



University of Tehran  
College of Engineering  
School of Electrical and Computer Engineering



# Intelligent Systems

Dr.Hosseini

## Final Project

Navid Dehban

810198390

Soroush Mesforush

810198472

Nika Emami

810198356

Reza Jahani

810198377

Bahman 01

## Contents

<b>1</b>	<b>Question 1 : Image reconstruction from feature maps</b>	<b>4</b>
1.1	Part A . . . . .	4
<b>2</b>	<b>Question 2 : Object detection and recognition</b>	<b>10</b>
2.1	History . . . . .	10
2.1.1	Part 1 . . . . .	10
2.1.2	Part 2 . . . . .	16
2.2	Implementing YOLO . . . . .	17
2.2.1	Familiarity with the model . . . . .	17
2.2.2	Understanding the theory of the problem . . . . .	17
2.2.3	Implementation steps . . . . .	23

**Abstract**

In this project we shall study machine vision in two sections, first we shall perform image reconstruction, we do this using the **AlexNet** model.

In the second section we shall perform object detection using the **YOLOv3** model.

# 1 Question 1 : Image reconstruction from feature maps

## 1.1 Part A

Here we want to design neural networks using pretrained **AlexNet** weights to recognize images taken from different layers of the **AlexNet** network.

First of all we take a look at the **AlexNet** architecture.

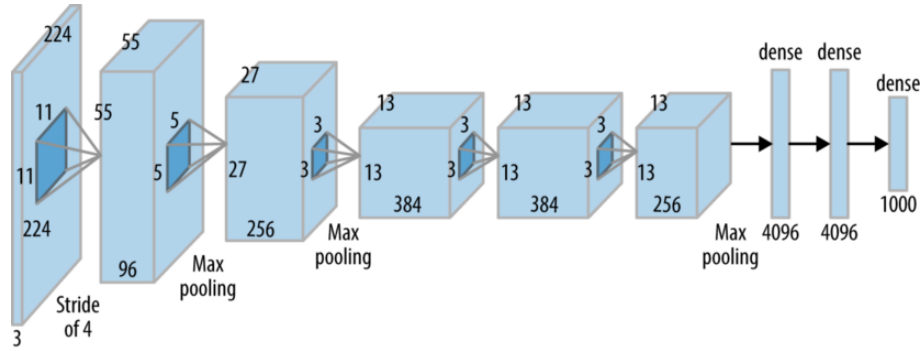


Figure 1: Alexnet architecture

This diagram matches the description of the **AlexNet** paper as follows.

layer	CONV1		CONV2		CONV3	CONV4	CONV5		FC6		FC7		FC8
processing steps	conv1 relu1	mpool1 norm1	conv2 relu2	mpool2 norm2	conv3 relu3	conv4 relu4	conv5 relu5	mpool5	fc6 relu6	drop6	fc7 relu7	drop7	fc8
out size	55	27	27	13	13	13	13	6	1	1	1	1	1
out channels	96	96	256	256	384	384	256	256	4096	4096	4096	4096	1000

Figure 2: Alexnet architecture

Now we go on to get the output of each layer and design the proper neural networks for each one accordingly.

```
1 def imshow(img):
2     npimg = img.numpy()
3     plt.imshow(np.transpose(npimg, (1, 2, 0)))
4     plt.show()
5
6 def show_batch(dataloader):
7     dataiter = iter(dataloader)
8     images, labels = dataiter.next()
9     imshow(make_grid(images))
10
11 def show_image(dataloader):
12     dataiter = iter(dataloader)
13     images, labels = dataiter.next()
14     random_num = randint(0, len(images)-1)
15     imshow(images[random_num])
16     label = labels[random_num]
17     print(f'Label: {label}, Shape: {images[random_num].shape}')
18
```

