

USING OBJECT DETECTION TECHNIQUES TO DETECT AND CLASSIFY ACNE LESIONS

Submitted by

**ASTHA GUPTA [Reg No: RA1611003010228]
CHINMAY AGRAWAL [Reg No: RA1611003011199]**

Under the guidance of
(, Department of)

*in partial fulfillment for the award of the degree
of*

BACHELOR OF TECHNOLOGY

in

COMPUTER SCIENCE & ENGINEERING

of

FACULTY OF ENGINEERING AND TECHNOLOGY



S.R.M. Nagar, Kattankulathur, Kancheepuram District

JUNE 2020

SRM UNIVERSITY

(Under Section 3 of UGC Act, 1956)

BONAFIDE CERTIFICATE

Certified that this project report titled “**USING OBJECT DETECTION TECHNIQUES TO DETECT AND CLASSIFY ACNE LESIONS**” is the bonafide work of “ **ASTHA GUPTA [Reg No: RA1611003010228], CHINMAY AGRAWAL [Reg No: RA1611003011199]**”, who carried out the project work under my supervision. Certified further, that to the best of my knowledge the work reported herein does not form any other project report or dissertation on the basis of which a degree or award was conferred on an earlier occasion on this or any other candidate.

SIGNATURE

GUIDE

Dept. of

Signature of the Internal Examiner

SIGNATURE

Dr. B. AMUTHA
HEAD OF THE DEPARTMENT
Dept. of Computer Science & Engineering

Signature of the External Examiner

OWN WORK DECLARATION



Department of Computer Science and Engineering

SRM Institute of Science & Technology

Own Work* Declaration Form

This sheet must be filled in (each box ticked to show that the condition has been met). It must be signed and dated along with your student registration number and included with all assignments you submit – work will not be marked unless this is done.

To be completed by the student for all assessments

Degree/ Course : BTECH

Student Name : ASTHA GUPTA , SUMANTA BANERJEE

Registration Number : RA1611003010228 , RA1611003010052

Title of Work : CLASSIFICATION OF SKIN LESIONS USING TRANSFER LEARNING APPROACH

I / We hereby certify that this assessment complies with the University's Rules and Regulations relating to Academic misconduct and plagiarism**, as listed in the University Website, Regulations, and the Education Committee guidelines.

I / We confirm that all the work contained in this assessment is my / our own except where indicated, and that I / W have met the following conditions:

- Clearly references / listed all sources as appropriate
- Referenced and put in inverted commas all quoted text (from books, web, etc)
- Given the sources of all pictures, data etc. that are not my own
- Not made any use of the report(s) or essay(s) of any other student(s) either past or present
- Acknowledged in appropriate places any help that I have received from others (e.g. fellow students, technicians, statisticians, external sources)
- Compiled with any other plagiarism criteria specified in the Course handbook / University website

I understand that any false claim for this work will be penalised in accordance with the University policies and regulations.

DECLARATION:

I am aware of and understand the University's policy on Academic misconduct and plagiarism and I certify that this assessment is my / our own work, except where indicated by referring, and that I have followed the good academic practices noted above.

If you are working in a group, please write your registration numbers and sign with the date for every student in your group.

ACKNOWLEDGEMENTS

We express our humble gratitude to Dr. Sandeep Sancheti, Vice Chancellor, SRM Institute of Science and Technology, for the facilities extended for the project work and his continued support.

We extend our sincere thanks to Dr. C. Muthamizhchelvan, Director, Faculty of Engineering and Technology, SRM Institute of Science and Technology, for his invaluable support.

We wish to thank Dr. B. Amutha, Professor Head, Department of Computer Science and Engineering, SRM Institute of Science and Technology, for her valuable suggestions and encouragement throughout the period of the project work.

We are extremely grateful to our Academic Advisor Dr.R.Annie Uthra , Associate Professor, Department of Computer Science and Engineering, SRM Institute of Science and Technology, for her great support at all the stages of project work.

We register our immeasurable thanks to our Faculty Advisor, Mr A.M.J Muthukumar, Associate Professor and Dr Revathi Venkataraman, Professor, Department of Computer Science and Engineering, SRM Institute of Science and Technology, for leading and helping us to complete our course.

Our inexpressible respect and thanks to my guide, Mr U.M Prakash, Associate Professor, Department of Computer Science and Engineering, SRM Institute of

Science and Technology, for providing me an opportunity to pursue my project under his/her mentorship. He provided me the freedom and support to explore the research topics of my interest. His passion for solving the real problems and making a difference in the world has always been inspiring.

We sincerely thank staff and students of the Computer Science and Engineering Department, SRM Institute of Science and Technology, for their help during my research. Finally, we would like to thank my parents, our family members and our friends for their unconditional love, constant support and encouragement.

Sumanta Banerjee

ABSTRACT

The problem of acne is one that affects people all around the world. Mild forms of acne resolve on their own within a few weeks. However, if left untreated, chronic forms of acne lesions could turn into more severe forms resulting in inflammation, swelling, pain and infection. In order to determine the proper medical treatment plan for a patient, a dermatologist has to not only ascertain the severity of acne, but also keep in mind the type of acne lesions developed. This results in a significant amount of time being devoted to the manual task of identifying and classifying the lesions for every patient. In order to help with this issue, an automatic technique could be devised for the purpose of detecting and classifying acne lesions in an image. Through this project, we aim to devise a deep learning-based system to count the number of acne lesions in an image and classify each of them into two categories namely inflammatory and non-inflammatory acne (whiteheads and blackheads). The data set comprises 469 images distributed among both classes. The proposed method relies on the use of widely known object detection algorithms to localise and detect the spots followed by a classification network to predict the labels for each spot. We believe that this system would prove to be very helpful to help to reduce errors and manual effort, while also working as an efficient safety net to help with a patient's diagnosis.

TABLE OF CONTENTS

OWN WORK DECLARATION	iii
ACKNOWLEDGEMENTS	iv
ABSTRACT	vi
LIST OF TABLES	x
LIST OF FIGURES	xi
ABBREVIATIONS	xii
1 INTRODUCTION	1
1.1 Acne	1
1.2 Problem Statement	3
1.3 Deep Learning	4
2 LITERATURE SURVEY	5
2.1 Methods Used	5
2.2 Conclusion	7
3 System Analysis	8
3.1 Dataset	8
3.2 Use of Object Detection Based Methods	9
3.3 Faster- RCNN	10
3.4 You Only Look Once	12
3.5 Single Shot Detector	13

3.6	Hardware Dependencies	15
3.7	Software Dependencies	15
4	System Design	17
4.1	YOLO Modules	19
4.1.1	Pre-processing	19
4.1.2	Backbone Architecture	19
4.1.3	Training	19
4.1.4	Evaluation	20
4.2	SSD Modules	21
4.2.1	Pre-processing	21
4.2.2	Backbone Architecture	21
4.2.3	Training	21
4.3	FR-CNN Modules	23
4.3.1	Pre-processing	23
4.3.2	Backbone Architecture	23
4.3.3	Training	24
4.3.4	Evaluation	24
5	RESULTS	25
5.1	YOLO	25
5.2	SSD	28
5.3	FR-CNN	32
6	CONCLUSION	36
7	Future Enhancements	38
7.1	Dynamic Region of Interest (ROI) Extraction	38
7.2	Sub-Classification of Inflammatory Class Acne	38
7.3	Severity and classification	39

7.4	Forming a benchmark dataset	39
7.5	Mobile application development	39
A	Coding	40
A.1	Pre-processing	40
A.2	YOLO Training	41
A.2.1	Darknet	41
A.2.2	Custom CFG File	46
A.3	SSD Training	67
A.4	FR-CNN Training	73
A.5	Evaluation	86
	REFERENCES	97

LIST OF TABLES

LIST OF FIGURES

3.1	Architecture Diagram for FR-CNN [8]	11
3.2	Architecture Diagram for YOLO [9]	13
3.3	Architecture Diagram for SSD [10]	15
4.1	Architecture Diagram for all three modules	18
5.1	Mean Average Precision for YOLO	26
5.2	Precision vs. Recall Curves for YOLO	27
5.3	An image sample of the detection results	29
5.4	Mean Average Precision for SSD	30
5.5	Precision vs. Recall Curves for SSD	31
5.6	Sample output image for FR-CNN	33
5.7	Mean Average Precision for SSD	34
5.8	Precision vs. Recall Curves for FR-CNN	35

ABBREVIATIONS

CNN	Convolutional Neural Network
ResNet	Residual Network
FR-CNN	Faster Regional Convolutional Neural Network
YOLO	You Only Look Once
SSD	Single Shot Detector
mAP	mean Average Precision
RMS	Root Mean Squared

CHAPTER 1

INTRODUCTION

1.1 Acne

Diagnosis of acne involves a medical examination of skin by a dermatologist. First step in diagnosis of acne is to determine the severity of acne by grading them in the range of one to four (one being mild and four being severe), then the doctor may examine and determine the types of acne on the patient's skin. The area of skin affected is also taken into account during the diagnosis process. After the skin examination, some questions are asked about the frequency, occurrence, if acne runs in the family.

While doctors are adept at the diagnosis process, there's always a scope for human error during this process. The part of the diagnosis most affected by the human error is the skin examination. Sometimes the doctors are required to quantify the number of acne in an area and a ballpark is taken into account. This is where the intervention of Artificial Intelligence can come handy, to aid the process of diagnosis and consequently resulting in better treatment.

Although the problem of acne is mostly present during the adolescent age period, it can be a chronic problem for many people. Women, especially those suffering from hormonal disorders, can end up with more severe forms of this disease which can lead to permanent scarring. Even though acne is not a life-threatening disease, it can affect one's mental health and lead to more serious mental health issues. Patients have to deal with the social stigma that revolves

around this issue and end up constantly being worried about the cosmetic disadvantages that acne problems can result in.

Acne are of several types and they possess certain characteristics which make the types discernible to human beings. For example, plugged pores can be whiteheads or blackheads depending on if they are open or closed. There are two main classes of acne- inflammatory and non-inflammatory.

Just like the name suggests, the former class consists of those acne types which result in the formation of inflammation around the acne lesion and are marked by pain and redness of the skin. They consist of cysts, nodules and pustules and papules. Red bumps which are tender and small can be papules, papules with pus are called pustules or pimples, nodules are larger in size and firm lumps under the surface of skin and pus filled nodules can be classified as cysts. While pustules and papules can be treated with over-the-counter remedies, the former two types require the care of a dermatologist. The second class of acne lesions consist of whiteheads and blackheads. They are usually self-treatable using readily available treatment options.

1.2 Problem Statement

Chronic acne, while affecting the cosmetic appearance of a person, can eventually lead to underlying mental health conditions especially during the adolescent years. It becomes imperative to take some action for the treatment of acne even if it is of the non-inflammatory type.

Furthermore, the transformation of non-inflammatory types to inflammatory types of acne can result in quite painful lesions which would require medical treatment, often in the form of topical drugs. In order to assess the progress of the treatment as well as to study the effects of the prescribed drugs, it is essential for the dermatologists to regularly review and monitor the patient regularly. This can lead to a substantial amount of time and money being involved throughout the treatment on both sides.

By developing a method to help quantify and classify the acne lesions on a facial image, we aim to provide patients with a way to monitor and assess their lesions. Not only would it help them save time and act as a safety net, it could also help dermatologists save time and effort.

Through our project, we wish to develop a deep learning-based object detection and classification technique in the field of dermatology. Using transfer learning to extract weights followed by building a classifier to train the model, the technique would allow the model to detect acne lesions in an image and classify them into one of three classes namely inflammatory, blackheads and whiteheads.

1.3 Deep Learning

A subset of Artificial Intelligence, the field of deep learning involves the use of neural networks in order to build models. Over the years, deep learning has been used majorly in the area of computer vision to detect, localise and classify images. Mimicking the task of neurons in the human body, the deep neural networks consist of multiple layers of nodes which process information as it passes through them.

Traditional techniques for image detection often rely on a domain expert to hand-craft features which are considered to be important for the task at hand. Often in the field of medical imaging, the variability between different classes of images might not be discernable by the human eye. This dependence can eventually result in inefficient results when tested over a large scale of images. In order to deal with this, deep learning offers an automated method to detect features and extract them for the model at hand.

Deep learning techniques involve pre-processing the data before being fed into a neural network. The network, then, learns features from the input images and assigns weights to the selected features. This is followed by a series of training steps for batches of images which help to adjust the weights and reduce the loss of detection. Finally, the trained model is used to detect objects in new images which form the test set.

CHAPTER 2

LITERATURE SURVEY

2.1 Methods Used

In recent years, many dermatologists have worked with computer scientists to help solve the problem of tracking and counting acne lesions. With the advent of computer vision methods, numerous techniques have been proposed to detect and/or classify acne lesions. The studies [1] to [7] use image processing and computer vision techniques to detect the acne lesions in facial images.

Fartash Vasefi et al. propose a color-based detection technique wherein they convert the original images into their CIE L*a*b versions and use a Gaussian filter applied on the a* plane (which depicts redness) to detect the acne spots using Otsu thresholding [1]. Then, they classify that spot as a papule or pustule by detecting if the segmented binary area has a pus-filled center or not, by using the Euler number generated for that particular lesion.

Another acne detection technique has been proposed in [2] wherein the RGB images are converted into their HSV and normalised grayscale versions and the subtraction between the V plane of the former and the latter is performed and the resulting image is converted into its binary threshold-ed form. This has then been used to find the suitable acne lesions.

The study by Thanapha Chantharaphaichi et al. focussed on the reduction of errors in the detection of acne spots uses a Bayesian classifier on top of a blob

detector via segmentation, with a histogram-based feature extraction [3]. The training is supervised and the testing is unsupervised. The main idea of this paper is to reduce the error after the extraction phase using a probabilistic method.

