

CSC420

Final Project

Xiao Xuan Wang (wangx576) 1003422124

Benli Wang(wangbenl) 1004156637

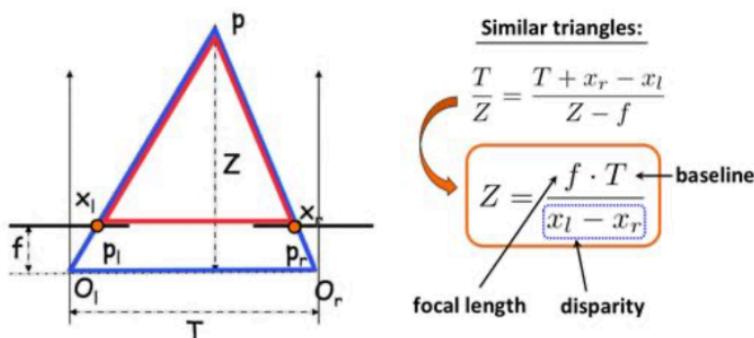
Road and object detection

Question 1

Compute disparity between the two stereo images.

We used two method for this question. Patch matching and SGBM.

For patch matching. Since we have two parallel cameras (call result image for left camera to be left image, result for right camera to be right image), thus in order to find the disparity of a patch on left image, we can check the corresponding rows of patch and find a best match in that row. After finding a best match, we can calculate the disparity and depth by compute



And the code I wrote to calculate the disparity is

```
def detect_match(img, template):
    """
    with help of documentation of matchTemplate in opencv
    """
    res = cv.matchTemplate(img,template, cv.TM_CCORR_NORMED)
    min_val, max_val, min_loc, max_loc = cv.minMaxLoc(res)
    top_left = max_loc[0] + (template.shape[1])/2
    return top_left

def compute_disparity(left_img, right_img, patch_size = 5):
    radius = (patch_size-1)/2
    result = np.zeros(left_img.shape[:2])
    row,col,color = left_img.shape
    for cr in range(row):
        for cc in range(col):
            lr = int(max(0,cr-radius))
            rr = int(min(row,cr+radius+1))
            uc = int(max(0,cc-radius))
            dc = int(min(col,cc+radius+1))
            temp = left_img[lr:rr, uc:dc,:]
            xr = detect_match(right_img[int(max(0,lr-radius)):int(min(row,rr+radius)), :, :], temp)
            result[cr,cc] = cc-xr
    print(result.shape)
    return result
```

```

lpath = 'sync/image02/data/'
rpath = 'sync/image03/data/'

for i in range(107):
    imgnl = 'sync/image_02/data/{}.png'.format(str(i).zfill(10))
    imgnr = 'sync/image_03/data/{}.png'.format(str(i).zfill(10))
    left_img = cv.imread(imgnl)
    right_img = cv.imread(imgnr)
    disparity = compute_disparity(left_img, right_img, patch_size = 5)
    #disparity = gaussian_filter(disparity, sigma=1)
    disparity = (disparity // 10)*10
    plt.imsave(str(i).zfill(6)+".png", disparity, cmap='gray')

```

But however, the result is very noisy and inaccurate.

When left image is



I got result like



The second method we used is SGBM (Stereo Processing by Semiglobal Matching and Mutual Information) algorithm, the basic idea of this algorithm is making a disparity map by select the disparity value for each pixel. Then Set a energy function as following. Then try to minimize the E(D).

$$E(D) = \sum_p \left(C(p, D_p) + \sum_{q \in N_p} P_1 I[|D_p - D_q| = 1] + \sum_{q \in N_p} P_2 I[|D_p - D_q| > 1] \right)$$

p and q are the pixels in the image, N_p stands for 8 pixels around p, $C(p, D_p)$ is the cost when the disparity value for p is D_p , P_1 is the penalty coefficient for pixels around p with difference of disparity value less and equal to 1, P_2 is the penalty coefficient for pixels around p with difference of disparity value higher than 1, $I[]$ will return 1 if and only if the statement inside is True. To optimize the solution is a NP-complete problem. However, it can be approximately solved by separate to many one-dimension problems and use dynamic programming to solve them. Thus, there are 8 one dimensional problems to solve.