Figure 3: Helper functions

```
1 def generate_dataloader(data, name, transform):
2     if data is None:
3         return None
4     if transform is None:
5         dataset = datasets.ImageFolder(data, transform=T.ToTensor())
6     else:
7         dataset = datasets.ImageFolder(data, transform=transform)
8     # idx = (dataset.targets==1) | (dataset.targets==2) | (dataset.targets==3) | (dataset.targets==4)
9     # dataset.targets = dataset.targets[idx]
10    # dataset.data = dataset.data[idx]/
11    if use_cuda:
12        kwargs = {"pin_memory": True, "num_workers": 1}
13    else:
14        kwargs = {}
15    dataloader = DataLoader(dataset, batch_size=batch_size,
16        shuffle=(name=="train"),
17        **kwargs)
18    return dataloader
19
```

Figure 4: Dataloader

```
1 preprocess_transform_pretrain = T.Compose([
2     T.Resize(256),
3     T.CenterCrop(224),
4     T.RandomHorizontalFlip(),
5     T.ToTensor(),
6     T.Normalize(mean=[0.485, 0.456, 0.406],
7                 std=[0.229, 0.224, 0.225]))
8
```

Figure 5: Preprocessing

```
1 class AlexNetConv2(nn.Module):
2     def init(self):
3         super(AlexNetConv2, self).init()
4         self.features = nn.Sequential(*list(model_alexnet.features.children())
5                                       [:-7])
6     def forward(self, x):
7         x = self.features(x)
8         return x
9
10 class AlexNetConv5(nn.Module):
11     def __init__(self):
12         super(AlexNetConv5, self).__init__()
13         self.featuresdfd = nn.Sequential(*list(model_alexnet.features.children())
14                                         [:-4])
15     def forward(self, x):
16         x = self.featuresdfd(x)
17         return x
18
19 class AlexNetFc8(nn.Module):
20     def init(self):
21         super(AlexNetFc8, self).init()
22         self.features = nn.Sequential(*list(model_alexnet.features.children())
23                                       [:-1])
24     def forward(self, x):
25         x = self.features(x)
26         return x
27
```

Figure 6: Layer Outputs

```
1  class MyModel_conv5(nn.Module):
2      def __init__(self, model):
3          super(MyModel_conv5, self).__init__()
4          self.pre_trained = model
5          self.seq = nn.Sequential(
6              nn.Conv2d(256, 256, 3, 1),
7              nn.Conv2d(256, 256, 3, 1),
8              nn.Conv2d(256, 256, 3, 1),
9              nn.ConvTranspose2d(256, 256, 5, 2),
10             nn.ConvTranspose2d(256, 128, 5, 2),
11             nn.ConvTranspose2d(128, 64, 5, 2),
12             nn.ConvTranspose2d(64, 32, 5, 2),
13             nn.ConvTranspose2d(32, 3, 5, 2),
14             nn.MaxPool2d((94, 94), stride=(1, 1)),
15         )
16
17     def forward(self, x):
18         x = self.pre_trained(x)
19         x = self.seq(x)
20     return x
21
22 pre_trained = AlexNetConv5()
23 model = MyModel_conv5(pre_trained)
24
```

Figure 7: Neural network

```
1 def train(dataloader, model, loss_fn, optimizer):
2     size = len(dataloader.dataset)
3     loss_values = []
4     model.train()
5     running_loss = 0
6     for batch, (X, y) in enumerate(dataloader):
7         X, y = X.to(device), y.to(device)
8         pred = model(X)
9         loss = loss_fn(pred, X)
10        optimizer.zero_grad()
11        loss.backward()
12        optimizer.step()
13        running_loss += loss.item()
14        loss_values.append(running_loss)
15        if batch % 100 == 0:
16            loss, current = loss.item(), batch * len(X)
17            print(f"loss: {loss:>7f} [{current:>5d}/{size:>5d}]")
18    plt.plot(np.array(loss_values), 'r')
```

Figure 8: Train

```
1 def test(dataloader, model, loss_fn):
2     loss_values = []
3     size = len(dataloader.dataset)
4     num_batches = len(dataloader)
5     model.eval()
6     test_loss, correct = 0, 0
7     with torch.no_grad():
8         for X, y in dataloader:
9             X, y = X.to(device), y.to(device)
10            pred = model(X)
11            test_loss += loss_fn(pred, X).item()
12            correct += (pred.argmax(1) == y).type(torch.float).sum().item()
13            test_loss /= num_batches
14            correct /= size
15        print(f"Test Error: \n Accuracy: {(100*correct):>0.1f}%, Avg loss: {test_loss:>8f} \n")
```

