth_input: test_ground_truth})
13 tf.logging.info('Final test accuracy = %.1f%% (N=%d)' %
14                 (test_accuracy * 100, len(test_bottlenecks)))

```

Results

The final test accuracy for this retrained model was 54.8%.

For example, an image of pizza Figure4.3, was fed into the model with the followng results:

- pizza 0.925
- pancakes 0.008
- nachos 0.007
- beef carpaccio 0.006
- tiramisu 0.004



Figure 4.3: Pizza

In contrast, Figure 4.4 was not classified as a pizza.

Empirical Analysis

These poor results are not that surprising. This is because we have many classes to train for, 101, and no parameter tuning has been carried out on the running of this code.

Figure 4.4 was not classified correctly, this is most likely due to the fact that the pizza does not take up much of the image.

4.2 Retrain with Extended Dataset

Overview

As the Food-101 dataset mostly consisted of meals Bossard, Guillaumin, and Van Gool, (2014), it was decided to extend the dataset slightly by including some single foods such as:

- cheese
- grapes



Figure 4.4: Pizza not classified correctly by the model

- banana
- apple
- orange
- spaghetti
- roll

In order to collect these images, the ImageNet repository was utilised to search for these foods individually and then download the subset of images to be included with the Food 101 dataset Deng et al., (2009). The `retrain.py` script was run on the extended dataset as in 4.1.

Network Architecture

The Inception V3 Model was used for this experiment.

Dataset

In this experiment the extended Food-101 dataset was used.

Libraries

Tensorflow and numpy are used in the retrain.py script.

Script

As seen in 4.2.

Results

For this model, an accuracy of 55.3% was achieved.

For example, an image of a banana (Figure 4.5) was fed into the model with the following results:

- banana 0.9962
- orange 0.0009
- cheese 0.0003
- frozen yoghurt carpaccio 0.0002
- churros 0.0001

Empirical Analysis

The slight increase in accuracy in 4.1, from 54.8% to 55.3%, makes sense. Since the model was pre-trained using the ImageNet dataset and all the new images used were from ImageNet, we would expect a higher classification accuracy on the new additions to the dataset. This would overall increase the average classification accuracy.



Figure 4.5: Banana

4.3 Retrain with Parameter Tuning

Overview

Within the `retrain.py` script *How to Retrain Inception's Final Layer for New Categories* (n.d.), as mentioned in previous experiments, there are various parameters that can be set and changed. Various combinations of these parameters were changed to see if it would increase the test accuracy of the model.

Network Architecture

The Inception V3 architecture was used for this model.

Dataset

The Food-101 dataset Bossard, Guillaumin, and Van Gool, (2014) with additional classes, as per 4.1, was used for this model.

Libraries

The libraries in use for this experiment are tensorflow and numpy.

Script

Script as seen in 4.1 but with some additions to calculate Top 5 accuracy as seen below:

```
python retrain_top5.py \ --image_dir  
~/dataset_directory \ --how_many_training_steps 4000 \ --learning_rate  
--testing_percentage 10 \ --validation_percentage 10
```

Top 5 accuracy of the model was also calculated using the code below.

```
1 print('Top 5 Evaluation')  
2 #Variables used to store a list of all classes , the amount of images  
3 # tested,  
4 # the amount of images with a top 5 accuracy and the sum of the highest  
5 # probabilities  
6 classes = list(image_lists.keys())  
7 image_count = 0  
8 top_5_count=0  
9 total_of_top1_probs = 0  
10 i = 0  
11 class_counter = 0  
12  
13 #Loops through all test images, selecting 10 images per class  
14 #and running them through the method in 'label_image.py'.  
15 while(i < len(test_bottlenecks)):  
16     class_counter += 1  
17     if class_counter > 10:  
18         i = int(round((i + len(test_bottlenecks))/len(classes)) - 10))  
19         class_counter = 0
```

```

19     else :
20         image_count += 1
21         results_from_classifier = label_image.runModel(test_filenames[i])
22         results = results_from_classifier [0]
23         probabilities = results_from_classifier [1]
24
25         if classes [test_ground_truth[i]] in results :
26             top_5_count += 1
27         else :
28             print("Expected: " + classes [test_ground_truth[i]])
29             for result in results :
30                 print("Classes: " + result)
31             print("")
32
33         total_of_top1_probs += max(probabilities)
34         i = i + 1
35     print(str(top_5_count))
36     average_probabilities = total_of_top1_probs/(len(classes)*10))
37
38 #Prints out the amount of test images used, the top 5 accuracy
39 # and the average probability of predictions .
40 print("Amount of test images: " + str(image_count))
41 print("Top 5 Accuracy: " + str((top_5_count/image_count)*100))
42 print("Average probability: " + str(average_probabilities ))

```

Method from label_image.py.

```

1 def runModel(file_name):
2     model_file = \
3         "/tmp/output_graph.pb"
4     label_file = "/tmp/output_labels.txt"
5     input_height = 299
6     input_width = 299

```

```

7     input_mean = 128
8     input_std = 128
9     input_layer = "Mul"
10    output_layer = " final_result "
11
12    graph = load_graph(model_file)
13    t = read_tensor_from_image_file(file_name,
14                                    input_height=input_height,
15                                    input_width=input_width,
16                                    input_mean=input_mean,
17                                    input_std=input_std)
18
19    input_name = "import/" + input_layer
20    output_name = "import/" + output_layer
21    input_operation = graph.get_operation_by_name(input_name)
22    output_operation = graph.get_operation_by_name(output_name)
23
24    with tf.Session(graph=graph) as sess:
25        results = sess.run(output_operation.outputs[0],
26                            {input_operation.outputs[0]: t})
27        results = np.squeeze(results)
28
29        top_k = results.argsort() [-5:][::-1]
30        labels = load_labels( label_file )
31
32        setIndex = False
33
34        top5_results = [None] * 5
35        index = 0
36        for i in top_k:
37            top5_results [index] = labels[i]
38            index += 1

```

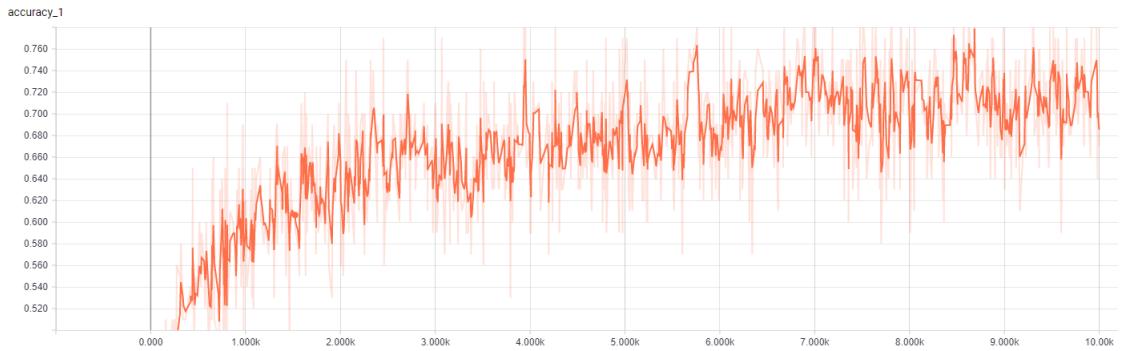


Figure 4.6: Graph of accuracy of the test dataset during training

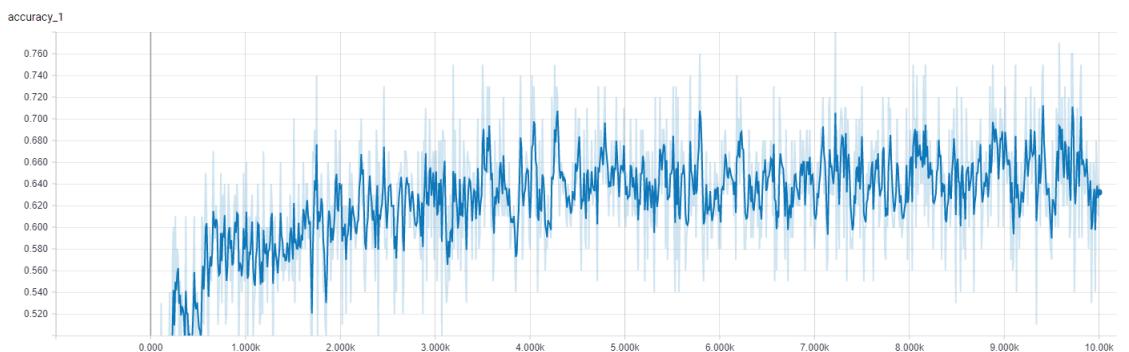


Figure 4.7: Graph of accuracy of the validation dataset during training

```

39
40     final_results = [top5_results, results]
41     return final_results

```

Some further parameters could be set such as:

- `-flip_left_right`
- `-random_crop`
- `-random_scale`
- `-random_brightness`

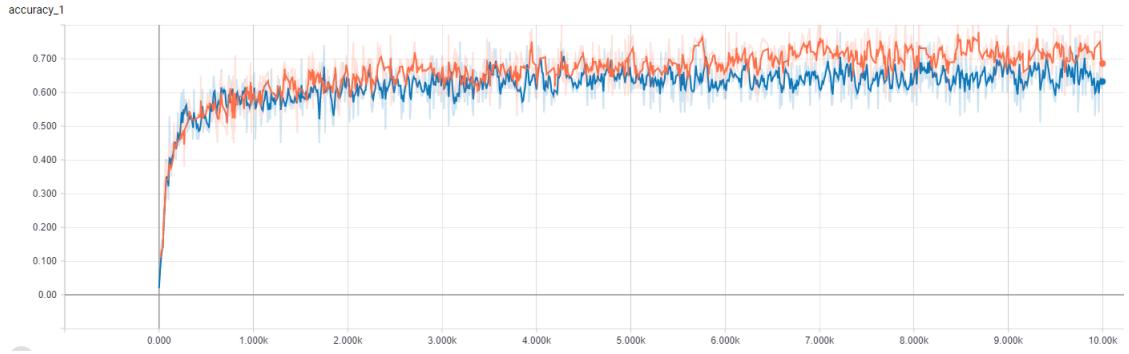


Figure 4.8: Comparison of accuracy

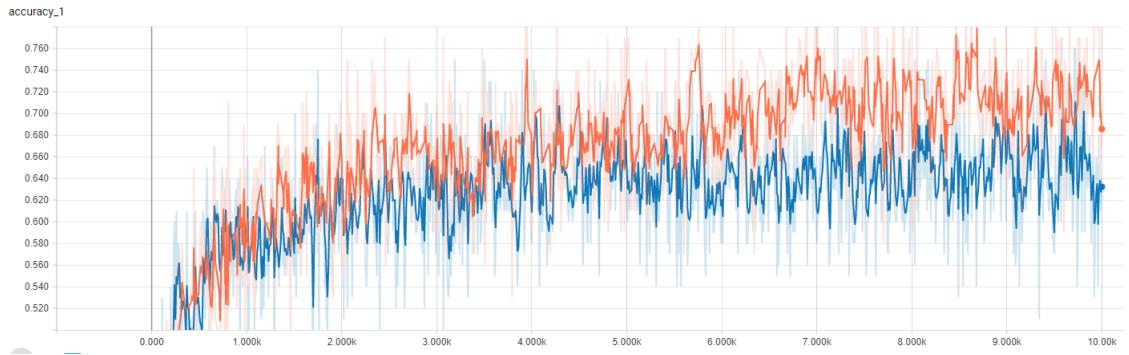


Figure 4.9: Comparison of accuracy

Table 4.1: Comparison of parameters

Parameter Tuning	Steps	Learning Rate	Test %	Validation %	Results
Configuration 1	8000	default	default	default	59.1%
Configuration 2	8000	0.1	default	default	65.8%
Configuration 3	10000	0.1	default	default	66.3%
Configuration 4	12000	0.1	default	default	66.6%
Configuration 5	10000	0.2	default	default	66.0%
Configuration 6	10000	0.1	15	15	66.3%

Results

The results of each set of parameters can be seen in Table 4.1. Using the parameters of 10,000 steps, and a learning rate of 0.1 and therefore a final test accuracy of 66.3%, the model achieved a Top 5 accuracy of 85.96%. This figure was calculated from 1090 images in the test dataset and the average probability of the predictions were at 0.62.

Empirical Analysis

The set of parameters that seem to be the most effective are 10000 steps with a 0.1 learning rate. Graphs of this model can be see in Figures 4.7 and 4.6. These are based on the validation set and then the test set respectively. A side by side comparison can also be seen in Figures 4.8 and 4.9 where orange is for during training and blue for the validation set.

There were two separate factors that each increased classification accuracy of about 5% each. These were training steps and learning rate.

Training steps are related to the number of images so before, when the training steps were at 4000, not all of our training images were being used. As the steps were increased twofold we saw a 3.8% increase in accuracy.

Another parameter that increased accuracy significantly was learning rate. The default learning rate is 0.01 which was increased to 0.1. This resulted in an increase of 6.7%. This is most likely due to the fact that since we are only looking at the last layer, we can afford to change the weights more significantly.

The Top 5 accuracy of the model was quite good but it was not run on the same number of images as the final test accuracy. This is due to the fact that tensorflow does not have an API for calculating Top 5 accuracy. As a result, it had to be calculated manually and this is very time intensive. The average probability of 0.62 is interestingly close to the overall Top 1 accuracy.

4.4 MobileNet

Overview

Due to the fact that the end goal for this project is to have a smartphone application that a user can use to keep track of their calorie measurement, there are a couple of options in how to achieve this. Firstly, an image can be taken on the phone and sent to a server to run a classification algorithm. Secondly, a model can be stored on the phone for computation. Transfer learning was once again used but on a different, smaller architecture called MobileNet Howard et al., (2017).

Network Architecture

The network architecture used for this experiment is MobileNet Howard et al., (ibid.). A parameter can be set to retrain this model in the code, as mentioned in previous experiments, instead of inception V-3 *How to Retrain Inception's Final Layer for New Categories* (n.d.). This architecture is designed to be smaller and much more efficient so that it can be used on smartphones which have less powerful resources available.

Mobilenet is designed for various different use cases on mobile phones such as image classification, landmark recognition, object detection and face attributes. It is based on "depthwise separable convolutions" whose purpose is the split a convolutional layer into a depthwise convolution and a pointwise convolution (1x1 convolution) Howard et al., (2017). The purpose of the former is to apply "a single filter to each input channel" Howard et al., (ibid.) while the latter combines these outputs. This in turn reduces the computational time. The first layer of the network is a standard convolutional layer and does not get split up.

Unfortunatley there is a significant decrease in accuracy in using Mobilenet as it is a much thinner model.

Dataset

The Food 101 dataset Bossard, Guillaumin, and Van Gool, (2014) with added classes was used for this experiment.

Libraries

Tensorflow and numpy.

Script

The retrain.py script *How to Retrain Inception’s Final Layer for New Categories* (n.d.) was used, with a different command parameter.

