

Vehicle Detection



Vehicle Detection

1. Overview
2. Project Introduction
3. Project Pipeline
 - 3.1 Read images
 - 3.2 YOLO
 - 3.2.1 Model details
 - 3.2.2 Filtering with threshold on class scores
 - 3.2.3 Non-max suppression
 - 3.2.4 Wrapping up the filtering
 - 3.3 Test on images
 - 3.3.1 Defining classes, anchors and image shape
 - 3.3.2 Loading pertained model
 - 3.3.3 Convert output of the model to usable bounding box tensors
 - 3.3.4 Filtering boxes
 - 3.3.5 Run on an image
 4. Discussion

1. Overview

The project is to detect vehicle on the road using **YOLO** model. Many ideas of this project are described in the two **YOLO** papers: Redmon et al., 2016 (<https://arxiv.org/abs/1506.02640>) and Redmon and Farhadi, 2016 (<https://arxiv.org/abs/1612.08242>), and my project of *Coursera/Deep Learning/Convolutional Neural Networks/Week 3/Car Detection for Autonomous Driving* with the certificate number `wY7JJSS7Q9B9`

2. Project Introduction

Comparing with the class text, the **YOLO** model is relatively easier but needs higher calculate ability. This model doesn't need a lot of image preprocess, and can directly feed the model with *RGB* images, and the output is classified object and the boxes in the image.

3. Project Pipeline

The pipeline is below:

1. read the test images

2. load the anchor boxes and the classification
3. define YOLO model and load pre-trained *h5* file
4. predict the image
5. define the pipeline and generate the video

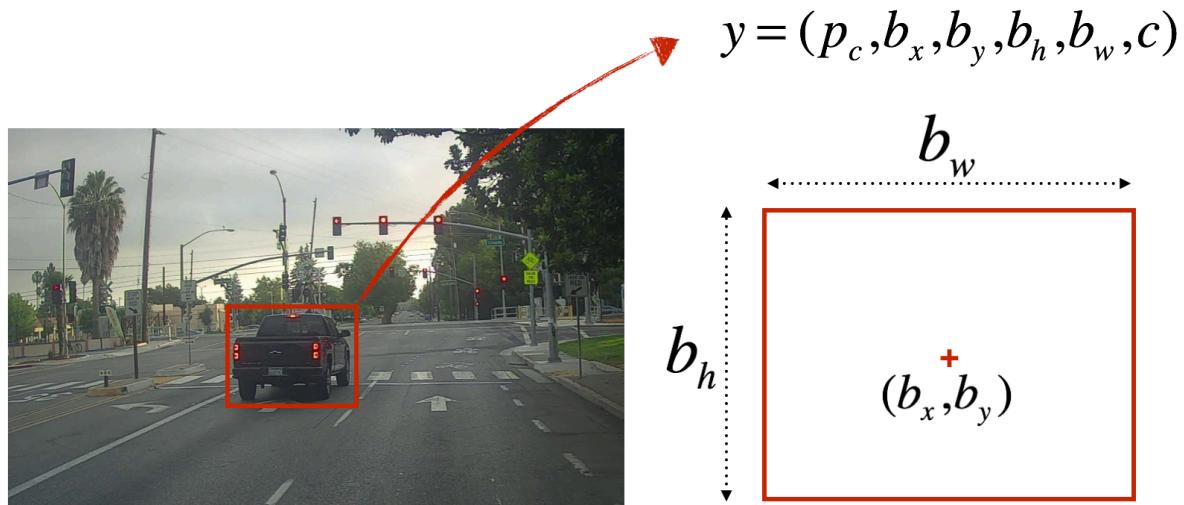
3.1 Read images

The read image function shows the original test images.



3.2 YOLO

YOLO ("you only look once") is a popular algorithm because it achieves high accuracy while also being able to run in real-time. This algorithm "only looks once" at the image in the sense that it requires only one forward propagation pass through the network to make predictions. After non-max suppression, it then outputs recognized objects together with the bounding boxes.



$p_c = 1$: confidence of an object being present in the bounding box
 $c = 3$: class of the object being detected (here 3 for “car”)

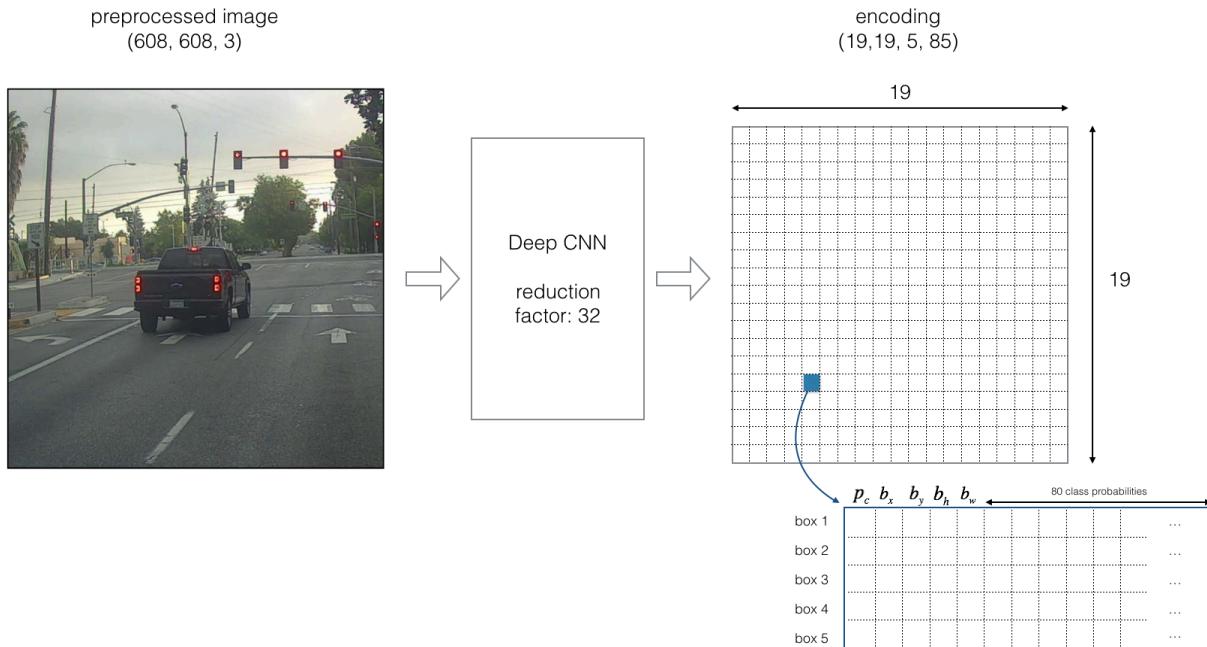
3.2.1 Model details

First things to know:

- The **input** is a batch of images of shape (m, 608, 608, 3)
- The **output** is a list of bounding boxes along with the recognized classes. Each bounding box is represented by 6 numbers ($p_c, b_x, b_y, b_h, b_w, c$) as explained above. If you expand c into an 80-dimensional vector, each bounding box is then represented by 85 numbers.

I am using 5 anchor boxes. So the YOLO architecture is as the following: IMAGE (m, 608, 608, 3) -> DEEP CNN -> ENCODING (m, 19, 19, 5, 85).