Figure 9: Test



## Part B&C

In this part we shall run the following code.

```

1 epochs = 5
2 for t in range(epochs):
3     print(f"Epoch {t+1}\n-----")
4     train(train_dataloader, model, loss_func, optimizer)
5     test(test_dataloader, model, loss_func)
6     print("Done!")
7

```

Figure 10: Test

The obtained results are as follows.

```

1 loss: 1.257729 [ 0/100000]
2 loss: 6.333055 [ 6400/100000]
3 loss: 1.987759 [12800/100000]
4 loss: 1.413734 [19200/100000]
5 loss: 1.384520 [25600/100000]
6 loss: 1.442773 [32000/100000]
7 loss: 1.484911 [38400/100000]
8 loss: 1.264908 [44800/100000]
9 loss: 1.304636 [51200/100000]
10 loss: 1.251667 [57600/100000]
11 loss: 1.091855 [64000/100000]
12 loss: 1.137109 [70400/100000]
13 loss: 1.184907 [76800/100000]
14 loss: 1.285305 [83200/100000]
15 loss: 1.299947 [89600/100000]
16 loss: 1.248004 [96000/100000]
17

```

Figure 11: Test

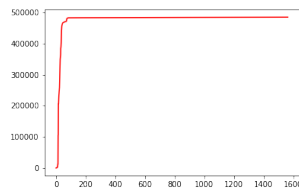


Figure 12: IoU

## Part D

We can deduce that neural network layers are responsible for extracting and recognizing a feature of the input data, for example convolutional layers are responsible for finding feature for error minimization, the first edges attempt to find simple and low level features such as the edges, by deepening our network we can extract and recognize more high level networks that are a bit abstract. Other than this, it is trivial to assume that these layers discard redundant features.

## 2 Question 2 : Object detection and recognition

### 2.1 History

#### 2.1.1 Part 1

Here we shall give some preliminary explanations about the idea behind each algorithm and model.

##### **Hog Algorithm**

The Histogram of oriented gradients algorithm also known as the **HOG** algorithm is a popular feature descriptor widely used in image processing and computer vision.

This algorithm is a simple and powerful feature descriptor employing linear **SVM** machine learning algorithm. By feature descriptors we mean the representation of an image that simply extracts the useful information and disregards the unnecessary information from the image. HOG is not only used for face detection but also it is widely used for object detection like cars, pets, and fruits. **HOG** is very robust for object detection because object shapes are characterized using the local intensity gradient distribution and edge direction.

Now we shall explain the functionality of the algorithm. The algorithm consists of three steps.

- **Step 1**

Dividing the photo into smaller cells

- **Step 2**

Calculation of the histogram for each cell. This histogram is formed by the calculated gradients of each cell. For this approach two kernels are employed to calculate the gradient. Then the magnitude and the direction of the gradients are calculated via the following formulas.

$$g = \sqrt{g_x^2 + g_y^2}, \quad \theta = \arctan\left(\frac{g_y}{g_x}\right)$$

One of the key points in this section is that the magnitude of the gradient increases when sharp changes in intensity occurs, this is handy in figuring out the edges.

-1	0	1	-1
			0
			1

Figure 13: The Kernels

- **Step 3**

After completing step two, we merge all the histograms to form one histogram which is used as a feature vector.

To understand the concept behind **HOG**, we must pay attention to the fact that the role of feature descriptor is used to define an image by its color intensities of pixels. The feature descriptor tries to capture all edges or curves within an image in order to know outline information. Using the first order directive kernels, you can calculate gradient magnitude and then obtain the orientation. Each orientation is assigned to different bins to obtain a histogram. In each bin, gradient magnitudes are stored according to their orientation angles.

Now we shall take a look at the pros and cons of the **HOG** algorithm.

