

March 24, 2023

In this lab, we are given images and need to predict multiple instances of an object class and the bounding boxes of those objects. The COCO dataset was split into training and validation. Below is an output of 3 images from each of the three classes. The plot purposefully chose images in the dataset that had multiple object instances in an image.

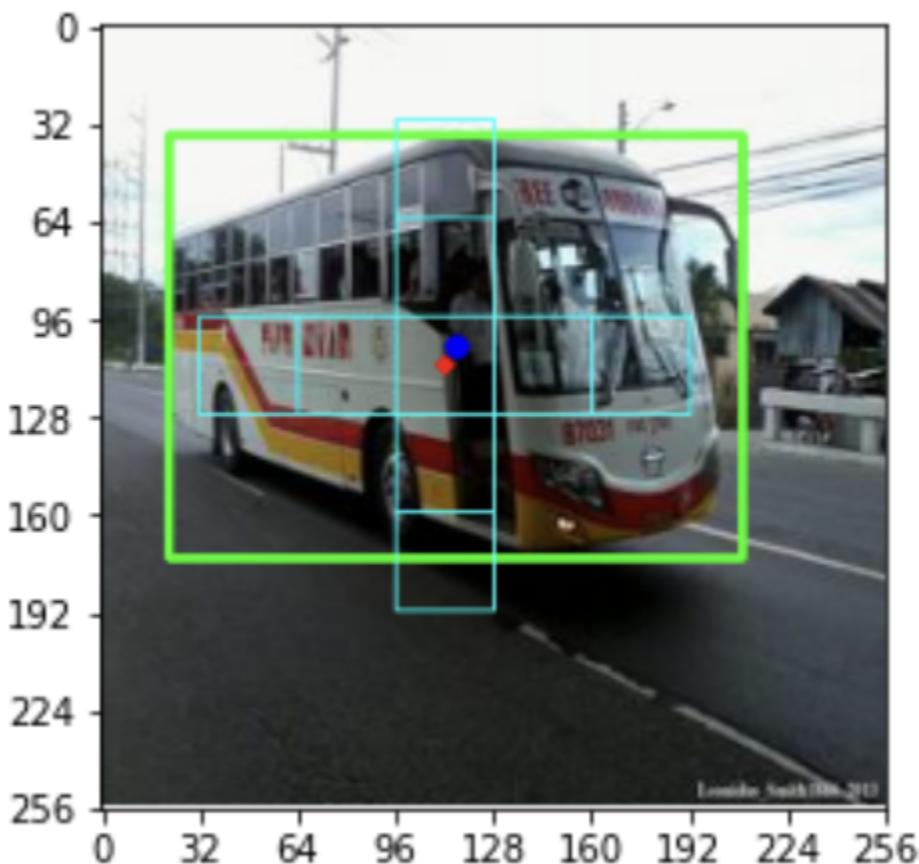


Building the Network:

This network used a ground truth yolo tensor of the shape $8 \times 8 \times 5 \times 8$. Specifically, that means that there were 8 yolo grid cells in the horizontal direction, 8 yolo grid cells in the vertical direction. 8 grid cells in each direction was chosen because it was easily divisible by a 256×256 image. It would make the size of each grid cell 32×32 pixels. Furthermore, for each yolo grid cell, there were 5 possible anchor boxes, which is what the 5 in the shape above

March 24, 2023

represents. Anchor boxes have ratios of 5×1 , 3×1 , 1×1 , 1×3 , and 1×5 . If the 0th anchor box was chosen in the tensor above, it would mean the 5×1 anchor box was chosen. A 1×1 anchor box has a height and width of the yolo interval, which was set to 32 in this lab. To make this more concrete, below is an example of a ground truth image with the possible anchor box choices:



Note that the red dot in the middle of the anchor boxes is the center of the yolo grid cell. The blue dot is the center of the bounding box. In this case, we can visually see that the object instance's bounding box center belongs to the 4th yolo grid cell in the horizontal direction and the 4th yolo grid cell in the vertical direction.

The final dimension, 8, is meant to represent the yolo vector within the yolo tensor. The first element represents "objectness" - whether an object's center is present in that yolo grid and anchor box combination. The second and third elements are the displacement of the center of the yolo grid (in pixel terms) compared to the center of the bounding box in the x and y directions, respectively. The fourth element is the log of the bounding box width minus the log of the anchor box width. The fifth element is the log of the bounding box height minus the log of

March 24, 2023

the anchor box height. The sixth, seventh, and eighth elements are a one-hot encoded vector of the class of the object represented in that specific yolo vector.

The network architecture mainly consisted of downsampled skip blocks which were then put through linear operators to a final output vector size of 2,560 (which is the size of the yolo tensor). Using this output, we could compute the loss of the output yolo tensor against the ground truth yolo tensor.

In the following pages are the code blocks for the network. Note that due to screenshot dimensions, the code blocks are going across multiple pages.