In our project, we used the opencv package to implement the SGBM and generate the more accurate disparity images by following codes. The output will have a black area on the left side and we filled up with the most left non-zero value column.

```
# SGBM Parameters -----
window_size = 3
# wsize default 3; 5; 7 for SGBM reduced size image; 15 for SGBM full size image (1300px and above); 5 Works nicely

left_matcher = cv2.StereoSGBM_create(
    minDisparity=0,
    numDisparities=112,
    blockSize=5,
    P1=8 * 3 * window_size ** 2,
    P2=32 * 3 * window_size ** 2,
    disp12MaxDiff=1,
    uniquenessRatio=15,
    speckleWindowSize=0,
    speckleRange=2,
    preFilterCap=63,
    mode=cv2.STEREO_SGBM_MODE_SGBM_3WAY
)

right_matcher = cv2.ximgproc.createRightMatcher(left_matcher)

# FILTER Parameters
lmbda = 80000
sigma = 1.2

wls_filter = cv2.ximgproc.createDisparityWLSFilter(matcher_left=left_matcher)
wls_filter.setLambda(lmbda)
wls_filter.setSigmaColor(sigma)

print('computing disparity...', index)
displ = left_matcher.compute(imgL, imgR) # .astype(np.float32)/16
dispr = right_matcher.compute(imgR, imgL) # .astype(np.float32)/16
displ = np.int16(displ)
dispr = np.int16(dispr)
filteredImg = wls_filter.filter(displ, imgL, None, dispr) # important to put "imgL" here!!!

filteredImg = cv2.normalize(src=filteredImg, dst=filteredImg, beta=0, alpha=255, norm_type=cv2.NORM_MINMAX)
filteredImg = np.uint8(filteredImg)
cv2.imshow('Disparity Map', filteredImg)
filteredImg[:,112] = np.tile(filteredImg[:,112], (112,1)).T
filteredImg = filteredImg / 2
cv2.imwrite('./data_road/disparity/depth_{}.png'.format(str(index).zfill(6)), filteredImg)
```

Here are some examples we get

Original image:



Disparity map:



Original image:



Disparity map:



Question 2

Compute depth of each pixel. Compute 3D location of each pixel.

Using method mentioned in patch match method in previous question, I can get the depth of a pixel using disparity by $depth = f * T / disparity$ where f is focal length, T is baseline.

As mentioned in set up, $T = 0.54$, and in calib_cam_to_cam.txt, the focal length f is 959.791.

```
15
16     def compute_depth(disparity, f, T):
17         disparity = disparity
18         disparity = cv.GaussianBlur(disparity, (21, 21), 1)
19         disparity = np.maximum(1, disparity)
20         print(disparity)
21         result = f * T / disparity
22         return result
```

To calculate 3D location, we need the information from calibration matrix of camera also from calib_cam_to_cam.txt. $px = 696.0217$ and $py = 224.1806$.

```
27
28     def calculate_3d(img, depth_map):
29         scale_y = py / (depth_map.shape[0] / 2)
30         scale_x = px / (depth_map.shape[1] / 2)
31         Z = depth_map
32         X = Z * (img[:, :, 1] * scale_x - px) / f
33         Y = Z * ((depth_map.shape[0] - img[:, :, 0]) * scale_y - py) / f
34         return np.stack((X, Y, Z), axis=-1)
35
```

```
depth_list = []
three_d_list = []
r, c = cv.imread('regular_disparity/' + str(0).zfill(6) + ".png", cv.IMREAD_GRAYSCALE).shape
img = np.zeros((r, c, 2))
for i in range(r):
    for j in range(c):
        img[i][j] = [i, j]
for i in range(100):
    disp_name = 'regular_disparity/' + str(i).zfill(6) + ".png"
    disparity = cv.imread(disp_name, cv.IMREAD_GRAYSCALE)
    depth = compute_depth(disparity, f, T)
    three_d = calculate_3d(img, depth)
    depth_list.append(depth)
    three_d_list.append(three_d)
```

Since from previous question, the result of disparity using SGBM is much better than patch matching, so in this question and later question we will use SGBM to generate disparity.

Question 3

**Train a road classifier on a set of annotated images,
and compute road pixels in your**

image. Which features would you use? Try to use both 2D and 3D features.

We tried two different machine learning approaches to solve this problem, SVM (support vector machine) and neural network. But, for both methods, we feed with same features. Here are nine features we give:

Red, Green, Blue value, Hue, Saturation and Value for the image, using the disparity map to calculate the 3D location (X, Y, Z) of each pixel.

Using RGB value because we want to distinguish the road and other object by color, it is a intuitive idea since human recognize the object by color first. Then applying HSV value will also help us to differentiate object in the image. The feature of 3D location is also added because we know in the ideal case, the y value for road will always be consistent. Besides, for all the training image we get, the road is in the center of the image, thus, 3D location should be an efficacious feature.

Here are our codes and some examples.