```
python tensorflow/examples/image_retraining/retrain.py \ --image_dir  
~/dataset_directory \ --architecture mobilenet_1.0_224 \  
--how_many_training_steps 10000 \ --learning_rate 0.1
```

Results

The final test accuracy of this model came to 50.2%.

Empirical Analysis

There was a decrease of 16.1% in this model to the highest accuracy from 4.4. This is due to the smaller architecture which is aimed to be faster and smaller with an expected decrease in accuracy.

4.5 Food 101 subset

Overview

In many of the papers that have been researched where food image classification was carried out, they attempted to classify a lot less than 108 food types as has been the case for experiments previously shown. Pouladzadeh, Villalobos, et al., (2012) used 12 classes, Pouladzadeh, Shirmohammadi, and Al-Maghribi, (2014) had 15, Abbirami.R.S et al., (2015) attempted to classify 11 classes of food, 20 classes were used in N. Chen et al., (2010) and Zhang et al., (2015) predicted 15 classes. Due to the lower number of classes in these papers, it was decided to retrain inception on a subset of the food-101 extended dataset to benchmark results. 13 classes were selected from food-101 for training.

Network Architecture

Retrained Inception model.

Dataset

A subset of the Food 101 dataset was used for this experiment Bossard, Guil- laumin, and Van Gool, (2014).

Results

A Final test Accuracy of 92.6% was recorded for this experiment which performs quite high in comparison to the data in Table 4.2. A Top 5 accuracy of 100% was calculated with an average prediction probability of 0.89 using 130 images. This was calculated using the script defined in 4.3.

Table 4.2: Accuracy of other studies

Reference	Classes	Accuracy
Novel SVM	12	92.6%
Measuring Calorie and Nutrition	15	90.41%
Large Scale Learning	11	78%
Toward Dietary Assessment	20	~91.67%
Snap-n-eat	15	85%

Empirical Analysis

It would make sense the accuracy of our model would increase when the number of classes are reduced as the margin of error is decreased. The Top 5 accuracy of the model was very successful. There were only 10 images per class tested though, totaling at 130 images so if the number of images increased, this accuracy would probably decrease slightly. If the code is run again, the same results cannot even be replicated for sure. On another run of this code a Top 5 accuracy of 96% was recorded. While this is still a very good accuracy, it does not compare to the first run of the script.

Chapter 5

Analysing the Trained Model

The purpose of this chapter is to analyse the models trained previously to gain understanding and applicability to the problem statement.

5.1 Sliding Window

Overview

In the previous experiments, the one-shot approach to food image classification has been looked at. That is, the model will give a prediction of the most likely food item in that image. This is a problem when there is a composite image i.e are multiple foods in an image, see Figure 5.1. There are a few options to combat this problem. Firstly, one could detect objects in the image, segment the image according to these objects and then run each segment through the model. A simple approach to this would be to segment the image into a number of sections and then run each section through the model. In order to follow the latter approach, a sliding window approach was used. This sliding window would move across the image and classify the segment of the image in the window. There are three options for window sliding shape as defined by a command line argument.

```
python sliding_window.py --image=~/image_dir --window_shape grid
```



Figure 5.1: Bowl of fruit

There are three options for window shape:

- Grid based window as per Figure 5.2
- Row based window as per Figure 5.3
- Column based window as per Figure 5.4

Libraries

For this experiment Tensorflow provided the classification of each segment while also helping with resizing along with Numpy. OpenCv was used to implement the sliding window as per *Sliding Windows for Object Detection with Python and OpenCV* (n.d.)

Script

There were four main elements to the script. Firstly, extracting a window to be classified. Secondly, resizing the image to be compatible with the Tensorflow

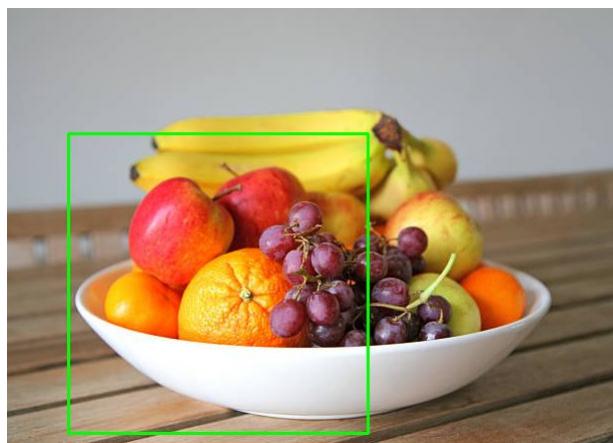


Figure 5.2: Grid based window

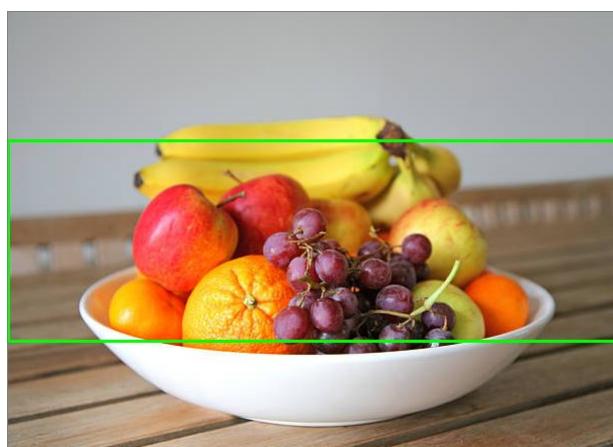


Figure 5.3: Row based window

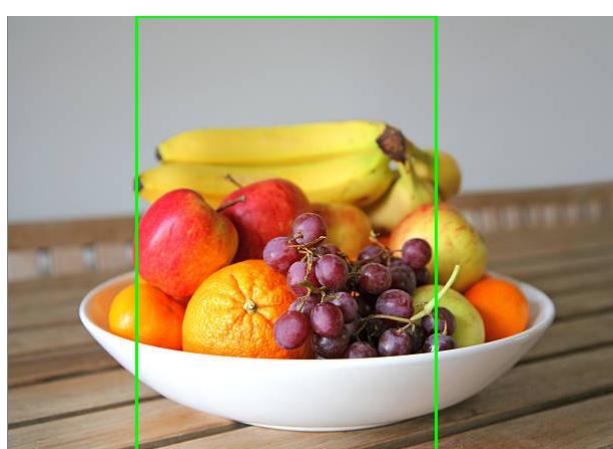


Figure 5.4: Column Based Window

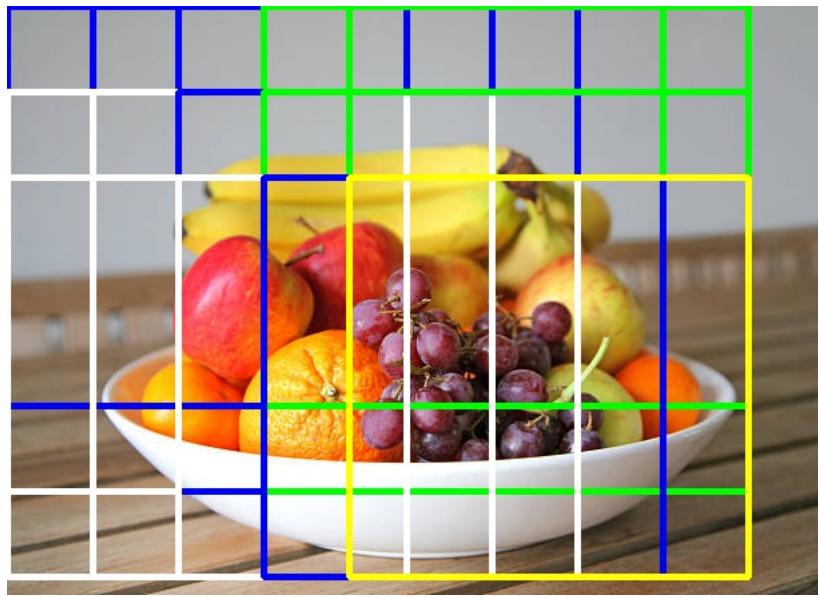


Figure 5.5: Fruit with Color Overlay

model. Thirdly, running the window through the Tensorflow model and finally saving a new image with a coloured overlay of classifications.

Extracting the window from the image

```

1  for (x, y, window) in sliding_window(resized, stepSize=64,
2      windowHeight=(winW, winH)):
3      # if the window does not meet our desired window size, ignore it
4      if window.shape[0] != winH or window.shape[1] != winW:
5          continue
6
7  def sliding_window(image, stepSize, windowHeight):
8      # slide a window across the image
9      for y in xrange(0, image.shape[0], stepSize):
10         for x in xrange(0, image.shape[1], stepSize):
11             # yield the current window
12             yield (x, y, image[y:y + windowHeight[1], x:x + windowHeight[0]])

```

Resizing the window

```
1 window = cv2.resize(window, (299, 299))

1 resized_image = tf.reshape(image, [1, input_height, input_width, 3])
2 resized = tf.image.resize_area(resized_image, [input_height, input_width])
3 normalized = tf.divide(tf.subtract(resized, [input_mean]), [input_std])
```

Running the Tensorflow model

```
1 with tf.Session() as sess:
2     numpy_image = sess.run(normalized)
3
4 with tf.Session(graph=graph) as sess:
5     results = sess.run(output_operation.outputs[0],
6                         {input_operation.outputs[0]: numpy_image})
7     probabilities = np.squeeze(results)
```

Saving the image with colour overlay

As seen in Figure 5.5, each square represents a window and each colour is for a different classification. Blue is for an apple, yellow for banana, green for grape, white for orange and black if an unexpected prediction is made.

```
1 cv2.rectangle(display_image, (x, y), (x + winW, y + winH),
2                 colour_dict.get(top1,
3 (0,0,0)), 4)
```

Results

Grid based window

The grid based window resulted in fifteen separate classification. As seen in Figure 5.1, there are multiple fruits in the image. Of these fruits, our model

Table 5.1: Grid Based Sliding Window Results

Food type	No. of Top-1 Classifications
Apple	5
Banana	1
Grape	4
Orange	5

Table 5.2: Row Based Sliding Window Results

Food type	No. of Top-1 Classifications
Apple	1
Banana	0
Grape	0
Orange	0
Other	3

is trained on four, apple, banana, orange and grapes. This method classified all four to Top-1 accuracy at least once each. This method took 42.8 seconds to run.

Row based window

The row based method resulted in four predictions as follows in Figure 5.2. Out of these four predictions, only one classified a known fruit at Top-1 accuracy,

Table 5.3: Column Based Sliding Window Results

Food type	No. of Top-1 Classifications
Apple	3
Banana	1
Grape	0
Orange	0
Other	1

an apple. An apple was also predicted to Top-5 accuracy in another instance. The runtime of this method was 16.1 seconds.

Column based window

The column based window approach had five total predictions and ran for a total of 13 seconds. As seen in 5.3, two out of four known fruits were classified to a Top-1 accuracy with all other fruits predicted to Top-5 accuracy. Only one Top-1 prediction did not contain a correct fruit.

Empirical Analysis

These results are very interesting because while a banana was only predicted to Top-1 accuracy once in grid based, once in column based and zero times in row based, if the whole image is ran through the model, a banana is at the Top-1 accuracy.

5.2 Recursive Refinement

Overview

After the sliding window code was run on Figure 5.1 in 5.1, it was observed that a sliding window was predicting grapes correctly in regions that contained a bunch of grapes. Since it would make sense that the model would be able detect an individual grape, it was decided that recursive refinement would be ran on a window that contained a grape. Due to the model requiring a 299 x 299 image size, the window could only be refined once as very small segments could not be resized up to 299 x 299. A window of 70 x 70 size was used.

Script

As you would think with recursive refinement, a recursive function would be used, but this was unnecessary due to image size restrictions. Instead, a conditional

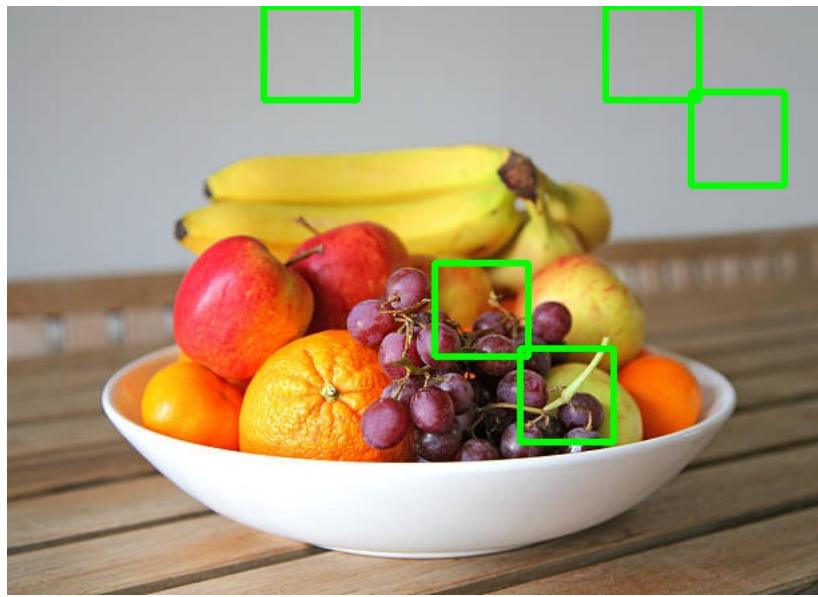


Figure 5.6: Recursive refinement 1

for loop was added to the existing code.