Below is a code block for the network architecture:

```
3 class ResBlock(nn.Module):
4     def __init__(self, in_ch, out_ch):
5         super(ResBlock, self).__init__() #make sure inherited classes are instantiated
6         self.in_ch = in_ch
7         self.out_ch = out_ch
8         self.convo1 = nn.Conv2d(in_ch, out_ch, kernel_size = 3, stride = 1, padding = 1) #creates instance of
9             #Conv2d that preserves shape of input!
10        self.convo2 = nn.Conv2d(in_ch, out_ch, kernel_size = 3, stride = 1, padding = 1) #preserve shape
11        self.convo3 = nn.Conv2d(in_ch, out_ch, kernel_size = 3, stride = 1, padding = 1) #preserve shape
12        self.bn1 = nn.BatchNorm2d(out_ch) #batch norm to improve performance
13        self.bn2 = nn.BatchNorm2d(out_ch) #batch norm to improve performance
14        self.bn3 = nn.BatchNorm2d(out_ch) #batch norm to improve performance
15
16    def forward(self, x):
17        identity = x #save input so we can add it to the output
18        #run three convolutions that preserve the shape, does BN, and applies an activation function
19        out = F.relu(self.bn1(self.convo1(x)))
20        out = F.relu(self.bn2(self.convo2(x)))
21        out = F.relu(self.bn3(self.convo3(x)))
22        # add input to output to give the network more of the original signal
23        out = out + identity
24
25
26
27
28 class HW6Net(nn.Module):
29     def __init__(self, input_nc, output_nc, ngf=8, n_blocks = 4): #ngf = num of conv filters in first convo layer
30         #n_blocks is number of ResNet blocks
31         assert(n_blocks >=0)
32         super(HW6Net, self).__init__()
33         # first conv layer
34         model = [nn.ReflectionPad2d(3),
35                 nn.Conv2d(input_nc, ngf, kernel_size = 7, padding = 0),
36                 nn.BatchNorm2d(ngf),
37                 nn.ReLU(True)]
38         #add downsampling layers
39         n_downsampling = 4
40         for i in range(n_downsampling):
41             mult = 2 ** i
42             model = model + [nn.Conv2d(ngf * mult, ngf * mult * 2, kernel_size = 3, stride = 2, padding = 1),
43                             nn.BatchNorm2d(ngf * mult * 2),
44                             nn.ReLU(True)]
45         #add own Skip blocks
```

March 24, 2023

```

44             nn.ReLU(True)]
45     #add own Skip blocks
46     mult = 2 ** n_downsampling
47     for i in range(n_blocks):
48         model = model + [ResBlock(64, 64)] # do 4 skip connections with 64 in_ch and 64 out_ch
49     self.model = nn.Sequential(*model)
50     ### YOLO head ####
51
52     yolo_head = [nn.Linear(16384, 10000), nn.ReLU(inplace=True), nn.Linear(10000, 2560)]
53     self.yolo_head = nn.Sequential(*yolo_head)
54
55     # class_head = [nn.Linear(16384, 3)] #classification task just takes the skip connection's output and
56     #applies linear function
57     self.class_head = nn.Sequential(*class_head)
58     ### bounding box regression
59
60     bbox_head_conv = [nn.Conv2d(64, 64, kernel_size = 3, padding = 1), nn.BatchNorm2d(64),
61     nn.ReLU(inplace = True), nn.Conv2d(64, 64, 3 , 1),
62     nn.BatchNorm2d(64), nn.ReLU(inplace = True)]
63     self.bbox_head_conv = nn.Sequential(*bbox_head_conv)
64     bbox_head_fc = [nn.Linear(12544, 1024), nn.ReLU(inplace = True), nn.Linear(1024,512),
65     nn.ReLU(inplace = True),
66     nn.Linear(512,4), nn.Sigmoid()] #sigmoid at the end to keep values between 0-1
67     self.bbox_head_fc = nn.Sequential(*bbox_head_fc)
68     # for i in range(n_blocks):
69     def forward(self, x):
70         #skip blocks:
71         ft = self.model(x)
72         #classification task:
73         ft_r = ft.view(ft.shape[0], -1) #reshape output of skipblock to vector, then convert it to
74         # a vector (from original shape)
75         yolo_tensor = self.yolo_head(ft_r) #run classification task, return the 3 node output
76         #bbox regression task
77         bbox_temp = self.bbox_head_conv(ft)
78         bbox_temp = bbox_temp.view(bbox_temp.shape[0], -1)
79         bbox = self.bbox_head_fc(bbox_temp) #return the bounding box.
80
81     return yolo_tensor
82

```

Logic behind dataloading, training, and evaluation:

Dataloader:

In order to write the dataloader, a dictionary was created that stores each image's bounding box annotations, original image height and original image width. Within the dataloader, the bounding box is rescaled to account for using a 256 x 256 image. The dataloader also implements creation of the yolo tensor. Specifically, 5 anchor boxes were used along with a yolo interval of 32 pixels in each direction. Therefore, there are a total of 64 x 5, or 320 possible yolo vectors in each image. Since the dataloader calls for only passing an image and the ground truth, the ground truth yolo vector was created by finding each bounding box's center and the yolo cell that that specific bounding box center belongs to. Then, the IOU for all the possible anchor box configurations was computed, and the anchor box that gave the largest IOU was chosen as the anchor box for that specific yolo cell. Now that the yolo cell and anchor box was chosen, it is now possible to fill out the 8 elements of the yolo vector (as described earlier in the report).

Training:

In order to run training, the image and ground truth yolo tensor was returned from the dataloader in a batch of 20. The images were fed into the network, and the predicted yolo tensor was returned. Since the output is a 20 x 2,560 matrix, it was reshaped back to 20 x 8 x 8 x 5 x 8, which represents a yolo tensor for each of the 20 images. In order to find the correct parts of the

March 24, 2023

yolo tensor to compute loss without looping through all the images, I used the function `torch.nonzeros()` on the ground truth yolo tensor's objectness portion of the yolo vector. Specifically, this means passing in the 0th element of every yolo vector into the nonzero function. This can be done by indexing into `[:, :, :, :, 0]` of the ground truth batch of 20. Passing that into `torch.nonzeros()` gets us the indices of the ground truth yolo tensor that contain objects.