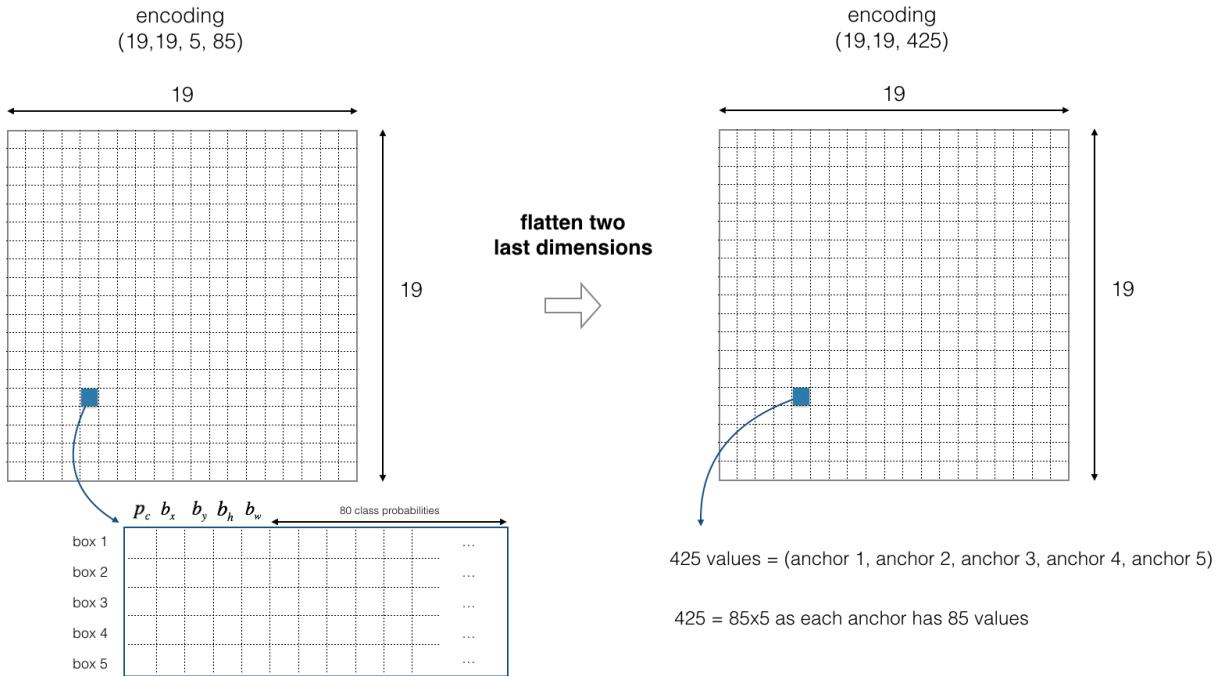
Lets look in greater detail at what this encoding represents.



If the center/midpoint of an object falls into a grid cell, that grid cell is responsible for detecting that object.

Since we are using 5 anchor boxes, each of the 19x19 cells thus encodes information about 5 boxes. Anchor boxes are defined only by their width and height.

For simplicity, we will flatten the last two last dimensions of the shape (19, 19, 5, 85) encoding. So the output of the Deep CNN is (19, 19, 425).



Now, for each box (of each cell) we will compute the following elementwise product and extract a probability that the box contains a certain class.

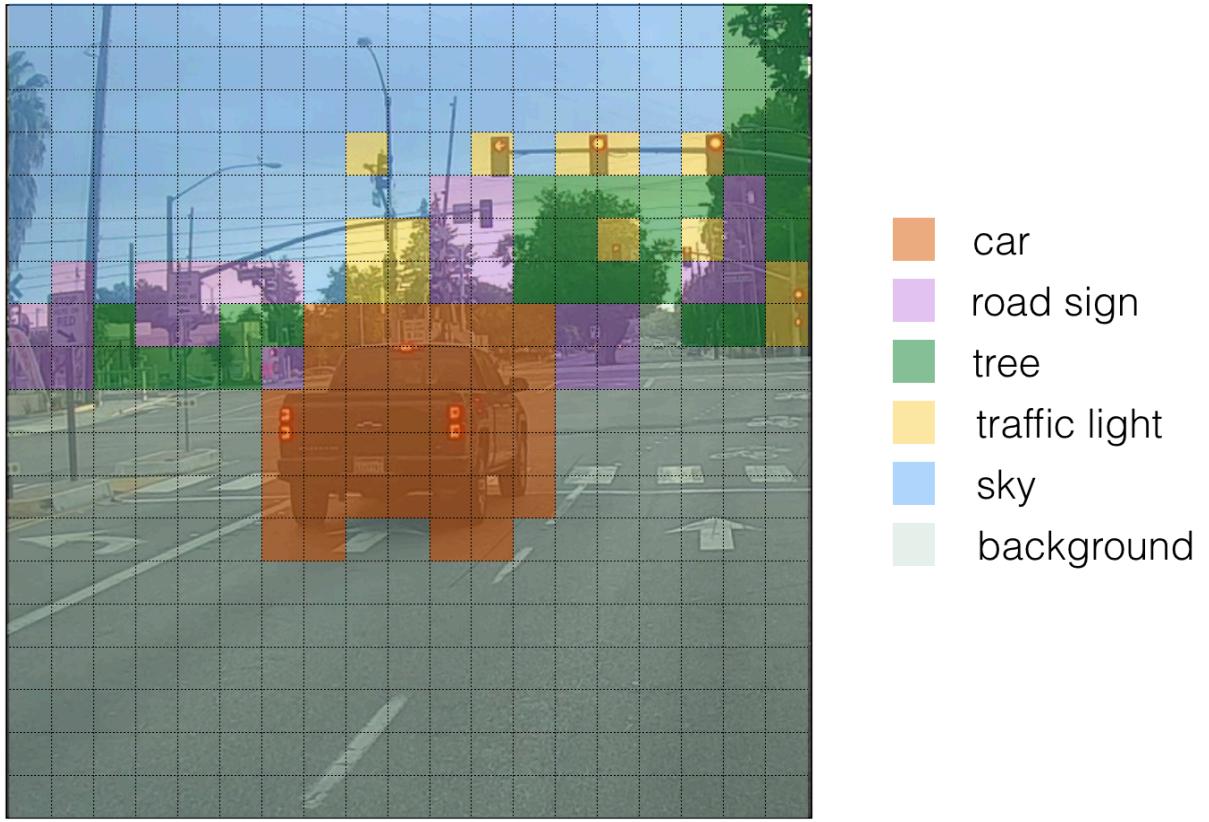
$$\begin{aligned}
 & \text{box 1} \quad [p_c \ | \ b_x \ | \ b_y \ | \ b_h \ | \ b_w \ | \ c_1 \ | \ c_2 \ | \ c_3 \ | \ c_4 \ | \ c_5 \ | \ \dots \ | \ c_{76} \ | \ c_{77} \ | \ c_{78} \ | \ c_{79} \ | \ c_{80}] \\
 & \text{scores} = p_c * \begin{pmatrix} c_1 \\ c_2 \\ c_3 \\ \vdots \\ c_{78} \\ c_{79} \\ c_{80} \end{pmatrix} = \begin{pmatrix} p_c c_1 \\ p_c c_2 \\ p_c c_3 \\ \vdots \\ p_c c_{78} \\ p_c c_{79} \\ p_c c_{80} \end{pmatrix} = \begin{pmatrix} 0.12 \\ 0.13 \\ 0.44 \\ \vdots \\ 0.07 \\ 0.01 \\ 0.09 \end{pmatrix} \xrightarrow{\text{find the max}} \text{score: } 0.44 \\
 & \text{box: } (b_x, b_y, b_h, b_w) \\
 & \text{class: } c = 3 (\text{"car"})
 \end{aligned}$$

the box (b_x, b_y, b_h, b_w) has detected $c = 3$ ("car") with probability score: 0.44

Here's one way to visualize what YOLO is predicting on an image:

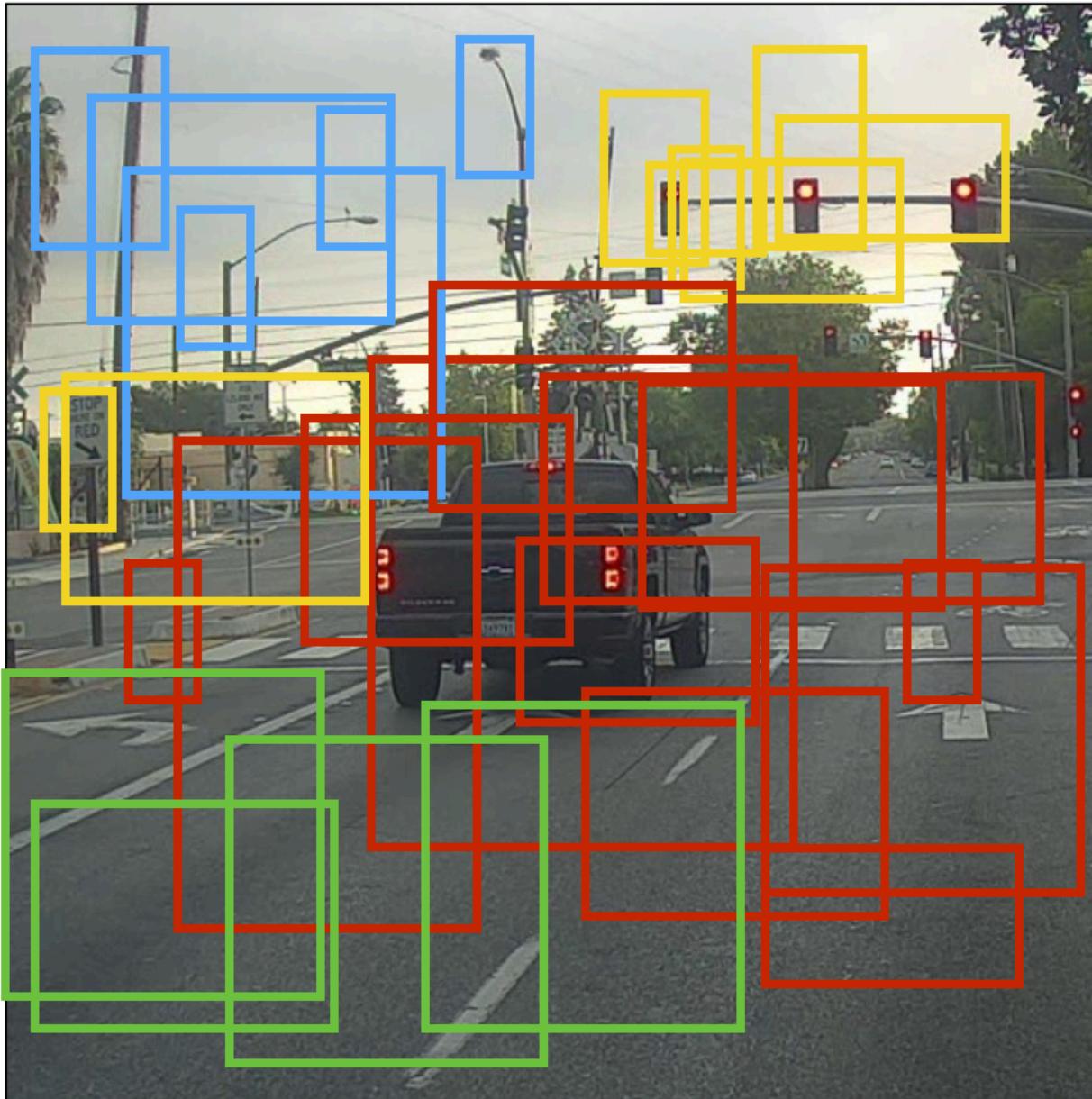
- For each of the 19×19 grid cells, find the maximum of the probability scores (taking a max across both the 5 anchor boxes and across different classes).
- Color that grid cell according to what object that grid cell considers the most likely.

Doing this results in this picture:



Note that this visualization isn't a core part of the YOLO algorithm itself for making predictions; it's just a nice way of visualizing an intermediate result of the algorithm.

Another way to visualize YOLO's output is to plot the bounding boxes that it outputs. Doing that results in a visualization like this:



In the figure above, we plotted only boxes that the model had assigned a high probability to, but this is still too many boxes. You'd like to filter the algorithm's output down to a much smaller number of detected objects. To do so, you'll use non-max suppression. Specifically, you'll carry out these steps:

- Get rid of boxes with a low score (meaning, the box is not very confident about detecting a class)
- Select only one box when several boxes overlap with each other and detect the same object.

3.2.2 Filtering with threshold on class scores

The next step is to get rid of any box for which the class "score" is less than the threshold.

The model gives you a total of $19 \times 19 \times 5 \times 85$ numbers, with each box described by 85 numbers. It'll be convenient to rearrange the $(19, 19, 5, 85)$ (or $(19, 19, 425)$) dimensional tensor into the following variables:

- `box_confidence` : tensor of shape $(19 \times 19, 5, 1)$ containing p_c (confidence probability)

that there's some objects) for each of the 5 boxes predicted in each of the 19×19 cells.

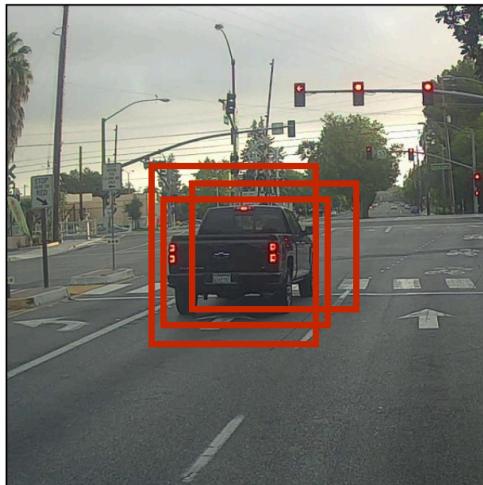
- `boxes` : tensor of shape $(19 \times 19, 5, 4)$ containing (b_x, b_y, b_h, b_w) for each of the 5 boxes per cell.
- `box_class_probs` : tensor of shape $(19 \times 19, 5, 80)$ containing the detection probabilities $(c_1, c_2, \dots, c_{80})$ for each of 80 classes for each of the 5 boxes per cell.

```
1 def yolo_filter_boxes(box_confidence, boxes, box_class_probs,
2     threshold=.6):
3     """Filters YOLO boxes by thresholding on object and class
4     confidence.
5
6     Arguments:
7         box_confidence -- tensor of shape (19, 19, 5, 1)
8         boxes -- tensor of shape (19, 19, 5, 4)
9         box_class_probs -- tensor of shape (19, 19, 5, 80)
10        threshold -- real value, if [ highest class probability score <
11                      threshold], then get
12                          rid of the corresponding box
13
14    Returns:
15        scores -- tensor of shape (None,), containing the class probability
16        score for
17            selected boxes
18        boxes -- tensor of shape (None, 4), containing (b_x, b_y, b_h, b_w)
19        coordinates of
20            selected boxes
21        classes -- tensor of shape (None,), containing the index of the
22        class detected by
23            the selected boxes
24
25        Note: "None" is here because you don't know the exact number of
26        selected boxes, as
27            it depends on the threshold.
28            For example, the actual output size of scores would be (10,) if
29            there are 10 boxes.
30
31    """
32
33    ### Code Implementation ###
34
35
36    return scores, boxes, classes
```

3.2.3 Non-max suppression

Even after filtering by thresholding over the classes scores, you still end up a lot of overlapping boxes. A second filter for selecting the right boxes is called non-maximum suppression (NMS).

Before non-max suppression



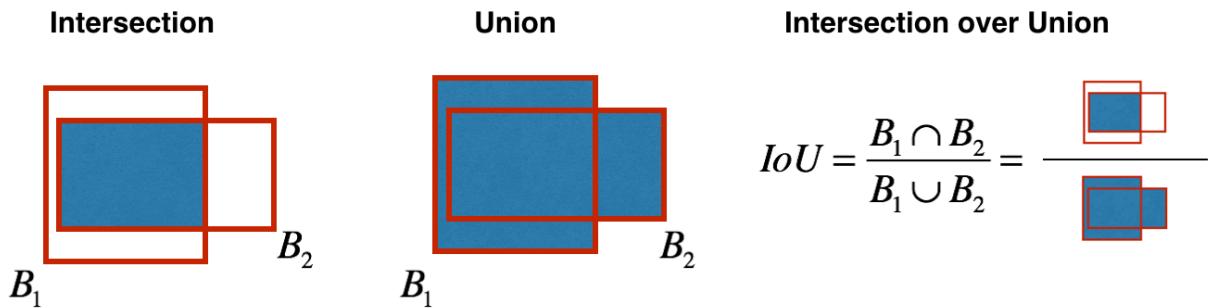
After non-max suppression



Non-Max Suppression



Non-max suppression uses the very important function called "**Intersection over Union**", or IoU.



```
1 def iou(box1, box2):
2     """Implement the intersection over union (IoU) between box1 and box2
3
4     Arguments:
5     box1 -- first box, list object with coordinates (x1, y1, x2, y2)
6     box2 -- second box, list object with coordinates (x1, y1, x2, y2)
7     """
8
9     ### Code Implementation ####
10
11    return iou
```

Now to implement non-max suppression function, the steps are:

1. Select the box that has the highest score.
2. Compute its overlap with all other boxes, and remove boxes that overlap it more than `iou_threshold`.
3. Go back to step 1 and iterate until there's no more boxes with a lower score than the current selected box.