Recently, machine learning and deep learning have been used by researchers for the task at hand. A SVM-based classifier has been proposed by C. Chang et al. [4], where the region-of-interest extraction is followed by defect extraction. It involves various stages including the extraction of the region of interest from the input image. The three classes in this study include normal skin, acne lesions and spots.

In [5], Lim ZV et al. uses deep learning to develop a classification model based on severity for three classes. The pre-processed images are augmented and fed to the CNN architecture. A comparison between the three widely known models namely, ResNet, DenseNet and Inception, is performed.

The method proposed by Zhao, Tingting et al. [6] uses five levels of acne severity as the basis of classification. Making the use of transfer learning, wherein the CNN architecture uses features extracted from a base data set, which are then calibrated to suit the domain-focussed data set. A pre-trained ResNet-52 architecture has been chosen as the CNN architecture. The results of the model are compared with those of clinicians to ascertain its performance.

Finally, Shen, X. at al. used a template matching technique with a sliding window to first perform a binary classification between skin and non skin patches, followed by classification into seven different categories of acne type [7]. This technique divides the pre-processed images into patches of fixed sizes and uses the window to classify each of them first as skin or non-skin, followed by the type of acne.

2.2 Conclusion

Having surveyed the literature on the various techniques used for the problem statement, we can conclude that all those studies fell broadly into one of two categories in terms of the methodology used.

The traditional Computer Vision techniques focus on the use of image filtering, edge detection and threshold-based segmentation. The feature extraction is performed on the basis of handcrafted features specified by a domain expert. The papers that use traditional methods perform detection mostly on the basis of a colour-based segmentation. The main problem with using these methods for our problem is that all the different types of acne have a complex set of similar as well as dissimilar features. While both blackheads and whiteheads do not have any redness around them, they have opposite colours in the RGB spectrum. Thus, the use of traditional methods will not prove to be very effective because it will suffer from low-quality feature extraction.

With deep learning, the task follows an end-to-end approach. The task of pattern detection is performed by detecting similarities between the image of a class on a pixel level. The nodes of the network learn the similarities of images in the same class and use them to perform a probability-based classification. The feature extraction is performed by the algorithm with a focus on even the most minute of details. When it comes to the task of classification of images, deep learning methods have been established to perform much better than traditional techniques.

CHAPTER 3

SYSTEM ANALYSIS

3.1 Dataset

The acne lesion dataset consists of a total of x images, x of which are inflammatory acne images and y are non-inflammatory acne lesions. The non-inflammatory are divided into z images of blackheads and w images of whiteheads. The dataset was obtained by scraping images from Dermnet, the largest independent source for dermatology images. Several scripts were written in Python using the BeautifulSoup framework to scrape the images from the website and save them locally on the workstation.

After obtaining the images, they were fed into a preprocessing pipeline which was required before passing it to various object detection algorithms. One of the major parts of the pipeline is annotating the dataset to conform with the implementation of the multiple object detection algorithms. The mostly used data annotation format is PASCAL Visual Object Classes (VOC) format. The dataset was annotated using an open-source data annotation tool known as LabelImg.

3.2 Use of Object Detection Based Methods

Given that the problem statement specifies that we need to build a system which would be able to not only count the acne lesions but also classify it into one of two categories, the decision to use object detection and classification algorithms was not a difficult one. One of the main reasons to not choose traditional methods like image segmentation was due to the variability of the different classes. It would have been very complex to decide for a single parameter of difference between all the classes.

For instance, the inflammatory types are marked by redness around them, but that could not be used for determining the threshold because within the non-inflammatory types the blackheads and the whiteheads are of completely opposite colors in the RGB spectrum. Thus, color-based segmentation was out of the question if we wanted the model to not be biased.

Another important consideration that we kept in mind while choosing deep learning was the high complexity of the data set. The variability between the skin color and the size of the lesion would make it difficult for us to define a fixed set of features. Given that deep learning based methods would break down the image pixel-by-pixel and extract defining features intuitively made it to be a better choice for the problem at hand. It is our belief that while traditional methods do work well in many cases, this is one where convolutional networks would perform better.

3.3 Faster- RCNN

One of the state of the art object detection algorithms used for the particular problem statement was Faster Regional Convolutional Neural Networks. It is the third iteration of Regional Convolutional Neural Networks, Fast Regional Convolutional Neural Networks being the second. It is famously trained using the PASCAL VOC format, hence the dataset was labeled in the PASCAL VOC format.

The typical pipeline of Faster-RCNN includes annotating the dataset after resizing the image maintaining the aspect ratio. After this, the Faster-RCNN model is trained on a base network with data augmentation occurring during the training process. The checkpoints of the model are saved after every fixed set of epochs making sure that the loss is converging.

The architecture of Faster-RCNN consists of two major networks. The first network is the Region Proposal Network (RPN) which outputs proposals for the main networks using a base network and convolution function instead of fully connected layers. This ensures that images of different sizes can be used to train the model. The main network of Faster-RCNN outputs the final coordinates of the objects with a classification score via a softmax function.

Being the third successor of Regional Convolutional Neural Networks, it can test images faster and also do it more accurately when compared to other models in terms of mean average precision(mAP). The advantages of this complex architecture includes being able to detect objects of various sizes in the test image, accepting images of multiple sizes as the input and also the speed at which it tests the images.

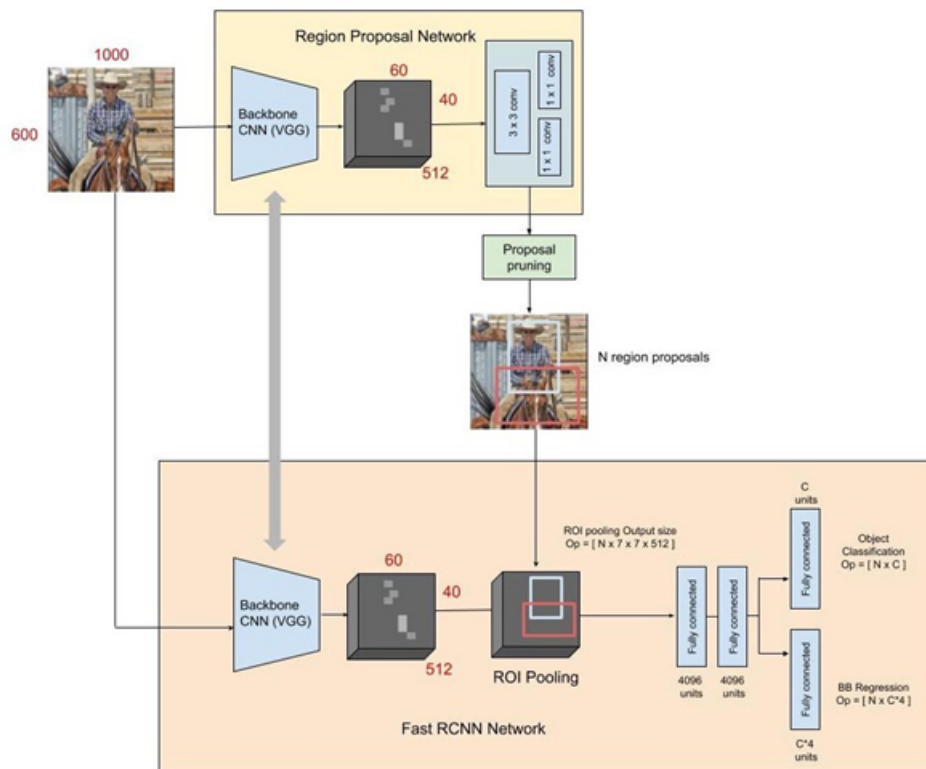


Figure 3.1: Architecture Diagram for FR-CNN [8]

3.4 You Only Look Once

Most object detection algorithms rely on the use of a regression task followed by a classification task. Even though it is a robust technique, it takes up a large amount of time and the underlying process is often complex. To deal with this problem, the YOLO architecture for object detection and classification was proposed.

In YOLO, each image is passed through the convolution layers only once. The image is divided into grids and anchor points for each grid are specified. The task for each grid is to predict if any object is present in its area with its associated probability. The grid responsible for predicting the object is the one in which the centre of the detected object lies. Relying on the probability for the class, YOLO sees the task only as a regression problem.

Over the years, four versions of YOLO have been released. The major benefit of this architecture is the speed it offers, making it the go-to option for real time object detection. Furthermore, unlike region-based detectors, it looks at the whole picture rather than specific regions. Another great advantage that YOLO offers is that it learns not only from the object-specific information, but also considers contextual information about the images. When tested on the COCO dataset for 80 classes, it was evaluated to have an mean Average Precision value of approximately 60 percent, making it at par with its contenders.

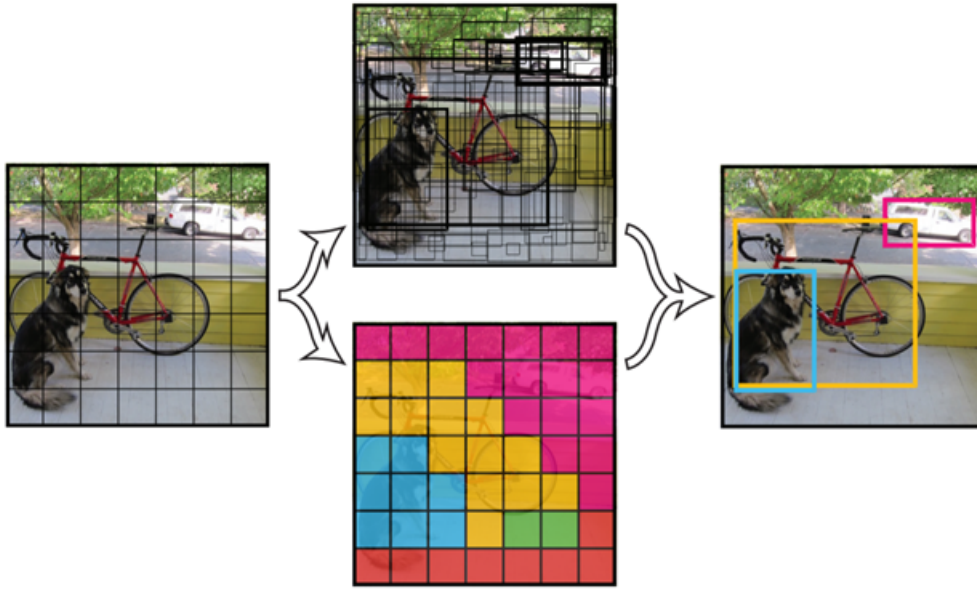


Figure 3.2: Architecture Diagram for YOLO [9]

3.5 Single Shot Detector

Single Shot Detector (SSD) is an object detection algorithm which uses a backbone ConvNet which is pre-trained to extract features. Fully connected final classification layer is removed from the ConvNet, and more convolution layers added to the backbone and the output is received as bounding boxes and their classes. Backbones can be any deep ConvNet trained for image classification, some popular architectures make use of Inception and ResNet networks as backbones.

SSD does not need a sliding window and multiple passes to detect objects, it can detect objects in a single pass, using grid cells. The image is first divided into grids and the resulting cells are responsible for detecting the object inside that cell. Every cell can detect the form and class of the object in their region. Dimension of the grid can also be decided as per the setting requirement.

Anchor boxes are the pre-defined rectangular boxes which are assigned to every cell and are responsible for detecting different sizes and forms of objects. Every anchor box has different dimension ratios and multiple anchor boxes can be assigned to a single cell. These anchor boxes have a specific aspect ratio and zoom level to define them. The size of the object within a cell may be smaller or larger, and to take care of that, zoom can be defined to scale the box with respect to the cell. For example a giraffe and a dog are of different sizes and both would require their bounding boxes to be of different ratios. A giraffe would require a vertically taller box with less width and a dog would have a horizontal box with enough height. The ground truth boxes which are annotated manually before training, are matched with anchor boxes to find out the best overlap during the training. The size of the object within a cell may be smaller or larger, and to take care of that, zoom can be defined to scale the box with respect to the cell.

The most important attribute of SSD is the use of receptive fields. Receptive field is the region which is being used by CNN's feature at that time in the input space. It is known that the different sizes of the region are the result of features at different layers, due to convolution operations. The deeper layers have larger feature size and the size gets smaller as we move to the upper layers. These variable feature maps help the SSD to detect objects of different sizes and give out an accurate bounding box. SSD has more number of convolution layers at the backbone which result in a separate object detection output. Due to having different sizes of receptive fields, earlier layers are able to detect small objects.

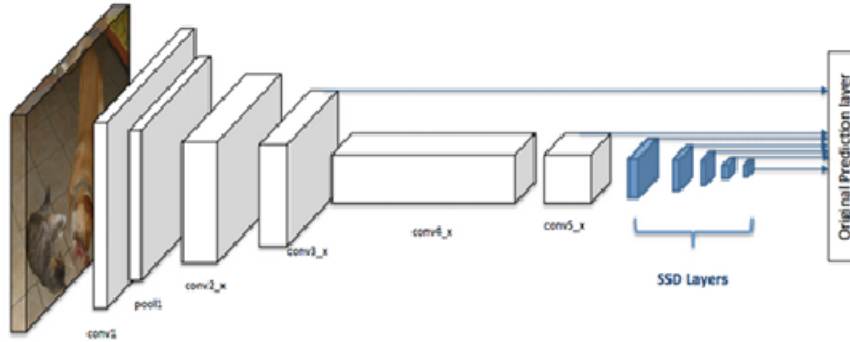


Figure 3.3: Architecture Diagram for SSD [10]

3.6 Hardware Dependencies

The major hardware dependencies include a Graphics Processing Unit consisting of 1xTesla K80 , compute 3.7 having 2496 CUDA cores and 12GB GDDR5 VRAM. It also includes a CPU consisting of 1xsingle core hyper threaded Xeon Processors @2.3Ghz i.e (1 core, 2 threads). These are the configurations of an instance of a Google Colab notebook using which the object detection models are trained.