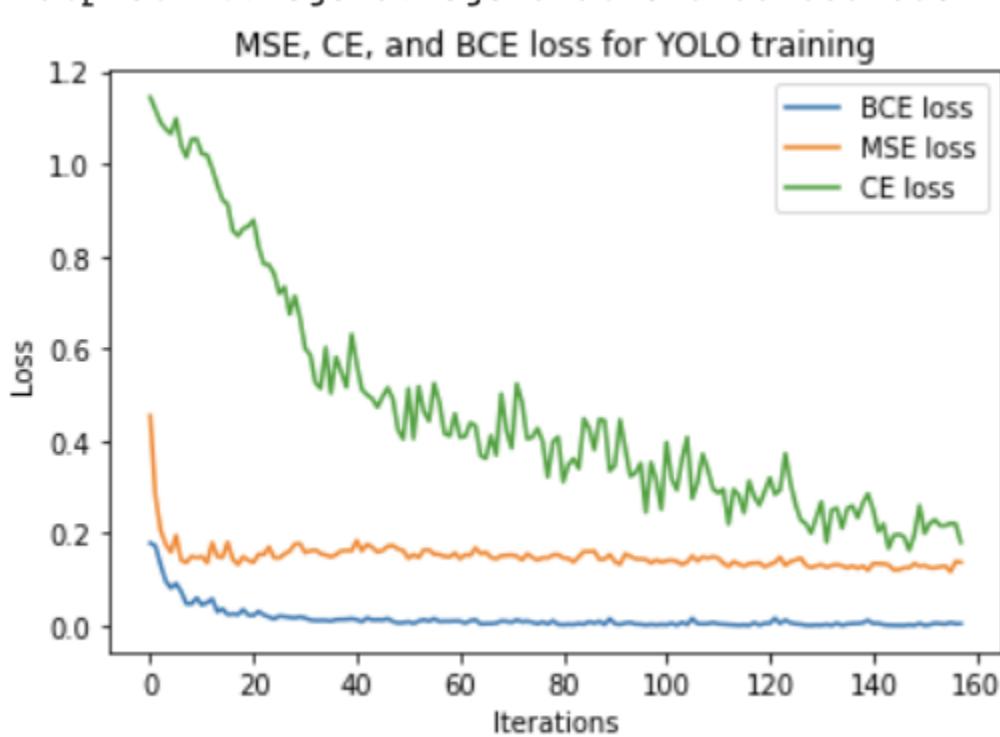
I then passed each index returned from the `torch.nonzeros` function (about 20 - 30 per batch of 20) and was easily able to access the “objectness”, predicted bounding box parameters, and class probabilities vector.

In order to compute the loss, I found the length of the array returned from the `nonzeros` function, and called that the number of objects in the batch. For the BCE loss, since we are just comparing 2 distinct values, the predicted objectness was a vector of number of object x 1. Using the indices from the `torch.nonzero()` function, I was able to get the first element of every ground truth and predicted yolo vector and computed the BCE loss.

Similarly, MSE loss was computed except it required creating a tensor of number of objects x 4, since there are 4 values for each bounding box tensor. The prediction and ground truth were then passed into the MSE loss function.

CE loss was also computed in a similar way (the tensor holding all of the predicted classes was number of objects in batch x 3). The 2 `number_of_objects` x 3 tensors (one representing predictions, one representing ground truth), were passed into the CE loss function.

Training was run for 10 epochs total. Below is the output of the loss curves for BCE, MSE, and CE loss.



Evaluation:

To run evaluation, the batch size was updated to 1. I capped the number of total possible object instances to 5. I then captured the values of the 5 highest “objectness” probabilities, assuming that they were above a threshold of 0.99 in the first place, and stored the indices of the resized yolo tensor that those probabilities belonged to (as part of an $8 \times 8 \times 5 \times 8$ tensor). I then accessed each of those stored indices, checked if the predicted object was actually larger than 4,096 pixels (since that was the minimum size in the dataset), and returned the corresponding bounding boxes represented in the correct yolo vectors.

Results



Results Discussion:

The detector returns many instances of the correct classification. For example, pizzas, cats, and buses all seemed to classify correctly but with multiple bounding boxes over a single object instance. This inability to return one single instance of a bounding box seems to be a pitfall of the detector - the cap of objects was set to 5 and the threshold was also quite high, but, there were still more instances of objects returned than actual objects in the image. The bounding boxes also seem a bit inaccurate - at an eyeball's glance there seems to be roughly a 50% IOU of the predicted bounding box with the ground truth.

March 24, 2023

Source Code:**Script to prepare data from COCO:**

```

1 #Write a script that filters through images and annotations to generate your training and testing dataset
2 # such that any image meets the following criteria: 1) contains at least one foreground object from bus, cat, pizza.
3 # must be at least 64 x 64 = 4096 pixels. There can be multiple foreground objects, but, if there are none, we
4 # can discard it. Once we have the images downloaded, resize them to 256 x 256 and also rescale the bbox params
5 # accordingly.
6 class_list = ['bus', 'cat', 'pizza']
7
8 ### PREPROCESSING TO GENERATE IMAGE UNIQUE IDS ###
9 catIds = coco.getCatIds(catNms = class_list)
10 imgIds = coco.getImgIds(catIds = catIds)
11 imgIds_list = []
12 for catId in catIds:
13     imgIds_list.append(coco.getImgIds(catIds = catId))
14 concatenated_list = []
15 for x in imgIds_list:
16     for y in x:
17         concatenated_list.append(y)
18
19 np_li = np.array(concatenated_list)
20 unique_imgIds = np.unique(np_li) #make the list unique - there might be duplicate image ids
21
22 unique_imgIds = list(unique_imgIds) #convert numpy array back to list
23
24 for n, i in enumerate(unique_imgIds):
25     unique_imgIds[n] = int(i) #convert each element of list to integer, so it can easily be interpreted by
26     #coco api
27 print(len(unique_imgIds)) #7799

