

Object Detection using TensorFlow

Menita Koonani, Renu Hadke

{hadke.r, koonani.m}@husky.neu.edu

INFO 7390, Spring 2018, Northeastern University

Abstract -- Object detection is a computer technology connected to computer vision and image processing that deals with detecting instances of objects of a certain class in digital images and videos. It has applications in many areas of computer vision such as face detection, people counting, vehicle detection, anomaly detection and for classifying images online. Our objective is to identify and classify objects spotted on images and real-time video using TensorFlow and to determine the accuracy of each identification. There are various models of ConvNets that can be taken into consideration for making comparisons. In this paper, we are considering three models namely, SSD with MobileNet, SSD Inception V2 model and Faster RCNN Inception model to compare the accuracy and size of the models and speed of the model.

I. Introduction

Object recognition has been an area of widespread research for a long time. Object recognition and detection is one of the most important focus for Computer Vision. Constant research is in progress in order to find methods to have a recognition or detection system as powerful as the humans. Given an image, we want a system capable of answering certain questions like: What is main subject present in the image? Is there chair or a table in it? Is there a person in it? These questions can be so trivial to human beings but, are extremely difficult for a computer system, although the same computer is way quicker and efficient than human beings in terms of sequential executions and calculations. This paper focusses on considering three models of ConvNets and comparing their detection accuracy and speed of detection.

II. Dataset

The paper implements three different models of TensorFlow Object Detection API and compares them in terms of speed and accuracy. The dataset used is a COCO Dataset that has been pre-trained using Convolution Neural Network to detect various objects in a common set of 90 images. It is downloaded from https://github.com/tensorflow/models/blob/master/research/object_detection/g3doc/detection_model_zoo.md and it has 200K labelled images categorized into 90 classes.

III. Convolution Neural Network

Convolution Neural Network is a class of deep, feed-forward artificial neural networks that has been applied to analysing visual imagery.¹ These are a special class of Multilayer perceptron which are well suited for pattern classification. It is specifically designed to recognize 2D shapes with a high level of invariance, skewing and scaling. They are made up of neurons that have learnable weights and biases. Each neuron receives some input, performs a dot product and optionally follows it with a non-linearity. The whole network still expresses a single differentiable score function: from the raw image pixels on one end to class scores at the other. A simple ConvNet is a sequence of layers, and every layer of a ConvNet transforms one volume of activations to another through a differentiable function. There are three main types of layers to build ConvNet architectures: Convolutional Layer, Pooling Layer, and Fully-Connected Layer

Layers in CNN:²

- INPUT layer will hold the raw pixel values of the image, of width 32, height 32, and with three color channels Red, Green, Blue.
- CONV layer will compute the output of neurons that are connected to local regions in the input, each computing a dot product between their weights and a small region they are connected to in the input volume.
- RELU layer will apply an elementwise activation function, such as the max (0, x) thresholding at zero.
- POOL layer will perform a down sampling operation along the spatial dimensions (width, height)
- FC (i.e. fully-connected) layer will compute the class scores.

IV. Various models using ConvNets:

SSD with MobileNet:

Out of the many detection models, we chose to work with the combination of Single Shot Detectors(SSDs) and MobileNets architecture as they are fast, efficient and do not require huge computational capability to fulfill the Object Detection task. The SSD approach is based on a feed-forward convolutional neural network which produces a fixed-size collection of bounding boxes and scores for the object class instances present in those boxes. The main difference between a “traditional” CNN’s and the MobileNet architecture is instead of a single 3x3 convolution layer followed by batch norm and ReLU, MobileNets split the convolution into a 3x3 depthwise conv and a 1x1 pointwise conv.

SSD Inception V2 Model:

Given an input image and a set of ground truth labels, SSD does the following:

- It passes the image through a series of convolutional layers, providing several sets of feature maps at different scales.
- For each location in each of these feature maps, a 3x3 convolutional filter is used to evaluate a small set of default bounding boxes.
- For each box, it simultaneously predicts the bounding box offset and the probabilities of each class.

- During training, it matches the ground truth box with these predicted boxes based on IoU (Intersection over Union). The best predicted box is labelled a “positive” along with all the other boxes having an IoU with the truth greater than 0.5.

Faster RCNN Inception Model:

The main insight of Faster R-CNN was to replace the slow selective search algorithm that was used in the R-CNN (Region-based Convolutional Neural Network), with a fast-neural net.³

The R-CNN basically works in three simple steps:

- Scans the input image for possible objects using a Selective Search algorithm thus generating ~2000 region proposals.
- Runs a convolutional neural net (CNN) on top of each of these region proposals
- Takes the output of each CNN and feed it into an SVM (Support Vector Machine) to classify the region and a linear regressor to tighten the bounding box of the object, if such an object exists.