This will remove all boxes that have a large overlap with the selected boxes. Only the "best" boxes remain.

```

1 def yolo_non_max_suppression(scores, boxes, classes, max_boxes=10,
2 iou_threshold=0.5)
3     """
4         Applies Non-max suppression (NMS) to set of boxes
5
6     Arguments:
7         scores -- tensor of shape (None,), output of yolo_filter_boxes()
8         boxes -- tensor of shape (None, 4), output of yolo_filter_boxes()
9         that have been
10            scaled to the image size (see later)
11         classes -- tensor of shape (None,), output of yolo_filter_boxes()
12         max_boxes -- integer, maximum number of predicted boxes you'd like
13         iou_threshold -- real value, "intersection over union" threshold
14         used for NMS
15             filtering
16
17     Returns:
18         scores -- tensor of shape (, None), predicted score for each box
19         boxes -- tensor of shape (4, None), predicted box coordinates
20         classes -- tensor of shape (, None), predicted class for each box
21
22     Note: The "None" dimension of the output tensors has obviously to be
23 less than
24     max_boxes. Note also that this function will transpose the shapes of
25 scores, boxes,
26 classes. This is made for convenience.
27 """
28
29
30     ### Code Implementation ###
31
32
33     return scores, boxes, classes

```

3.2.4 Wrapping up the filtering

This function is to take the output of the $(19 \times 19 \times 5 \times 85)$ dimensional deep CNN and filter through all the boxes using the function just implemented.

The function `yolo_eval()` takes the output of the **YOLO** encoding and filters the boxes using score threshold and NMS.

There are a few ways of representing boxes, such as via their corners or via their midpoint and height/width. YOLO converts between a few such formats at different times, using the following functions (defined in the `yolo_utils.py`):

```

1 boxes = yolo_boxes_to_corners(box_xy, box_wh)

```

which converts the yolo box coordinates (x,y,w,h) to box corners' coordinates (x1, y1, x2, y2) to fit the input of `yolo_filter_boxes`

```
1 boxes = scale_boxes(boxes, image_shape)
```

YOLO's network was trained to run on 608x608 images. If you are testing this data on a different size image--for example, the car detection dataset had 720x1280 images--this step rescales the boxes so that they can be plotted on top of the original 720x1280 image.

```
1 def yolo_eval(yolo_outputs, image_shape = (720., 1280.), max_boxes=10,
2               score_threshold=.5, iou_threshold=.5):
3     """
4         Converts the output of YOLO encoding (a lot of boxes) to your
5         predicted boxes along
6         with their scores, box coordinates and classes.
7
8     Arguments:
9         yolo_outputs -- output of the encoding model (for image_shape of
10            (608, 608, 3)),
11            contains 4 tensors:
12                box_confidence: tensor of shape (None, 19, 19, 5, 1)
13                box_xy: tensor of shape (None, 19, 19, 5, 2)
14                box_wh: tensor of shape (None, 19, 19, 5, 2)
15                box_class_probs: tensor of shape (None, 19, 19, 5,
16
16            80)
17            image_shape -- tensor of shape (2,) containing the input shape, in
18            this notebook we
19            use (608., 608.) (has to be float32 dtype)
20            max_boxes -- integer, maximum number of predicted boxes you'd like
21            score_threshold -- real value, if [ highest class probability score
22            < threshold],
23                            then get rid of the corresponding box
24            iou_threshold -- real value, "intersection over union" threshold
25            used for NMS
26
27            filtering
28
29        Returns:
30            scores -- tensor of shape (None, ), predicted score for each box
31            boxes -- tensor of shape (None, 4), predicted box coordinates
32            classes -- tensor of shape (None,), predicted class for each box
33        """
34
35    ### Code Implementation ####
36
37
38    return scores, boxes, classes
```

Conclusion:

- Input image (608, 608, 3)
- The input image goes through a CNN, resulting in a (19,19,5,85) dimensional output.

- After flattening the last two dimensions, the output is a volume of shape (19, 19, 425):
 - Each cell in a 19x19 grid over the input image gives 425 numbers.
 - $425 = 5 \times 85$ because each cell contains predictions for 5 boxes, corresponding to 5 anchor boxes, as seen in lecture.
 - $85 = 5 + 80$ where 5 is because $(p_c, b_x, b_y, b_h, b_w)$ has 5 numbers, and 80 is the number of classes we'd like to detect
- Select only few boxes based on:
 - Score-thresholding: throw away boxes that have detected a class with a score less than the threshold
 - Non-max suppression: Compute the Intersection over Union and avoid selecting overlapping boxes

3.3 Test on images

```
1 sess = K.get_session()
```

3.3.1 Defining classes, anchors and image shape

```
1 # TXT includes 80 classes information
2 class_names = read_classes("model_data/coco_classes.txt")
3 # TXT includes 5 boxes
4 anchors = read_anchors("model_data/yolo_anchors.txt")
5 # Define image shape
6 image_shape = (720., 1280.)
```

The `coco_classes.txt` includes 80 classes that we want YOLO to recognize, the class label either as an integer from 1 to 80, or as an 80-dimensional vector (with 80 numbers) one component of which is 1 and the rest of which are 0.

3.3.2 Loading pertained model

```
1 yolo_model = load_model("model_data/yolo.h5")
```

This loads the weights of a trained YOLO model. Here's a summary of the layers your model contains.

```
1 yolo_model.summary()
```