```

Create Training Data Dictionary for Dataloader:

```

1 data_dict_train = {}
2
3 #os.chdir("/Users/alim/Documents/ECE60146/hw6/")
4 for n, imgId in enumerate(unique_imgIds):
5     annId = coco.getAnnIds(imgId, catIds=catIds, iscrowd=False)
6     #now that we have the annotation id, we want to check if there is a dominant object. In order to get
7     #the details of that annotation id, we need to load the annotation using that id:
8     annId_metadata = coco.loadAnns(annId)
9     key = -1 #temp value for key that doesn't exist
10    for ann in annId_metadata:
11        [x, y, w, h] = ann['bbox']
12        #print(ann)
13        if ((ann['area'] <= 4096) or (w <= 64) or (h <= 64)):
14            continue #do not consider looking at annotation details for that annId
15        else:
16            #store this annotation:
17            imgId_metadata = coco.loadImgs(imgId)
18            key = imgId_metadata[0]['file_name']
19            if (key in data_dict_train):
20                #we need to append to existing dictionary
21                data_dict_train[key].append({'ann_id': ann['id'], 'bbox': ann['bbox'], 'cat_id': ann['category_id'], 'url':imgId_metadata[0]['coco_url'],
22                                            'img_id' : imgId_metadata[0]['id'], 'img_h' : imgId_metadata[0]['height'],
23                                            'img_w' : imgId_metadata[0]['width']})
24
25        else:
26            #we need to create new dictionary key
27            data_dict_train[key] = [{'ann_id': ann['id'], 'bbox': ann['bbox'], 'cat_id': ann['category_id'], 'url':imgId_metadata[0]['coco_url'],
28                                    'img_id' : imgId_metadata[0]['id'], 'img_h' : imgId_metadata[0]['height'],
29                                    'img_w' : imgId_metadata[0]['width']}]
30
31        #download this image:
32        #if (key != -1):
33        #    coco.download(tarDir = 'train_orig/', imgIds = [data_dict_train[key][0]['img_id']])
34
35
36 ###ORGANIZATION OF data_dict_train:
37 #key is file_name
38 #value is a list of elements.
39 #element 0: annotation id
40 #element 1: bbox
41 #element 2: category id
42 #element 3: img_id coco_url
43 #element 4: img_id
44 #element 5: img height
45 #element 6: img width
46 #print("imgid", imgId)
47 full_len
48 print(len(data_dict_train)) #6712 images in training dataset
49 for k in data_dict_train:
50     temp_len = len(data_dict_train[k])
51     full_len += temp_len
52 print(full_len) #8343 annotations total

```

Create Validation Data Dictionary for Dataloader:

```
 1 data_dict_val = {}
 2
 3 #os.chdir("./Users/alim/Documents/ECE60146/hw6/")
 4 for n, imgId in enumerate(unique_imgIds_test):
 5     annId = coco.getAnnIds(imgId, catIds=catIds, iscrowd=False)
 6     #now that we have the annotation id, we want to check if there is a dominant object. In order to get
 7     #the details of that annotation id, we need to load the annotation using that id:
 8     annId_metadata = coco.loadAnns(annId)
 9     key = -1 #temp value for key that doesn't exist
10     for ann in annId_metadata:
11         [x, y, w, h] = ann['bbox']
12         #print(ann)
13         if ((ann['area'] <= 4096) or (w <= 64) or (h <= 64)):
14             continue #do not consider looking at annotation details for that annId
15         else:
16             #store this annotation:
17             imgId_metadata = coco.loadImgs(imgId)
18             key = imgId_metadata[0]['file_name']
19             if (key in data_dict_val):
20                 #we need to append to existing dictionary
21                 data_dict_val[key].append({'ann_id' : ann['id'], 'bbox' : ann['bbox'], 'cat_id': ann['category_id'],
22                                         'url':imgId_metadata[0]['coco_url'],
23                                         'img_id' : imgId_metadata[0]['id'], 'img_h' : imgId_metadata[0]['height'],
24                                         'img_w' : imgId_metadata[0]['width']})
25
26             else:
27                 #we need to create new dictionary key
28                 data_dict_val[key] = [{"ann_id" : ann['id'], 'bbox' : ann['bbox'], 'cat_id': ann['category_id'], 'url':imgId_metadata[0]['coco_url'],
29                                         'img_id' : imgId_metadata[0]['id'], 'img_h' : imgId_metadata[0]['height'],
30                                         'img_w' : imgId_metadata[0]['width']}]
31             #download this image:
32             #if (key != -1):
33             #    coco.download(tarDir = 'val_orig/', imgIds = [data_dict_val[key][0]['img_id']])
34
35 full_len = 0
36 print(len(data_dict_val)) #3427 images in val dataset
37 for k in data_dict_val:
38     temp_len = len(data_dict_val[k])
39     full_len += temp_len
40 print(full_len) #4198 annotations total in val dataset
```

