

# Social distancing detection with computer vision techniques

Lee Hoang

BSc Computer Science  
City University, London  
Date

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Output Summary</b>	<b>3</b>
<b>3</b>	<b>Background Research</b>	<b>4</b>
3.1	What is computer vision? . . . . .	4
3.2	Deep learning . . . . .	4
3.2.1	Neural networks (NN) . . . . .	4
3.2.2	Convolutional neural networks (CNN) . . . . .	5
3.3	WILDTRACK dataset . . . . .	5
3.3.1	Research papers with WILDTRACK dataset . . . . .	6
3.4	Existing systems . . . . .	6
3.4.1	Example (Y. C. Hou, 2020) . . . . .	6
<b>4</b>	<b>Method</b>	<b>7</b>
4.1	Agile development . . . . .	7
4.2	Management tools . . . . .	7
4.2.1	Version control . . . . .	7
4.3	Deep learning architectures for object detection . . . . .	7
4.3.1	YOLO (You Only Look Once) . . . . .	7
4.3.2	SSD (Single shot detection) . . . . .	8
4.4	Pre-trained models (With COCO dataset) . . . . .	8
4.5	Programming language and libraries . . . . .	8
4.5.1	Python . . . . .	8
4.5.2	NumPy . . . . .	8
4.5.3	OpenCV . . . . .	9
4.5.4	matplotlib . . . . .	9
<b>5</b>	<b>Results</b>	<b>10</b>
5.1	Detecting objects with 'yolov3' (Analysis) . . . . .	10
5.1.1	Filtering classification methods . . . . .	10

# Chapter 1

## Introduction

## Chapter 2

### Output Summary

# Chapter 3

## Background Research

### 3.1 What is computer vision?

Computer vision (CV) is a scientific discipline that studies how computer can efficiently perceive, process, and understand information from visual data such as images and videos.

As humans, we can classify three-dimensional objects with ease, whether the pictures are the same object with different colours or angles, we are good at determining the object we are classifying. Computer vision has been developed to detect edges from a pixelated image, face detection, and has been used to develop 3D models from a snapshot yet the technology we have today could be compared to a young child's biological vision.

Computer vision is used in various real world application such as traffic surveillance or medical imaging (SZELISKI, 2020), where people are now able to utilize magnetic resonance imaging (MRI) to safely analysis the heart wall motion where the end result is a 3d model of the heart pumping (Metaxas, 1997).

In recent years, computer vision has been adapting deep learning algorithms to efficiently classify unseen objects within pixel images and videos.

### 3.2 Deep learning

Deep learning uses artificial intelligence (AI) to try and simulate the choices that a human brain will make. Problems that have regression or classification outputs can be solved by passing data/inputs through artificial neurons which were previously tweaked for the specific problem by training data. There are many different variants of deep learning algorithms such as Artificial Neural Networks (ANN) or Long Short-Term Memory (LSTM) Networks (Hochreiter, 1997) which build onto each other.

#### 3.2.1 Neural networks (NN)

Neural networks is a network formed of interconnected perceptrons which each carry weights and biases. The weights ( $w$ ) and biases ( $b$ ) are each represented by a float value which are used to multiply and add to the input respectively. The outcome is then put into an activation function (e.g. ReLU or sigmoid) which determines if the

neuron should be activated. These activations chain together to output a value of what the neural network thinks the solution is.

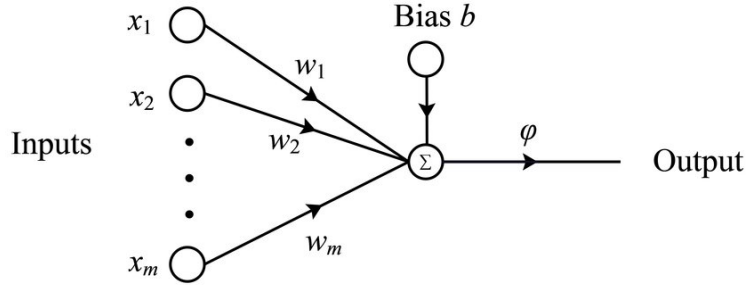


Figure 1: A single perceptron with inputs, weights, and biases

The values of weights and biases are updated when training data is pass through the network using back propagation, adjusting their values depending on the actual and network’s solution to the problem. A network can have multiple layers of several perceptrons in order to solve more complex problems, but has the down side of overfitting the accuracy of the output.