```
1 if top1 == "grape" and window_shape == "grid" and rr_grape:  
2     for (x_grape, y_grape, grape_window) in  
3         sliding_window(window_resized, stepSize=64, windowSize=(70,  
4             70)):  
5             #reshape to square  
6             grape_window_resized = cv2.resize(grape_window, (299, 299))  
7             top1_grape = subSample.classify(grape_window_resized,  
8                 window_shape)  
9             if top1_grape == "grape":  
10                 cv2.rectangle(display_image, (x_grape + x, y_grape + y),  
11                     (x_grape + x + 70, y_grape + y + 70),  
12                     colour_dict.get(top1, (0,0,0)), 4)  
13                 #cv2.imshow("Window", grape_window_resized)  
14                 cv2.waitKey(1)  
15             time.sleep(0.025)
```

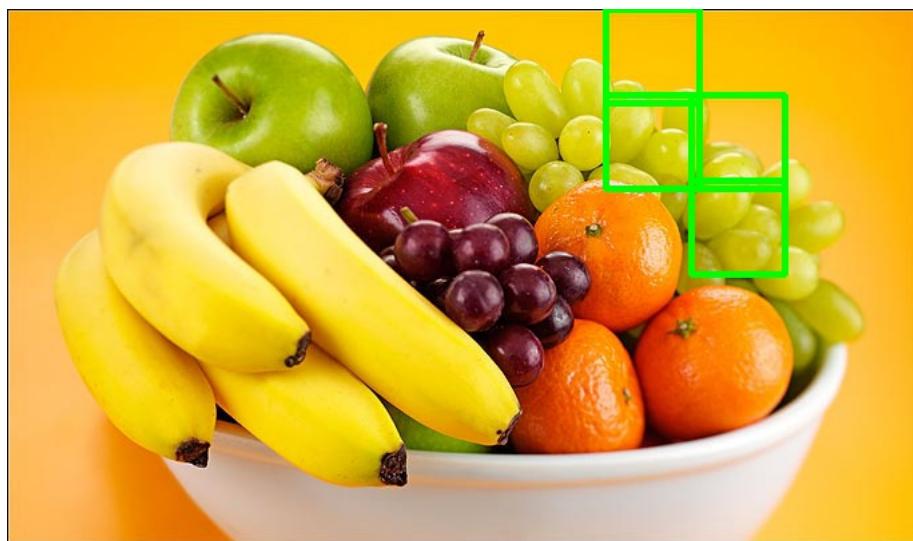


Figure 5.7: Recursive refinement 2

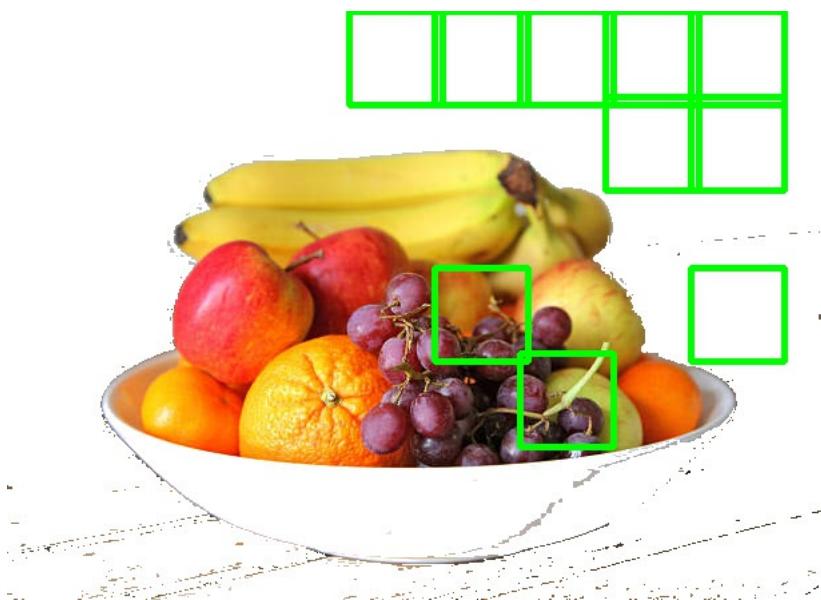


Figure 5.8: Recursive refinement 3



Figure 5.9: Bowl of fruit with background removed

Results

Some very interesting results were recorded on three separate images. Two new images are seen here which we will explore in future experiments. In all three images, while we are getting some expected predictions, grapes are being classified in locations that have nothing resembling a grape. These can be viewed in Figures 5.6, 5.7 and 5.8.

Empirical Analysis

The instances where false positives were recorded are mind boggling. An analysis to why this may be occurring is outlined in the section 5.5.

Table 5.4: Comparison of fruit image sliding window results with and without background

Food type	Grid	Row	Column	White Bg Grid	White Bg Row	White Bg Column
Apple	5	1	3	10	0	4
Banana	1	0	1	0	0	0
Grape	4	0	0	1	0	1
Orange	5	0	0	3	0	0
Other	0	3	1	1	4	0

5.3 Impact of Background

Overview

As we can see in 5.1, using sliding windows to classify many sections of an image, there were some cases where some unexpected predictions were made. Due to this, the decision was made to analyse the effect the background of the image on its classification. The sliding window code was then ran on a new image. This new image was the same fruit bowl as used previously but the background was filled in as white as per Figure 5.9.

Results

Grid

For grid based sliding window approach, the results turned out to be less successful than with the background. In this experiment, fourteen out of fifteen of top-1 classification were of an expected food type rather than fifteen out of fifteen with the background present. We expected the food types of apple, orange, grape and banana to appear in this image but while a banana was detected to a top-5 accuracy on a few occasions it was never predicted to a top-1 accuracy. The contrast between the image results can be seen in Table 5.4.

Table 5.5: Comparison of fruit bowl images

Food type	Grid	Row	Column	New Grid	New Row	New Column
Apple	5	1	3	4	1	0
Banana	1	0	1	5	0	5
Grape	4	0	0	2	1	0
Orange	5	0	0	0	0	0
Other	0	3	1	1	2	1

Row

The row based sliding window again had worse result than its counterpart, with zero out of four correct classifications as opposed to one. In this case, an orange appeared at top-5 accuracy once. The most common prediction was ice-cream which appeared at top-1 accuracy in three out of four instances.

Column

In contrast to our previous two methods of sliding window, this method outperformed its counterpart with correct predictions of all five windows while before we only had four out of five. In this experiment, an apple was predicted four times and a grape once, with all correct fruits appearing to top-5 accuracy.

Empirical Analysis

It is quite interesting that removing the background to the image reduced our accuracy overall. Many white foods were classified instead which makes sense due to the impact of colour expected.



Figure 5.10: Alternative Bowl of fruit

5.4 Alternative Test Image

Overview

In our previous sliding window oriented experiments, we had only used a single image. In order to see whether this image had biases unknown to us, another fruit bowl image had to be tested. This image was selected as fruit took up a larger portion of the image as seen in Figure 5.10.

Results

Grid

The performance of this experiment was slightly worse than with the previously used image. When the grid based sliding window was executed on Figure 5.10, fourteen out of fifteen predictions had an expected value. Out of the fourteen predictions orange was not predicted to top-1 accuracy at all. This can be seen, in comparison to previously used image, in Table 5.5.

Row

In the column based window for the new fruit image, the results were not very successful as has been the trend for most row based classification. Two out of four predictions had an expected value at top-1 accuracy.

Column

The column based approach had a similar result to its counterpart in that only one of its predictions was unexpected. Although, due to the size of the new image, another column was created and thus has a better overall accuracy.

Empirical Analysis

A possible reason that an orange was not classified in any of these images is because in Figure 5.10, a more mandarin food is displayed.

5.5 Analysing Results of Recursive Refinement Further

Overview

In 5.2, some interesting results were obtained when attempting to run recursive refinement on an image with grapes in it. Originally it was meant to see if individual grapes could be recognised. On evaluation of the images used for training, it was found that bunches of grapes were used for training, not individual grapes. Therefore the results expected were never feasibly going to be obtained. Even though this was the case, the script resulted in finding grapes in portions of the background, as seen in Figure 5.6. The purpose of this experiment is to look into why this is occurring. A new image of a fruit bowl (Figure 5.11), taken on a mobile phone was selected for this experiment. The image was used to run the script as defined in 5.2 but on two separate

models. The model with tuned parameters, trained on 108 classes, with a final test accuracy of 66.3% and the model trained on 13 classes, with a final accuracy of 92.6%.

Results

The outputs of the script can be seen in Figure 5.12 and Figure 5.13. The output using the model trained on 13 classes had less false positive predictions. As the script ran (on both models), the segments predicted as grapes were saved to disk and then were each classified using 'label_image.py'. The probabilities of the predictions were recorded with results as seen in Table 5.6.

Empirical Analysis

Comparing the two models

As we can see evidence of in Table 5.6, along with Figure 5.12 and Figure 5.13, the model trained on 13 classes predicted only four false positives while the larger model predicted 6. This is most likely due to the larger model not being able to separate grapes as well as the smaller model. It is worth mentioning that the grape class has the lowest number of training images of all the other classes.

Analysing Probabilities

The probabilities of predictions for grapes in these false positive classifications are all low, in both models. In contrast, some of the correctly classified segments were run through the model and the results all had probabilities in the high nineties. When these background segments are run through the model, some prediction has to be made and none of these false predictions have a probability of over 0.4. The average probability of these segments being grapes is in fact 0.257 (rounded to three decimal places). A counter measure to this prob-



Figure 5.11: Fruit image taken on a mobile phone

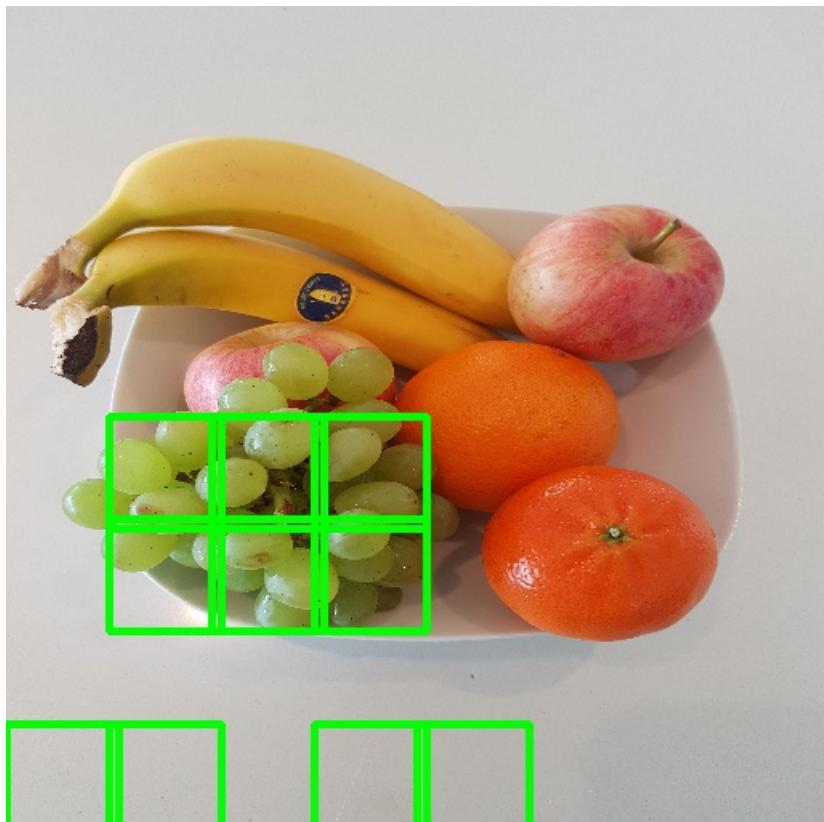


Figure 5.12: Image after sliding window - 13 classes

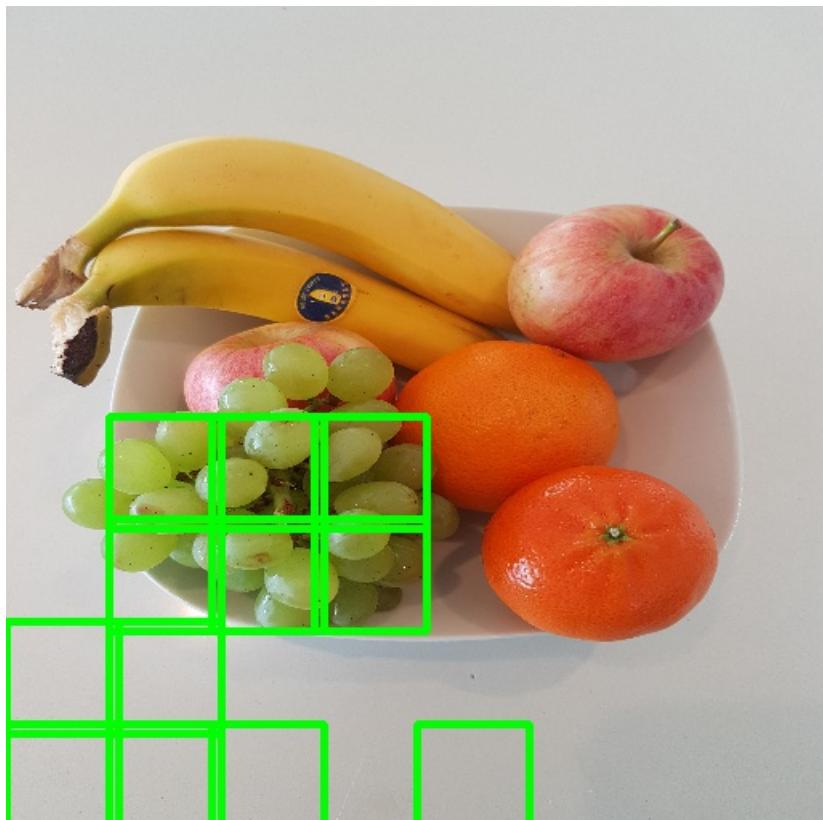


Figure 5.13: Image after sliding window - 108 classes

Table 5.6: Results of recursive refinement segment classifications using two models

Fruitbowl Segment	13 class model	108 class model
False Positive 1	grape 0.31476045 apple 0.25866747 orange 0.11228518 apple pie 0.10234257 chocolate cake 0.09062083	grape 0.31820437 panna cotta 0.091400646 macarons 0.08560496 roll 0.08188102 apple 0.03981942
False Positive 2	grape 0.29391274 apple 0.2371384 orange 0.16196558 chocolate cake 0.10039616 apple pie 0.096121624	panna cotta 0.21783036 grape 0.105718106 apple 0.087030284 macarons 0.07753955 orange 0.037592527
False Positive 3	grape 0.22402291 chocolate cake 0.16412543 apple 0.14523567 orange 0.12525927 apple pie 0.099948086	grape 0.2835378 panna cotta 0.088551655 macarons 0.06166248 orange 0.041896477 roll 0.040348493
False Positive 4	apple 0.37422368 grape 0.28592423 orange 0.10230055 chocolate cake 0.085367255 apple pie 0.053731006	grape 0.20476209 macarons 0.17529227 panna cotta 0.16139281 roll 0.04359601 orange 0.036991216
False Positive 5	N/A	grape 0.22574392 macarons 0.14840269 panna cotta 0.08670294 roll 0.08650105 orange 0.046900786
False Positive 6	N/A	grape 0.31369162 panna cotta 0.11833602 macarons 0.09010604 apple 0.057512555 roll 0.047116004