Create Dataloader:

```
1 ##### build the dataloader.
2
3 #root = '/Users/alim/Documents/ECE60146/hw6/'
4 folders = ['train_resized/', 'val_resized/']
5
6 class MyDataset(torch.utils.data.Dataset):
7     def __init__(self, folder):
8         super(MyDataset).__init__()
9         #self.root = root
10        self.folder = folder
11        if (self.folder == 'train_resized/'):
12            self.data_dict = data_dict_train
13        if (self.folder == 'val_resized/'):
14            self.data_dict = data_dict_val
15
16        self.mapping = {6: 0, 17: 1, 59: 2}
17        self.one_hot_encoding = {0: torch.tensor(np.array([1, 0, 0])),
18                                1: torch.tensor(np.array([0, 1, 0])),
19                                2: torch.tensor(np.array([0, 0, 1]))}
20        self.images = os.listdir(folder) #create list of files in the train or val directory. We will use
21        # this list to get bbox params, read image files, etc.
22
23        for img in self.images:
24            if (img == ".DS_Store"):
25                self.images.remove(".DS_Store") #handle case when image isn't an image. Just remove it from the
26                #image list.
27
28        self.to_Tensor_and_Norm = tvt.Compose([tvt.ToTensor(), tvt.Normalize([0], [1])])
29
30        #find the image with the most objects, get that number of objects so we can create an empty list
31        # with that number of elements:
32        self.max_obj_instances = 0 #this becomes 14 for train, 13 for val
33        for key in self.data_dict:
34            if len(self.data_dict[key]) > self.max_obj_instances:
35                self.max_obj_instances = len(self.data_dict[key])
36        self.cell_and_anchors_tracker = np.zeros((8, self.max_obj_instances)) -99 #instantiate to -99
37        #tracker is a 3 x 14 matrix. Each col will hold data of an obj instance
38        #row 1 and row 2 will be x,y idx of assigned cell
39        #row 3 will be the assigned anchor box. The assignments will be made in the __getitem__ function
40        #rows 4 to 7 will be the x,y,w,h (normalized to a 256 to 256 image)
41        #row 8 is the class in terms of 0, 1, 2
42
43
44
45    def __len__(self):
46        return len(self.data_dict)
47
48    def __getitem__(self, index):
49        #prepare image:
50        PIL_img = Image.open(self.folder + self.images[index])
51        torch_img = self.to_Tensor_and_Norm(PIL_img)
52
53        #prepare class labels
54        li_class_labels = []
```

```

54     li_class_labels = []
55     for i in range(len(self.data_dict[self.images[index]])):
56         #print(self.data_dict[self.images[index]])
57         class_label = self.data_dict[self.images[index]][i]['cat_id']
58         class_label = self.mapping(class_label)
59         class_label = torch.tensor(class_label)
60         li_class_labels.append(class_label)
61         #for each object in the image, figure out which grid cell it belongs (i.e, cell which is closest)
62         #to the center of the bbox.
63         #will have a 8 by 8 grid cell, so, the yolo interval will be 256 / 8 = 32
64         #get bbox center:
65         yolo_interval = 32
66         bbox_anchor_list = []
67         yolo_tensor = torch.zeros(8,8,5,8)
68         dict_assigned_yolo_vectors = {}
69
70         for i in range(len(self.data_dict[self.images[index]])): #iterate through each object instance
71             if i == 0: #start to print image for sanity check
72                 np_img = np.uint8(PIL_img)
73                 ticks = np.linspace(start = 0, stop = 256, num=9)
74
75
76             [x, y, w, h] = self.data_dict[self.images[index]][i]['bbox']
77             img_w = self.data_dict[self.images[index]][i]['img_w']
78             img_h = self.data_dict[self.images[index]][i]['img_h']
79             #scale x y w h in terms of a 256 x 256 image.
80             scaling_factor_x = 256.0 / img_w
81             scaling_factor_y = 256.0 / img_h
82             x = scaling_factor_x * x
83             y = scaling_factor_y * y
84             w = scaling_factor_x * w
85             h = scaling_factor_y * h
86             #done scaling original x, y, w, h parameters
87
88             #find center of bbox
89             cx_bbox = (x + x + w) / 2 #center of bbox (x coord)
90             cy_bbox = (y + y + h) / 2 #center of bbox (y coord)
91
92             #find the yolo cell that the bbox belongs to:
93             yolo_grid_x_idx = cx_bbox // yolo_interval
94             yolo_grid_y_idx = cy_bbox // yolo_interval
95             #find the exact yolo cell center coords. We will need this to compute dx, dy, the difference
96             #between center of bbox and center of yolo grid
97             yolo_grid_x_center = yolo_grid_x_idx * yolo_interval + (yolo_interval / 2)
98             yolo_grid_y_center = yolo_grid_y_idx * yolo_interval + (yolo_interval / 2)
99
100            #compute displacements between the center of the cell and the center of the bounding box in a
101            #training image for a given object instance (dx and dy)
102            dx = cx_bbox - yolo_grid_x_center
103            dy = cy_bbox - yolo_grid_y_center
104
105            #plot bbox, center of yolo grid cell that was assigned, and center of bbox
106            np_img = cv2.rectangle(np_img, (int(x), int(y)),
107                                  (int((x+w)), int((y+h))), (0,255,0), 2)
108            #image, center_coordinates, radius, color, thickness

```

```

108     #image, center_coordinates, radius, color, thickness
109     np_img = cv2.circle(np_img, (int(cx_bbox), int(cy_bbox)), 1,
110                         (0,0,255), 5)
111     np_img = cv2.circle(np_img, (int(yolo_grid_x_center), int(yolo_grid_y_center)), 1, (255, 0, 0), 4)
112
113     #pick the correct anchor box. There are 5 anchor boxes with different aspect ratio:
114     # 5 x 1, 3 x 1, 1 x 1, 1 x 3, 1 x 5
115     #iou_box function in pytorch expects that coords are in [x1 y1 x2 y2] format.
116     #we already know the bbox:
117     GT_bbox_for_iou = torch.tensor([x, y, x+w, y+h]).view(1,4) #view reshapes to 1,4
118     IOU_li = []
119     anchor_box_options = [[5,1], [3,1], [1,1], [1,3], [1, 5]]
120     for n, ab in enumerate(anchor_box_options):
121
122         bbox_anchor_w = yolo_interval * ab[0]
123         bbox_anchor_h = yolo_interval * ab[1]
124         bbox_anchor = torch.tensor([yolo_grid_x_center, yolo_grid_y_center, bbox_anchor_w, bbox_anchor_h]).view(1,4)
125         #above bbox_anchor is in cx,cy,w,h fmt
126         bbox_anchor = ops.box_convert(bbox_anchor, in_fmt = 'cxcywh', out_fmt = 'xyxy')
127         bbox_anchor = torch.clamp(bbox_anchor, 0, 256)
128         bbox_anchor_list.append(bbox_anchor)
129
130     np_img = cv2.rectangle(np_img, (int(bbox_anchor_list[n+(5*i)][0][0]), int(bbox_anchor_list[n+(5*i)][0][1])), (int(bbox_anchor_list[n+(5*i)][0][2]),
131                           int(bbox_anchor_list[n+(5*i)][0][3])), (0, 255, 255), 1)
132
133     IOU = ops.box_iou(GT_bbox_for_iou, bbox_anchor)
134     IOU_li.append(IOU)
135     if n == 4:
136         temp_arr = np.empty(5)
137         for k, iou in enumerate(IOU_li):
138             temp_arr[k] = iou.view(-1).item() #put computed IOUs into nparray format so we can perform argmax
139             max_IOU_idx = np.argmax(temp_arr)
140             anchor_mapping = {0: [5,1], 1: [3, 1], 2: [1, 1], 3: [1, 3], 4: [1,5]} #gets from index to h/w of anchor box
141             #create yolo vector's first 3 elements:
142             yolo_vector = torch.empty((1,8))
143             yolo_vector[0,0] = 1
144             yolo_vector[0,1] = dx / yolo_interval
145             yolo_vector[0,2] = dy / yolo_interval
146             #compute sigma_w, sigma_h, the ratio between GT and anchor w and ratio between GT and anchor h
147             best_anchor_w = anchor_mapping[max_IOU_idx][0]
148             best_anchor_h = anchor_mapping[max_IOU_idx][1]
149             w_GT = torch.log(torch.tensor(w)) - torch.log(torch.tensor(best_anchor_w * yolo_interval))
150             h_GT = torch.log(torch.tensor(h)) - torch.log(torch.tensor(best_anchor_h * yolo_interval))
151
152             yolo_vector[0,3] = w_GT
153             yolo_vector[0,4] = h_GT
154             #assign class label as one hot vector:
155             one_hot_vector = self.one_hot_encoding[int(li_class_labels[i])]
156             numerical_class = torch.argmax(one_hot_vector)
157             #print("ONE HOT VEC", one_hot_vector)
158             yolo_vector[0,5:] = one_hot_vector
159             #print("FINAL YOLO VECTOR", yolo_vector.view(-1))
160             #print("anchor box chosen was:", anchor_mapping[max_IOU_idx], "where anchbox is W, H")
161             #print("yolo cell was:", yolo_grid_x_idx, yolo_grid_y_idx)
162