1 Layer (type)	Output Shape	Param #
Connected to		
=====	=====	=====
=====	=====	=====

```
3 input_1 (InputLayer)           (None, 608, 608, 3) 0
4
5 conv2d_1 (Conv2D)             (None, 608, 608, 32) 864
6 input_1[0][0]
7
8 batch_normalization_1 (BatchNor (None, 608, 608, 32) 128
9 conv2d_1[0][0]
10
11 leaky_re_lu_1 (LeakyReLU)     (None, 608, 608, 32) 0
12 batch_normalization_1[0][0]
13
14 max_pooling2d_1 (MaxPooling2D) (None, 304, 304, 32) 0
15 leaky_re_lu_1[0][0]
16
17 conv2d_2 (Conv2D)             (None, 304, 304, 64) 18432
18 max_pooling2d_1[0][0]
19
20 batch_normalization_2 (BatchNor (None, 304, 304, 64) 256
21 conv2d_2[0][0]
22
23 leaky_re_lu_2 (LeakyReLU)     (None, 304, 304, 64) 0
24 batch_normalization_2[0][0]
25
26 max_pooling2d_2 (MaxPooling2D) (None, 152, 152, 64) 0
27 leaky_re_lu_2[0][0]
28
29 conv2d_3 (Conv2D)             (None, 152, 152, 128) 73728
30 max_pooling2d_2[0][0]
31
32 batch_normalization_3 (BatchNor (None, 152, 152, 128) 512
33 conv2d_3[0][0]
34
35 leaky_re_lu_3 (LeakyReLU)     (None, 152, 152, 128) 0
36 batch_normalization_3[0][0]
```

```
27 conv2d_4 (Conv2D)           (None, 152, 152, 64) 8192
28 leaky_re_lu_3[0][0]
29
30
31 leaky_re_lu_4 (LeakyReLU)    (None, 152, 152, 64) 0
32 batch_normalization_4 (BatchNor (None, 152, 152, 64) 256
33 conv2d_4[0][0]
34
35
36
37 leaky_re_lu_5 (LeakyReLU)    (None, 152, 152, 128 0
38 batch_normalization_5 (BatchNor (None, 152, 152, 128 512
39 conv2d_5[0][0]
40
41
42
43 leaky_re_lu_5[0][0]
44
45
46
47 leaky_re_lu_5 (LeakyReLU)    (None, 76, 76, 128) 0
48 max_pooling2d_3 (MaxPooling2D) (None, 76, 76, 128) 0
49 max_pooling2d_3[0][0]
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
```

```
51 leaky_re_lu_7 (LeakyReLU)      (None, 76, 76, 128)  0
batch_normalization_7[0][0]
52
53 conv2d_8 (Conv2D)            (None, 76, 76, 256)  294912
leaky_re_lu_7[0][0]
54
55 batch_normalization_8 (BatchNor (None, 76, 76, 256)  1024
conv2d_8[0][0]
56
57 leaky_re_lu_8 (LeakyReLU)      (None, 76, 76, 256)  0
batch_normalization_8[0][0]
58
59 max_pooling2d_4 (MaxPooling2D) (None, 38, 38, 256)  0
leaky_re_lu_8[0][0]
60
61 conv2d_9 (Conv2D)            (None, 38, 38, 512)  1179648
max_pooling2d_4[0][0]
62
63 batch_normalization_9 (BatchNor (None, 38, 38, 512)  2048
conv2d_9[0][0]
64
65 leaky_re_lu_9 (LeakyReLU)      (None, 38, 38, 512)  0
batch_normalization_9[0][0]
66
67 conv2d_10 (Conv2D)            (None, 38, 38, 256)  131072
leaky_re_lu_9[0][0]
68
69 batch_normalization_10 (BatchNo (None, 38, 38, 256)  1024
conv2d_10[0][0]
70
71 leaky_re_lu_10 (LeakyReLU)     (None, 38, 38, 256)  0
batch_normalization_10[0][0]
72
73 conv2d_11 (Conv2D)            (None, 38, 38, 512)  1179648
leaky_re_lu_10[0][0]
74
```

```
75 batch_normalization_11 (BatchNo (None, 38, 38, 512) 2048
76 conv2d_11[0][0]
77 _____
78 leaky_re_lu_11 (LeakyReLU)      (None, 38, 38, 512) 0
79 batch_normalization_11[0][0]
80 _____
81 conv2d_12 (Conv2D)            (None, 38, 38, 256) 131072
82 leaky_re_lu_11[0][0]
83 _____
84 batch_normalization_12 (BatchNo (None, 38, 38, 256) 1024
85 conv2d_12[0][0]
86 _____
87 leaky_re_lu_12 (LeakyReLU)      (None, 38, 38, 256) 0
88 batch_normalization_12[0][0]
89 _____
90 conv2d_13 (Conv2D)            (None, 38, 38, 512) 1179648
91 leaky_re_lu_12[0][0]
92 _____
93 batch_normalization_13 (BatchNo (None, 38, 38, 512) 2048
94 conv2d_13[0][0]
95 _____
96 leaky_re_lu_13 (LeakyReLU)      (None, 38, 38, 512) 0
97 batch_normalization_13[0][0]
98 _____
```

```
99 conv2d_15 (Conv2D)           (None, 19, 19, 512) 524288
100 leaky_re_lu_14[0][0]
101
102 batch_normalization_15 (BatchNo (None, 19, 19, 512) 2048
103 conv2d_15[0][0]
104
105 leaky_re_lu_15 (LeakyReLU)   (None, 19, 19, 512) 0
106 batch_normalization_15[0][0]
107
108 conv2d_16 (Conv2D)           (None, 19, 19, 1024) 4718592
109 leaky_re_lu_15[0][0]
110
111 batch_normalization_16 (BatchNo (None, 19, 19, 1024) 4096
112 conv2d_16[0][0]
113
114 leaky_re_lu_16 (LeakyReLU)   (None, 19, 19, 1024) 0
115 batch_normalization_16[0][0]
116
117 conv2d_17 (Conv2D)           (None, 19, 19, 512) 524288
118 leaky_re_lu_16[0][0]
119
120 batch_normalization_17 (BatchNo (None, 19, 19, 512) 2048
121 conv2d_17[0][0]
122
123 leaky_re_lu_17 (LeakyReLU)   (None, 19, 19, 512) 0
124 batch_normalization_17[0][0]
125
126 conv2d_18 (Conv2D)           (None, 19, 19, 1024) 4718592
127 leaky_re_lu_17[0][0]
128
129 batch_normalization_18 (BatchNo (None, 19, 19, 1024) 4096
130 conv2d_18[0][0]
131
132 leaky_re_lu_18 (LeakyReLU)   (None, 19, 19, 1024) 0
133 batch_normalization_18[0][0]
134
```

```
123 conv2d_19 (Conv2D)           (None, 19, 19, 1024) 9437184
    leaky_re_lu_18[0][0]
124
125 batch_normalization_19 (BatchNo (None, 19, 19, 1024) 4096
    conv2d_19[0][0]
126
127 conv2d_21 (Conv2D)           (None, 38, 38, 64)   32768
    leaky_re_lu_13[0][0]
128
129 leaky_re_lu_19 (LeakyReLU)   (None, 19, 19, 1024) 0
    batch_normalization_19[0][0]
130
131 batch_normalization_21 (BatchNo (None, 38, 38, 64) 256
    conv2d_21[0][0]
132
133 conv2d_20 (Conv2D)           (None, 19, 19, 1024) 9437184
    leaky_re_lu_19[0][0]
134
135 leaky_re_lu_21 (LeakyReLU)   (None, 38, 38, 64)   0
    batch_normalization_21[0][0]
136
137 batch_normalization_20 (BatchNo (None, 19, 19, 1024) 4096
    conv2d_20[0][0]
138
139 space_to_depth_x2 (Lambda)   (None, 19, 19, 256)   0
    leaky_re_lu_21[0][0]
140
141 leaky_re_lu_20 (LeakyReLU)   (None, 19, 19, 1024) 0
    batch_normalization_20[0][0]
142
143 concatenate_1 (Concatenate)  (None, 19, 19, 1280) 0
    space_to_depth_x2[0][0]
144
145 leaky_re_lu_20[0][0]
146 conv2d_22 (Conv2D)           (None, 19, 19, 1024) 11796480
    concatenate_1[0][0]
```

147

148 batch_normalization_22 (BatchNorm (None, 19, 19, 1024) 4096
conv2d_22[0][0]

149

150 leaky_re_lu_22 (LeakyReLU) (None, 19, 19, 1024) 0
batch_normalization_22[0][0]

151

152 conv2d_23 (Conv2D) (None, 19, 19, 425) 435625
leaky_re_lu_22[0][0]