- **Pros**

- Very useful for representing structural objects with low form variation, this means things that have the same common form such as buildings, people etc...
- Powerful descriptor
- Accurate for object classification

**Cons**

- May result in large feature vectors
- Not a very fast descriptor
- Accurate for object classification
- If we have objects with structural variation such as rotation, the algorithm shall be faulty.

**Viola-Jones**

The **Viola-Jones** algorithm, is a framework to perform facial detection. This algorithm is of a cascaded object form, it uses **Haar**-like features and **Adaboost** learning to detect faces in photos. The algorithm starts with a set of simple rectangular features, which are evaluated on a scale-invariant image pyramid to produce a series of feature responses. These responses are then passed through a cascaded set of simple classifiers to quickly reject non-face regions, while retaining face regions for further analysis. The final output is a set of face regions in the input image.

It is important to note that the algorithm works with grayscale images, given an image, the algorithm looks at many smaller subregions and tries to find a face by looking for specific features in each subregion. It needs to check many different positions and scales because an image can contain many faces of various sizes. Viola and Jones used **Haar**-like features to detect faces in this algorithm.

Now we shall take a look at the steps of this algorithm.

- **Selecting Haar-like features**

- **Creating an integral image**
- **Running AdaBoost training**
- **Creating classifier cascades**

### **Haar-Like functions**

**Haar**-like features are digital image features utilized in object recognition. We know that faces share some properties that are the same for everybody, of the human face like the eyes region is darker than its neighbor pixels, and the nose region is brighter than the eye region.

A simple way to find out which region is lighter or darker is to sum up the pixel values of both regions and compare them. The sum of pixel values in the darker region will be smaller than the sum of pixels in the lighter region. If one side is lighter than the other, it may be an edge of an eyebrow or sometimes the middle portion may be shinier than the surrounding boxes, which can be interpreted as a nose. This can be accomplished using **Haar**-like features and with the help of them, we can interpret the different parts of a face.

There are 3 types of Haar-like features that Viola and Jones used in their research.

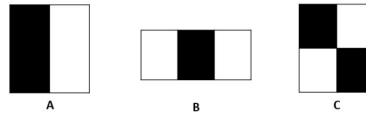
- **Edge features**
- **Four-sided features**
- **Line-features**

Edge features and Line features are useful for detecting edges and lines respectively. The four-sided features are used for finding diagonal features.

The value of the feature is calculated as a single number: the sum of pixel values in the black area minus the sum of pixel values in the white area. The value is zero for a plain surface in which all the pixels have the same value, and thus, provide no useful information.

Since our faces are of complex shapes with darker and brighter spots, a **Haar**-like feature gives you a large number when the areas in the black and white rectangles are very different. Using this value, we get a piece of valid information out of the image.

To be useful, a **Haar**-like feature needs to give you a large number, meaning that the areas in the black and white rectangles are very different.



### Integral Images

Calculations to perform different mathematical computations on the pixels of photos are time consuming and have a great cost, so we must use integral images also known as summed-area tables, this allows us to perform calculations quickly so we can understand whether a feature fits the criteria or not.

### AdaBoost

There are many features used in face identification, but all of them aren't optimal and very useful in performing this task, so we use the **AdaBoost** algorithm which is a machine learning algorithm to extract the best features in facial recognition.

**AdaBoost** checks the performance of all classifiers with the features we support to it, after this it chooses the best feature (type and size of the feature) based on the performance.

Now we shall give some complementary explanation about how we can calculate and evaluate the performance of a classifier. We must evaluate each classifier on the subregions used for training, some shall produce strong responses and some shall not, those that are classified as positive means that it has successfully been classified as a human face, so all in all classifiers that performed well are better and should be given a higher weight, we class this classifier a **boosted classifier**.

All in all, he algorithm is setting a minimum threshold to determine whether something can be classified as a useful feature or not.