```

```
150     h_GT = torch.log(torch.tensor(h)) - torch.log(torch.tensor(best_anchor_h * yolo_interval))
151
152
153     yolo_vector[0,3] = w_GT
154     yolo_vector[0,4] = h_GT
155     #assign class label as one hot vector:
156     one_hot_vector = self.one_hot_encoding[int(li_class_labels[i])]
157     numerical_class = torch.argmax(one_hot_vector)
158     #print("ONE HOT VEC", one_hot_vector)
159     yolo_vector[0,5:] = one_hot_vector
160     #print("FINAL YOLO VECTOR", yolo_vector.view(-1))
161     #print("anchor box chosen was:", anchor_mapping[max_IOU_idx], "where anchbox is W, H")
162     #print("yolo cell was:", yolo_grid_x_idx, yolo_grid_y_idx)
163     #assign yolo_vector to the correct index of yolo_tensor:
164     #print('assigned to:',int(yolo_grid_x_idx), int(yolo_grid_y_idx), max_IOU_idx)
165
166     yolo_tensor[int(yolo_grid_x_idx), int(yolo_grid_y_idx),max_IOU_idx,:] = yolo_vector
167     dict_assigned_yolo_vectors[(int(yolo_grid_x_idx), int(yolo_grid_y_idx),max_IOU_idx)] = \
168     {'grid_x' : int(yolo_grid_x_idx), 'grid_y' : int(yolo_grid_y_idx), 'anchor_box_idx' : max_IOU_idx}
169     self.cell_and_anchors_tracker[:,i] = [yolo_grid_x_idx, yolo_grid_y_idx, max_IOU_idx, x, y, w, h, numerical_class]
170     #print("yolo_vector", yolo_vector)
171     # if i == (len(self.data_dict[self.images[index]]) - 1):
172     #     plt.imshow(np_img)
173     #     plt.xticks(ticks)
174     #     plt.yticks(ticks)
175
176     return torch_img, yolo_tensor, self.cell_and_anchors_tracker
177
178 train_dataset = MyDataset( folders[0])
179 index = 6
180 train_dataset[index][1].shape
181 batch_size = 1
182 #create the dataloader:
183 train_dataloader = torch.utils.data.DataLoader(train_dataset, batch_size=batch_size, num_workers = 0, drop_last=True, shuffle=True)
184 # for n, i in enumerate(train_dataloader):
185 #     print(n)
186 #     print(i[2])
187
188 val_dataset = MyDataset(folders[1])
189 val_dataloader = torch.utils.data.DataLoader(val_dataset, batch_size = 1, num_workers = 0, drop_last = True, shuffle = True)
190 # for n, i in enumerate(val_dataloader):
191 #     print(n)
192
```