Apart from the specifications of the workstation, the end product will be developed into a mobile application. The initial development of the application will be for the ios platform requiring an iPhone and will later be developed for android operating systems as well.

3.7 Software Dependencies

The software dependencies includes the programming languages, platform used to train the object detection algorithms and the tools used for data annotation. The programming language used for the project was Python 3.7. Several frame-

works in Python were used to solve individual parts of the project. The deep learning frameworks include TensorFlow and Keras. The framework to scrape the acne images from Dermnet was BeautifulSoup. The frameworks for numerical computations include NumPy and Pandas.

The platform to train the object detection algorithms was Google Colab with the hardware specifications as mentioned in the heading above. Google Colab is a platform offered by Google Research. It is used to run code on the cloud in the form of a Jupyter Notebook. It is primarily used to train machine learning and deep learning models, data analysis, etc. The main advantage of using Google Colab is the free resources that are provided by Google which includes free GPU usage.

An open-source tool was used for data annotation known as LabelImg. It is an image annotation tool with a Graphical User Interface. It is used to make bounding boxes across the objects inside the image. It saves the annotations in the form of an XML file with regards to the PASCAL VOC format which was used by Imagenet. It also includes support for YOLO.

CHAPTER 4

SYSTEM DESIGN

A common architecture was decided for all the algorithms chosen. The basic idea was to perform a common pre-processing step for all and use the same evaluation metrics to adjudge their performances. Then the results were compared to determine the best one.

The first step was the extraction of annotations from the dataset of images. This was done in the PASCAL VOC format for FR-CNN and SSD, and a specific format for YOLO algorithm. The total number of images was 496. Data augmentation was also performed while training.

The next step was to train the model using pre-defined parameters initially. The parameters included the position of anchors, the size of the image batches, the number of epochs (iterations for YOLO) and the optimisers used. The loss function was also selected according to the specifications of the model. This step also included the use of pre-trained models for feature extraction.

Finally, the models were tested on the test cohort and the mean Average Precision (mAP) was calculated using the Precision vs. Recall Curve for each algorithm.

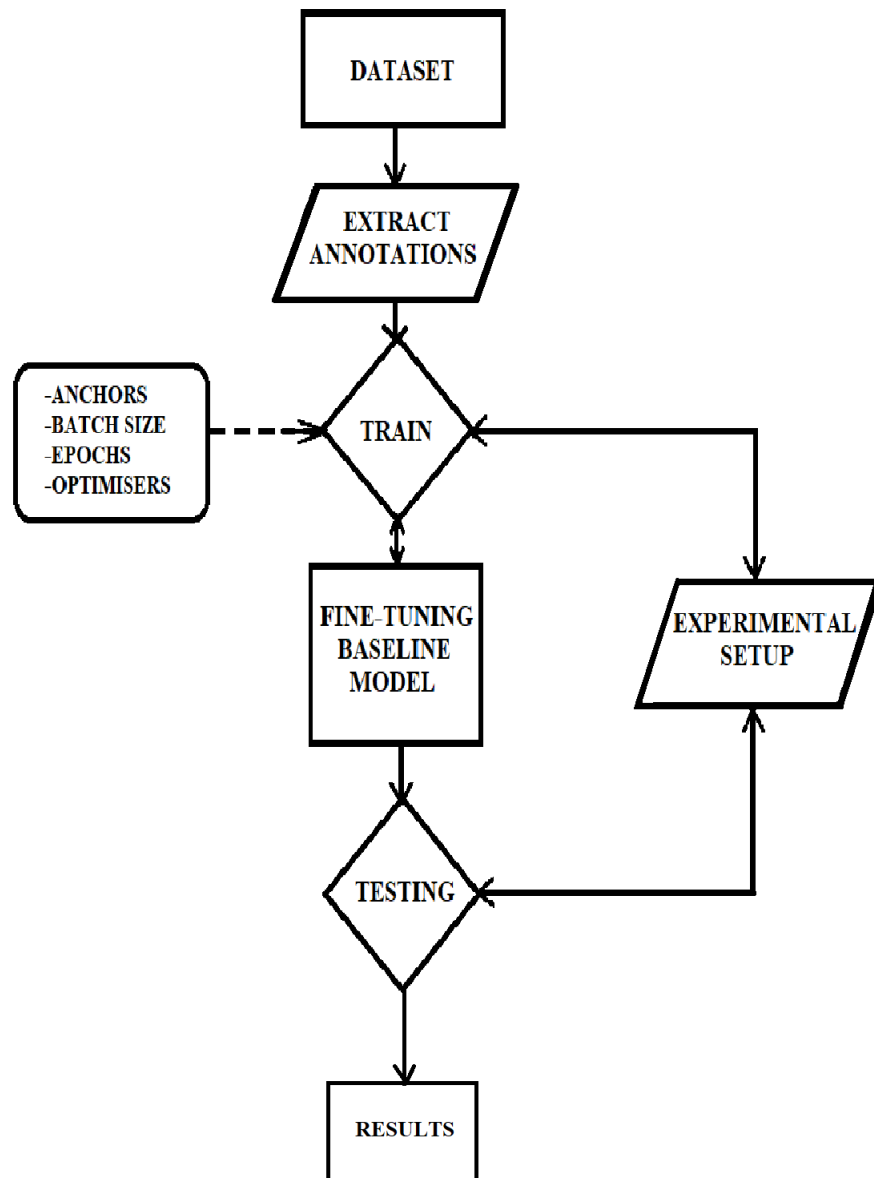


Figure 4.1: Architecture Diagram for all three modules

4.1 YOLO Modules

4.1.1 Pre-processing

YOLO requires a specific form of data annotation that is different from the PASCAL VOC format. The bounding box information for it needs to be in the form of a text file with five features for each object in an image. These five values are the class, xmin, ymin, xmax and ymax values. Furthermore, a separate txt file for each image needs to be created. Additionally, after splitting into train, test and validation sets, a text file needs to be created for each subset which would contain the location of each image file present in the subset. Finally, a yolo.names file and a yolo.data files need to be created in order to define the train, validation and class files. A location to save the last weights during training is also created.

4.1.2 Backbone Architecture

Yolo relies on the use of DarkNet as its base architecture. By means of transfer learning, weights for the base architecture are obtained by training it on the COCO dataset. These weights of the pre-trained model are used as the starting weights for training the YOLO model. DarkNet consists of a network of fully connected layers. Based in C, this architecture provides a high speed alternative to the other models. <diagram>

4.1.3 Training

The YOLO model was trained using initial weights extracted by training the Darknet architecture on the COCO dataset. The adjusted weights were saved

after every 1000 iterations and the final weights were saved in the backup folder. For the custom training, the classes were set as three and the filter size was set as 24. The model was trained for 6000 iterations (2000 iterations per class) with a learning rate of 0.001 and a batch size of 64.

4.1.4 Evaluation

In order to evaluate the model, the metric used was the Mean Average Precision. The maximum validation mAP was calculated as 20.04% with 6000 iterations. The confidence threshold was set as 0.5

4.2 SSD Modules

4.2.1 Pre-processing

Preprocessing steps for SSD are similar to FR-CNN. The standard annotation format for object detection algorithms is the PASCAL VOC format. The dataset, after being scraped from the DermNet Server, is contained inside a repository. Using the data annotation tool LabelImg, the images are manually labelled and annotated in the XML format.

The format majorly retains the coordinates of the bounding boxes that are marked on each object inside the image. Other options include marking whether the object in the image is difficult to locate and also if the object is truncated or not. To train using the Tensorflow object detection API, XML annotations need to be converted to TF Records format, for which the scripts are provided as open source by Tensorflow.

4.2.2 Backbone Architecture

Backbone architecture is the pre-trained ConvNet with the final fully connected layer removed. It is used for feature extraction, which helps the model to train on custom image data-set. For the purpose of this experiment we have used the Inception v2 based SSD model, trained on COCO dataset.

4.2.3 Training

For training 90% of images from the data-set were used and model was trained on the following configuration. The batch size was set to 8 with an RMS Prop

optimiser. The initial learning rate was set as 0.004 and the decay factor was set as 0.95. The input size of the images was set as 300x300.

4.3 FR-CNN Modules

4.3.1 Pre-processing

The standard annotation format for object detection algorithms is the PASCAL VOC format which is used by Faster-Regional Convolution Neural Networks. The dataset, after being scraped from the DermNet Server, is contained inside a repository. Using the data annotation tool LabelImg, the images are manually labelled and annotated in the XML format. The format majorly retains the coordinates of the bounding boxes that are marked on each object inside the image. Other options include marking whether the object in the image is difficult to locate and also if the object is truncated or not. One of the advantages of Faster-RCNN is being able to train on images with different sizes but ensuring that the smallest side in all the images are of the same size.

4.3.2 Backbone Architecture

The backbone architecture for Regional Convolutional Neural Networks included VGG16. It included the entire network pre-trained on ImageNet without the fully connected layers. The features that were extracted from the last convolutional layer which is conv5, are used as the base for training RCNN. One of the prominent features of object detection algorithms is transfer learning, using the features/weights from one architecture to generalise the model of another architecture. For Faster-RCNN, ResNet(50, 150, 200) are also used as the backbone architecture. ResNet or Residual Networks are also pretrained on ImageNet with deeper layers and includes skip connections which help in extracting better features when compared to VGG16.

4.3.3 Training

Faster-RCNN was trained on VGG16 and ResNet50 with pretrained weights for the backbone architecture. For custom training, three classes were used: inflammatory, whiteheads and blackheads. The weights were saved after every epoch only if the loss reduced after the corresponding epoch. The model was trained with a batch size of 64. The number of epochs were 100 with adam optimiser and a learning rate of 0.0001.

4.3.4 Evaluation

In order to evaluate the model, the metric used was the Mean Average Precision. The maximum validation mAP was calculated as 16.05% with 100 Epochs. The confidence threshold was set as 0.5.

CHAPTER 5

RESULTS

5.1 YOLO

The YOLO algorithm was trained on a dataset of 496 images with a test-train split of 10-90 percent. The backbone architecture used was DarkNet and the batch size was set as the default value of 64.

Data augmentation was used during the training process and all images were resized to 416x416. This was done automatically by the system. The confidence threshold was set at 0.5 for the calculation of precision and recall.

Accuracy was not chosen as the evaluation metric due to the problem of class imbalance. The mean Average Precision value is a more reliable metric for such a scenario. Overall, the model achieved an mAP value of 46.53%.

Certain drawbacks of YOLO are to be responsible for these results. However, the training and testing time for YOLO was lower than its counterparts.

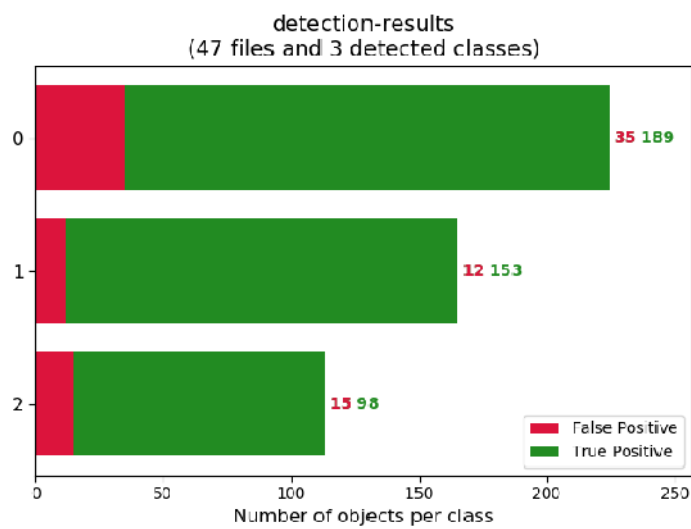
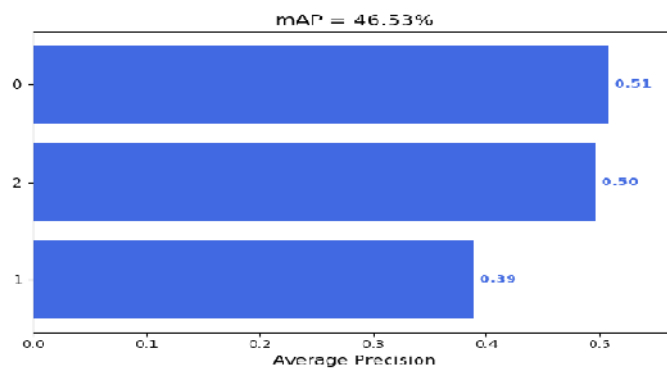


Figure 5.1: Mean Average Precision for YOLO

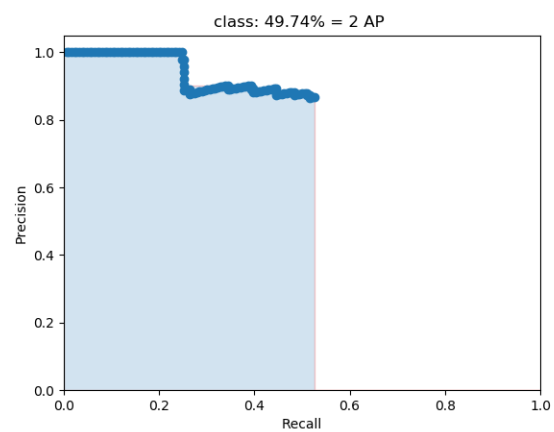
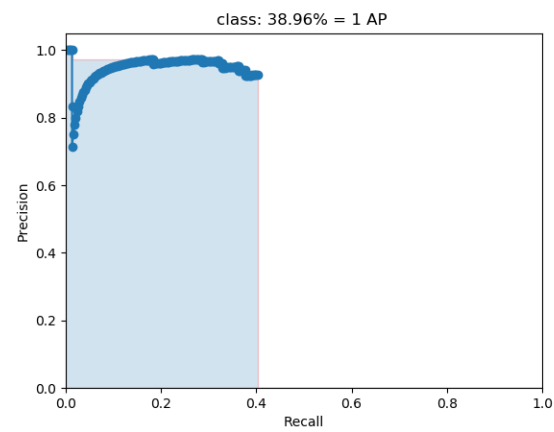
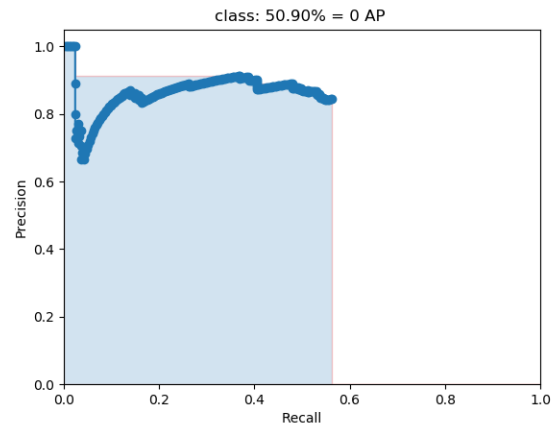


Figure 5.2: Precision vs. Recall Curves for YOLO

5.2 SSD

The SSD algorithm was trained on an NVIDIA 2080Ti workstation with a batch size of 8 for 341 epochs. The base network used was Inception V2. The optimiser used was the Root Mean Squared (RMS) prop optimizer.