### 2.1.2 Part 2

Here we have implemented the **Violo-Jones** algorithm on a photo of [David Beckham](#). The codes and output are as follows.

```
1  def detect_faces(image_path, cascade_Path):
2
3      faceCascade = cv2.CascadeClassifier(cascade_Path)
4
5      frame = cv2.imread(image_path)
6      frame = cv2.cvtColor(frame, cv2.COLOR_BGR2RGB)
7      gray = cv2.cvtColor(frame, cv2.COLOR_BGR2GRAY)
8      faces = faceCascade.detectMultiScale(gray, flags=cv2.CASCADE_SCALE_IMAGE)
9
10     return faces, frame
11
```

Figure 14: Face detection function

```
1  def show_detected_faces(faces, frame):
2      for (x, y, w, h) in faces:
3          cv2.rectangle(frame, (x, y), (x+w, y+h), (255, 0, 0), 3)
4      plt.imshow(frame)
5      plt.show()
6
```

Figure 15: Face showing function

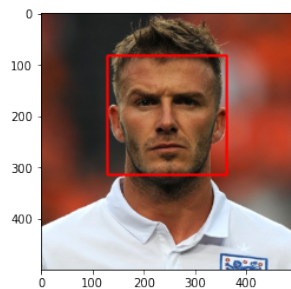


Figure 16: Face detection output

## 2.2 Implementing YOLO

### 2.2.1 Familiarity with the model

The **YOLOv3** model has been trained on the **COCO** dataset which has more than 30000 photos and 80 different classes.

In this project we have used the **YOLOv3-320** model, this model has 320 pixels in width and height.

### 2.2.2 Understanding the theory of the problem

#### Part 2.1

Here we shall firstly explain what a grid cell is in general.

##### Grid Cells

Grid cells have a biologic interpretation which is interesting, it is implied that grid cells are neurons which become active at multiple locations in an environment.

The idea behind **YOLO** is to output a number for the dimension which is large enough for all the objects to fit in it. so the original picture is cropped and divided into an  $n \times n$  grid, each cell in this grid is interpreted as a grid cell.

As implied in the original **YOLO** paper, each grid cell predicts two bounding boxes and one class, this is a limitation of the **YOLO** model.

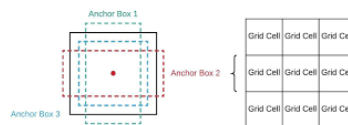


Figure 17: Grid Cell

#### Part 2.2

Here we shall explain the concept of a bounding box, obviously bounding boxes are used in object detection and are generally boxes in the shape of rectangles that surround an object which is to be detected, essentially it is the border's of the



coordinates that enclose the object, these boxes are used to identify the position and type of the object, a bounding box in an image is shown below.

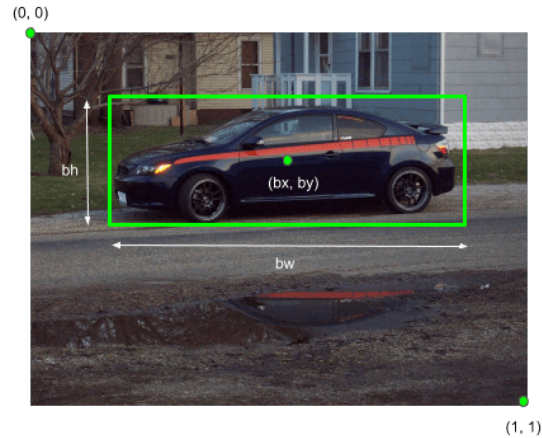


Figure 18: Bounding box

To procure a bounding box we need the coordinates of its edges and the coordinates of the middle is also appreciated. It is also important to imply that the class of the object inside the bounding box is important. The parameters needed to define a bounding box are listed below.

- **Class**
- **Coordinates of edges**
- **Coordinates of center**
- **Height and Width**
- **Confidence**

This represents the possibility of the object being in the box.

### Part 2.3

The tensor output of the YOLO model contains elements described as below:

$(Pc, x, y, h, w, c_1, c_2, c_3, \dots, c_{80})$

- **Pc**:Probability of existence of an object within the cell
- **x**:x-axis center of the bounding box
- **y**:y-axis center of the bounding box
- **h**:height of the bounding box
- **w**:width of the bounding box
- **c<sub>i</sub>**:1 if class i is in the box, 0 if class i is not the box

$$Outputsize = 5 + Num(classes)$$

### Part 2.4

Here we shall explain how the confidence parameter is obtained. First of all it is important to know what the confidence score indicates, this score shows the certainty that the box contains an object and how accurate the box performs its prediction.