SVM

```

def train_detector():
    print(validate[0], validate[0][0])
    # train_image, validate_image = train[0], validate[0]
    img_lst = []
    val_lst = []
    depth_lst = []
    for i in range(60):
        im_1 = cv.imread("data_road/data_left/training/image_2/im_{:}.png".format(str(i).zfill(6)))
        im_2 = cv.imread("data_road/data_right_left/training/pt_image_2/im_{:}.png".format(str(i).zfill(6)))
        im_3 = cv.imread("data_road/disparity/depth_{:}.png".format(str(i).zfill(6)))

        img_lst.append(cv.blur(im_1, (5,5)))
        val_lst.append(cv.blur(im_2, (5,5)))
        depth_lst.append(cv.blur(im_3, (5,5)))
    train_image = img_lst[0]

    validate_image = val_lst[0]
    depth = depth_lst[0]

    M, N, _ = train_image.shape
    m, n, _ = validate_image.shape
    locs = get_3d_location(depth)

    bayes = np.zeros((M,N,1))
    for val_im in val_lst:
        bayes = bayes + val_im[:, :, 0].reshape(M,N,1) / 255

    hsv = cv.cvtColor(train_image, cv.COLOR_BGR2HSV)
    hsv = np.insert(hsv, 3, train_image[:, :, 0], axis=2)
    hsv = np.insert(hsv, 4, train_image[:, :, 1], axis=2)
    hsv = np.insert(hsv, 5, train_image[:, :, 2], axis=2)

    X = np.reshape(hsv, (M * N, 6))
    X = X / 255

    X = np.reshape(validate_image, (M * N, 3))

    y = np.zeros([0])
    y[y > 125] = 255
    y[y <= 125] = 0

```

```

for i in range(1, 80):
    train_image, validate_image, depth = img_lst[i], val_lst[i], depth_lst[i]
M, N, _ = train_image.shape
m, n, _ = validate_image.shape
hsv = cv.cvtColor(train_image, cv.COLOR_BGR2HSV)
hsv = np.insert(hsv, 3, train_image[:, :, 0], axis=2)
hsv = np.insert(hsv, 4, train_image[:, :, 1], axis=2)
hsv = np.insert(hsv, 5, train_image[:, :, 2], axis=2)
new_x = np.reshape(hsv, (M * N, 6))
new_x = new_x / 255

locs = get_3d_location(depth)

# normalize the 3d location
loc_flat1 = locs[:, 0].flatten()
loc_flat1.sort()
max_loc_1 = loc_flat1[int(0.9 * loc_flat1.shape[0])]
loc_flat2 = locs[:, 1].flatten()
loc_flat2.sort()
max_loc_2 = loc_flat2[int(0.9 * loc_flat2.shape[0])]
loc_flat3 = locs[:, 2].flatten()
loc_flat3.sort()
max_loc_3 = loc_flat3[int(0.9 * loc_flat3.shape[0])]

new_x = np.insert(new_x, 6, locs[:, 0] / max_loc_1, axis=1)
new_x = np.insert(new_x, 7, locs[:, 1] / max_loc_2, axis=1)
new_x = np.insert(new_x, 8, locs[:, 2] / max_loc_3, axis=1)
new_y = np.reshape(validate_image, (m * n, 3))
new_y = np.reshape(validate_image, (m * n, 3))
new_y[new_y > 125] = 0
new_y[new_y <= 125] = 1
X = np.vstack((X, new_x))
y = np.concatenate((y, new_y))
print(i)
return X, y

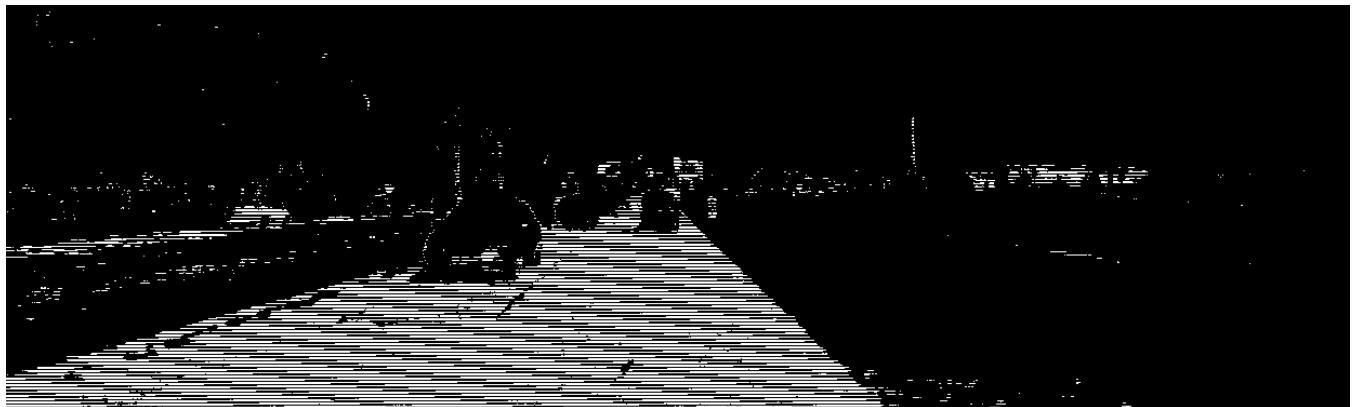
```

```

if __name__ == "__main__":
    X, y = train_detector()
    clf = LinearSVC(random_state=0, tol=1e-05)
    print('start')
    clf.fit(X, y)
    pkl_filename = "road_model.pkl"
    with open(pkl_filename, 'wb') as file:
        pickle.dump(clf, file)