Code for Model was given above

Code for Training Loop:

```
1 #create training loop:
2 torch.autograd.set_detect_anomaly(True)
3
4 test_net = HW6Net(input_nc = 3, output_nc = 3, ngf = 4, n_blocks=4)
5 test_net = test_net.to(torch.device("cuda:0"))
6 total_loss_list = []
7 bce_loss_list = []
8 ce_loss_list = []
9 mse_loss_list = []
10 criterion_CE = nn.CrossEntropyLoss()
11 criterion_MSE = nn.MSELoss()
12 criterion_BCE = nn.BCELoss()
13 optimizer = torch.optim.Adam(test_net.parameters(), lr = 1e-3, betas = (0.9, 0.99))
14 epochs = 10
15 num_layers = len(list(test_net.parameters()))
16 print(num_layers)
17
18 for epoch in range(epochs):
19     running_loss_CE = 0.0
20     running_loss_MSE = 0.0
21     running_loss_BCE = 0.0
22     print('in epoch', epoch)
23     for n, data in enumerate(train_dataloader):
24         imgs, GT_yolo_tensor, grid_info = data
25         imgs = imgs.to(torch.device("cuda:0"))
26         GT_yolo_tensor = GT_yolo_tensor.to(torch.device("cuda:0"))
27         optimizer.zero_grad()
28         output_yolo_tensor = test_net(imgs) #output is 20 x 2560
29
30         #compute BCE loss for whether there is an object present or not in specific anchor box:
31         GT_yolo_tensor = torch.reshape(GT_yolo_tensor, (batch_size, 8, 8, 5, 8))
32         output_yolo_tensor = torch.reshape(output_yolo_tensor, (batch_size, 8, 8, 5, 8))
33         #print(GT_yolo_tensor.shape)
34         #find idx of non-zeros:
35         non_zero_idx = torch.nonzero(GT_yolo_tensor[:,:,:,:,0])
36
37
38         #compute BCE loss for object vs not object:
39         num_objects = non_zero_idx.shape[0] #find the total number of objects in the GT.
40         #for each yolo vector, we will return 6 non-zero values (2 of the 8 are 0 because of the 3 class
41         #one hot encoding). We want to store the predictions of every 6th index returned from the non-zero
42         #indices as those will contain the "predicted objectness", or the first element of the yolo vector.
43         #Since we want to store every 6th index, we need to divide the total length of the non-zero indices
44         #by 6 so we have the correct length vector of predictions to feed into BCE loss.
45         #prepare arrays for pred/GT objectness:
46         pred_objectness = torch.zeros(num_objects)
47         GT_objectness = torch.ones(num_objects) #note that 3 is the batch size
48
49         #prepare arrays for pred/GT bbox:
50         pred_bbox = torch.zeros(num_objects , 4)
51         GT_bbox = torch.zeros(num_objects, 4)
52
53         #prepare arrays for pred/GT classification:
```

```

53     #prepare arrays for pred/GT classification:
54     pred_class_prob = torch.zeros(num_objects, 3)
55     GT_classes = torch.zeros(num_objects).type(torch.LongTensor) #torch expects GT categories to be
56     # "long tensors", which are really just integers in torch. So, we convert this to long tensor.
57     #compute CE loss for classification:
58
59     for idx, i in enumerate(non_zero_idx):
60         #the following line gets pred objectness for one of the inputs into BCE loss.
61         pred_objectness[idx] = nn.Sigmoid()(output_yolo_tensor[i[0], i[1], i[2], i[3], 0])
62
63         #these following 2 lines will be dx, dy, w ratio, h ratio for pred and GT for input into MSE loss
64         pred_bbox[idx, :] = output_yolo_tensor[i[0], i[1], i[2], i[3], 1:5]
65         GT_bbox[idx, :] = GT_yolo_tensor[i[0], i[1], i[2], i[3], 1:5]
66
67         #the following 2 lines get the pred/GT classifications for input into CE loss
68         pred_class_prob[idx,:] = output_yolo_tensor[i[0], i[1], i[2], i[3], 5:]
69         GT_classes[idx] = torch.argmax(GT_yolo_tensor[i[0], i[1], i[2], i[3], 5:])
70
71
72     loss_BCE = criterion_BCE(pred_objectness, GT_objectness) ##remember - BCE is not symmetric! Need
73     #to have the first arg be the prediction and the second arg be the target/GT.
74     loss_BCE.backward(retain_graph = True)
75     #total_loss = total_loss + loss_BCE
76
77     #convert bbox to xyxy fmt:
78     pred_bbox_MSE = ops.box_convert(pred_bbox, in_fmt = 'cxcywh', out_fmt = 'xyxy')
79     GT_bbox_MSE = ops.box_convert(GT_bbox, in_fmt = 'cxcywh', out_fmt = 'xyxy')
80
81     #compute MSE loss for bbox.
82     loss_MSE = criterion_MSE(pred_bbox_MSE, GT_bbox_MSE)
83     loss_MSE.backward(retain_graph = True)
84
85     #compute CE loss for classification:
86     loss_CE = criterion_CE(pred_class_prob, GT_classes )
87     loss_CE.backward()
88     running_loss_BCE = running_loss_BCE + loss_BCE.item()
89     running_loss_CE = running_loss_CE + loss_CE.item()
90     running_loss_MSE = running_loss_MSE + loss_MSE.item()
91
92     optimizer.step()
93     if (n+1) % 20 == 0:
94         print("[epoch: %d, batch: %5d] loss: %3f" % (epoch + 1, n + 1, running_loss_CE + running_loss_MSE +
95                                                 running_loss_BCE / 20))
96         print("running_loss_BCE:", running_loss_BCE)
97         print("running_loss_CE:", running_loss_CE)
98         print("running_loss_MSE:", running_loss_MSE)
99         bce_loss_list.append(running_loss_BCE / 20)
100        running_loss_BCE = 0.0
101        ce_loss_list.append(running_loss_CE / 20)
102        running_loss_CE = 0.0
103        mse_loss_list.append(running_loss_MSE / 20)
104        running_loss_MSE = 0.0
105        total_loss_list.append(running_loss_CE + running_loss_MSE +
106                               running_loss_BCE / 20)
107