We know the the YOLO divides the image into an  $n \times n$  grid, the confidence score shall be calculated using the following formula.

$$Confidence = Pr [Obj] \times IoU$$

Where  $Pr [Obj]$  is the probability of the object existing in the bounding box and **IoU** is the intersection over union between the predicted box and the real thing.

#### **IoU(intersection over union)**

This metric is used to evaluate the performance of object detection, this is done by comparing the ground truth bounding box to the predicted bounding box. A graphical description is shown below.

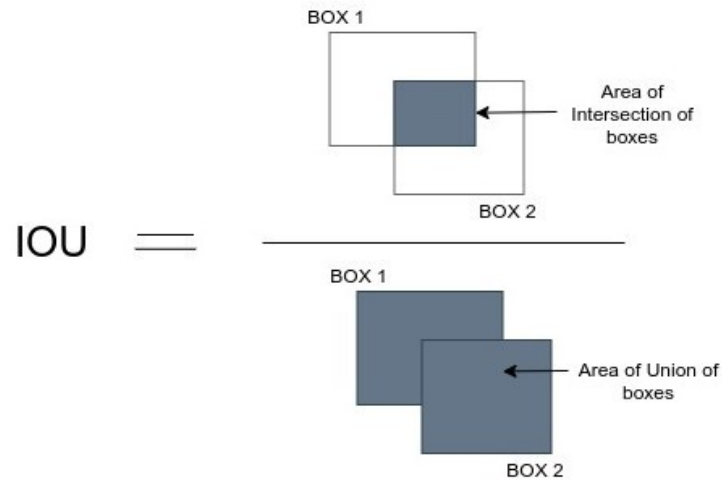


Figure 19: IoU

$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$

### Part 2.5

As pointed out in the **YOLO** paper, the sum-squared error due to ease in optimization, but this doesn't work out as we want it to, we aim to maximize the average precision, many grid cells may not contain objects, so the confidence decreases considerably, hence the model might become unstable. This is remedied by using an appropriate loss function.

The **YOLO** loss function is broken into three parts as follows.

- **Bounding box coordinates**
- **Bounding box score prediction**
- **Class score prediction**

**YOLO** predicts multiple bounding boxes per grid cell. At training time we only want one bounding box predictor to be responsible for each object. We assign

one predictor to be “responsible” for predicting an object based on which prediction has the highest current IOU with the ground truth. This leads to specialization between the bounding box predictors. Each predictor gets better at predicting certain sizes, aspect ratios, or classes of object, improving overall recall.

The loss function is as follows.

$$\begin{aligned}
& \lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{obj} [(x - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2] \\
& + \lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{obj} \left[ \left( \sqrt{\omega_i} - \sqrt{\hat{\omega}_i} \right)^2 + \left( \sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right] \\
& + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{obj} (C_i - \hat{C}_i)^2 + \lambda_{noobj} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{noobj} (C_i - \hat{C}_i)^2 \\
& + \sum_{i=0}^{S^2} \mathbb{1}_i^{obj} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2, \quad \lambda_{coord} = 5, \lambda_{noobj} = 0.5
\end{aligned}$$

In the following notation,  $\mathbb{1}_i^{obj}$  means that if the object appears in cell  $i$  and  $\mathbb{1}_{ij}^{obj}$  means that the  $j$ th bounding box predictor in cell  $i$  is responsible for it.

To show which part of the loss function belongs to which we have :

$$\begin{aligned}
\text{Bbox coords : } & \lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{obj} [(x - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2] \\
& + \lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{obj} \left[ \left( \sqrt{\omega_i} - \sqrt{\hat{\omega}_i} \right)^2 + \left( \sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right] \\
\text{Bbox score : } & \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{obj} (C_i - \hat{C}_i)^2 + \lambda_{noobj} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{noobj} (C_i - \hat{C}_i)^2 \\
\text{Class score : } & \sum_{i=0}^{S^2} \mathbb{1}_i^{obj} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2
\end{aligned}$$

The parameters effective on the loss function are shown above.