Faster R-CNN is similar to the original R-CNN but is improved on its detection speed through two augmentations:

- It performs feature extraction over the image even before proposing regions, thus running only one CNN over the entire image instead of running 2000 CNN’s across 2000 overlapping regions.
- It replaces the SVM with a softmax layer, thus extending the network for predictions instead of creating a new model.

V. TensorFlow

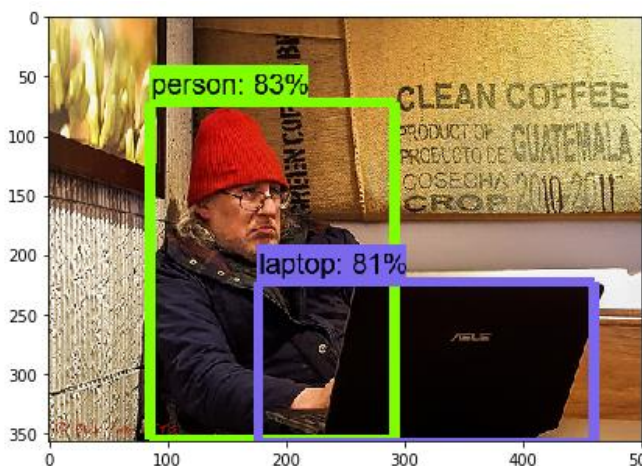
TensorFlow is an open source software library for high performance numerical computation. Its flexible architecture allows easy deployment of computation across a variety of platforms (CPUs, GPUs, TPUs), and from desktops to clusters of servers to mobile and edge devices. Originally developed by researchers and engineers from the Google Brain team within Google’s AI organization, it comes with strong support for machine learning and deep learning and the flexible numerical computation core is used across many other scientific domains.⁴

VI. OpenCV

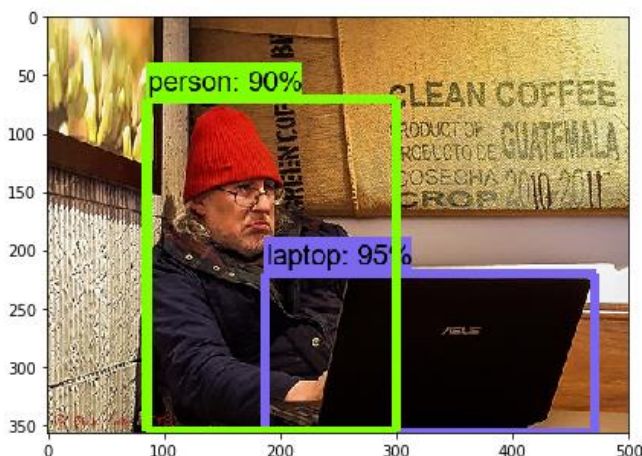
OpenCV (Open Source Computer Vision) is a library of programming functions mainly aimed at real-time computer vision. The C++ API provides a class 'videocapture' for capturing video from cameras or for reading video files and image sequences. It is basically used to access the Webcam of our computer to capture real-time videos. The frames from the videos can be read by creating an object of the VideoCapture class. This object which handles everything related to opening and closing of the webcam.

VII. Results

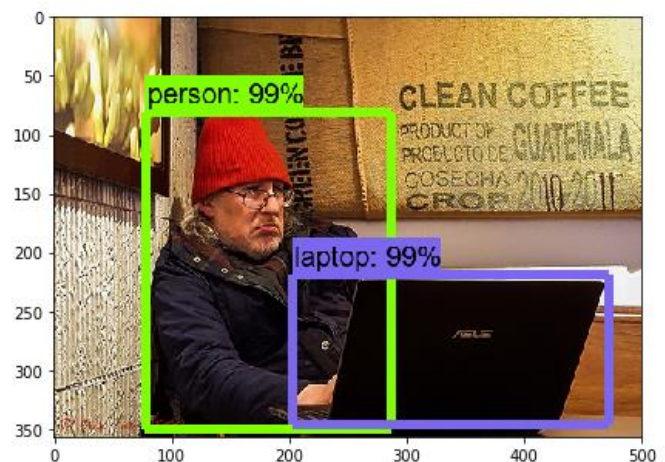
For SSD with MobileNet, the accuracy of object detection in image is 83% for person and 81% for laptop. This model worked fast though had the least accuracy.