```

Code for Evaluation:

```

1 #eval code:
2 #iterate over validation dataset, get predictions with "objectness" and plot the corresponding pred classification
3 #and bbox and the GT pred/bbox.
4 count_objects = 0
5 class_mapping = {0: 'bus', 1: 'cat', 2: 'pizza'}
6 fig, ax = plt.subplots(2,5, figsize=(20,20))
7 with torch.no_grad():
8     for n, data in enumerate(val_dataloader):
9         if n < 10:
10             max_obj_count = 0
11             print("in val dataloader iteration", n)
12             #print("STARTING EVAL CODE")
13             images, GT_yolo_tensor, grid_info = data
14             images = images.to(torch.device("cuda:0"))
15             GT_yolo_tensor = GT_yolo_tensor.to(torch.device("cuda:0"))
16
17             #bbox_label = ops.box_convert(bbox_label, in_fmt = 'xyxy', out_fmt = 'cxcywh')
18             pred_yolo_tensor = test_net(images)
19             pred_yolo_tensor = torch.reshape(pred_yolo_tensor, (8, 8, 5, 8))
20             GT_yolo_tensor = torch.reshape(GT_yolo_tensor, (8, 8, 5, 8))
21
22             #open image, put into numpy fmt
23             img_np = np.array(images.squeeze(dim = 0).cpu())
24             img_np = img_np.transpose(1,2,0)
25             img_PIL = Image.fromarray(np.uint8(img_np * 255)) #this was originally a tensor, so, need to unnormalize!!
26             img_np = np.uint8(img_PIL)
27
28
29             #get maximum probabilities of "objectness of prediction" and once we have found those indices, plot their corresponding bboxes:
30
31             max_objectness = np.zeros((5))
32             max_objectness_pos = np.zeros((5,4))
33             for i in range(8):
34                 for j in range(8):
35                     for k in range(5):
36                         temp_objectness = np.array(nn.Sigmoid()(pred_yolo_tensor[i,j,k,0].cpu()))
37                         if (temp_objectness > 0.999):
38                             if np.min(max_objectness) < temp_objectness:
39                                 idx_temp = np.argmax(max_objectness)
40                                 max_objectness[idx_temp] = temp_objectness #store probability of object if current probability is higher than the lowest probability in the top 5 objects
41                                 max_objectness_pos[idx_temp] = [i,j,k,0]
42
43             if (i == 7) and (j == 7) and (k == 4):
44                 for num, x in enumerate(max_objectness_pos):
45                     #print("MAX objectness is", max_objectness)
46                     if np.array_equal(np.zeros((5)), max_objectness):
47                         continue #didn't find any objects above threshold!
48                     #print("X IS ", x)
49                     #print("the max objectness pos is", max_objectness_pos)
50                     #print("max_objectness is", max_objectness)
51                     #we know this has a high probability of being an object. Compute h/w to make sure it is a reasonable size.
52                     pred_bbox = pred_yolo_tensor[int(x[0]), int(x[1]), int(x[2]), 1:5]
53
54                     #we know this has a high probability of being an object. Compute h/w to make sure it is a reasonable size.
55                     pred_bbox = pred_yolo_tensor[int(x[0]), int(x[1]), int(x[2]), 1:5]
56
57                     center_x = x[0] * 32 + pred_bbox[0] * 32 #predicted yolo cell_x + offset_x from predicted yolo center gives us the predicted center of bbox
58                     center_y = x[1] * 32 + pred_bbox[1] * 32 #predicted yolo cell_y + offset_y from predicted yolo center gives us the predicted center of bbox
59                     anchor_mapping = {0: [5,1], 1: [3,1], 2: [1,1], 3: [1,3], 4: [1,5]} #gets from index to w/h of anchor box
60                     pred_bbox_w = torch.exp(pred_bbox[2]) * anchor_mapping[x[2]][0] * 32 #this value represents which anchor box we are looking at in the prediction.
61                     pred_bbox_h = torch.exp(pred_bbox[3]) * anchor_mapping[x[3]][1] * 32
62                     #print("pred_bbox_h", pred_bbox_h)
63                     if ((pred_bbox_h * pred_bbox_w) < 4096):
64                         #print("found an obj too small")
65                         continue #object is too small for dataset.
66                     max_obj_count += 1
67
68                     # we have passed test for object size. Convert it to xyxy fmt so bbox can be displayed on image:
69                     pred_bbox_real_coords = torch.tensor([center_x, center_y, pred_bbox_w, pred_bbox_h])
70                     pred_bbox_xyxy = ops.box_convert(pred_bbox_real_coords, in_fmt = 'cxcywh', out_fmt = 'xyxy')
71                     pred_bbox_xyxy = torch.clamp(pred_bbox_xyxy, 0, 256)
72                     x1 = int(pred_bbox_xyxy[0])
73                     y1 = int(pred_bbox_xyxy[1])
74                     x2 = int(pred_bbox_xyxy[2])
75                     y2 = int(pred_bbox_xyxy[3])
76                     #add bbox to image:
77                     img_np = cv2.rectangle(img_np, (x1, y1), (x2, y2), (0,255,0), 2)
78
79                     #get object classification:
80                     #print("pred_yolo_tensor[i, j, k, 5:]")
81                     pred_class_vector = int(torch.argmax(pred_yolo_tensor[int(x[0)], int(x[1]), int(x[2]), 5:]).cpu())
82                     #print("pred_class_vector", pred_class_vector)
83                     class_mapping = {0: 'bus', 1: 'cat', 2: 'pizza'}
84                     pred_classification = class_mapping[pred_class_vector]
85                     #print(pred_classification)
86                     img_np = cv2.putText(img_np, pred_classification, (x1, y1 + 25), cv2.FONT_HERSHEY_SIMPLEX, 0.8, (0,255, 0), 2)
87
88
89             ## Plot GT bboxes
90
91             GT_yolo_tens = grid_info.squeeze
92             non_zeros = non_zero_idx <built-in method squeeze of Tensor object at 0x7f7dfd821ef0>
93             grid_info = grid_info.squeeze(dim=0)
94             ticks = np.linspace(start = 0, stop = 256, num=9)
95             for num, i in enumerate(non_zero_idx):
96                 #print("grid_info", grid_info[:,num])
97                 x1, y1 = grid_info[3:5,num]
98                 w, h = grid_info[5:7, num]
99                 class_num = grid_info[-1, num]
100
101                 x2 = int(x1 + w)
102                 y2 = int(y1 + h)
103                 x1 = int(x1)
104                 y1 = int(y1)
105                 img_np = cv2.rectangle(img_np, (x1, y1), (x2, y2), (255, 0, 0), 2)

```

```
104     x1 = int(x1)
105     y1 = int(y1)
106     img_np = cv2.rectangle(img_np, (x1, y1), (x2, y2), (255, 0, 0), 2)
107     GT_bbox = GT_yolo_tensor[i[0], i[1], i[2], 1:5]
108     anchor_mapping = {0: [5,1], 1: [3, 1], 2: [1, 1], 3: [1, 3], 4: [1,5]} #gets from index to w/h of anchor box
109     GT_bbox_w = GT_bbox[2] * anchor_mapping[int(i[2])][0] * 32 #this value represents which anchor box we are looking at in the prediction.
110     #the above line multiplies the ratio of the GT/anch width by the anchor width by the yolo interval to get the predicted w
111     GT_bbox_h = GT_bbox[3] * anchor_mapping[int(i[2])][1] * 32
112
113
114
115     img_np = cv2.putText(img_np, class_mapping[int(class_num)], (x1, y1 + 25),
116                           cv2.FONT_HERSHEY_SIMPLEX, 0.8, (255,0, 0),2)
117     #print("about to plot image!!! when n is: ", n)
118     if n < 5:
119         ax[0,n].set_xticks(ticks)
120         ax[0,n].set_yticks(ticks)
121         ax[0,n].imshow(img_np)
122     if n >4:
123         ax[1,n-5].set_xticks(ticks)
124         ax[1,n-5].set_yticks(ticks)
125         ax[1,n-5].imshow(img_np)
126
127
128
129     if n > 10:
130         break
131 #print(count)
```