### Part 2.6

Here we shall explain the concept of **Non Max Suppression**.

To make it simple **Non Max Suppression** (NMS) is a technique used in computer vision, it is a series of algorithms to select a bounding box between overlapping boxes, this is somewhat a form of whittling down the boxes to a few which are needed for object recognition.

This algorithm considers two things as shown below.

- **IoU**
- **Confidence**

Considering these two parameters the best boxes are chosen.

### Part 2.7

To solve this, we use multiple bounding boxes. First, we predefine different shapes of bounding boxes. So we can have multiple predictions with multiple bounding boxes. Knowing that output of each cell would obey the format of  $(P_c, x, y, h, w, c_1, c_2, c_3)$ , when more than 1 object exists in a grid cell, the output would be the concatenated version of this type of output. For example, in the case of two bounding boxes and three classes the output will be vector of size  $(16, 1)$  and the vector contains,  $(P_{c_1}, x_1, y_1, h_1, w_1, c_1, c_2, c_3, P_{c_2}, x_2, y_2, h_2, w_2, c_1, c_2, c_3)$ .

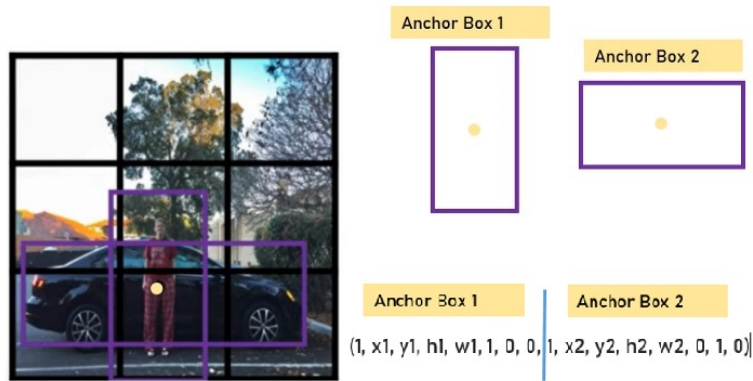


Figure 20: IoU

### 2.2.3 Implementation steps

Here we want to detect three objects, **Car**, **Bus** and **Person** from the wanted video.

#### Part 3.a & 3.b

As explained in part 3.a, to be able to perform object detection in a video, we must split the video into consecutive frames and then perform the desired task on each frame, we do this with the help of **open cv**. We know that in each frame we must extract the predicted outputs, so we design a **find\_objects** function as follows.

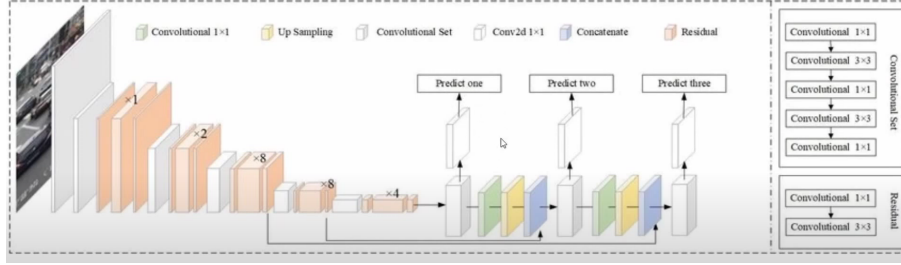


Figure 21: YOLO architecture

As observed in **Figure 16** final layers consist of three output layers. Each layer's output returns a tensor including different detections. Each of these detections is a tensor of size  $1 \times 85$ , below we present a short description of this tensor:

$$d = (P_c, x, y, h, w, c_1, c_2, \dots, c_{80})$$

- **P<sub>c</sub>**:Probability of existence of an object within the cell
- **x**:x-axis center of the bounding box
- **y**:y-axis center of the bounding box
- **h**:height of the bounding box