Train test split was consistent with other algorithms used in this study. Since the batch size is small, augmentation is not necessary in this model. The performance was measured based on the precision and recall curves for the three classes. It achieved an mAP of 52.04%.



Figure 5.3: An image sample of the detection results

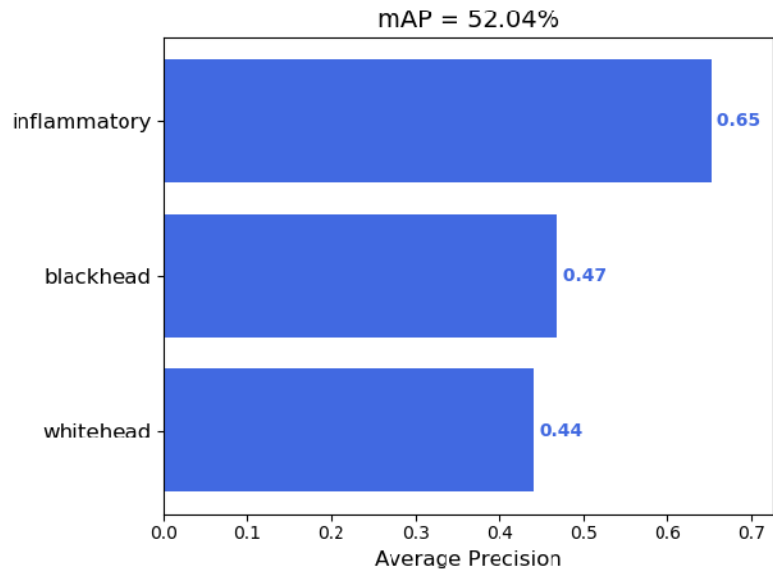


Figure 5.4: Mean Average Precision for SSD

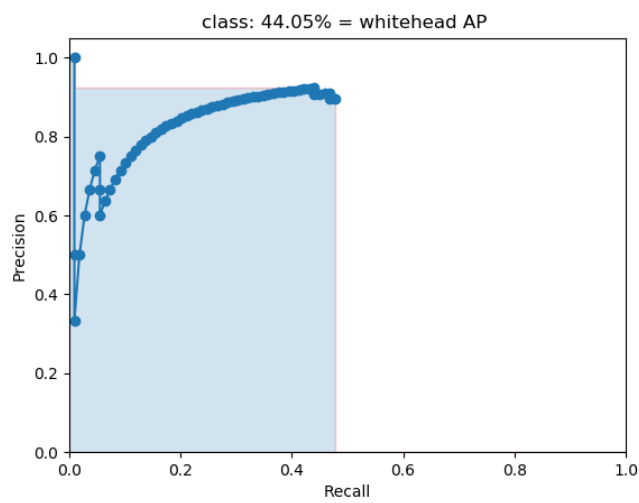
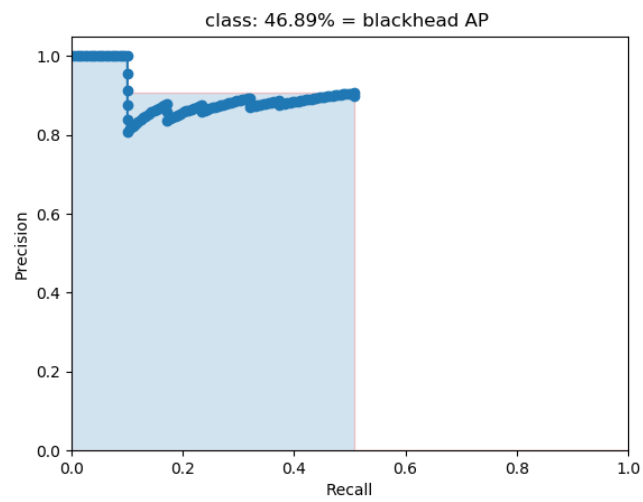
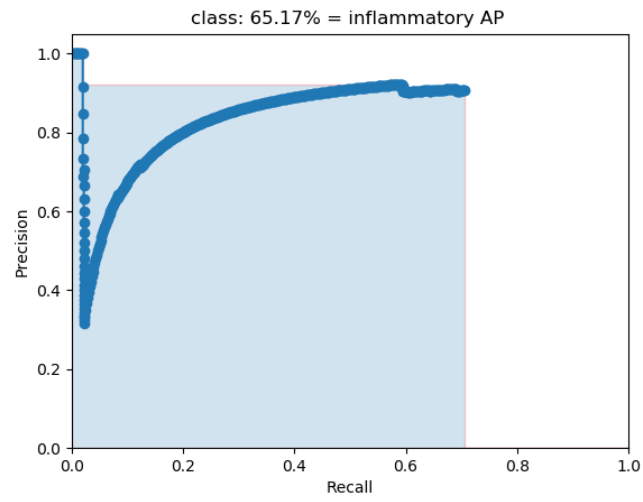


Figure 5.5: Precision vs. Recall Curves for SSD

5.3 FR-CNN

The Faster-RCNN algorithm was trained on an NVIDIA 2080Ti workstation with a batch size of 64 for 100 epochs. The base network used was ResNet50. The optimiser used was Adam with binary cross-entropy loss.

The train-test split was 90:10. There was no validation set associated with the model as the size of the dataset was really low. The test set included an equal number of images in which there were a majority of single class lesions.

To tackle the problem of class imbalance, we performed data augmentation while training the algorithms to ensure synthetic data is generated accordingly and can neutralise the problem as much as possible.

As shown in the graphs, the highest mAP achieved was 86.61% with Faster-RCNN. The individual Precision-Recall Curves for each class are also shown below.

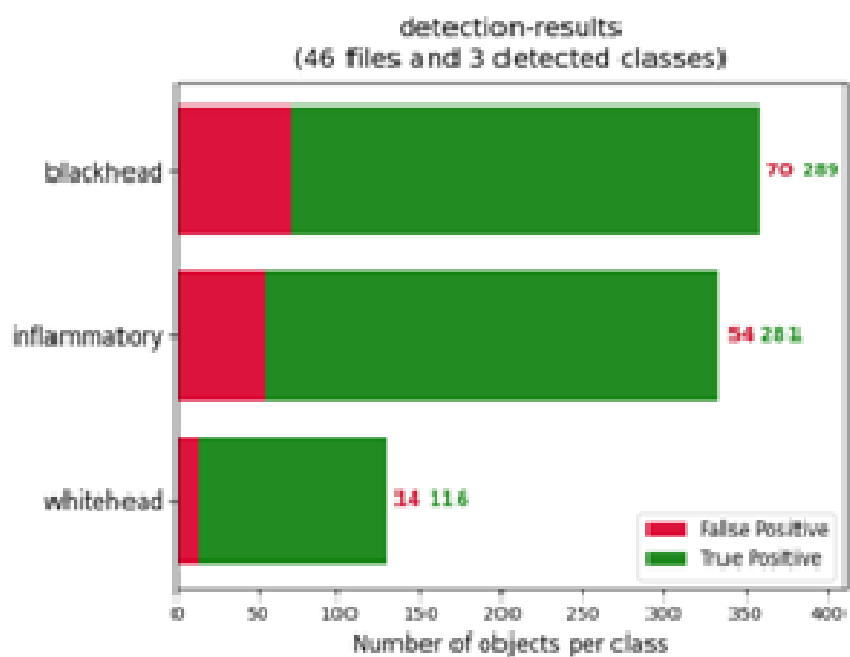
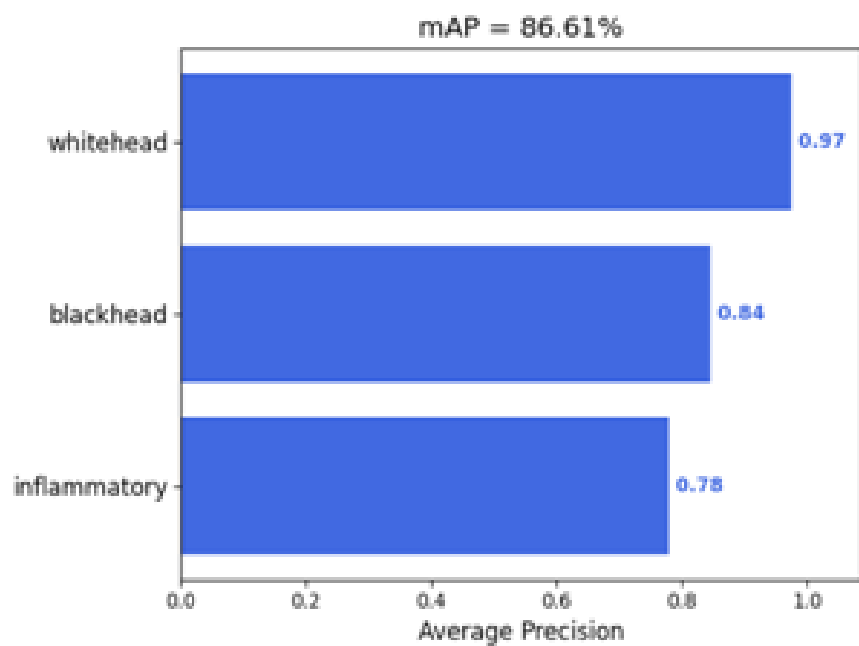


Figure 5.7: Mean Average Precision for SSD

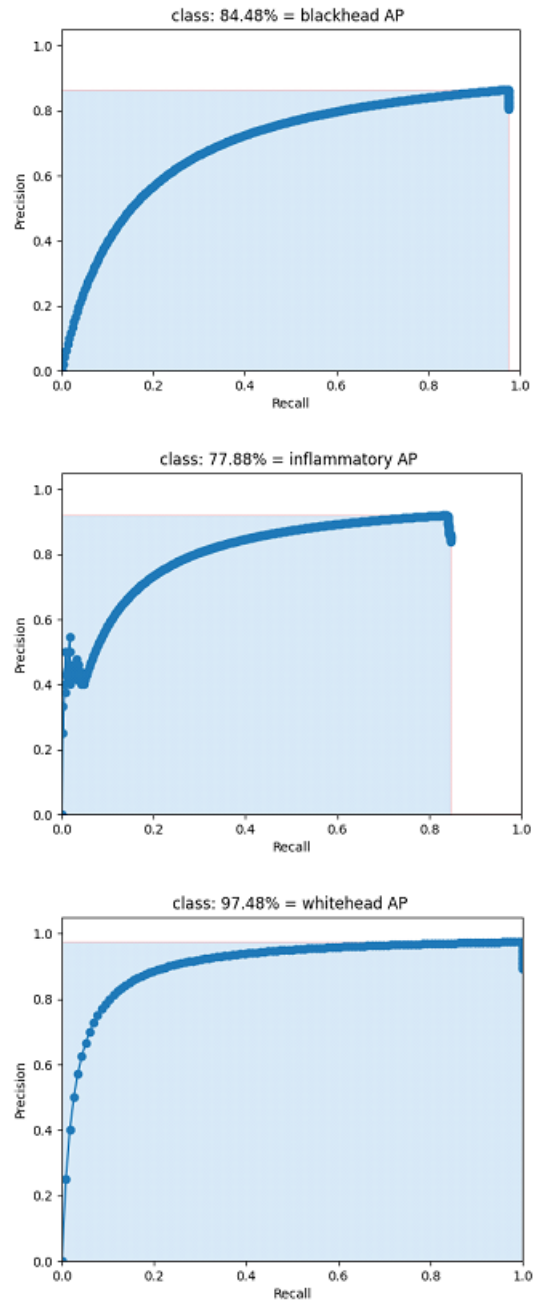


Figure 5.8: Precision vs. Recall Curves for FR-CNN

CHAPTER 6

CONCLUSION

Acne detection methods that specifically localise or classify acne lesions either based on severity or type are being used by dermatologists in day to day lives. A novel deep learning system that can both localise and classify acne lesions has been studied and developed in this project.

As we compare various object detection algorithms that were used to detect acne lesions, we come to the conclusion that Faster-RCNN performs better than YOLO and SSD Multibox in terms of Mean Average Precision obtained by calculating the area under the precision-recall curve for all the classes. One of the reasons for this is the two-network structure of Faster-RCNN. The model first identifies all the objects in the image and then moves further ahead to localise and classify them accurately. By changing the anchor scales in Faster-RCNN from the norm, we were able to successfully detect small lesions such as blackheads and whiteheads as well.