Figure 5.14: Banana Image Pre-Resolution

lem could be to disregard any predictions with a probability under a certain threshold, perhaps 0.4.

5.6 Scaling Down Images

Overview

In this experiment, it was decided to look into how scaling down the size of the image would influence its classification. Three images were scaled down to 10% of their size for this purpose, Figure 5.14, Figure 5.15 and Figure 5.16. Each of these images were run through the model before and after image resizing and the predictions were recorded.

Network Architecture

Inception V-3 architecture. Parameter tuned model was used for this.



Figure 5.15: Apple Pie Image Pre-Resolution



Figure 5.16: Pizza Image Pre-Resolution

Table 5.7: Results of Down Scaled Images

Food Image	Pre-Scaled Image	Scaled Image
Banana	Top 1 : 0.9971	Top 1 : 0.9995
Apple Pie	Top 1 : 0.7986	Top 5 : 0.0836
Pizza	Top 1 : 0.9753	Top 1 : 0.9705

Dataset

Food-101 extended dataset.

Results

The results of the images before and after were somewhat contrasting as seen in Table 5.7. The 'Top 1' notation in the table suggests that the food was predicted as the number one prediction. The decimal value is the probability of that food type being the correct classification.

Empirical Analysis

As we can see in Table 5.7, the banana (Figure 5.14) has a very similar prediction and probability. This would indicate that the classifier does not need clear quality images for bananas but maybe shape, texture and colour are important factors.

The apple pie (Figure 5.15) has drastically different results in Table 5.7. Originally the apple pie was correctly predicted with a probability of 0.7986 but after downscaling, the apple pie was the fifth food predicted with a likelihood of 0.0836.

The image of the pizza (Figure 5.16), had very similar results before and after resizing. This might indicate that pizza is determined more by shape, texture and colour than fine quality.



Figure 5.17: Banana Image Greyscale

5.7 Effect of Colour

Overview

As seen in the last experiment, it is possible that colour plays a significant part in the overall classification of some food types, along with shape and texture. Due to this, what would happen if colour was removed from the image? The three images, Figure 5.14, Figure 5.15 and Figure 5.16, were all converted to greyscale to test this and the resulting images can be seen in Figure 5.17, Figure 5.18 and Figure 5.19. Each of these images were ran through the model before and after converting to greyscale and the results recorded.

Network Architecture

Inception V-3 architecture. Parameter tuned model was used for this experiment.



Figure 5.18: Apple Pie Image Greyscale

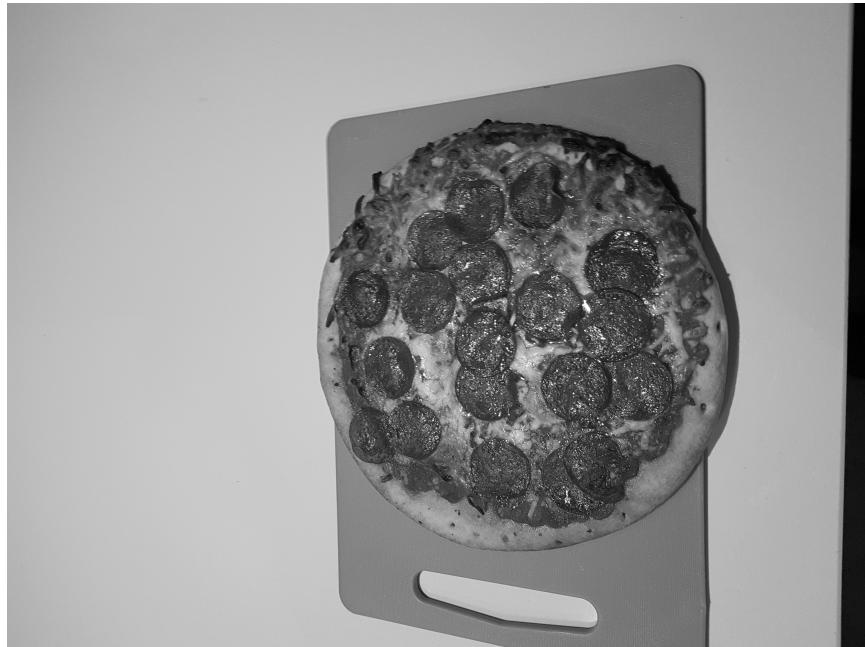


Figure 5.19: Pizza Image Greyscale

Table 5.8: Effect of Colour

Food Image	Pre-Scaled Image	Greyscale
Banana	Top 1 : 0.9971	Top 1 : 0.9986
Apple Pie	Top 1 : 0.7986	Top 1 : 0.3462
Pizza	Top 1 : 0.9753	Top 2 : 0.1299

Dataset

Food-101 extended dataset.

Results

The results for this experiment can be seen in Table 5.8. The 'Top 1' notation in the table suggests that the food was predicted as the number one prediction. The decimal value is the probability of that food type being the correct classification.

Empirical Analysis

Colour

In regards to the images of the banana (Figure 5.14 and Figure 5.17), there is very little difference in classification. They were both classified to a Top 1 accuracy and their probabilities differing by only 0.0015 with the grey scale image being of a higher probability interestingly enough. This would indicate that colour is not an important factor for classifying bananas and the coloured image may even have noise. This can be seen in Table 5.8.

The apple pie images (Figure 5.15 and Figure 5.18) were also both classified correctly but with a large gap in the likelihood of that classification being correct. As per Table 5.8, the coloured image had a probability of 0.7986 as opposed to one of 0.3462. This would indicate that while colour is important,

there is enough unique data from the rest of the image to result in correct classification.

Finally, the pizza images (Figure 5.16 and Figure 5.19) were the most contrasting in their results. While the original image was classified to Top 1 accuracy with a likelihood of 0.9753, the coloured image was classified to Top 2 accuracy while a probability of 0.1299. From this we can deduce that colour is vital for the classification of pizza.

Since this experiment only compared three images, we cannot say for sure if this analysis is biased or not.

Colour vs Scale

In the last experiment, it was seen that the pizza and banana images (Figure 5.16 and Figure 5.14) were not really affected by image quality. While this seems to be true, the image of apple pie, Figure 5.15, was greatly influenced by only being classified to a Top 5 accuracy when down scaled.

For the greyscale images, bananas were not effected greatly and neither was an apple pie but a pizza was greatly influenced.

From this, it is clear that the prominent unique features for each of these food types are different. The banana (Figure 5.14) may have focus on shape and texture, the pie (Figure 5.15) on image quality and some influence of colour while the pizza (Figure 5.16) may have a focus on shape, texture and colour.

5.8 Visualising Images Through the Network

Overview

Network Architecture

Dataset

API's

Script

Results

Empirical Analysis

Chapter 6

Prototype Application

In order to demonstrate practical use of the models trained in Chapter 4, a prototype application was developed. This application was developed for a smartphone running on the Android operating system in the Java programming language. A backend service was created to receive an image from the smartphone application, process the image using the CNN previously developed and return a response to the application in the form of a prediction.

6.1 Requirements

6.2 Design

The design process for this prototype application consisted of user interface design, android application system design and the activity of choosing resources for a backend host.

User Interface

The user interface mockups were created using (Fluid? insert reference). The aim of this design was to provide a simple, easy to use interface so that a user could use the application with minimal effort. The main benefit of using

computer vision for this type of application is to reduce the effort of the user or dietitian keeping track of food intake. Therefore, this application had to be very quick and easy to use.

Three main activities were needed for this application. The first activity, as seen in Figure 6.1, would consist of options to either take an image of food or to view the food logs of the user. If the user decides to take an image of food, they will be brought to the second activity which can be seen in Figure 6.2. The user is required to take a picture before reaching the activity in Figure 6.2. Once the user presses the send button the image is classified and the content of the activity changes to display the classification and calorie count as in Figure 6.3. Alternatively, the user can cancel the process and return to the first activity. The user would then submit the food classification for logging and automatically return to the first activity. The final activity is displayed when a user views their food logs. This activity has the ability to view a list of the food logs taken by the user by day, week or month. The calorie count of the selected time frame is to be displayed on this activity also.

System Architecture

Package Diagram

Architectural Patterns

The Model-View-Presenter (MVP)(reference) architectural pattern was adapted for this application. This pattern is very similar to the Model-View-Controller (MVC) architecture which is quite popular in the software development industry. The main difference between MVC and MVP is that whereas in MVC the controller is responsible for which view is used, in MVP the presenter is called through the view. This is due to the architecture of Android applications in general, where the view takes primary control.

(Talk more about MPV).

As illustrated in Figure 6.5, the views in the application have a dependency only on their corresponding presenter. The presenters contain the business

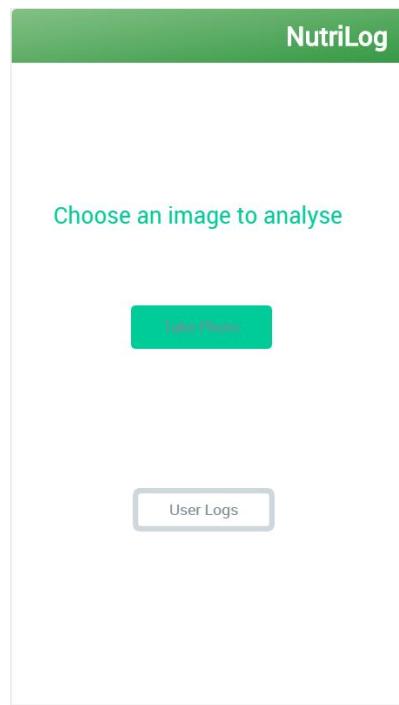


Figure 6.1: Landing Activity

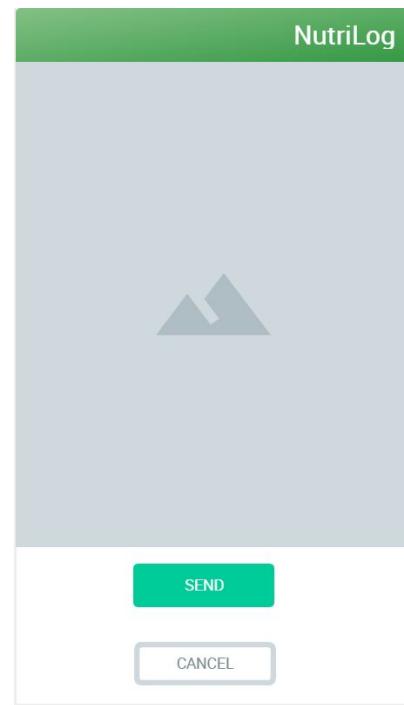


Figure 6.2: Image Submission Activity

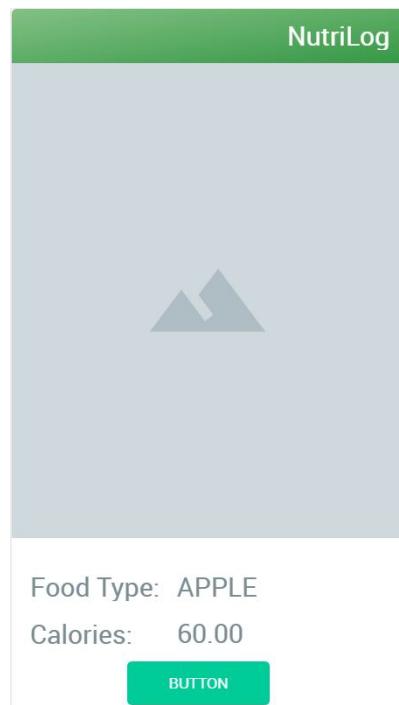


Figure 6.3: Classification Activity

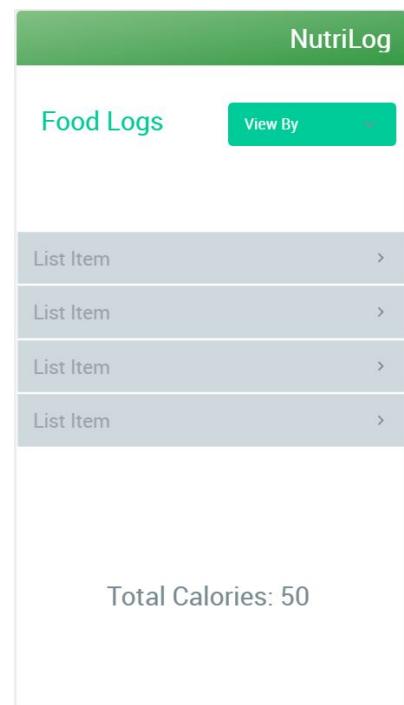


Figure 6.4: Food Logs Activity

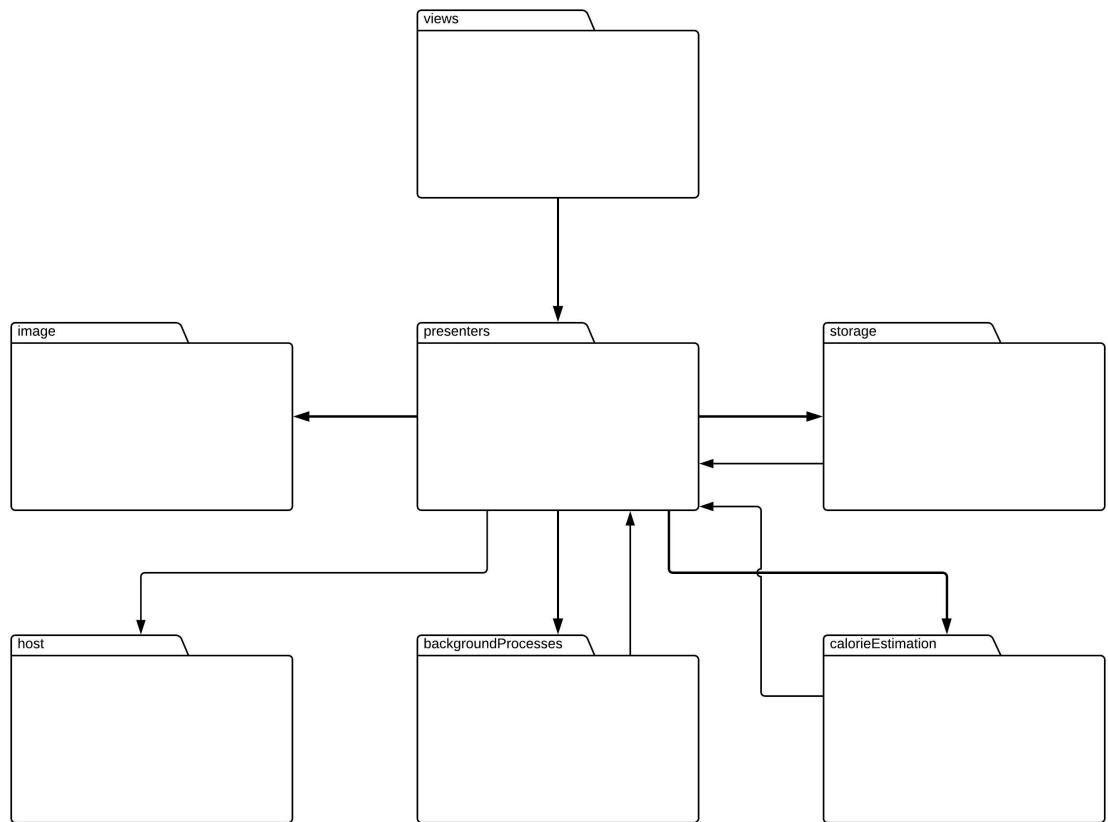


Figure 6.5: Package Diagram

logic of the application and have a dependency on all other packages. The packages of backgroundProcesses, calorieEstimation and storage have a dependency on the presenters package for callback purposes and also database creation in Android requires context of the activity information.

AWS Architecture

6.3 Implementation

Send image to server - AsyncTask