- **w**:width of the bounding box
- **c\_<sub>i</sub>**:1 if class i is in the box, 0 if class i is not the box

```

1  def find_objects(conf_Thresh, outputs, img):
2      height, width, channel = img.shape
3      bounding_box = []
4      label_Id = []
5      bb_conf = []
6
7      for output in outputs:
8          for detection in output:
9              scores = detection[5:]
10             label = np.argmax(scores)
11             confidence = scores[label]
12             if confidence > conf_Thresh:
13                 w,h = int(detection[2]* width), int(detection[3]*height)
14                 x,y = int((detection[0]*width)-w/2), int((detection[1]*height)-h/2)
15                 bounding_box.append([x,y,w,h])
16                 label_Id.append(label)
17                 bb_conf.append(float(confidence))
18  return bounding_box, label_Id, bb_conf, height, width
19

```

Figure 22: find\_objects

### Part 3.c

Here we design a function to show the detected objects and their labels on the input video, we shall also create a new video featuring labels on cars, people and buses.

```

1  def show_detected_object(labels, desired_labels, bounding_box, label_Id,
2                          bb_conf, conf_Thresh, NMS_Thresh):
3
4      indices = cv2.dnn.NMSBoxes(bounding_box, bb_conf, conf_Thresh, NMS_Thresh)
5      for i in indices:
6          box = bounding_box[i]
7          x,y,w,h = box[0], box[1], box[2], box[3]
8          if labels[label_Id[i]] in desired_labels:
9              cv2.rectangle(img, (x,y), (x+w,y+h), (255,0,255), 2)
10             cv2.putText(img, f'{labels[label_Id[i]].upper()} {int(bb_conf[i]*100)}%',
11                         (x,y-10), cv2.FONT_HERSHEY_SIMPLEX, 0.6, (255,0,255), 2)
12  return img

```

Figure 23: show\_detected\_object

```
1 def display_video(video):
2     fig = plt.figure(figsize=(10,10))
3
4     mov = []
5     for i in range(len(video)):
6         img = plt.imshow(video[i], animated=True)
7         plt.axis('off')
8         mov.append([img])
9
10    anime = animation.ArtistAnimation(fig, mov, interval=50, repeat_delay=1000)
11
12    plt.close()
13    return anime
14
```

Figure 24: display\_video

```
1 cap = cv2.VideoCapture(video_dir)
2 yolo_size = 320
3 conf_Thresh = 0.5
4 NMS_Thresh = 0.3
5 video = []
6
7 width = int(cap.get(3))
8 height = int(cap.get(4))
9 fps = cap.get(cv2.CAP_PROP_FPS)
10 out = cv2.VideoWriter('output.mp4', cv2.VideoWriter_fourcc('M', 'J', 'P', 'G'), fps, (width, height))
11
12 while True:
13     success, img = cap.read()
14     if success == False:
15         break
16     blob = cv2.dnn.blobFromImage(img, 1/255, (yolo_size, yolo_size), [0,0,0], crop=False)
17     model.setInput(blob)
18     layerNames = model.getLayerNames()
19     outputNames = [layerNames[i]-1 for i in model.getUnconnectedOutLayers()]
20     outputs = model.forward(outputNames)
21     bounding_box, label_Id, bb_conf, height, width = find_objects(conf_Thresh, outputs, img)
22     image = show_detected_object(labels, desired_labels, bounding_box, label_Id, bb_conf, conf_Thresh,
23                                 NMS_Thresh)
24     out.write(image)
25     video.append(image)
26
27 cap.release()
28 out.release()
29 cv2.destroyAllWindows()

```

Figure 25: Video Exporter



## References

- [1] [Reshad Hosseini](#), *Intelligent Systems Lecture Notes, Fall 01*
- [2] [Hands-On Convolutional Neural Networks with TensorFlow](#)
- [3] [Intersection over Union \(IoU\) for object detection](#)
- [4] [A Comprehensive Guide To Object Detection Using YOLO Framework](#)
- [5] [Learn about bounding boxes in image processing](#)
- [6] [A Gentle Introduction Into The Histogram Of Oriented Gradients](#)
- [7] [The Viola-Jones Face Detection Algorithm](#)
- [8] [Get middle layer in pytorch](#)