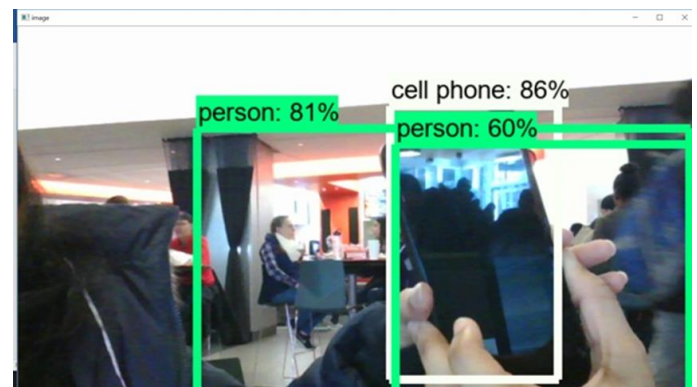
For SSD Inception V2 model, the accuracy of the objects detected in the images is 90% for the person and 95% for the laptop.



For Faster RCNN Inception Model, the accuracy of object detection in image is 99% for person and 99% for laptop



In addition, we were also successful in accessing the webcam of our system using OpenCV to detect real-time objects. The model used here is SSD with MobileNets as it produces much faster results as compared to the other two models.



VIII. Discussion

The principal difference between the two models is that Faster RCNN Inception V2 is optimized for accuracy, while the MobileNets are optimized to be small and efficient, at the cost of some accuracy. The SSD with MobileNets detects objects in only a single shot with just two components in its architecture namely, Feature Extraction and Detection Generator, while the Faster R-CNN consists of three components- Feature Extraction, Proposal Generation and Box Classifier.

In the Faster R-CNN Inception model, a region proposal network is used to generate regions of interest and then either fully-connected layers or position-sensitive convolutional layers to classify those regions. SSD does the two in a "single

shot,” simultaneously predicting the bounding box and the class as it processes the image.

Model name	Speed (ms)	COCO mAP[^1]
ssd_mobilenet_v1_coco	30	21
ssd_mobilenet_v2_coco	31	22
ssd_inception_v2_coco	42	24
faster_rcnn_inception_v2_coco	58	28

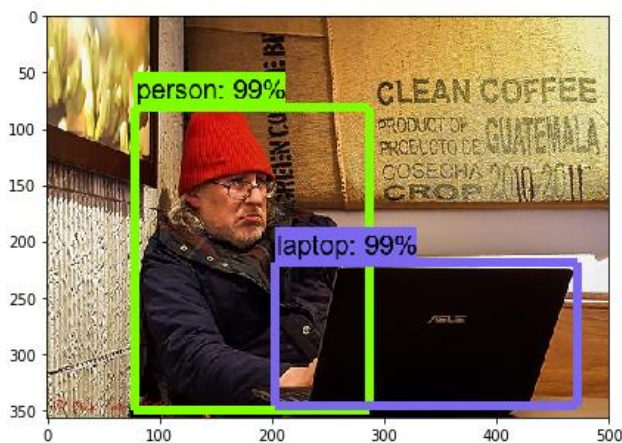
The higher the mAp (minimum average precision), the better the model. Based on the observations, SSD with MobileNets provided much better results in terms of speed but the Faster RCNN Inception Model provided a higher accuracy with some compromise on the speed. The complete code with Documentation can be referred in this link:

<https://github.com/renuHadke/ADS>

IX. Conclusion

SSD with MobileNet model accuracy:
 Person 83% and laptop 81%
 SSD Inception V2 model accuracy:
 Person 90% and laptop 95%
 Faster RCNN Inception Model accuracy:
 Person 99% and laptop 99%

As we can see above, Faster RCNN Inception Model gives the highest accuracy and SSD with MobileNet gives the lowest. But MobileNets are optimized to be small and efficient, at the cost of some accuracy.



X. Acknowledgement

We are thankful to our Prof. Nik Bear Brown, Assistant Professor, Northeastern University, Boston, MA for his valuable guidance, encouragement and co-operation during the implementation of this project.

XI. References

- [1]https://en.wikipedia.org/wiki/Convolutional_neural_network
- [2]<http://cs231n.github.io/convolutional-networks/>
- [3]<https://towardsdatascience.com/deep-learning-for-object-detection-a-comprehensive-review-73930816d8d9>
- [4] <https://www.tensorflow.org/>
- [5]https://github.com/tensorflow/models/tree/master/research/object_detection
- [6]<https://codelabs.developers.google.com/codelabs/tensorflow-for-poets/#0>
- [7]<https://hackernoon.com/creating-insanely-fast-image-classifiers-with-mobilenet-in-tensorflow-f030ce0a2991>
- [8]<https://pdfs.semanticscholar.org/3ac3/6574d1b1af029974cef6d6709feee2502194.pdf>
- [9]<https://towardsdatascience.com/deep-learning-for-object-detection-a-comprehensive-review-73930816d8d9>