There were several limitations that were encountered along the way which are mentioned in the report. Some of these limitations were resolved using data augmentation techniques and optimisations in the model architecture. These changes helped achieve a better model to detect acne lesions in terms of precision.

This project is aimed to set a benchmark for acne detection models in the future. The images used in the future to develop better models must include high quality images with proper labeling. The study also shows that there is a dire

need to maintain a centralised dataset for acne lesion detection to build stronger and more powerful models in the future.

CHAPTER 7

FUTURE ENHANCEMENTS

The research in applications of Artificial Intelligence in the field of healthcare is fast growing and starting to make its presence. Based on the results and work provided in this work, there's a large scope for further development and improvement of current methods.

7.1 Dynamic Region of Interest (ROI) Extraction

Currently, any region-of-interest is extracted manually before being fed into the model. However, in a future enhancement it would be possible to develop a way to extract the regions automatically through image processing techniques. The location of the pupil could serve as a reference point using which the system would be able to obtain cheek, nose and forehead areas to be treated as separate images for the training process. This would remove the background outliers leading to an increase in efficiency and performance.

7.2 Sub-Classification of Inflammatory Class Acne

There are several sub classifications of inflammatory acne such as inflamed comedones, papules, pustules, nodules and cysts. These types can be diagnosed using object detection algorithms if there is enough data available. Improvement of models to detect a wide range of classes can make the model applicable in the real life healthcare setting.

7.3 Severity and classification

Current research work related to acne, focuses on specific problems such as calculating the severity of acne or detecting the classes using image data. After gaining sufficient maturity of work in specific fields, these can be clubbed to provide broader results.

7.4 Forming a benchmark dataset

Currently the research being done on the acne image data, is using multiple sources to extract and collect images from multiple sources. Most of the researchers who have worked on problems related to acne images, have not made the dataset open source. The need of benchmark datasets for AI applications in healthcare is critical for future research and development.

7.5 Mobile application development

To aid dermatologists in the real time diagnosis process, a mobile application can be really helpful. Therefore the development and optimization of lighter object detection models can immensely improve the mobile diagnosis process. These models can run tests on the mobile backend to analyse input data and give the output instantly.

APPENDIX A

CODING

A.1 Pre-processing

```
1 from os import listdir
2 from os.path import isfile, join
3 import os
4 import cv2
5 from PIL import Image
6 import imutils
7 from xml.dom import minidom
8 import xml.etree.ElementTree as ET
9
10 # change this to your path
11 annotations_path = "/Users/ronald/Desktop/work/sengupta/
    acne_deep_learning/keras-frcnn/ACNE/Annotations/"
12 images_path = "/Users/ronald/Desktop/work/sengupta/
    acne_deep_learning/keras-frcnn/ACNE/JPEGImages/"
13
14 # get all xml files in a list
15 annotations_files = sorted(
16     [
17         join(annotations_path, f)
18         for f in listdir(annotations_path)
19         if isfile(join(annotations_path, f))
20     ]
21 )
22
23 # change <path> content to the filename.
```

```

24 for f in annotations_files:
25     tree = ET.parse(f)
26     root = tree.getroot()
27     for item in root.iter("path"):
28         item.text = f
29     tree.write(f)
30     print("File -", f, "written.")

```

A.2 YOLO Training

A.2.1 Darknet

```

1 from ctypes import *
2 import math
3 import random
4
5 def sample(probs):
6     s = sum(probs)
7     probs = [a/s for a in probs]
8     r = random.uniform(0, 1)
9     for i in range(len(probs)):
10         r = r - probs[i]
11         if r <= 0:
12             return i
13     return len(probs)-1
14
15 def c_array(ctype, values):
16     arr = (ctype*len(values))()
17     arr[:] = values
18     return arr
19
20 class BOX(Structure):
21     _fields_ = [("x", c_float),
22                 ("y", c_float),

```

```

23         ("w", c_float),
24         ("h", c_float)]
25
26 class DETECTION(Structure):
27     _fields_ = [("bbox", BOX),
28                 ("classes", c_int),
29                 ("prob", POINTER(c_float)),
30                 ("mask", POINTER(c_float)),
31                 ("objectness", c_float),
32                 ("sort_class", c_int)]
33
34
35 class IMAGE(Structure):
36     _fields_ = [("w", c_int),
37                 ("h", c_int),
38                 ("c", c_int),
39                 ("data", POINTER(c_float))]
40
41 class METADATA(Structure):
42     _fields_ = [("classes", c_int),
43                 ("names", POINTER(c_char_p))]
44
45
46
47 #lib = CDLL("/home/pjreddie/documents/darknet/libdarknet.so",
48             RTLD_GLOBAL)
49 lib = CDLL("libdarknet.so", RTLD_GLOBAL)
50 lib.network_width.argtypes = [c_void_p]
51 lib.network_width.restype = c_int
52 lib.network_height.argtypes = [c_void_p]
53 lib.network_height.restype = c_int
54
55 predict = lib.network_predict
56 predict.argtypes = [c_void_p, POINTER(c_float)]
57 predict.restype = POINTER(c_float)

```

```

57
58 set_gpu = lib.cuda_set_device
59 set_gpu.argtypes = [c_int]
60
61 make_image = lib.make_image
62 make_image.argtypes = [c_int, c_int, c_int]
63 make_image.restype = IMAGE
64
65 get_network_boxes = lib.get_network_boxes
66 get_network_boxes.argtypes = [c_void_p, c_int, c_int, c_float,
    c_float, POINTER(c_int), c_int, POINTER(c_int)]
67 get_network_boxes.restype = POINTER(DETECTION)
68
69 make_network_boxes = lib.make_network_boxes
70 make_network_boxes.argtypes = [c_void_p]
71 make_network_boxes.restype = POINTER(DETECTION)
72
73 free_detections = lib.free_detections
74 free_detections.argtypes = [POINTER(DETECTION), c_int]
75
76 free_ptrs = lib.free_ptrs
77 free_ptrs.argtypes = [POINTER(c_void_p), c_int]
78
79 network_predict = lib.network_predict
80 network_predict.argtypes = [c_void_p, POINTER(c_float)]
81
82 reset_rnn = lib.reset_rnn
83 reset_rnn.argtypes = [c_void_p]
84
85 load_net = lib.load_network
86 load_net.argtypes = [c_char_p, c_char_p, c_int]
87 load_net.restype = c_void_p
88
89 do_nms_obj = lib.do_nms_obj
90 do_nms_obj.argtypes = [POINTER(DETECTION), c_int, c_int,

```

```

    c_float]
91
92 do_nms_sort = lib.do_nms_sort
93 do_nms_sort.argtypes = [POINTER(DETECTION), c_int, c_int,
    c_float]
94
95 free_image = lib.free_image
96 free_image.argtypes = [IMAGE]
97
98 letterbox_image = lib.letterbox_image
99 letterbox_image.argtypes = [IMAGE, c_int, c_int]
100 letterbox_image.restype = IMAGE
101
102 load_meta = lib.get_metadata
103 lib.get_metadata.argtypes = [c_char_p]
104 lib.get_metadata.restype = METADATA
105
106 load_image = lib.load_image_color
107 load_image.argtypes = [c_char_p, c_int, c_int]
108 load_image.restype = IMAGE
109
110 rgbgr_image = lib.rgbgr_image
111 rgbgr_image.argtypes = [IMAGE]
112
113 predict_image = lib.network_predict_image
114 predict_image.argtypes = [c_void_p, IMAGE]
115 predict_image.restype = POINTER(c_float)
116
117 def classify(net, meta, im):
118     out = predict_image(net, im)
119     res = []
120     for i in range(meta.classes):
121         res.append((meta.names[i], out[i]))
122     res = sorted(res, key=lambda x: -x[1])
123     return res

```



```

124
125 def detect(net, meta, image, thresh=.5, hier_thresh=.5, nms
    =.45):
126     im = load_image(image, 0, 0)
127     num = c_int(0)
128     pnum = pointer(num)
129     predict_image(net, im)
130     dets = get_network_boxes(net, im.w, im.h, thresh,
    hier_thresh, None, 0, pnum)
131     num = pnum[0]
132     if (nms): do_nms_obj(dets, num, meta.classes, nms);
133
134     res = []
135     for j in range(num):
136         for i in range(meta.classes):
137             if dets[j].prob[i] > 0:
138                 b = dets[j].bbox
139                 res.append((meta.names[i], dets[j].prob[i], (b.
    x, b.y, b.w, b.h)))
140     res = sorted(res, key=lambda x: -x[1])
141     free_image(im)
142     free_detections(dets, num)
143     return res
144
145 if __name__ == "__main__":
146     #net = load_net("cfg/densenet201.cfg", "/home/pjreddie/
    trained/densenet201.weights", 0)
147     #im = load_image("data/wolf.jpg", 0, 0)
148     #meta = load_meta("cfg/imagenet1k.data")
149     #r = classify(net, meta, im)
150     #print r[:10]
151     net = load_net("cfg/tiny-yolo.cfg", "tiny-yolo.weights", 0)
152     meta = load_meta("cfg/coco.data")
153     r = detect(net, meta, "data/dog.jpg")
154     print r