```
1     @Override
2     protected String doInBackground(Void... params) {
3         String result = "";
4         OkHttpClient client = new OkHttpClient();
5         String imageToSend = image;
6         RequestBody requestBody = new
7             MultipartBody.Builder()
8                 .setType(MultipartBody.FORM)
9                 .addFormDataPart("image",
10                   imageToSend)
11                 .build();
12
13
14         Request request = new
15             Request.Builder().url(host.getUrl())
16                 .post(requestBody).build();
17
18
19         Response response = null;
20         try {
21             response =
22                 client.newCall(request).execute();
23             result = response.body().string();
24             response.body().close();
25         } catch (IOException e) {
26             e.printStackTrace();
27         }
28     }
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
179
180
181
182
183
184
185
186
187
188
189
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
209
210
211
212
213
214
215
216
217
218
218
219
219
220
221
222
223
224
225
226
227
227
228
228
229
229
230
231
232
233
234
235
235
236
236
237
237
238
238
239
239
240
240
241
241
242
242
243
243
244
244
245
245
246
246
247
247
248
248
249
249
250
250
251
251
252
252
253
253
254
254
255
255
256
256
257
257
258
258
259
259
260
260
261
261
262
262
263
263
264
264
265
265
266
266
267
267
268
268
269
269
270
270
271
271
272
272
273
273
274
274
275
275
276
276
277
277
278
278
279
279
280
280
281
281
282
282
283
283
284
284
285
285
286
286
287
287
288
288
289
289
290
290
291
291
292
292
293
293
294
294
295
295
296
296
297
297
298
298
299
299
300
300
301
301
302
302
303
303
304
304
305
305
306
306
307
307
308
308
309
309
310
310
311
311
312
312
313
313
314
314
315
315
316
316
317
317
318
318
319
319
320
320
321
321
322
322
323
323
324
324
325
325
326
326
327
327
328
328
329
329
330
330
331
331
332
332
333
333
334
334
335
335
336
336
337
337
338
338
339
339
340
340
341
341
342
342
343
343
344
344
345
345
346
346
347
347
348
348
349
349
350
350
351
351
352
352
353
353
354
354
355
355
356
356
357
357
358
358
359
359
360
360
361
361
362
362
363
363
364
364
365
365
366
366
367
367
368
368
369
369
370
370
371
371
372
372
373
373
374
374
375
375
376
376
377
377
378
378
379
379
380
380
381
381
382
382
383
383
384
384
385
385
386
386
387
387
388
388
389
389
390
390
391
391
392
392
393
393
394
394
395
395
396
396
397
397
398
398
399
399
400
400
401
401
402
402
403
403
404
404
405
405
406
406
407
407
408
408
409
409
410
410
411
411
412
412
413
413
414
414
415
415
416
416
417
417
418
418
419
419
420
420
421
421
422
422
423
423
424
424
425
425
426
426
427
427
428
428
429
429
430
430
431
431
432
432
433
433
434
434
435
435
436
436
437
437
438
438
439
439
440
440
441
441
442
442
443
443
444
444
445
445
446
446
447
447
448
448
449
449
450
450
451
451
452
452
453
453
454
454
455
455
456
456
457
457
458
458
459
459
460
460
461
461
462
462
463
463
464
464
465
465
466
466
467
467
468
468
469
469
470
470
471
471
472
472
473
473
474
474
475
475
476
476
477
477
478
478
479
479
480
480
481
481
482
482
483
483
484
484
485
485
486
486
487
487
488
488
489
489
490
490
491
491
492
492
493
493
494
494
495
495
496
496
497
497
498
498
499
499
500
500
501
501
502
502
503
503
504
504
505
505
506
506
507
507
508
508
509
509
510
510
511
511
512
512
513
513
514
514
515
515
516
516
517
517
518
518
519
519
520
520
521
521
522
522
523
523
524
524
525
525
526
526
527
527
528
528
529
529
530
530
531
531
532
532
533
533
534
534
535
535
536
536
537
537
538
538
539
539
540
540
541
541
542
542
543
543
544
544
545
545
546
546
547
547
548
548
549
549
550
550
551
551
552
552
553
553
554
554
555
555
556
556
557
557
558
558
559
559
560
560
561
561
562
562
563
563
564
564
565
565
566
566
567
567
568
568
569
569
570
570
571
571
572
572
573
573
574
574
575
575
576
576
577
577
578
578
579
579
580
580
581
581
582
582
583
583
584
584
585
585
586
586
587
587
588
588
589
589
590
590
591
591
592
592
593
593
594
594
595
595
596
596
597
597
598
598
599
599
600
600
601
601
602
602
603
603
604
604
605
605
606
606
607
607
608
608
609
609
610
610
611
611
612
612
613
613
614
614
615
615
616
616
617
617
618
618
619
619
620
620
621
621
622
622
623
623
624
624
625
625
626
626
627
627
628
628
629
629
630
630
631
631
632
632
633
633
634
634
635
635
636
636
637
637
638
638
639
639
640
640
641
641
642
642
643
643
644
644
645
645
646
646
647
647
648
648
649
649
650
650
651
651
652
652
653
653
654
654
655
655
656
656
657
657
658
658
659
659
660
660
661
661
662
662
663
663
664
664
665
665
666
666
667
667
668
668
669
669
670
670
671
671
672
672
673
673
674
674
675
675
676
676
677
677
678
678
679
679
680
680
681
681
682
682
683
683
684
684
685
685
686
686
687
687
688
688
689
689
690
690
691
691
692
692
693
693
694
694
695
695
696
696
697
697
698
698
699
699
700
700
701
701
702
702
703
703
704
704
705
705
706
706
707
707
708
708
709
709
710
710
711
711
712
712
713
713
714
714
715
715
716
716
717
717
718
718
719
719
720
720
721
721
722
722
723
723
724
724
725
725
726
726
727
727
728
728
729
729
730
730
731
731
732
732
733
733
734
734
735
735
736
736
737
737
738
738
739
739
740
740
741
741
742
742
743
743
744
744
745
745
746
746
747
747
748
748
749
749
750
750
751
751
752
752
753
753
754
754
755
755
756
756
757
757
758
758
759
759
760
760
761
761
762
762
763
763
764
764
765
765
766
766
767
767
768
768
769
769
770
770
771
771
772
772
773
773
774
774
775
775
776
776
777
777
778
778
779
779
780
780
781
781
782
782
783
783
784
784
785
785
786
786
787
787
788
788
789
789
790
790
791
791
792
792
793
793
794
794
795
795
796
796
797
797
798
798
799
799
800
800
801
801
802
802
803
803
804
804
805
805
806
806
807
807
808
808
809
809
810
810
811
811
812
812
813
813
814
814
815
815
816
816
817
817
818
818
819
819
820
820
821
821
822
822
823
823
824
824
825
825
826
826
827
827
828
828
829
829
830
830
831
831
832
832
833
833
834
834
835
835
836
836
837
837
838
838
839
839
840
840
841
841
842
842
843
843
844
844
845
845
846
846
847
847
848
848
849
849
850
850
851
851
852
852
853
853
854
854
855
855
856
856
857
857
858
858
859
859
860
860
861
861
862
862
863
863
864
864
865
865
866
866
867
867
868
868
869
869
870
870
871
871
872
872
873
873
874
874
875
875
876
876
877
877
878
878
879
879
880
880
881
881
882
882
883
883
884
884
885
885
886
886
887
887
888
888
889
889
890
890
891
891
892
892
893
893
894
894
895
895
896
896
897
897
898
898
899
899
900
900
901
901
902
902
903
903
904
904
905
905
906
906
907
907
908
908
909
909
910
910
911
911
912
912
913
913
914
914
915
915
916
916
917
917
918
918
919
919
920
920
921
921
922
922
923
923
924
924
925
925
926
926
927
927
928
928
929
929
930
930
931
931
932
932
933
933
934
934
935
935
936
936
937
937
938
938
939
939
940
940
941
941
942
942
943
943
944
944
945
945
946
946
947
947
948
948
949
949
950
950
951
951
952
952
953
953
954
954
955
955
956
956
957
957
958
958
959
959
960
960
961
961
962
962
963
963
964
964
965
965
966
966
967
967
968
968
969
969
970
970
971
971
972
972
973
973
974
974
975
975
976
976
977
977
978
978
979
979
980
980
981
981
982
982
983
983
984
984
985
985
986
986
987
987
988
988
989
989
990
990
991
991
992
992
993
993
994
994
995
995
996
996
997
997
998
998
999
999
1000
1000
1001
1001
1002
1002
1003
1003
1004
1004
1005
1005
1006
1006
1007
1007
1008
1008
1009
1009
1010
1010
1011
1011
1012
1012
1013
1013
1014
1014
1015
1015
1016
1016
1017
1017
1018
1018
1019
1019
1020
1020
1021
1021
1022
1022
1023
1023
1024
1024
1025
1025
1026
1026
1027
1027
1028
1028
1029
1029
1030
1030
1031
1031
1032
1032
1033
1033
1034
1034
1035
1035
1036
1036
1037
1037
1038
1038
1039
1039
1040
1040
1041
1041
1042
1042
1043
1043
1044
1044
1045
1045
1046
1046
1047
1047
1048
1048
1049
1049
1050
1050
1051
1051
1052
1052
1053
1053
1054
1054
1055
1055
1056
1056
1057
1057
1058
1058
1059
1059
1060
1060
1061
1061
1062
1062
1063
1063
1064
1064
1065
1065
1066
1066
1067
1067
1068
1068
1069
1069
1070
1070
1071
1071
1072
1072
1073
1073
1074
1074
1075
1075
1076
1076
1077
1077
1078
1078
1079
1079
1080
1080
1081
1081
1082
1082
1083
1083
1084
1084
1085
1085
1086
1086
1087
1087
1088
1088
1089
1089
1090
1090
1091
1091
1092
1092
1093
1093
1094
1094
1095
1095
1096
1096
1097
1097
1098
1098
1099
1099
1100
1100
1101
1101
1102
1102
1103
1103
1104
1104
1105
1105
1106
1106
1107
1107
1108
1108
1109
1109
1110
1110
1111
1111
1112
1112
1113
1113
1114
1114
1115
1115
1116
1116
1117
1117
1118
1118
1119
1119
1120
1120
1121
1121
1122
1122
1123
1123
1124
1124
1125
1125
1126
1126
1127
1127
1128
1128
1129
1129
1130
1130
1131
1131
1132
1132
1133
1133
1134
1134
1135
1135
1136
1136
1137
1137
1138
1138
1139
1139
1140
1140
1141
1141
1142
1142
1143
1143
1144
1144
1145
1145
1146
1146
1147
1147
1148
1148
1149
1149
1150
1150
1151
1151
1152
1152
1153
1153
1154
1154
1155
1155
1156
1156
1157
1157
1158
1158
1159
1159
1160
1160
1161
1161
1162
1162
1163
1163
1164
1164
1165
1165
1166
1166
1167
1167
1168
1168
1169
1169
1170
1170
1171
1171
1172
1172
1173
1173
1174
1174
1175
1175
1176
1176
1177
1177
1178
1178
1179
1179
1180
1180
1181
1181
1182
1182
1183
1183
1184
1184
1185
1185
1186
1186
1187
1187
1188
1188
1189
1189
1190
1190
1191
1191
1192
1192
1193
1193
1194
1194
1195
1195
1196
1196
1197
1197
1198
1198
1199
1199
1200
1200
1201
1201
1202
1202
1203
1203
1204
1204
1205
1205
1206
1206
1207
1207
1208
1208
1209
1209
1210
1210
1211
1211
1212
1212
1213
1213
1214
1214
1215
1215
1216
1216
1217
1217
1218
1218
1219
1219
1220
1220
1221
1221
1222
1222
1223
1223
1224
1224
1225
1225
1226
1226
1227
1227
1228
1228
1229
1229
1230
1230
1231
1231
1232
1232
1233
1233
1234
1234
1235
1235
1236
1236
1237
1237
1238
1238
1239
1239
1240
1240
1241
1241
1242
1242
1243
1243
1244
1244
1245
1245
1246
1246
1247
1247
1248
1248
1249
1249
1250
1250
1251
1251
1252
1252
1253
1253
1254
1254
1255
1255
1256
1256
1257
1257
1258
1258
1259
1259
1260
1260
1261
1261
1262
1262
1263
1263
1264
1264
1265
1265
1266
1266
1267
1267
1268
1268
1269
1269
1270
1270
1271
1271
1272
1272
1273
1273
1274
1274
1275
1275
1276
1276
1277
1277
1278
1278
1279
1279
1280
1280
1281
1281
1282
1282
1283
1283
1284
1284
1285
1285
1286
1286
1287
1287
1288
1288
1289
1289
1290
1290
1291
1291
1292
1292
1293
1293
1294
1294
1295
1295
1296
1296
1297
1297
1298
1298
1299
1299
1300
1300
1301
1301
1302
1302
1303
1303
1304
1304
1305
1305
1306
1306
1307
1307
1308
1308
1309
1309
1310
1310
1311
1311
1312
1312
1313
1313
1314
1314
1315
1315
1316
1316
1317
1317
1318
1318
1319
1319
1320
1320
1321
1321
1322
1322
1323
1323
1324
1324
1325
1325
1326
1326
1327
1327
1328
1328
1329
1329
1330
1330
1331
1331
1332
1332
1333
1333
1334
1334
1335
1335
1336
1336
1337
1337
1338
1338
1339
1339
1340
1340
1341
1341
1342
1342
1343
1343
1344
1344
1345
1345
1346
1346
1347
1347
1348
1348
1349
1349
1350
1350
1351
1351
1352
1352
1353
1353
1354
1354
1355
1355
1356
1356
1357
1357
1358
1358
1359
1359
1360
1360
1361
1361
1362
1362
1363
1363
1364
1364
1365
1365
1366
1366
1367
1367
1368
1368
1369
1369
1370
1370
1371
1371
1372
1372
1373
1373
1374
1374
1375
1375
1376
1376
1377
1377
1378
1378
1379
1379
1380
1380
1381
1381
1382
1382
1383
1383
1384
1384

```