153 =====

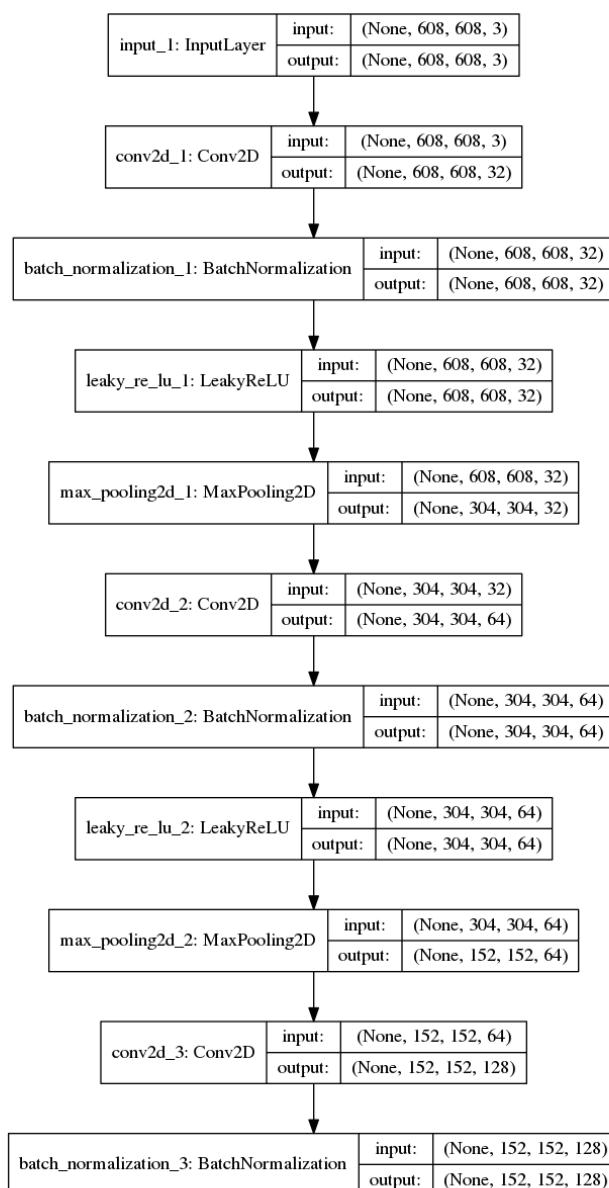
=====

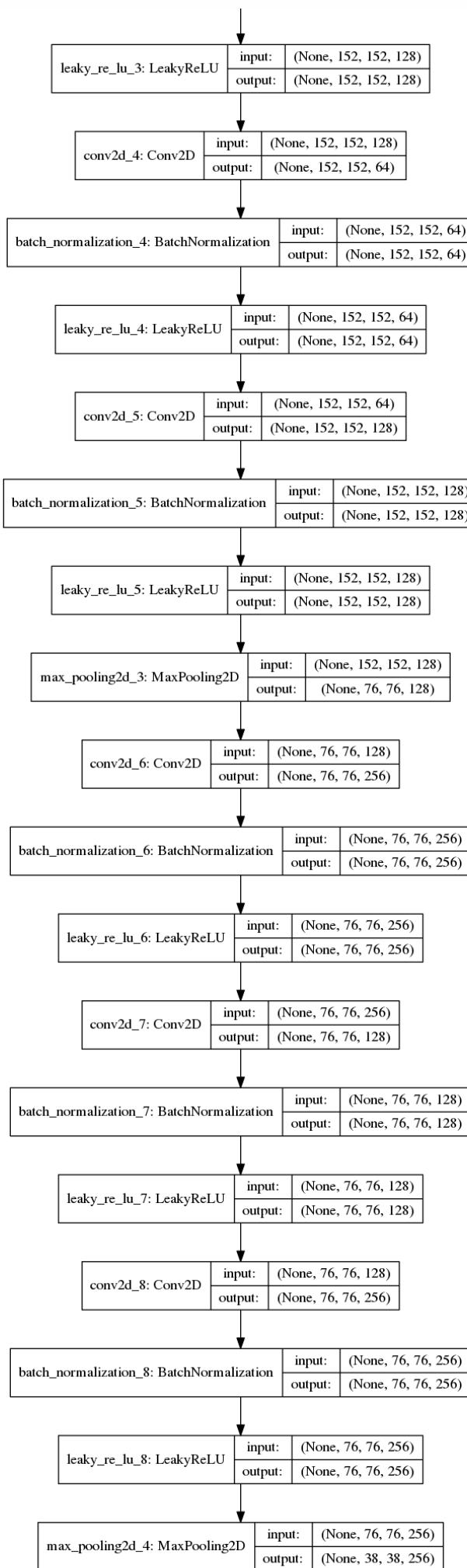
154 Total params: 50,983,561

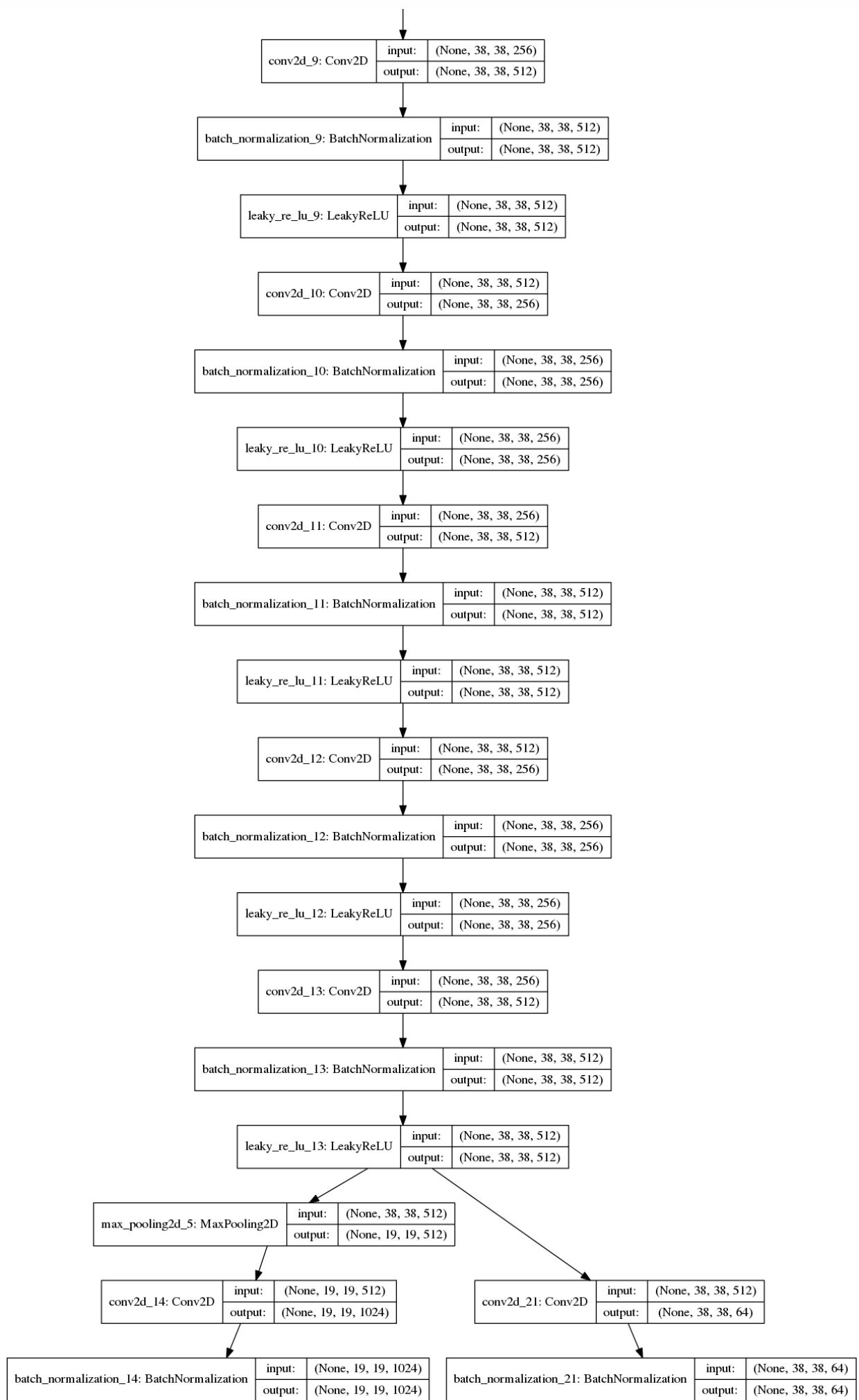
155 Trainable params: 50,962,889

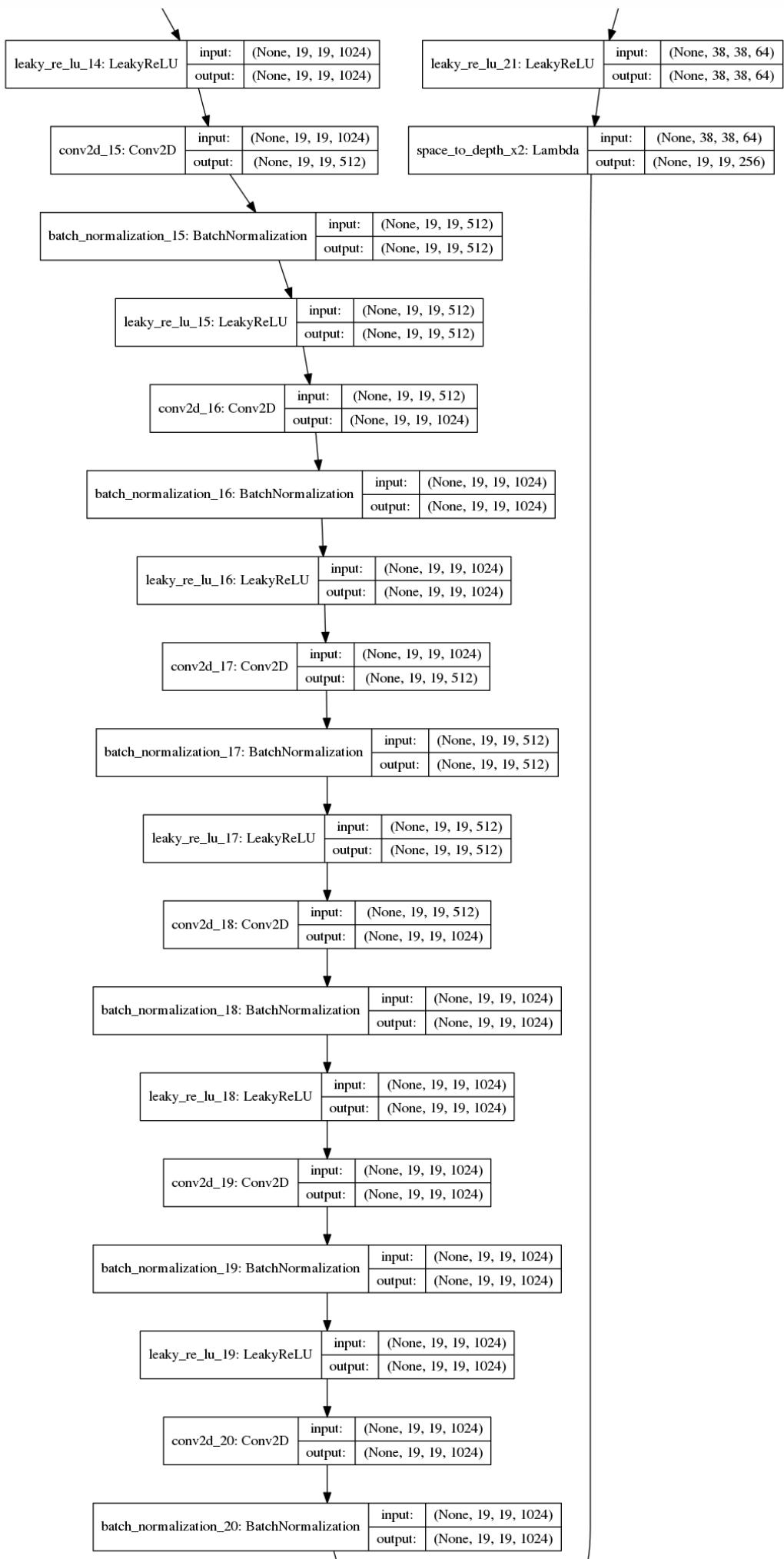
156 Non-trainable params: 20,672

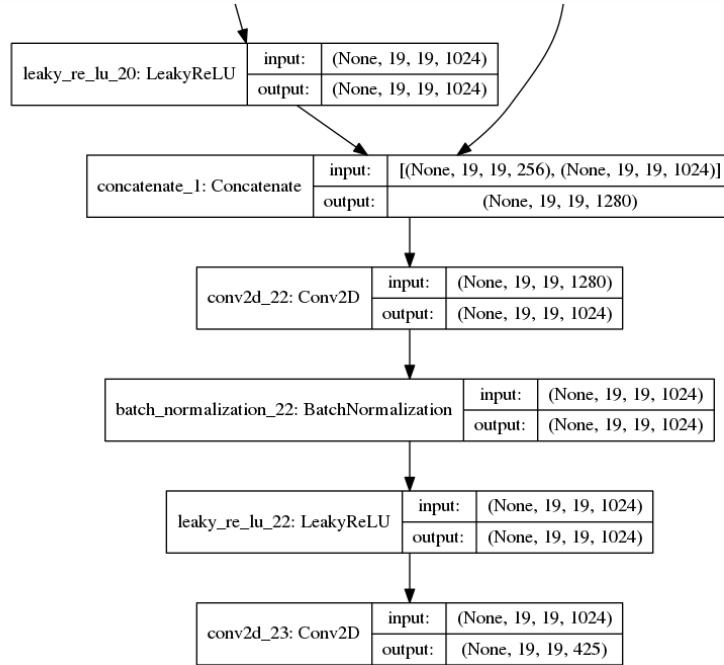
157











3.3.3 Convert output of the model to usable bounding box tensors

```
1 yolo_outputs = yolo_head(yolo_model.output, anchors, len(class_names))
```

3.3.4 Filtering boxes

`yolo_outputs` provided all the predicted boxes of `yolo_model` in the correct format. Now is ready to perform filtering and select only the best boxes. Lets now call `yolo_eval`, which you had previously implemented, to do this.

```
1 scores, boxes, classes = yolo_eval(yolo_outputs, image_shape)
```

3.3.5 Run on an image

Let the fun begin. A `sess` graph is created that can be summarized as follows:

1. is given to `yolo_model`. The model is used to compute the output
2. is processed by `yolo_head`. It gives you
3. goes through a filtering function, `yolo_eval`. It outputs your predictions.

The code below also uses the following function:

```
1 image, image_data = preprocess_image("images/" + image_file,  
model_image_size = (608, 608))
```

which outputs:

- image: a python (PIL) representation of your image used for drawing boxes. You won't need to use it.