```

A.2.2 Custom CFG File

```
1
2 [net]
3 #Testing
4 #batch=1
5 #subdivisions=1
6 #Training
7 batch=64
8 subdivisions=64
9 width=416
10 height=416
11 channels=3
12 momentum=0.9
13 decay=0.0005
14 angle=0
15 saturation = 1.5
16 exposure = 1.5
17 hue=.1
18
19 learning_rate=0.001
20 burn_in=100
21 max_batches = 6000
22 policy=steps
23 steps=5400
24 scales=.1,.1
25
26 [convolutional]
27 batch_normalize=1
28 filters=32
29 size=3
30 stride=1
31 pad=1
32 activation=leaky
```

```
33
34 # Downsample
35
36 [convolutional]
37 batch_normalize=1
38 filters=64
39 size=3
40 stride=2
41 pad=1
42 activation=leaky
43
44 [convolutional]
45 batch_normalize=1
46 filters=32
47 size=1
48 stride=1
49 pad=1
50 activation=leaky
51
52 [convolutional]
53 batch_normalize=1
54 filters=64
55 size=3
56 stride=1
57 pad=1
58 activation=leaky
59
60 [shortcut]
61 from=-3
62 activation=linear
63
64 # Downsample
65
66 [convolutional]
67 batch_normalize=1
```

```
68 filters=128
69 size=3
70 stride=2
71 pad=1
72 activation=leaky
73
74 [convolutional]
75 batch_normalize=1
76 filters=64
77 size=1
78 stride=1
79 pad=1
80 activation=leaky
81
82 [convolutional]
83 batch_normalize=1
84 filters=128
85 size=3
86 stride=1
87 pad=1
88 activation=leaky
89
90 [shortcut]
91 from=-3
92 activation=linear
93
94 [convolutional]
95 batch_normalize=1
96 filters=64
97 size=1
98 stride=1
99 pad=1
100 activation=leaky
101
102 [convolutional]
```

```
103 batch_normalize=1
104 filters=128
105 size=3
106 stride=1
107 pad=1
108 activation=leaky
109
110 [shortcut]
111 from=-3
112 activation=linear
113
114 # Downsample
115
116 [convolutional]
117 batch_normalize=1
118 filters=256
119 size=3
120 stride=2
121 pad=1
122 activation=leaky
123
124 [convolutional]
125 batch_normalize=1
126 filters=128
127 size=1
128 stride=1
129 pad=1
130 activation=leaky
131
132 [convolutional]
133 batch_normalize=1
134 filters=256
135 size=3
136 stride=1
137 pad=1
```

```
138 activation=leaky
139
140 [shortcut]
141 from=-3
142 activation=linear
143
144 [convolutional]
145 batch_normalize=1
146 filters=128
147 size=1
148 stride=1
149 pad=1
150 activation=leaky
151
152 [convolutional]
153 batch_normalize=1
154 filters=256
155 size=3
156 stride=1
157 pad=1
158 activation=leaky
159
160 [shortcut]
161 from=-3
162 activation=linear
163
164 [convolutional]
165 batch_normalize=1
166 filters=128
167 size=1
168 stride=1
169 pad=1
170 activation=leaky
171
172 [convolutional]
```

```
173 batch_normalize=1
174 filters=256
175 size=3
176 stride=1
177 pad=1
178 activation=leaky
179
180 [shortcut]
181 from=-3
182 activation=linear
183
184 [convolutional]
185 batch_normalize=1
186 filters=128
187 size=1
188 stride=1
189 pad=1
190 activation=leaky
191
192 [convolutional]
193 batch_normalize=1
194 filters=256
195 size=3
196 stride=1
197 pad=1
198 activation=leaky
199
200 [shortcut]
201 from=-3
202 activation=linear
203
204
205 [convolutional]
206 batch_normalize=1
207 filters=128
```

```
208 size=1
209 stride=1
210 pad=1
211 activation=leaky
212
213 [convolutional]
214 batch_normalize=1
215 filters=256
216 size=3
217 stride=1
218 pad=1
219 activation=leaky
220
221 [shortcut]
222 from=-3
223 activation=linear
224
225 [convolutional]
226 batch_normalize=1
227 filters=128
228 size=1
229 stride=1
230 pad=1
231 activation=leaky
232
233 [convolutional]
234 batch_normalize=1
235 filters=256
236 size=3
237 stride=1
238 pad=1
239 activation=leaky
240
241 [shortcut]
242 from=-3
```



```
243 activation=linear
244
245 [convolutional]
246 batch_normalize=1
247 filters=128
248 size=1
249 stride=1
250 pad=1
251 activation=leaky
252
253 [convolutional]
254 batch_normalize=1
255 filters=256
256 size=3
257 stride=1
258 pad=1
259 activation=leaky
260
261 [shortcut]
262 from=-3
263 activation=linear
264
265 [convolutional]
266 batch_normalize=1
267 filters=128
268 size=1
269 stride=1
270 pad=1
271 activation=leaky
272
273 [convolutional]
274 batch_normalize=1
275 filters=256
276 size=3
277 stride=1
```

```
278 pad=1
279 activation=leaky
280
281 [shortcut]
282 from=-3
283 activation=linear
284
285 # Downsample
286
287 [convolutional]
288 batch_normalize=1
289 filters=512
290 size=3
291 stride=2
292 pad=1
293 activation=leaky
294
295 [convolutional]
296 batch_normalize=1
297 filters=256
298 size=1
299 stride=1
300 pad=1
301 activation=leaky
302
303 [convolutional]
304 batch_normalize=1
305 filters=512
306 size=3
307 stride=1
308 pad=1
309 activation=leaky
310
311 [shortcut]
312 from=-3
```

```
313 activation=linear
314
315
316 [convolutional]
317 batch_normalize=1
318 filters=256
319 size=1
320 stride=1
321 pad=1
322 activation=leaky
323
324 [convolutional]
325 batch_normalize=1
326 filters=512
327 size=3
328 stride=1
329 pad=1
330 activation=leaky
331
332 [shortcut]
333 from=-3
334 activation=linear
335
336
337 [convolutional]
338 batch_normalize=1
339 filters=256
340 size=1
341 stride=1
342 pad=1
343 activation=leaky
344
345 [convolutional]
346 batch_normalize=1
347 filters=512
```

```
348 size=3
349 stride=1
350 pad=1
351 activation=leaky
352
353 [shortcut]
354 from=-3
355 activation=linear
356
357
358 [convolutional]
359 batch_normalize=1
360 filters=256
361 size=1
362 stride=1
363 pad=1
364 activation=leaky
365
366 [convolutional]
367 batch_normalize=1
368 filters=512
369 size=3
370 stride=1
371 pad=1
372 activation=leaky
373
374 [shortcut]
375 from=-3
376 activation=linear
377
378 [convolutional]
379 batch_normalize=1
380 filters=256
381 size=1
382 stride=1
```

```
383 pad=1
384 activation=leaky
385
386 [convolutional]
387 batch_normalize=1
388 filters=512
389 size=3
390 stride=1
391 pad=1
392 activation=leaky
393
394 [shortcut]
395 from=-3
396 activation=linear
397
398
399 [convolutional]
400 batch_normalize=1
401 filters=256
402 size=1
403 stride=1
404 pad=1
405 activation=leaky
406
407 [convolutional]
408 batch_normalize=1
409 filters=512
410 size=3
411 stride=1
412 pad=1
413 activation=leaky
414
415 [shortcut]
416 from=-3
417 activation=linear
```

```
418
419
420 [convolutional]
421 batch_normalize=1
422 filters=256
423 size=1
424 stride=1
425 pad=1
426 activation=leaky
427
428 [convolutional]
429 batch_normalize=1
430 filters=512
431 size=3
432 stride=1
433 pad=1
434 activation=leaky
435
436 [shortcut]
437 from=-3
438 activation=linear
439
440 [convolutional]
441 batch_normalize=1
442 filters=256
443 size=1
444 stride=1
445 pad=1
446 activation=leaky
447
448 [convolutional]
449 batch_normalize=1
450 filters=512
451 size=3
452 stride=1
```

```
453 pad=1
454 activation=leaky
455
456 [shortcut]
457 from=-3
458 activation=linear
459
460 # Downsample
461
462 [convolutional]
463 batch_normalize=1
464 filters=1024
465 size=3
466 stride=2
467 pad=1
468 activation=leaky
469
470 [convolutional]
471 batch_normalize=1
472 filters=512
473 size=1
474 stride=1
475 pad=1
476 activation=leaky
477
478 [convolutional]
479 batch_normalize=1
480 filters=1024
481 size=3
482 stride=1
483 pad=1
484 activation=leaky
485
486 [shortcut]
487 from=-3
```

```
488 activation=linear
489
490 [convolutional]
491 batch_normalize=1
492 filters=512
493 size=1
494 stride=1
495 pad=1
496 activation=leaky
497
498 [convolutional]
499 batch_normalize=1
500 filters=1024
501 size=3
502 stride=1
503 pad=1
504 activation=leaky
505
506 [shortcut]
507 from=-3
508 activation=linear
509
510 [convolutional]
511 batch_normalize=1
512 filters=512
513 size=1
514 stride=1
515 pad=1
516 activation=leaky
517
518 [convolutional]
519 batch_normalize=1
520 filters=1024
521 size=3
522 stride=1
```



```

523 pad=1
524 activation=leaky
525
526 [shortcut]
527 from=-3
528 activation=linear
529
530 [convolutional]
531 batch_normalize=1
532 filters=512
533 size=1
534 stride=1
535 pad=1
536 activation=leaky
537
538 [convolutional]
539 batch_normalize=1
540 filters=1024
541 size=3
542 stride=1
543 pad=1
544 activation=leaky
545
546 [shortcut]
547 from=-3
548 activation=linear
549 #####
550 [convolutional]
551 batch_normalize=1
552 filters=512
553 size=1
554 stride=1
555 pad=1
556 activation=leaky
557 [convolutional]

```

```
558 batch_normalize=1
559 size=3
560 stride=1
561 pad=1
562 filters=1024
563 activation=leaky
564 [convolutional]
565 batch_normalize=1
566 filters=512
567 size=1
568 stride=1
569 pad=1
570 activation=leaky
571 [convolutional]
572 batch_normalize=1
573 size=3
574 stride=1
575 pad=1
576 filters=1024
577 activation=leaky
578 [convolutional]
579 batch_normalize=1
580 filters=512
581 size=1
582 stride=1
583 pad=1
584 activation=leaky
585 [convolutional]
586 batch_normalize=1
587 size=3
588 stride=1
589 pad=1
590 filters=1024
591 activation=leaky
592 [convolutional]
```

```

593 size=1
594 stride=1
595 pad=1
596 filters=24
597 activation=linear
598 [yolo]
599 mask = 6,7,8
600 anchors = 10,13, 16,30, 33,23, 30,61, 62,45, 59,119,
        116,90, 156,198, 373,326
601 classes=3
602 num=9
603 jitter=.3
604 ignore_thresh = .7
605 truth_thresh = 1
606 random=1
607 [route]
608 layers = -4
609 [convolutional]
610 batch_normalize=1
611 filters=256
612 size=1
613 stride=1
614 pad=1
615 activation=leaky
616 [upsample]
617 stride=2
618 [route]
619 layers = -1, 61
620 [convolutional]
621 batch_normalize=1
622 filters=256
623 size=1
624 stride=1
625 pad=1
626 activation=leaky

```

```
627 [convolutional]
628 batch_normalize=1
629 size=3
630 stride=1
631 pad=1
632 filters=512
633 activation=leaky
634 [convolutional]
635 batch_normalize=1
636 filters=256
637 size=1
638 stride=1
639 pad=1
640 activation=leaky
641 [convolutional]
642 batch_normalize=1
643 size=3
644 stride=1
645 pad=1
646 filters=512
647 activation=leaky
648 [convolutional]
649 batch_normalize=1
650 filters=256
651 size=1
652 stride=1
653 pad=1
654 activation=leaky
655 [convolutional]
656 batch_normalize=1
657 size=3
658 stride=1
659 pad=1
660 filters=512
661 activation=leaky
```

```

662 [convolutional]
663 size=1
664 stride=1
665 pad=1
666 filters=24
667 activation=linear
668 [yolo]
669 mask = 3,4,5
670 anchors = 10,13, 16,30, 33,23, 30,61, 62,45, 59,119,
        116,90, 156,198, 373,326
671 classes=3
672 num=9
673 jitter=.3
674 ignore_thresh = .7
675 truth_thresh = 1
676 random=1
677 [route]
678 layers = -4
679 [convolutional]
680 batch_normalize=1
681 filters=128
682 size=1
683 stride=1
684 pad=1
685 activation=leaky
686 [upsample]
687 stride=2
688 [route]
689 layers = -1, 36
690 [convolutional]
691 batch_normalize=1
692 filters=128
693 size=1
694 stride=1
695 pad=1

```

```
696 activation=leaky
697 [convolutional]
698 batch_normalize=1
699 size=3
700 stride=1
701 pad=1
702 filters=256
703 activation=leaky
704 [convolutional]
705 batch_normalize=1
706 filters=128
707 size=1
708 stride=1
709 pad=1
710 activation=leaky
711 [convolutional]
712 batch_normalize=1
713 size=3
714 stride=1
715 pad=1
716 filters=256
717 activation=leaky
718 [convolutional]
719 batch_normalize=1
720 filters=128
721 size=1
722 stride=1
723 pad=1
724 activation=leaky
725 [convolutional]
726 batch_normalize=1
727 size=3
728 stride=1
729 pad=1
730 filters=256
```

```

731 activation=leaky
732 [convolutional]
733 size=1
734 stride=1
735 pad=1
736 filters=24
737 activation=linear
738 [yolo]
739 mask = 0,1,2
740 anchors = 10,13, 16,30, 33,23, 30,61, 62,45, 59,119,
          116,90, 156,198, 373,326
741 classes=3
742 num=9
743 jitter=.3
744 ignore_thresh = .7
745 truth_thresh = 1
746 random=1

```

A.3 SSD Training

```

1 model {
2   ssd {
3     num_classes: 3
4     box_coder {
5       faster_rcnn_box_coder {
6         y_scale: 10.0
7         x_scale: 10.0
8         height_scale: 5.0
9         width_scale: 5.0
10    }
11  }
12  matcher {
13    argmax_matcher {
14      matched_threshold: 0.5

```

```

15     unmatched_threshold: 0.5
16     ignore_thresholds: false
17     negatives_lower_than_unmatched: true
18     force_match_for_each_row: true
19 }
20 }
21 similarity_calculator {
22     iou_similarity {
23     }
24 }
25 anchor_generator {
26     ssd_anchor_generator {
27         num_layers: 6
28         min_scale: 0.05
29         max_scale: 0.95
30         aspect_ratios: 1.0
31         aspect_ratios: 2.0
32         aspect_ratios: 0.5
33         #aspect_ratios: 3.0
34         #aspect_ratios: 0.3333
35         reduce_boxes_in_lowest_layer: true
36     }
37 }
38 image_resizer {
39     fixed_shape_resizer {
40         height: 300
41         width: 300
42     }
43 }
44 box_predictor {
45     convolutional_box_predictor {
46         min_depth: 0
47         max_depth: 0
48         num_layers_before_predictor: 0
49         use_dropout: false

```



```

50     dropout_keep_probability: 0.8
51     kernel_size: 3
52     box_code_size: 4
53     apply_sigmoid_to_scores: false
54     conv_hyperparams {
55         activation: RELU_6,
56         regularizer {
57             l2_regularizer {
58                 weight: 0.00004
59             }
60         }
61         initializer {
62             truncated_normal_initializer {
63                 stddev: 0.03
64                 mean: 0.0
65             }
66         }
67     }
68 }
69 }
70 feature_extractor {
71     type: 'ssd_inception_v2'
72     min_depth: 16
73     depth_multiplier: 1.0
74     conv_hyperparams {
75         activation: RELU_6,
76         regularizer {
77             l2_regularizer {
78                 weight: 0.00004
79             }
80         }
81         initializer {
82             truncated_normal_initializer {
83                 stddev: 0.03
84                 mean: 0.0

```

```

85         }
86     }
87     batch_norm {
88         train: true,
89         scale: true,
90         center: true,
91         decay: 0.9997,
92         epsilon: 0.001,
93     }
94 }
95 override_base_feature_extractor_hyperparams: true
96 }
97 loss {
98     classification_loss {
99         weighted_sigmoid {
100         }
101     }
102     localization_loss {
103         weighted_smooth_l1 {
104         }
105     }
106     hard_example_miner {
107         num_hard_examples: 3000
108         iou_threshold: 0.99
109         loss_type: CLASSIFICATION
110         max_negatives_per_positive: 3
111         min_negatives_per_image: 0
112     }
113     classification_weight: 1.0
114     localization_weight: 1.0
115 }
116 normalize_loss_by_num_matches: true
117 post_processing {
118     batch_non_max_suppression {
119         score_threshold: 1e-8

```

```

120     iou_threshold: 0.6
121     max_detections_per_class: 100
122     max_total_detections: 100
123   }
124   score_converter: SIGMOID
125 }
126 }
127 }
128
129 train_config: {
130   batch_size: 8
131   optimizer {
132     rms_prop_optimizer: {
133       learning_rate: {
134         exponential_decay_learning_rate {
135           initial_learning_rate: 0.004
136           decay_steps: 800720
137           decay_factor: 0.95
138         }
139       }
140       momentum_optimizer_value: 0.9
141       decay: 0.9
142       epsilon: 1.0
143     }
144   }
145   fine_tune_checkpoint: "pre-trained-model/model.ckpt"
146   from_detection_checkpoint: true
147   # Note: The below line limits the training process to 200K
148   # steps, which we
149   # empirically found to be sufficient enough to train the pets
150   # dataset. This
151   # effectively bypasses the learning rate schedule (the
152   # learning rate will
153   # never decay). Remove the below line to train indefinitely.
154   num_steps: 200000

```

```

152 data_augmentation_options {
153     random_horizontal_flip {
154     }
155 }
156 data_augmentation_options {
157     ssd_random_crop {
158     }
159 }
160 }
161
162 train_input_reader: {
163     tf_record_input_reader {
164         input_path: "annotations/train.record"
165     }
166     label_map_path: "annotations/label_map.pbtxt"
167 }
168
169 eval_config: {
170     #metrics_set: "coco_detection_metrics"
171     num_examples: 62
172     # Note: The below line limits the evaluation process to 10
173     # evaluations.
174     # Remove the below line to evaluate indefinitely.
175     max_evals: 100
176 }
177
178 eval_input_reader: {
179     tf_record_input_reader {
180         input_path: "annotations/test.record"
181     }
182     label_map_path: "annotations/label_map.pbtxt"
183     shuffle: false
184     num_readers: 1
185 }