```
19     } catch (IOException e) {
20         e.printStackTrace();
21     }
22
23     return result;
24 }
```

Get Calories from API - AsyncTask

```
1     @Override
2     protected String doInBackground(Void... params) {
3         String result = "";
4         OkHttpClient client = new OkHttpClient();
5
6         Request request = new
7             Request.Builder().url(url)
8                 .build();
9
10        Response response = null;
11        try {
12            response =
13                client.newCall(request).execute();
14            result = response.body().string();
15            response.body().close();
16        } catch (IOException e) {
17            e.printStackTrace();
18        }
19        try {
20            JSONObject jsonObject = new
21                JSONObject(result);
22            JSONArray jsonArray =
23                jsonObject.getJSONArray("hits");
24        }
```

```

20         JSONObject hits =
21             jsonArray.getJSONObject(0);
22         JSONObject fields =
23             hits.getJSONObject("fields");
24         String calories =
25             fields.getString("nf_calories");
26         result = calories;
27     } catch (JSONException e) {
28         e.printStackTrace();
29     }
30     return result;
31 }
```

Encode Bitmap to base64 String

```

1     public String toBase64() {
2         BitmapFactory.Options options = new
3             BitmapFactory.Options();
4         options.inSampleSize = 8;
5
6         Bitmap imageBitmap =
7             BitmapFactory.decodeFile(image.getAbsolutePath(),
8                 options);
9         Bitmap.createScaledBitmap(imageBitmap, 300,
10            400, false);
11        ByteArrayOutputStream byteArrayOutputStream
12            = new ByteArrayOutputStream();
13        imageBitmap.compress(Bitmap.CompressFormat.PNG,
14            100, byteArrayOutputStream);
15        byte[] byteArray =
16            byteArrayOutputStream.toByteArray();
17
18    }
```

```
11         return Base64.encodeToString(byteArray,
12             Base64.DEFAULT);
13     }
```

DAOFactory class

```
1 public class DAOFactory {
2
3     public DAO getDAO(Context context){
4         return SqlLiteDAO.getInstance(context);
5     }
6 }
```

DAO interface

```
1 public interface DAO {
2     void addFoodLog(FoodLog foodLog);
3     List<FoodLog> getLogsByDay(Date date);
4     List<FoodLog> getLogsByWeek(Date date);
5     List<FoodLog> getLogsByMonth(Date date);
6     void deleteFoodLogs(List<FoodLog> foodLogs);
7     void deleteDb();
8 }
```

HostFactory class

```
1 public class HostFactory {
2
3     public Host createHost() {
4         Host host = new
5             Host.HostBuilder("52.214.205.157")
```

```
5           .withDns("ec2-52-214-205-157.eu-west-1.compu
6           .withPort(5000)
7           .withRoute("/classifyImage/")
8           .build();
9
10      return host;
11  }
```

HostBuilder

```
1      public static class HostBuilder {
2
3          private final String ipv4;
4
5          private String dns;
6
7          private int port;
8
9          private String route;
10
11
12         public HostBuilder(String ipv4) {
13             this.ipv4 = ipv4;
14         }
15
16         public HostBuilder withDns(String dns) {
17             this.dns = dns;
18             return this;
19         }
20
21         public HostBuilder withPort(int port) {
22             this.port = port;
23             return this;
24         }
25
26         public HostBuilder withRoute(String route) {
27             this.route = route;
28         }
29
30     }
31
32     public static void main(String[] args) {
33         HostBuilder hostBuilder = new HostBuilder()
34             .withDns("ec2-52-214-205-157.eu-west-1.compu
35             .withPort(5000)
36             .withRoute("/classifyImage/")
37             .build();
38
39         System.out.println(hostBuilder);
40     }
41
42 }
```

```
23         return this;
24     }
25
26     public Host build() {
27         return new Host(this);
28     }
29
30 }
```

Singleton DAO Object

```
1     public static DAO getInstance(Context
2         context) {
3             if(instance == null) {
4                 instance = new SqlLiteDAO(context);
5             }
6             return instance;
7 }
```

Main Activity

```
1 public class MainActivity extends AppCompatActivity {
2
3     private MainPresenter mainPresenter;
4
5     @Override
6     protected void onCreate(Bundle
7         savedInstanceState) {
8         super.onCreate(savedInstanceState);
9         setContentView(R.layout.activity_main);
```

```
10         mainPresenter = new MainPresenter(this);
11     }
12
13     public void onClickTakePhoto(View view) {
14         mainPresenter.takePhoto();
15     }
16
17     public void onClickUserLogs(View view) {
18         mainPresenter.userLogs();
19     }
20 }
```

Presenter for Main Activity

```
1 public class MainPresenter {
2
3     private Intent intent;
4     private Context context;
5
6     public MainPresenter(Context context) {
7         this.context=context;
8     }
9
10    public void takePhoto() {
11        intent = new Intent(context,
12                            CaptureImageActivity.class);
13        context.startActivity(intent);
14    }
15
16    public void userLogs() {
17        intent = new Intent(context,
18                            FoodLogsActivity.class);
```

```

17         context.startActivity(intent);
18     }
19
20 }

```

6.4 Testing

Junit

Mockito

Unit Tests

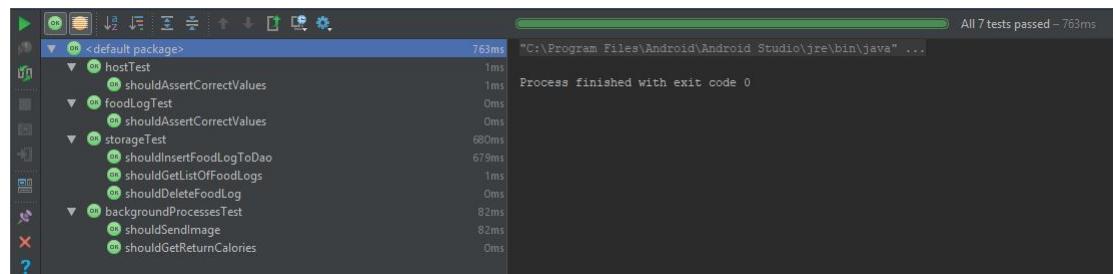


Figure 6.6: Automated Unit Tests

Test to ensure FoodLog object stores information correctly