- image_data: a numpy-array representing the image. This will be the input to the CNN.

```

1 def predict(sess, image, printlabel, scores, boxes, classes,
2 class_names):
3     """
4         Runs the graph stored in "sess" to predict boxes for "image_file".
5         Prints and plots
6         the predictions.
7
8         Arguments:
9             sess -- your tensorflow/Keras session containing the YOLO graph
10            image_file -- name of an image stored in the "images" folder.
11
12            Returns:
13                out_scores -- tensor of shape (None, ), scores of the predicted
14                boxes
15                out_boxes -- tensor of shape (None, 4), coordinates of the predicted
16                boxes
17                out_classes -- tensor of shape (None, ), class index of the
18                predicted boxes
19
20                Note: "None" actually represents the number of predicted boxes, it
21                varies between 0
22                and max_boxes.
23
24
25    ### Code Implementation ###
26
27
28    return out_scores, out_boxes, out_classes, image

```

The output image shows below:

```

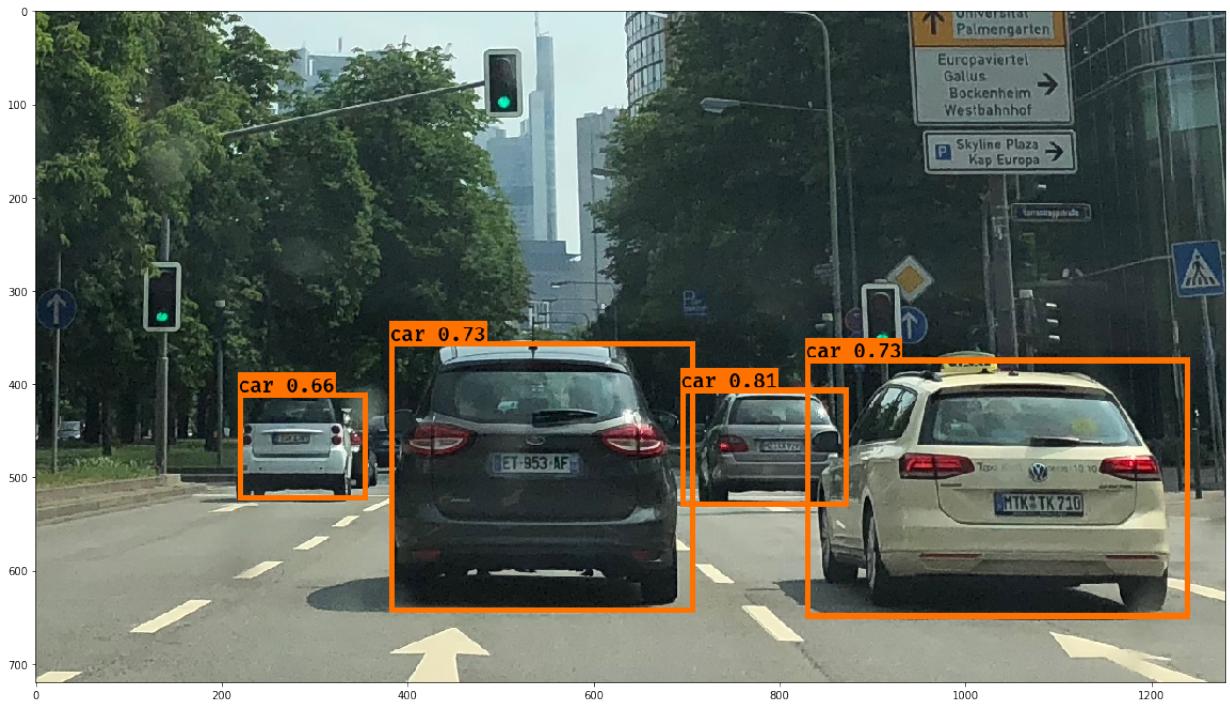
1     out_scores, out_boxes, out_classes, output_image = predict(sess,
2     image,True,scores, boxes, all_classes,class_names=classes)
3     print('Found {} boxes for {}'.format(len(out_boxes), f))
4     plt.imshow(output_image)