```

A.4 FR-CNN Training

```
1 from __future__ import division
2 import random
3 import pprint
4 import sys
5 import time
6 import numpy as np
7 from optparse import OptionParser
8 import pickle
9
10 from keras import backend as K
11 from keras.optimizers import Adam, SGD, RMSprop
12 from keras.layers import Input
13 from keras.models import Model
14 from keras_frcnn import config, data_generators
15 from keras_frcnn import losses as losses
16 import keras_frcnn.roi_helpers as roi_helpers
17 from keras.utils import generic_utils
18
19 sys.setrecursionlimit(40000)
20
21 parser = OptionParser()
22
23 parser.add_option("-p", "--path", dest="train_path", help="Path
    to training data.")
24 parser.add_option(
25     "-o",
26     "--parser",
27     dest="parser",
28     help="Parser to use. One of simple or pascal_voc",
29     default="pascal_voc",
30 )
31 parser.add_option(
```

```

32     "-n",
33     "--num_rois",
34     type="int",
35     dest="num_rois",
36     help="Number of RoIs to process at once.",
37     default=32,
38 )
39 parser.add_option(
40     "--network",
41     dest="network",
42     help="Base network to use. Supports vgg or resnet50.",
43     default="resnet50",
44 )
45 parser.add_option(
46     "--hf",
47     dest="horizontal_flips",
48     help="Augment with horizontal flips in training. (Default=
49 false).",
50     action="store_true",
51     default=False,
52 )
53 parser.add_option(
54     "--vf",
55     dest="vertical_flips",
56     help="Augment with vertical flips in training. (Default=
57 false).",
58     action="store_true",
59     default=False,
60 )
61 parser.add_option(
62     "--rot",
63     "--rot_90",
64     dest="rot_90",
65     help="Augment with 90 degree rotations in training. (
66     Default=false).",

```

```

64     action="store_true",
65     default=False,
66 )
67 parser.add_option(
68     "--num_epochs",
69     type="int",
70     dest="num_epochs",
71     help="Number of epochs.",
72     default=2000,
73 )
74 parser.add_option(
75     "--config_filename",
76     dest="config_filename",
77     help="Location to store all the metadata related to the
78     training (to be used when testing).",
79     default="config.pickle",
80 )
81 parser.add_option(
82     "--output_weight_path",
83     dest="output_weight_path",
84     help="Output path for weights.",
85     default="./model_frcnn.hdf5",
86 )
87 parser.add_option(
88     "--input_weight_path",
89     dest="input_weight_path",
90     help="Input path for weights. If not specified, will try to
91     load default weights provided by keras.",
92 )
93 (options, args) = parser.parse_args()
94 if not options.train_path: # if filename is not given
95     parser.error(
96         "Error: path to training data must be specified. Pass

```

```

--path to command line"
97     )
98
99 if options.parser == "pascal_voc":
100     from keras_frcnn.pascal_voc_parser import get_data
101 elif options.parser == "simple":
102     from keras_frcnn.simple_parser import get_data
103 else:
104     raise ValueError(
105         "Command line option parser must be one of 'pascal_voc'
106         or 'simple'"
107     )
108 # pass the settings from the command line, and persist them in
109 # the config object
110 C = config.Config()
111 C.use_horizontal_flips = bool(options.horizontal_flips)
112 C.use_vertical_flips = bool(options.vertical_flips)
113 C.rot_90 = bool(options.rot_90)
114
115 C.model_path = options.output_weight_path
116 C.num_rois = int(options.num_rois)
117
118 if options.network == "vgg":
119     C.network = "vgg"
120     from keras_frcnn import vgg as nn
121 elif options.network == "resnet50":
122     from keras_frcnn import resnet as nn
123
124     C.network = "resnet50"
125 else:
126     print("Not a valid model")
127     raise ValueError
128

```



```

129
130 # check if weight path was passed via command line
131 if options.input_weight_path:
132     C.base_net_weights = options.input_weight_path
133 else:
134     # set the path to weights based on backend and model
135     C.base_net_weights = nn.get_weight_path()
136
137 all_imgs, classes_count, class_mapping = get_data(options.
    train_path)
138
139 if "bg" not in classes_count:
140     classes_count["bg"] = 0
141     class_mapping["bg"] = len(class_mapping)
142
143 C.class_mapping = class_mapping
144
145 inv_map = {v: k for k, v in class_mapping.items()}
146
147 print("Training images per class:")
148 pprint.pprint(classes_count)
149 print("Num classes (including bg) = {}".format(len(
    classes_count)))
150
151 config_output_filename = options.config_filename
152
153 with open(config_output_filename, "wb") as config_f:
154     pickle.dump(C, config_f)
155     print(
156         "Config has been written to {}, and can be loaded when
    testing to ensure correct results".format(
157             config_output_filename
158         )
159     )
160

```

```

161 random.shuffle(all_imgs)
162
163 num_imgs = len(all_imgs)
164
165 train_imgs = [s for s in all_imgs if s["imageset"] == "trainval
    "]
166 val_imgs = [s for s in all_imgs if s["imageset"] == "test"]
167
168 print("Num train samples {}".format(len(train_imgs)))
169 print("Num val samples {}".format(len(val_imgs)))
170
171
172 data_gen_train = data_generators.get_anchor_gt(
173     train_imgs,
174     classes_count,
175     C,
176     nn.get_img_output_length,
177     K.image_dim_ordering(),
178     mode="train",
179 )
180 data_gen_val = data_generators.get_anchor_gt(
181     val_imgs,
182     classes_count,
183     C,
184     nn.get_img_output_length,
185     K.image_dim_ordering(),
186     mode="val",
187 )
188
189 if K.image_dim_ordering() == "th":
190     input_shape_img = (3, None, None)
191 else:
192     input_shape_img = (None, None, 3)
193
194 img_input = Input(shape=input_shape_img)

```

```

195 roi_input = Input(shape=(None, 4))
196
197 # define the base network (resnet here, can be VGG, Inception,
    etc)
198 shared_layers = nn.nn_base(img_input, trainable=True)
199
200 # define the RPN, built on the base layers
201 num_anchors = len(C.anchor_box_scales) * len(C.
    anchor_box_ratios)
202 rpn = nn.rpn(shared_layers, num_anchors)
203
204 classifier = nn.classifier(
205     shared_layers, roi_input, C.num_rois, nb_classes=len(
    classes_count), trainable=True
206 )
207
208 model_rpn = Model(img_input, rpn[:2])
209 model_classifier = Model([img_input, roi_input], classifier)
210
211 # this is a model that holds both the RPN and the classifier,
    used to load/save weights for the models
212 model_all = Model([img_input, roi_input], rpn[:2] + classifier)
213
214 try:
215     print("loading weights from {}".format(C.base_net_weights))
216     model_rpn.load_weights(C.base_net_weights, by_name=True)
217     model_classifier.load_weights(C.base_net_weights, by_name=
    True)
218 except:
219     print(
220         "Could not load pretrained model weights. Weights can
    be found in the keras application folder \
221     https://github.com/fchollet/keras/tree/master/keras/
    applications"
222     )

```

```

223
224 optimizer = Adam(lr=1e-5)
225 optimizer_classifier = Adam(lr=1e-5)
226 model_rpn.compile(
227     optimizer=optimizer,
228     loss=[losses.rpn_loss_cls(num_anchors), losses.
229         rpn_loss_regr(num_anchors)],
230 )
231 model_classifier.compile(
232     optimizer=optimizer_classifier,
233     loss=[losses.class_loss_cls, losses.class_loss_regr(len(
234         classes_count) - 1)],
235     metrics={"dense_class_{}".format(len(classes_count)): "
236         accuracy"},
237 )
238 model_all.compile(optimizer="sgd", loss="mae")
239
240
241 epoch_length = 1000
242 num_epochs = int(options.num_epochs)
243 iter_num = 0
244
245 losses = np.zeros((epoch_length, 5))
246 rpn_accuracy_rpn_monitor = []
247 rpn_accuracy_for_epoch = []
248 start_time = time.time()
249
250 best_loss = np.Inf
251
252 class_mapping_inv = {v: k for k, v in class_mapping.items()}
253 print("Starting training")
254
255 vis = True
256
257 for epoch_num in range(num_epochs):

```

```

255 progbar = generic_utils.Progbar(epoch_length)
256 print("Epoch {}/{}".format(epoch_num + 1, num_epochs))
257
258 while True:
259     try:
260
261         if len(rpn_accuracy_rpn_monitor) == epoch_length
and C.verbose:
262             mean_overlapping_bboxes = float(sum(
rpn_accuracy_rpn_monitor)) / len(
263                 rpn_accuracy_rpn_monitor
264             )
265             rpn_accuracy_rpn_monitor = []
266             print(
267                 "Average number of overlapping bounding
boxes from RPN = {} for {} previous iterations".format(
268                     mean_overlapping_bboxes, epoch_length
269                 )
270             )
271             if mean_overlapping_bboxes == 0:
272                 print(
273                     "RPN is not producing bounding boxes
that overlap the ground truth boxes. Check RPN settings or
keep training."
274                 )
275
276             X, Y, img_data = next(data_gen_train)
277
278             loss_rpn = model_rpn.train_on_batch(X, Y)
279
280             P_rpn = model_rpn.predict_on_batch(X)
281
282             R = roi_helpers.rpn_to_roi(
283                 P_rpn[0],
284                 P_rpn[1],

```

```

285         C,
286         K.image_dim_ordering(),
287         use_regr=True,
288         overlap_thresh=0.7,
289         max_boxes=300,
290     )
291     # note: calc_iou converts from (x1,y1,x2,y2) to (x,
y,w,h) format
292     X2, Y1, Y2, IouS = roi_helpers.calc_iou(R, img_data
, C, class_mapping)
293
294     if X2 is None:
295         rpn_accuracy_rpn_monitor.append(0)
296         rpn_accuracy_for_epoch.append(0)
297         continue
298
299     neg_samples = np.where(Y1[0, :, -1] == 1)
300     pos_samples = np.where(Y1[0, :, -1] == 0)
301
302     if len(neg_samples) > 0:
303         neg_samples = neg_samples[0]
304     else:
305         neg_samples = []
306
307     if len(pos_samples) > 0:
308         pos_samples = pos_samples[0]
309     else:
310         pos_samples = []
311
312     rpn_accuracy_rpn_monitor.append(len(pos_samples))
313     rpn_accuracy_for_epoch.append((len(pos_samples)))
314
315     if C.num_rois > 1:
316         if len(pos_samples) < C.num_rois // 2:
317             selected_pos_samples = pos_samples.tolist()

```

```

318         else:
319             selected_pos_samples = np.random.choice(
320                 pos_samples, C.num_rois // 2, replace=
False
321                 ).tolist()
322         try:
323             selected_neg_samples = np.random.choice(
324                 neg_samples,
325                 C.num_rois - len(selected_pos_samples),
326                 replace=False,
327                 ).tolist()
328         except:
329             selected_neg_samples = np.random.choice(
330                 neg_samples,
331                 C.num_rois - len(selected_pos_samples),
332                 replace=True,
333                 ).tolist()
334
335         sel_samples = selected_pos_samples +
selected_neg_samples
336         else:
337             # in the extreme case where num_rois = 1, we
pick a random pos or neg sample
338             selected_pos_samples = pos_samples.tolist()
339             selected_neg_samples = neg_samples.tolist()
340             if np.random.randint(0, 2):
341                 sel_samples = random.choice(neg_samples)
342             else:
343                 sel_samples = random.choice(pos_samples)
344
345         loss_class = model_classifier.train_on_batch(
346             [X, X2[:, sel_samples, :]],
347             [Y1[:, sel_samples, :], Y2[:, sel_samples, :]],
348         )
349

```

```

350         losses[iter_num, 0] = loss_rpn[1]
351         losses[iter_num, 1] = loss_rpn[2]
352
353         losses[iter_num, 2] = loss_class[1]
354         losses[iter_num, 3] = loss_class[2]
355         losses[iter_num, 4] = loss_class[3]
356
357         progbar.update(
358             iter_num + 1,
359             [
360                 ("rpn_cls", losses[iter_num, 0]),
361                 ("rpn_regr", losses[iter_num, 1]),
362                 ("detector_cls", losses[iter_num, 2]),
363                 ("detector_regr", losses[iter_num, 3]),
364             ],
365         )
366
367         iter_num += 1
368
369         if iter_num == epoch_length:
370             loss_rpn_cls = np.mean(losses[:, 0])
371             loss_rpn_regr = np.mean(losses[:, 1])
372             loss_class_cls = np.mean(losses[:, 2])
373             loss_class_regr = np.mean(losses[:, 3])
374             class_acc = np.mean(losses[:, 4])
375
376             mean_overlapping_bboxes = float(sum(
377                 rpn_accuracy_for_epoch)) / len(
378                     rpn_accuracy_for_epoch
379                 )
380             rpn_accuracy_for_epoch = []
381
382             if C.verbose:
383                 print(
384                     "Mean number of bounding boxes from RPN

```



```

overlapping ground truth boxes: {}".format(
384         mean_overlapping_bboxes
385     )
386 )
387 print(
388     "Classifier accuracy for bounding boxes
from RPN: {}".format(
389         class_acc
390     )
391 )
392 print("Loss RPN classifier: {}".format(
loss_rpn_cls))
393 print("Loss RPN regression: {}".format(
loss_rpn_regr))
394 print("Loss Detector classifier: {}".format(
(loss_class_cls))
395 print("Loss Detector regression: {}".format(
(loss_class_regr))
396 print("Elapsed time: {}".format(time.time()
- start_time))
397
398 curr_loss = (
399     loss_rpn_cls + loss_rpn_regr +
loss_class_cls + loss_class_regr
400 )
401 iter_num = 0
402 start_time = time.time()
403
404 if curr_loss < best_loss:
405     if C.verbose:
406         print(
407             "Total loss decreased from {} to
{}, saving weights".format(
408                 best_loss, curr_loss
409         )

```

```

410         )
411         best_loss = curr_loss
412         model_all.save_weights(C.model_path)
413
414         break
415
416     except Exception as e:
417         print("Exception: {}".format(e))
418         continue
419
420 print("Training complete, exiting.")