```

1      @Before
2      public void setUp() {
3          foodLog = new FoodLogImpl(0, "banana",
4              59.03, new Date());
5
6      @Test
7      public void shouldAssertCorrectValues() {

```

```
8         assert food.equals(foodLog.getFood());
9         assert calories == foodLog.getCalories();
10        assert
11            timestamp.equals(foodLog.getTimestamp());
12        assert id == foodLog.getId();
13    }
```

Test to ensure Host object stores information correctly

```
1     @Before
2
3     public void setUp() {
4         host = new
5             Host.HostBuilder("52.214.205.157")
6                 .withDns("ec2-52-214-205-157.eu-west-1.c
7
8         }
9
10    @Test
11
12    public void shouldAssertCorrectValues() {
13        assert this.ip.equals(host.getIpv4());
14        assert this.dns.equals(host.getDns());
15        assert this.port == host.getPort();
16        assert this.route.equals(host.getRoute());
17        assert this.urlString.equals(host.getUrl());
18    }
```

Test to ensure DAO object stores and removes data correctly

```
1     @Before
2
3     public void setup() {
```

```

3         dao = mock(SqlLiteDAO.class);
4         date = new Date();
5         foodLog = new FoodLogImpl(0, "test", 0.0,
6                         date);
7     }
8
9     @Test
10    public void shouldGetListOfFoodLogs() {
11        foodLogs = dao.getLogsByDay(date);
12        foodLogs = dao.getLogsByWeek(date);
13        foodLogs = dao.getLogsByMonth(date);
14    }
15
16    @Test
17    public void shouldInsertFoodLogToDao() {
18        dao.addFoodLog(foodLog);
19        setTestLog();
20        assert foodLogs.contains(testLog);
21    }
22
23    @Test
24    public void shouldDeleteFoodLog() {
25        List<FoodLog> foodLogs = new ArrayList<>();
26        foodLogs.add(testLog);
27        dao.deleteFoodLogs(foodLogs);
28        setTestLog();
29        assert testLog.equals(null);
30    }
31
32    private void setTestLog() {
33        foodLogs = dao.getLogsByDay(date);
34        foodLogs.forEach(f -> {

```

```
34         if(f.getFood().equals("test")) {
35             testLog = f;
36         }
37     });
38 }
```

6.5 Backend

Python Flask

Implementation

Backend service code

```
1 import base64
2 import label_image
3 import time
4
5 app = Flask(__name__)
6
7 @app.route('/classifyImage/', methods=['POST'])
8 def classify():
9     imgdata = base64.b64decode(request.form.get('image'))
10    filename = 'images/' + str(time.strftime("%Y%m%d-%H%M%S")) +
11        '.jpg'
12    with open(filename, 'wb') as f:
13        f.write(imgdata)
14
15    results = label_image.runModel(filename)
16    predictions = results[0]
17
18    resultString = """
```

```

18     for i in predictions :
19         result_String += str(i) + ","
20     return result_String
21
22 if __name__ == '__main__':
23     app.run(host='0.0.0.0', threaded=True)

```

Method called from 'label_image.py'

```

1 def runModel(file_name):
2     model_file = \
3         "models/output_graph13.pb"
4     label_file = "labels/output_labels13.txt"
5     input_height = 299
6     input_width = 299
7     input_mean = 128
8     input_std = 128
9     input_layer = "Mul"
10    output_layer = " final_result "
11
12    graph = load_graph(model_file)
13    t = read_tensor_from_image_file(file_name,
14                                    input_height=input_height,
15                                    input_width=input_width,
16                                    input_mean=input_mean,
17                                    input_std=input_std)
18
19    input_name = "import/" + input_layer
20    output_name = "import/" + output_layer
21    input_operation = graph.get_operation_by_name(input_name)
22    output_operation = graph.get_operation_by_name(output_name)
23
24    with tf.Session(graph=graph) as sess:

```