### 3.2.2 Convolutional neural networks (CNN)

CNNs builds upon the NN, specialising in working with grid-structured data such as images or videos. Unlike NN which takes in a 1 dimensional vector as its input, CNN uses tensors (high dimensional vector/matrix) for processing images. The convolutional layer of the network includes kernels/filters (a type of matrix) which efficiently detects features within an image by calculate the sum of the kernel multiplied by a subgrid of the image. These filters are trained to extract details within the images, and to make better predictions over the course of training. Different techniques can be applied such padding to reduce the lost of data when filtering.

The information about the deep learning architectures we have discussed were influenced by taking part in the ‘Introduction to AI’ (IN3062) and ‘Programming and Mathematics for Artificial Intelligence’ (IN3063) module at City University.

## 3.3 WILDTRACK dataset

The ‘WILDTRACK’ dataset provides footage of pedestrians from 7 different angle, each video having a length of 60 minutes and were filmed using three ‘GoPro Hero 4’ and four ‘GoPro Hero 3’ at high resolution (1920x1080). The location of the video took place near a public research university called ETH Zurich, Switzerland. The cameras used were calibrated to allow for precise calculations when wanting the distance between two objects.

While there are other datasets that offer similar features such as the ‘EPFL-RLC’ dataset, the videos themselves do not contain the same density of pedestrians relative to the ‘WILDTRACK’ dataset. Footage from other datasets showcase people who are more static, making it less challenging for the project.

### **3.3.1 Research papers with WILDTRACK dataset**

There are many published research papers that have used this dataset for their own project. For example, a research paper that was dedicated to detecting the same pedestrian using all camera footage and creating a shared top down view of the 'point of interest' which indicates the location of the pedestrians (López-Cifuentes A. 2018). What was interesting about this paper was that the author used another dataset that did not include calibrated cameras. When comparing the results at the end, the accuracy of the 'WILDTRACK' was marginally better as the other dataset had calibration errors, therefore the 'point of interest' were at different locations when looking at the shared top down view.

## **3.4 Existing systems**

Many different social distancing detections have been made ever since the outbreak of the corona virus. Most of the systems use deep learning architectures paired with the OpenCV library to help classify pedestrians within a video.

### **3.4.1 Example (Y. C. Hou, 2020)**

This system uses a combination of the YOLO (You Only Look Once) model with the COCO (Common Objects in Context) dataset to train their model. The goal of the project was to produce a top down view of the pedestrians, showing the distance between each person who were identified within the frame. The results of the system were very accurate, as they were able to make use of calibrated cameras. What was very interesting about the project is that some pedestrians were not classified due to hiding behind others. This showcases the limitation of the dataset used, as there were no overlapping footage of the field of view. Further improvements to the system were suggested such as mask and human body temperature detection. The system overall is similar to what this project will achieve with the difference of dataset.

# Chapter 4

## Method

### 4.1 Agile development

An agile approach with sprints was used during the development of the project. The sub objectives for this project required the previous sub objectives to be completed, so a sprint timetable seemed suitable. The approach gave time to reflect and adjust the work being done through out the weeks while allowing focus on the core functionalities of the software being produced.

The outcome of each sprint allowed for early prototypes of the system, easier analysis on code and product, and fixing any bugs within the code. The consequences of only partially completing a sprint's objective can be solved by adding more time for further development. Partially completion of the sprints are inevitable therefore spacing out the sprints instead of one after the other will help allocate more time.

### 4.2 Management tools

A diary is used to record sprints and work done throughout the week, recording the results produced and research done in order to produce code.

#### 4.2.1 Version control

GitHub is used to store previous versions of the code used for the project. This is in case of any mistakes made during development, backup versions of the project can be restored. GitHub can allow the user to backup their version no matter how small the changes are, providing flexibility throughout the project.