```

```

1 car 0.66 (218, 409) (357, 524)
2 car 0.73 (381, 354) (709, 645)
3 car 0.73 (828, 372) (1241, 651)
4 car 0.81 (694, 404) (874, 531)
5 Found 4 boxes for test_images/1.jpg

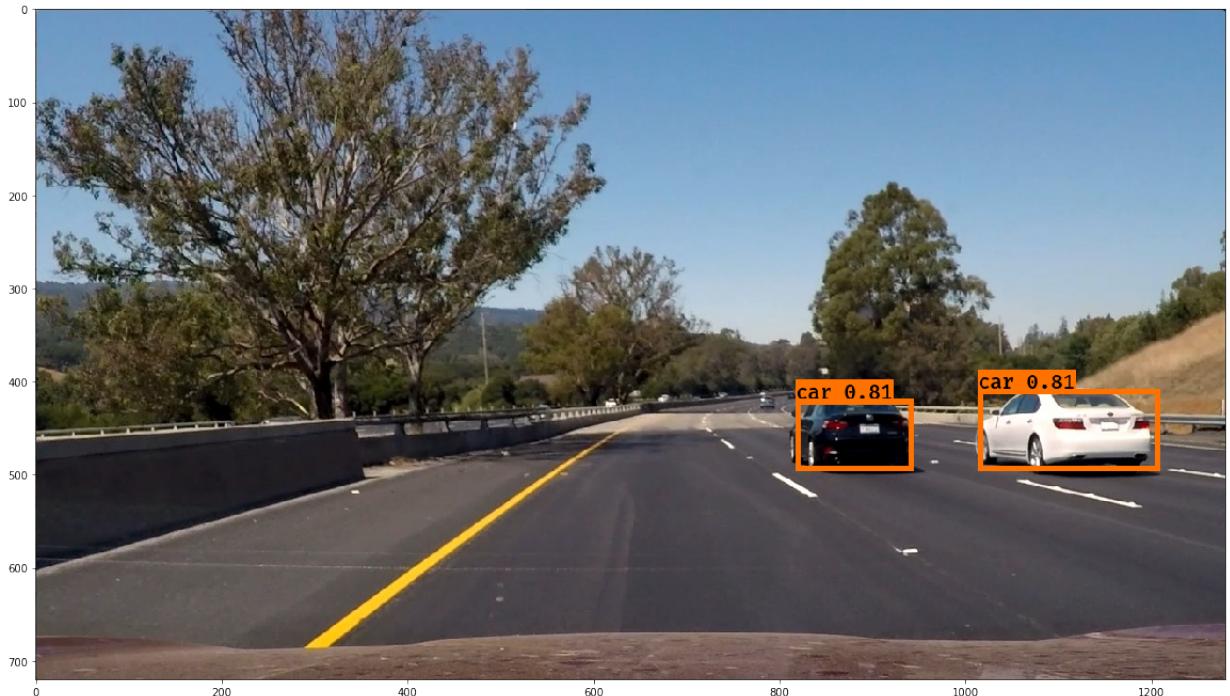
```



```

1 car 0.81 (1014, 408) (1209, 496)
2 car 0.81 (818, 419) (944, 496)
3 Found 2 boxes for test_images/test6.jpg

```



Conclusion:

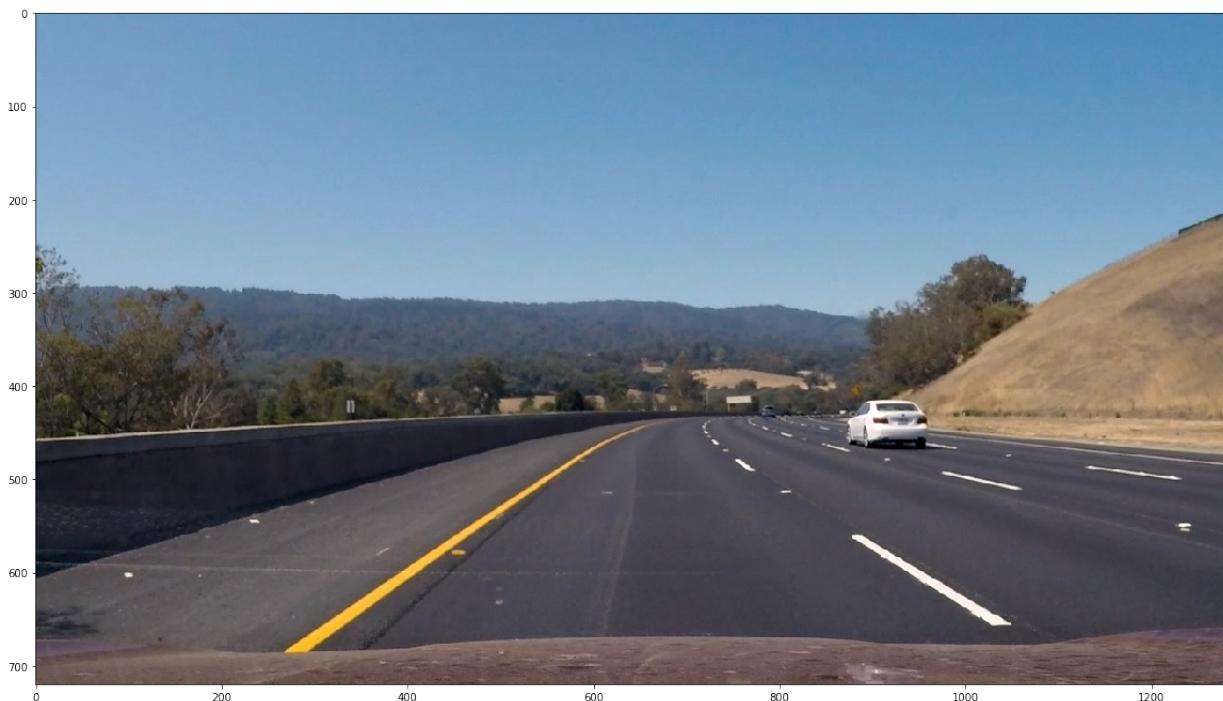
- YOLO is a state-of-the-art object detection model that is fast and accurate
- It runs an input image through a CNN which outputs a $19 \times 19 \times 5 \times 85$ dimensional volume.
- The encoding can be seen as a grid where each of the 19×19 cells contains information about 5 boxes.
- I filter through all the boxes using non-max suppression. Specifically:

- Score thresholding on the probability of detecting a class to keep only accurate (high probability) boxes
- Intersection over Union (IoU) thresholding to eliminate overlapping boxes
- Because training a YOLO model from randomly initialized weights is non-trivial and requires a large dataset as well as lot of computation, I used previously trained model parameters in this exercise.

4. Discussion

There is one image that the vehicle detection is failed.

```
1 | Found 0 boxes for test_images/test3.jpg
```



The car is very clear for human eyes, but the **YOLO** model in this case fails to detect it.

The reason may come from the loaded pertained model. The YOLO model is very computationally expensive to train, that's why I loaded from [The official YOLO website](#).

The way to generate `.h5` file is explained in this [github](#). What I did is as below:

1. Download `.cfg` and `.weights` of YOLOv2 608 × 608 from [The official YOLO website](#)
2. Open **Terminal** or **CMD** to run the following command (I defined the `.cfg` as `yolo.cfg`, same as `weights`):

```
1 | python3 yad2k.py yolo.cfg yolo.weights model_data/yolo.h5
```

3. Load `.h5` into the project

There are several models in the official website, some other models with input size adaption may give a better results.