```

25     results = sess.run(output_operation.outputs[0],
26                         {input_operation.outputs[0]: t})
27     results = np.squeeze(results)
28
29     top_k = results.argsort() [-5:][::-1]
30     labels = load_labels( label_file )
31
32     setIndex = False
33
34     top5_results = [None] * 5
35     index = 0
36     for i in top_k:
37         top5_results [index] = labels[i]
38         index += 1
39
40     final_results = [top5_results , results ]
41     return final_results

```

Deployment

6.6 Resources Utilised

Chapter 7

Discussion and Conclusion

Bibliography

- Abbirami.R.S et al. (2015). “Large Scale Learning for Food Image Classification”. In: *International Journal on Recent and Innovation Trends in Computing and Communication* 3.3, pp. 973–978. URL: <http://dx.doi.org/10.17762/ijritcc2321-8169.150317>.
- About - OpenCV library.* <https://opencv.org/about.html>. Accessed: 2017-02-18.
- Aurélien, Géron (2017). *Hands-On Machine Learning with Scikit-Learn & TensorFlow*.
- Bossard, Lukas, Matthieu Guillaumin, and Luc Van Gool (2014). “Food-101 – Mining Discriminative Components with Random Forests”. In: *European Conference on Computer Vision*.
- Build an Image Dataset in Tensorflow.* https://github.com/aymericdamien/TensorFlow-Examples/blob/master/examples/5_DataManagement/build_an_image_dataset.py. Accessed: 2017-02-15.
- Chen, Mei-Yun et al. (2012). “Automatic chinese food identification and quantity estimation”. In: *SIGGRAPH Asia 2012 Technical Briefs*. ACM, p. 29.
- Chen, Nicholas et al. (2010). “Toward dietary assessment via mobile phone video cameras”. In: *AMIA Annual Symposium Proceedings*. Vol. 2010. American Medical Informatics Association, p. 106.
- Complete Guide to Tensorflow for Deep Learning with Python.* <https://www.udemy.com/complete-guide-to-tensorflow-for-deep-learning-with-python/learn/v4/overview>. Accessed: 2017-02-15.
- Deng, J. et al. (2009). “ImageNet: A Large-Scale Hierarchical Image Database”. In: *CVPR09*.

- Donahue, Jeffrey (2017). “Transferrable Representations for Visual Recognition”. PhD thesis. EECS Department, University of California, Berkeley.
- Girshick, Ross (2015). “Fast R-CNN”. In: *Proceedings of the International Conference on Computer Vision (ICCV)*.
- Girshick, Ross et al. (2014). “Rich feature hierarchies for accurate object detection and semantic segmentation”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- He, Kaiming et al. (2017). “Mask R-CNN”. In: *arXiv preprint arXiv:1703.06870*.
- He, Ye et al. (2013). “Food image analysis: Segmentation, identification and weight estimation”. In: *Multimedia and Expo (ICME), 2013 IEEE International Conference on*. IEEE, pp. 1–6.
- Hoiem, Derek, Yodsawalai Chodpathumwan, and Qieyun Dai (2012). “Diagnosing Error in Object Detectors”. In: *Computer Vision – ECCV 2012: 12th European Conference on Computer Vision, Florence, Italy, October 7-13, 2012, Proceedings, Part III*. Ed. by Andrew Fitzgibbon et al. Berlin, Heidelberg: Springer Berlin Heidelberg, pp. 340–353.
- How to Retrain Inception’s Final Layer for New Categories.* https://www.tensorflow.org/tutorials/image_retraining. Accessed: 2018-01-24.
- Howard, Andrew G et al. (2017). “Mobileneets: Efficient convolutional neural networks for mobile vision applications”. In: *arXiv preprint arXiv:1704.04861*.
- Image Segmentation Using DIGITS 5.* <https://devblogs.nvidia.com/image-segmentation-using-digits-5/>. Accessed: 2017-03-06.
- Joutou, Taichi and Keiji Yanai (2009). “A food image recognition system with multiple kernel learning”. In: *Image Processing (ICIP), 2009 16th IEEE International Conference on*. IEEE, pp. 285–288.
- Kagaya, Hokuto, Kiyoharu Aizawa, and Makoto Ogawa (2014). “Food detection and recognition using convolutional neural network”. In: *Proceedings of the 22nd ACM international conference on Multimedia*. ACM, pp. 1085–1088.
- Kawano, Yoshiyuki and Keiji Yanai (2014). “Food image recognition with deep convolutional features”. In: *Proceedings of the 2014 ACM Interna-*

- tional Joint Conference on Pervasive and Ubiquitous Computing: Adjunct Publication.* ACM, pp. 589–593.
- Krizhevsky, Alex, Vinod Nair, and Geoffrey Hinton. “CIFAR-10 (Canadian Institute for Advanced Research)”. In: URL: <http://www.cs.toronto.edu/~kriz/cifar.html>.
- LeCun, Yann and Corinna Cortes (2010). “MNIST handwritten digit database”. In: URL: <http://yann.lecun.com/exdb/mnist/>.
- Liu, Chang et al. (2016). “Deepfood: Deep learning-based food image recognition for computer-aided dietary assessment”. In: *International Conference on Smart Homes and Health Telematics*. Springer, pp. 37–48.
- Luger, George F. (2005). *Artificial Intelligence. Structures and Strategies for Complex Problem Solving/Luger GF*. Pearson Education Limited.
- Mao, Dongyuan, Qian Yu, and Jingfan Wang. “Deep Learning Based Food Recognition”. In:
- Marsland, Stephen (2015). *Machine learning: an algorithmic perspective*. CRC press, pp. 71–110.
- Mezgec, Simon and Barbara Koroušić Seljak (2017). “Nutrinet: A deep learning food and drink image recognition system for dietary assessment”. In: *Nutrients* 9.7, p. 657.
- Mitchell, Tom M. (1997). *Machine Learning*. International Edition 1997. New York: The McGraw-Hill Companies, Inc., pp. 81–126.
- Pouladzadeh, Parisa, Shervin Shirmohammadi, and Rana Al-Maghribi (2014). “Measuring calorie and nutrition from food image”. In: *IEEE Transactions on Instrumentation and Measurement* 63.8, pp. 1947–1956.
- Pouladzadeh, Parisa, Shervin Shirmohammadi, and Abdulsalam Yassine (2014). “Using graph cut segmentation for food calorie measurement”. In: *Medical Measurements and Applications (MeMeA), 2014 IEEE International Symposium on*. IEEE, pp. 1–6.
- Pouladzadeh, Parisa, Gregorio Villalobos, et al. (2012). “A novel SVM based food recognition method for calorie measurement applications”. In: *Multi-media and Expo Workshops (ICMEW), 2012 IEEE International Conference on*. IEEE, pp. 495–498.

- Ren, Shaoqing et al. (2015). “Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks”. In: *Neural Information Processing Systems (NIPS)*.
- Schap, TE et al. (2014). “Merging dietary assessment with the adolescent lifestyle”. In: *Journal of human nutrition and dietetics* 27.s1, pp. 82–88.
- Shelhamer, Evan, Jonathan Long, and Trevor Darrell (2017). “Fully Convolutional Networks for Semantic Segmentation”. In: *IEEE Trans. Pattern Anal. Mach. Intell.* 39.4, pp. 640–651.
- Sliding Windows for Object Detection with Python and OpenCV*. <https://www.pyimagesearch.com/2015/03/23/sliding-windows-for-object-detection-with-python-and-opencv/>. Accessed: 2018-01-30.
- Support Vector Machines for Machine Learning*. <https://machinelearningmastery.com/support-vector-machines-for-machine-learning/>. Accessed: 2017-12-20.
- 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.
- Villalobos, Gregorio, Rana Almaghrabi, Behnoosh Hariri, et al. (2011). “A personal assistive system for nutrient intake monitoring”. In: *Proceedings of the 2011 international ACM workshop on Ubiquitous meta user interfaces*. ACM, pp. 17–22.
- Villalobos, Gregorio, Rana Almaghrabi, Parisa Pouladzadeh, et al. (2012). “An image procesing approach for calorie intake measurement”. In: *Medical Measurements and Applications Proceedings (MeMeA), 2012 IEEE International Symposium on*. IEEE, pp. 1–5.
- Yanai, Keiji and Yoshiyuki Kawano (2015). “Food image recognition using deep convolutional network with pre-training and fine-tuning”. In: *Multi-media & Expo Workshops (ICMEW), 2015 IEEE International Conference on*. IEEE, pp. 1–6.
- Zhang, Weiyu et al. (2015). “Snap n Eat - Food Recognition and Nutrition Estimation on a Smartphone”. In: *Journal of diabetes science and technology* 9.3, pp. 525–533.

Zhu, Fengqing et al. (2011). “Segmentation assisted food classification for dietary assessment”. In: *Computational Imaging IX*. Vol. 7873. International Society for Optics and Photonics, 78730B.

Appendix A

Appendix

A.1 Image Segmentation

Fully Convolutional Neural Networks for Semantic Segmentation

There has been a very interesting paper from UC Berkeley focused on using Full Convolutional Networks for semantic segmentation Shelhamer, Long, and Darrell, (2017). Fully Convolutional Networks (FCN) do not have any fully connected layers. They are replaced with more filtering layers. Nvidia Digits have a semantic segmentation implementation based off the work of this paper.

They took this approach because "feedforward computation and backpropagation are much more efficient when computer layer-by-layer over an entire image instead of independently patch-by-patch" Shelhamer, Long, and Darrell, (ibid.). This was also because they were focused on object detection. Normal classifiers do not work very well when they are to classify more than one subject in an image and image segmentation was a way to solve this.

There are a set of steps you can follow to turn a CNN into a FCN for semantic segmentation as follows ie. change to a convolutional layer from a fully connected one:

- The size of the filters must be set to the size of the input layers.
- For every neuron in the fully connected layer, have a filter.

Segmentation

Graph Based Segmentation

Graph cut segmentation has been used extensively in image segmentation. OpenCV has an implementation of a graph cut algorithm called grabcut which has been used to segment food on occasion Pouladzadeh, Shirmohammadi, and Yassine, (2014).

Table A.1: FCN Results Shelhamer, Long, and Darrell, (2017)

	FCN
VOC11 mean IU	62.7
VOC12 mean IU	62.2
PASCAL VOC10 pixel acc.	67.0
PASCAL VOC10 mean acc.	50.7
PASCAL VOC10 mean IU	37.8
PASCAL VOC10 f.w. IU	52.5

Table A.2: Results

	Single Food Portion	Non-mixed Food	Mixed Food
Color and Texture	92.21	N/A	N/A
Graph Based	95 (3% increase)	5% increase	15% increase

According to Pouladzadeh, Shirmohammadi, and Yassine, (ibid.), "Graph cut based method is well-known to be efficient, robust, and capable of finding the best contour of objects in n image, suggesting it to be a good method for separating food portions in a food image for calorie measurement". Along with the graph cut segmentation algorithm, this research team also used color and texture segmentation. Gabor filters were used to measure texture features Pouladzadeh, Shirmohammadi, and Yassine, (ibid.). When color and texture segmentation was applied, the method came into difficulty with mixed foods but by applying graph cut segmentation, clearer object boundaries were shown. In conclusion, the accuracy of the classification increased when using graph based segmentation rather than color and texture as seen in A.2.

Local Variation Framework

Another paper was published in which the research team attempted to create a food calorie estimation system Y. He et al., (2013). This system would compromise of three steps, image segmentation, image classification and weight estimation. For the segmentation module, a local variation approach to seg-

VOC 2010 test	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	motorbike	person	plant	shop	sofa	train	tv	mAP
DPM v0 [Girshick et al.] ⁸	49.2	53.8	13.1	15.3	35.5	51.4	49.7	27.0	17.2	28.8	14.7	17.8	46.4	51.2	47.7	10.8	34.2	20.7	43.8	38.3	31.4
UVA [Ullman et al. 2013]	56.2	42.4	15.3	12.6	21.8	49.3	36.8	46.1	12.9	32.1	30.0	36.5	41.5	52.9	32.9	15.3	41.1	31.8	47.0	44.8	35.1
Regionlets [Wang et al. 2013a]	65.0	48.9	25.9	24.6	24.5	50.1	54.5	51.2	17.0	28.9	30.2	35.8	40.2	55.7	43.5	14.3	43.9	32.6	54.0	45.9	39.7
SegDPM [Fidler et al. 2013] ⁹	61.4	53.4	25.6	25.2	35.5	51.7	50.6	50.8	19.3	31.8	26.8	40.4	48.3	54.4	47.1	14.8	38.7	35.0	52.8	41.1	40.4
R-CNN	67.1	64.1	46.7	32.0	30.5	56.4	57.2	65.9	27.0	47.3	40.9	66.6	57.8	65.9	53.6	26.7	56.5	38.1	52.8	50.2	50.2
R-CNN BB	71.8	65.8	53.0	36.8	35.9	59.7	60.0	69.9	27.9	50.6	41.4	70.0	62.0	69.0	58.1	29.5	59.4	39.3	61.2	52.4	53.7

Figure A.1: Detection Average Precision Donahue, (2017)

mentation was performed. Local variation is by which intensity differences between neighbouring pixels is measured. This is a type of graph based segmentation.

The team also carried out some segmentation refinement when the segmentation algorithm had been performed. This consisted of removed small segments (defined as less than 50 pixels) and trying to prevent over and under segmentation. After classification was performed on each segment, segments with low confidence values were removed Y. He et al., (ibid.).

Conclusion

Both of the above papers of Pouladzadeh, Shirmohammadi, and Yassine, (2014) and Y. He et al., (2013) used a graph based segmentation. The first paper used a more generic implementation while the second used a local variation framework. Both methods provided successful results in the image segmentation process.

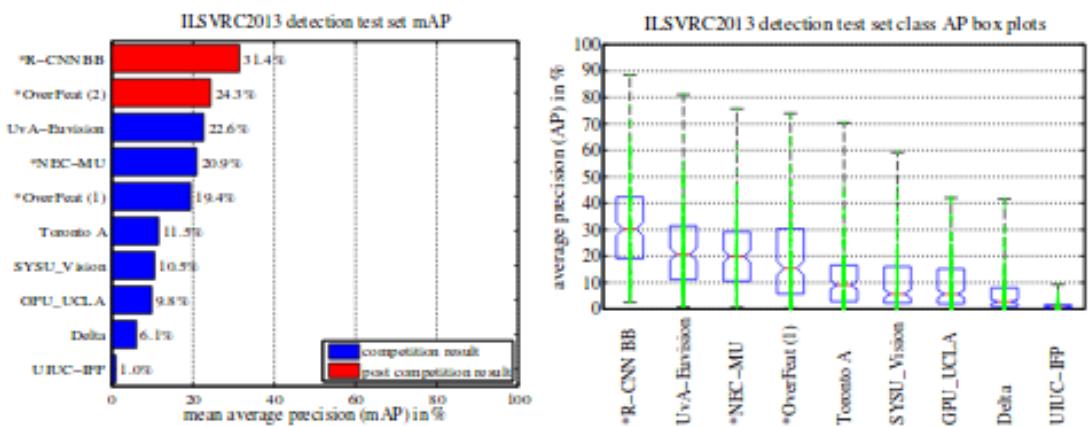


Figure A.2: Mean Average Precision Donahue, (2017)

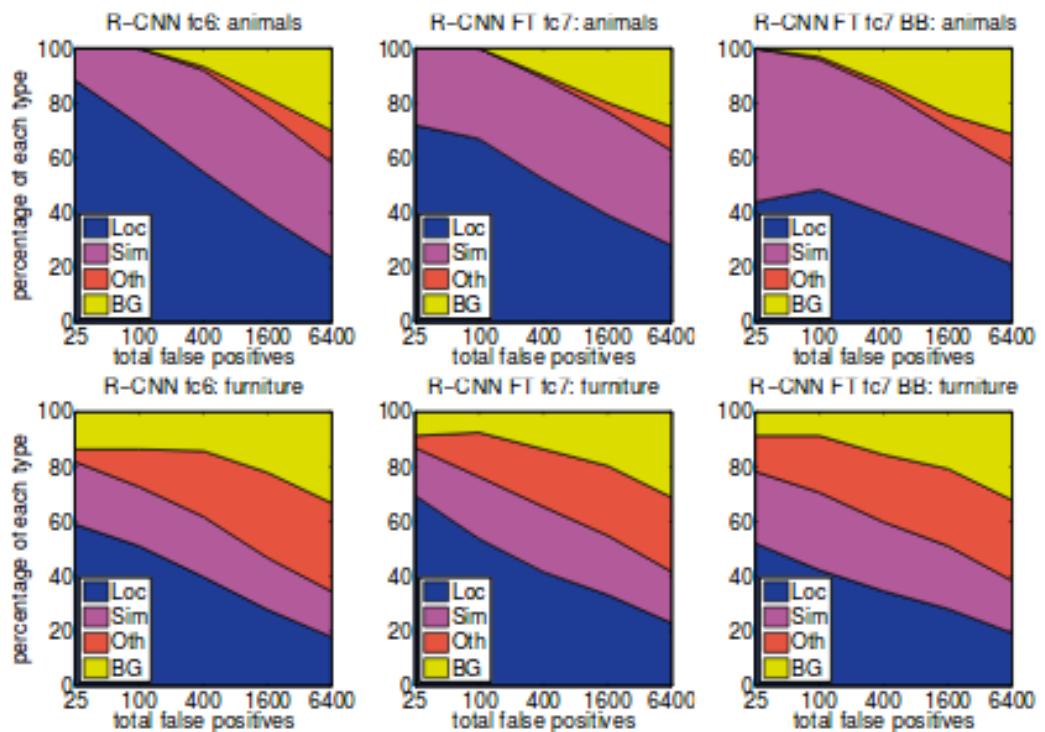


Figure A.3: Distribution of top-ranked false positives Donahue, (2017)

	<i>full</i> R-CNN		<i>fg</i> R-CNN		<i>full+fg</i> R-CNN	
O ₂ P (Carreira et al., 2012)	fc ₆	fc ₇	fc ₆	fc ₇	fc ₆	fc ₇
46.4	43.0	42.5	43.7	42.1	47.9	45.8

Figure A.4: Segmentation Mean Accuracy Donahue, (2017)

VOC 2011 val	bg	aer	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	milk	person	plant	sheep	sofa	train	tv	mean
O ₂ P (Carreira et al., 2012)	84.0	69.0	21.7	47.7	42.2	42.4	64.7	65.8	57.4	12.9	37.4	20.5	43.7	35.7	52.7	51.0	35.8	51.0	28.4	59.8	49.7	46.4
<i>full</i> R-CNN fc ₆	81.3	56.2	23.9	42.9	40.7	38.8	59.2	56.5	53.2	11.4	34.6	16.7	48.1	37.0	51.4	46.0	31.5	41.0	24.3	53.7	51.1	43.0
<i>full</i> R-CNN fc ₇	81.0	52.8	25.1	43.8	40.5	42.7	55.4	57.7	51.3	8.7	32.5	11.5	48.1	37.0	50.5	46.4	30.2	42.1	21.2	57.7	56.0	42.5
<i>fg</i> R-CNN fc ₆	81.4	54.1	21.1	40.6	38.7	63.6	59.9	57.2	52.5	9.1	36.5	23.6	46.4	38.1	53.2	51.3	32.2	38.7	29.0	53.0	47.5	41.7
<i>fg</i> R-CNN fc ₇	80.9	50.1	20.0	40.2	34.1	40.9	59.7	59.8	52.7	7.3	32.1	14.3	48.8	42.9	54.0	48.6	28.9	42.6	24.9	52.2	48.8	42.1
<i>full+fg</i> R-CNN fc ₆	83.1	60.4	23.2	48.4	47.3	52.6	61.6	60.6	69.1	10.8	45.8	20.9	57.7	43.3	57.4	52.9	34.7	48.7	28.1	60.0	48.6	47.9
<i>full+fg</i> R-CNN fc ₇	82.3	56.7	20.6	40.9	44.2	43.6	59.3	61.3	57.8	7.7	38.4	15.1	53.4	43.7	50.8	52.0	34.1	47.8	24.7	60.1	55.2	45.7

Figure A.5: Per-category segmentation accuracy Donahue, (2017)

class	AP	class	AP	class	AP	class	AP	class	AP	class	AP
accordion	50.8	centipede	30.4	hair spray	13.8	pencil box	11.4	snowplow	69.2		
airplane	50.0	chain saw	14.1	hamburger	34.2	pencil sharpener	9.0	soap dispenser	16.8		
ant	31.8	chair	19.5	hammer	9.9	perfume	32.8	soccer ball	43.7		
antelope	53.8	chime	24.6	hamster	46.0	person	41.7	sofa	16.3		
apple	30.9	cocktail shaker	46.2	harmonica	12.6	piano	20.5	spatula	6.8		
armadillo	54.0	coffee maker	21.5	harp	50.4	pineapple	22.6	squirrel	31.3		
artichoke	45.0	computer keyboard	39.6	hat with a wide brim	40.5	ping-pong ball	21.0	starfish	45.1		
axe	11.8	computer mouse	21.2	head cabbage	17.4	pitcher	19.2	stethoscope	18.3		
baby bed	42.0	corkscrew	24.2	helmet	33.4	pizza	43.7	stove	8.1		
backpack	2.8	cream	29.9	hippopotamus	38.0	plastic bag	6.4	strainer	9.9		
bagel	37.5	croquet ball	30.0	horizontal bar	7.0	plate rack	15.2	strawberry	26.8		
balance beam	32.6	crutch	23.7	horse	41.7	pomegranate	32.0	stretcher	13.2		
banana	21.9	cucumber	22.8	hotdog	28.7	popsicle	21.2	sunglasses	18.8		

Figure A.6: Per-class segmentation accuracy Donahue, (2017)

Table A.3: Project Plan

Task/Deliverable	Deadline
Interim Report	21/12/17
Select method and paper to replicate	10/01/18
Iteration 1	???
Iteration 2	???
Iteration 3	???
Draft Final Report	08/03/18
FYP Product	10/04/18
FYP Report	17/04/18