### 4.3 Deep learning architectures for object detection

Deep learning has been a foundation for modern computer vision, allowing object detection to be 'automatic' by training a CNN which tunes itself after each batch of data, then being further developed by implementing algorithms for object detection which only requires one pass through the network. This project will specifically be using yolov3 and SSD architectures.



### 4.3.1 YOLO (You Only Look Once)

YOLO is a object detection architecture that uses convolutional neural networks to divide the image/input into a grid. Each box in the grid is then associated with a high dimensional vector which record data such as: if there is an object within the grid, the predicted position of the bounding box, and the class id of the object. The vector can be expanded if there is more than one object within the box which is called anchor boxes. See appendix 1.

It is possible for the architecture to identify the same object twice within the same box creating redundancy. Since the predictions are based on probability, the architecture chooses the highest probability and uses there bounding box to identify the object. This feature is called 'Non-Max suppression'.

There are many versions of the 'yolo' such as 'yolov2' but during 2018, 'yolov3' was released (REDMON J. 2018). When comparing to previous versions, it offers an increase in speed and efficiency when computing while also providing a better backbone classifier (The core convolutional neural network).

There are also different models of 'yolov3' that perform better speed wise but at the cost of accuracy such as 'yolov3-tiny' which can compute images at 220 frames per second. This project specifically will use 'yolov3-320' which will detect objects at 45 frames per second while still having good accuracy. The reason for using this model is that realistically, cameras used in offices or public areas for surveillance will most likely be around 30 to 60 frames per second.

### 4.3.2 SSD (Single shot detection)

## 4.4 Pre-trained models (With COCO dataset)

The Microsoft 'COCO' (Common Objects in Context) dataset is a large-scale object detection dataset that contains over 330,000 images and more than 80 different object categories. This project used a pre-trained 'yolov3-320' model with the 'COCO' dataset to detect people. The reasoning for using this is that the dataset is very large, and will take several days to train the network. In comparison, downloading the pre-trained weights (Redmon, 2021) took a few minutes. Another reason for using 'COCO' is that it

## 4.5 Programming language and libraries

This section discusses the chosen programming language and libraries used in the project, giving examples of functionalities within the coding aspect.

### 4.5.1 Python

Python was chosen for flexibility when organising and presenting code through out the project. Python allows the user to call modules from different files that contained reusable class and functions, making the overall structure of the code less cluttered.

Python has compatibility with libraries needed for the project.

### 4.5.2 NumPy

NumPy specialises in matrix/array arithmetic, being able compute high dimensional matrix multiplication with ease. For this project, we use NumPy to concatenate matrices to generate a window for the top down view.

### 4.5.3 OpenCV

OpenCV is a cross-platform library which can be used to develop real-time computer vision applications. OpenCV can be used to utilise the pre-trained weights previously discussed, allowing to input an image and filter out the objects within the image.

Another functionality of the OpenCV library is image transformation, transforming an image to a top-down view.

### 4.5.4 matplotlib

Matplotlib is a tool that specialises in creating graphs.

# Chapter 5

## Results

This section discusses the results found throughout the project.

### 5.1 Detecting objects with 'yolov3' (Analysis)

OpenCV enabled the use of 'yolov3' with the functions of loading the weights of the neural network and forwarding an image to output the objects detected. The first test that was done was to showcase how well the model detected objects. The model was tested with different image provided by the dataset.

The first image that was forwarded was from the third camera angle (figure 2). The main focus of using this specific image was to look at how well the model detected people who were overlapping each other. The outcome of model was over 10,000 objects detected within the image, taking 0.54 seconds to compute. The amount of objects detected might seem unnecessary, but the model was trained to detect 80 different objects, therefore will provide a large number of classifications per image.

Each classification comes with a probability that the object detected is the correct classification. Therefore objects that are 'person' could also have a probability of being identified as a 'car'.

#### 5.1.1 Filtering classification methods

The probabilities can be filtered by tuning the number of objects that can be classified and the probability threshold.