```

A.5 Evaluation

```

1 import os
2 import cv2
3 import numpy as np
4 import sys
5 import pickle
6 from optparse import OptionParser
7 import time
8 from keras_frcnn import config
9 import keras_frcnn.resnet as nn
10 from keras import backend as K
11 from keras.layers import Input
12 from keras.models import Model
13 from keras_frcnn import roi_helpers
14 from keras_frcnn import data_generators
15 from sklearn.metrics import average_precision_score
16
17
18 def get_map(pred, gt, f):
19     T = {}
20     P = {}

```

```

21  fx, fy = f
22
23  for bbox in gt:
24      bbox['bbox_matched'] = False
25
26  pred_probs = np.array([s['prob'] for s in pred])
27  box_idx_sorted_by_prob = np.argsort(pred_probs)[::-1]
28
29  for box_idx in box_idx_sorted_by_prob:
30      pred_box = pred[box_idx]
31      pred_class = pred_box['class']
32      pred_x1 = pred_box['x1']
33      pred_x2 = pred_box['x2']
34      pred_y1 = pred_box['y1']
35      pred_y2 = pred_box['y2']
36      pred_prob = pred_box['prob']
37      if pred_class not in P:
38          P[pred_class] = []
39          T[pred_class] = []
40      P[pred_class].append(pred_prob)
41      found_match = False
42
43  for gt_box in gt:
44      gt_class = gt_box['class']
45      gt_x1 = gt_box['x1']/fx
46      gt_x2 = gt_box['x2']/fx
47      gt_y1 = gt_box['y1']/fy
48      gt_y2 = gt_box['y2']/fy
49      gt_seen = gt_box['bbox_matched']
50      if gt_class != pred_class:
51          continue
52      if gt_seen:
53          continue
54      iou = data_generators.iou((pred_x1, pred_y1, pred_x2,
pred_y2), (gt_x1, gt_y1, gt_x2, gt_y2))

```

```

55     if iou >= 0.5:
56         found_match = True
57         gt_box['bbox_matched'] = True
58         break
59     else:
60         continue
61
62     T[pred_class].append(int(found_match))
63
64 for gt_box in gt:
65     if not gt_box['bbox_matched'] and not gt_box['difficult']:
66         if gt_box['class'] not in P:
67             P[gt_box['class']] = []
68             T[gt_box['class']] = []
69
70             T[gt_box['class']].append(1)
71             P[gt_box['class']].append(0)
72
73     #import pdb
74     #pdb.set_trace()
75     return T, P
76
77 sys.setrecursionlimit(40000)
78
79 parser = OptionParser()
80
81 parser.add_option("-p", "--path", dest="test_path", help="Path
    to test data.")
82 parser.add_option("-n", "--num_rois", dest="num_rois",
83                 help="Number of ROIs per iteration. Higher means more
    memory use.", default=32)
84 parser.add_option("--config_filename", dest="config_filename",
    help=
85                 "Location to read the metadata related to the training
    (generated when training).",

```

```

86         default="config.pickle")
87 parser.add_option("-o", "--parser", dest="parser", help="Parser
    to use. One of simple or pascal_voc",
88         default="pascal_voc"),
89
90 (options, args) = parser.parse_args()
91
92 if not options.test_path:    # if filename is not given
93     parser.error('Error: path to test data must be specified.
    Pass --path to command line')
94
95
96 if options.parser == 'pascal_voc':
97     from keras_frcnn.pascal_voc_parser import get_data
98 elif options.parser == 'simple':
99     from keras_frcnn.simple_parser import get_data
100 else:
101     raise ValueError("Command line option parser must be one of '
    pascal_voc' or 'simple'")
102
103 config_output_filename = options.config_filename
104
105 with open(config_output_filename, 'r') as f_in:
106     C = pickle.load(f_in)
107
108 # turn off any data augmentation at test time
109 C.use_horizontal_flips = False
110 C.use_vertical_flips = False
111 C.rot_90 = False
112
113 img_path = options.test_path
114
115
116 def format_img(img, C):
117     img_min_side = float(C.im_size)

```

```

118 (height,width,_) = img.shape
119
120 if width <= height:
121     f = img_min_side/width
122     new_height = int(f * height)
123     new_width = int(img_min_side)
124 else:
125     f = img_min_side/height
126     new_width = int(f * width)
127     new_height = int(img_min_side)
128 fx = width/float(new_width)
129 fy = height/float(new_height)
130 img = cv2.resize(img, (new_width, new_height), interpolation=
    cv2.INTER_CUBIC)
131 img = img[:, :, (2, 1, 0)]
132 img = img.astype(np.float32)
133 img[:, :, 0] -= C.img_channel_mean[0]
134 img[:, :, 1] -= C.img_channel_mean[1]
135 img[:, :, 2] -= C.img_channel_mean[2]
136 img /= C.img_scaling_factor
137 img = np.transpose(img, (2, 0, 1))
138 img = np.expand_dims(img, axis=0)
139 return img, fx, fy
140
141
142 class_mapping = C.class_mapping
143
144 if 'bg' not in class_mapping:
145     class_mapping['bg'] = len(class_mapping)
146
147 class_mapping = {v: k for k, v in class_mapping.iteritems()}
148 print(class_mapping)
149 class_to_color = {class_mapping[v]: np.random.randint(0, 255,
    3) for v in class_mapping}
150 C.num_rois = int(options.num_rois)

```

```

151
152 if K.image_dim_ordering() == 'th':
153     input_shape_img = (3, None, None)
154     input_shape_features = (1024, None, None)
155 else:
156     input_shape_img = (None, None, 3)
157     input_shape_features = (None, None, 1024)
158
159
160 img_input = Input(shape=input_shape_img)
161 roi_input = Input(shape=(C.num_rois, 4))
162 feature_map_input = Input(shape=input_shape_features)
163
164 # define the base network (resnet here, can be VGG, Inception,
165     etc)
166
167 shared_layers = nn.nn_base(img_input, trainable=True)
168
169 # define the RPN, built on the base layers
170
171 num_anchors = len(C.anchor_box_scales) * len(C.
172     anchor_box_ratios)
173
174 rpn_layers = nn.rpn(shared_layers, num_anchors)
175
176 classifier = nn.classifier(feature_map_input, roi_input, C.
177     num_rois, nb_classes=len(class_mapping), trainable=True)
178
179
180 model_rpn = Model(img_input, rpn_layers)
181 model_classifier_only = Model([feature_map_input, roi_input],
182     classifier)
183
184 model_classifier = Model([feature_map_input, roi_input],
185     classifier)
186
187
188 model_rpn.load_weights(C.model_path, by_name=True)
189 model_classifier.load_weights(C.model_path, by_name=True)
190

```

```

181 model_rpn.compile(optimizer='sgd', loss='mse')
182 model_classifier.compile(optimizer='sgd', loss='mse')
183
184 all_imgs, _, _ = get_data(options.test_path)
185 test_imgs = [s for s in all_imgs if s['imageset'] == 'test']
186
187
188 T = {}
189 P = {}
190 for idx, img_data in enumerate(test_imgs):
191     print('{}{}/{}'.format(idx, len(test_imgs)))
192     st = time.time()
193     filepath = img_data['filepath']
194
195     img = cv2.imread(filepath)
196
197     X, fx, fy = format_img(img, C)
198
199     if K.image_dim_ordering() == 'tf':
200         X = np.transpose(X, (0, 2, 3, 1))
201
202     # get the feature maps and output from the RPN
203     [Y1, Y2, F] = model_rpn.predict(X)
204
205     R = roi_helpers.rpn_to_roi(Y1, Y2, C, K.image_dim_ordering(),
206                               overlap_thresh=0.7)
207
208     # convert from (x1,y1,x2,y2) to (x,y,w,h)
209     R[:, 2] -= R[:, 0]
210     R[:, 3] -= R[:, 1]
211
212     # apply the spatial pyramid pooling to the proposed regions
213     bboxes = {}
214     probs = {}

```



```

215 for jk in range(R.shape[0] // C.num_rois + 1):
216     ROIs = np.expand_dims(R[C.num_rois * jk:C.num_rois * (jk +
217     1), :], axis=0)
218     if ROIs.shape[1] == 0:
219         break
220
221 if jk == R.shape[0] // C.num_rois:
222     # pad R
223     curr_shape = ROIs.shape
224     target_shape = (curr_shape[0], C.num_rois, curr_shape[2])
225     ROIs_padded = np.zeros(target_shape).astype(ROIs.dtype)
226     ROIs_padded[:, :curr_shape[1], :] = ROIs
227     ROIs_padded[0, curr_shape[1]:, :] = ROIs[0, 0, :]
228     ROIs = ROIs_padded
229
230 [P_cls, P_regr] = model_classifier_only.predict([F, ROIs])
231
232 for ii in range(P_cls.shape[1]):
233     if np.argmax(P_cls[0, ii, :]) == (P_cls.shape[2] - 1):
234         continue
235
236     cls_name = class_mapping[np.argmax(P_cls[0, ii, :])]
237
238     if cls_name not in bboxes:
239         bboxes[cls_name] = []
240         probs[cls_name] = []
241
242     (x, y, w, h) = ROIs[0, ii, :]
243
244     cls_num = np.argmax(P_cls[0, ii, :])
245     try:
246         (tx, ty, tw, th) = P_regr[0, ii, 4 * cls_num:4 * (
247             cls_num + 1)]
248         tx /= C.classifier_regr_std[0]

```

```

248         ty /= C.classifier_regr_std[1]
249         tw /= C.classifier_regr_std[2]
250         th /= C.classifier_regr_std[3]
251         x, y, w, h = roi_helpers.apply_regr(x, y, w, h, tx, ty,
tw, th)
252     except:
253         pass
254     bboxes[cls_name].append([16 * x, 16 * y, 16 * (x + w), 16
* (y + h)])
255     probs[cls_name].append(np.max(P_cls[0, ii, :]))
256
257 all_dets = []
258
259 for key in bboxes:
260     bbox = np.array(bboxes[key])
261
262     new_boxes, new_probs = roi_helpers.non_max_suppression_fast
(bbox, np.array(probs[key]), overlap_thresh=0.5)
263     for jk in range(new_boxes.shape[0]):
264         (x1, y1, x2, y2) = new_boxes[jk, :]
265         det = {'x1': x1, 'x2': x2, 'y1': y1, 'y2': y2, 'class':
key, 'prob': new_probs[jk]}
266         all_dets.append(det)
267
268
269 print('Elapsed time = {}'.format(time.time() - st))
270 t, p = get_map(all_dets, img_data['bboxes'], (fx, fy))
271 for key in t.keys():
272     if key not in T:
273         T[key] = []
274         P[key] = []
275     T[key].extend(t[key])
276     P[key].extend(p[key])
277 all_aps = []
278 for key in T.keys():

```

```
279     ap = average_precision_score(T[key], P[key])
280     print('{} AP: {}'.format(key, ap))
281     all_aps.append(ap)
282     print('mAP = {}'.format(np.mean(np.array(all_aps))))
283     #print(T)
284     #print(P)
285 measure_map.py
286 Displaying train_frcnn.py.
```

REFERENCES

- [1] Vasefi, Fartash Kemp, William MacKinnon, Nicholas Amini, Mohammad Valdebran, Manuel Huang, Kevin Zhang, Haomiao. (2018). "Automated facial acne assessment from smartphone images." 22. 10.1117/12.2292506.
- [2] T. Chantharaphaichi, B. Uyyanonvara, C. Sinthanayothin and A. Nishihara, "Automatic acne detection for medical treatment," 2015 6th International Conference of Information and Communication Technology for Embedded Systems (IC-ICTES), Hua-Hin, 2015, pp. 1-6, doi: 10.1109/ICTEmSys.2015.7110813.5
- [3] Thanapha Chantharaphaichi, Bunyarit Uyyanonvara et al., "Automatic Acne Detection With Featured Bayesian Classifier For Medical Treatment", International MultiConference BIM 2015
- [4] C. Chang and H. Liao, "Automatic Facial Spots and Acnes Detection System," Journal of Cosmetics, Dermatological Sciences and Applications, Vol. 3 No. 1A, 2013, pp. 28-35. doi: 10.4236/jcdsa.2013.31A006.
- [5] Lim ZV, Akram F, Ngo CP, et al., "Automated grading of acne vulgaris by deep learning with convolutional neural networks.", Skin Res Technol. 2020;26(2):187-192. doi:10.1111/srt.12794
- [6] Zhao, Tingting Zhang, Hang Spoelstra, Jacob. (2019). "A Computer Vision Application for Assessing Facial Acne Severity from Selfie Images." arXiv:1907.07901
- [7] Shen, X., Zhang, J., Yan, C. et al. "An Automatic Diagnosis Method of Facial Acne Vulgaris Based on Convolutional Neural Network". Sci Rep 8, 5839 (2018). <https://doi.org/10.1038/s41598-018-24204-6>

- [8] Ren S., He K., Girshick R. and Sun J. 2015 Faster R-CNN: Towards real-time object detection with region proposal networks *Advances in neural information processing systems* 91-99
- [9] Redmon J., Divvala S., Girshick R. and Farhadi A. 2016 You only look once: Unified, real-time object detection *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* 779-788
- [10] Liu, Wei et al. "SSD: Single Shot MultiBox Detector." *Lecture Notes in Computer Science* (2016): 21–37. Crossref. Web.