```

Example



The output can recognize the road location but is not perfect. Thus, I applied some filter on the output image to get a more accurate result.

```
for i in range(3):
    image = cv.blur(image, (3,3))
    image[image[:, :, 0] > 160] = [255, 255, 255]
    image[image[:, :, 0] < 40] = [0, 0, 0]
    kernel = cv.getStructuringElement(cv.MORPH_ELLIPSE, (5, 5))
    image = cv.morphologyEx(image, cv.MORPH_OPEN, kernel)

image[image[:, :, 0] > 0] = [255, 255, 255]
return image
```

Then the new output is



Moreover, the pros by choosing SVM is it can solve high dimension problems well and get the relation of the data and features. However, the training takes longer time and the worst case will take too long if the data set is large.

Neural Network

```
if __name__ == "__main__":
    BATCH_SIZE = 10
    SHUFFLE_BUFFER_SIZE = 100
    print("Num GPUs Available: ", len(tf.config.experimental.list_physical_devices('GPU')))

    with tf.device('/gpu:1'):

        x_train, y_train = train_detector()
        test_x, test_y = get_test()
        train_dataset = tf.data.Dataset.from_tensor_slices((x_train, y_train))
        test_dataset = tf.data.Dataset.from_tensor_slices((test_x, test_y))
        # Building a 3-layer DNN with 50 units each.
        model = tf.keras.Sequential([
            tf.keras.layers.Flatten(),
            tf.keras.layers.conv2d(),
            tf.keras.layers.Dense(1024, activation='relu'),
            tf.keras.layers.Dense(108, activation='relu'),
            tf.keras.layers.Dense(2, activation='softmax')
        ])

        train_dataset = train_dataset.shuffle(SHUFFLE_BUFFER_SIZE).batch(BATCH_SIZE)
        test_dataset = test_dataset.batch(BATCH_SIZE)
        # Use the train data to train this classifier

        model.compile(optimizer=tf.keras.optimizers.RMSprop(),
                      loss=tf.keras.losses.SparseCategoricalCrossentropy(),
                      metrics=[tf.keras.metrics.SparseCategoricalAccuracy()])
        model.fit(train_dataset, epochs=2, verbose=1)
        print('testing')
        model.evaluate(test_dataset)
        model.save('saved_model', save_format='tf')
```



The labels on the road and the lines are not recognize as road might because when we trained model, the most road don't have the white color and make the model to mis predict the road

Question 4

Fit a plane in 3D to the road pixels by using the depth of the pixels. Make sure your algorithm is robust to outliers

By using the pixel labeled as rode and their 3D location calculated in question 2, I can perform RANSAC to find a optimal plane. Since we know three 3d point can find a unique plane, so I can randomly pick three point in the road, and fit a plane using three point, calculate the difference between other road pixel and the plane, count the outliers(pixel that has a difference greater than threshold), and find the optimal plane that has the smallest outliers.

```
def calculate_error(three_d, road, a, b, c, d, threshold=0.5):
    # get 3d point for road pixel
    three_d_road = three_d[road[:,0], road[:,1]]
    predict = (d - a*three_d_road[:,0]-b*three_d_road[:,1])/c
    diff = np.abs(predict - three_d_road[:,2])
    return (diff>1).sum()
```

```

label_name = 'data_road/training/gt_image_2/um_lane_000000.png'
label = cv.imread(label_name)
label = np.sum(label, axis = 2)/255
road = np.asarray(np.where(label == 2))
road = road.T
three_d = three_d_list[0]
min_outlier = 1000000
optimal = None
for i in range(1000):
    pick3 = np.random.choice(len(road), 3)
    p1 = road[pick3[0]]
    p2 = road[pick3[1]]
    p3 = road[pick3[2]]
    p1 = three_d[p1[0],p1[1]]
    p2 = three_d[p2[0],p2[1]]
    p3 = three_d[p3[0],p3[1]]
    v1 = p3 - p1
    v2 = p2 - p1
    cp = np.cross(v1, v2)
    a, b, c = cp
    v2 = p2 - p1
    d = np.dot(cp, p3)
    outlier = calculate_error(three_d, road, a, b, c, d)
    if outlier < min_outlier:
        min_outlier = outlier
        optimal = [a,b,c,d]
print(p1,p2,p3,min_outlier,optimal)

```

For this iteration, result is $0.2534022227414456x - 33.70033785527126y + 0.7452026089080129z = 42.45503035930558$

Set $x, z = 0$, we have $y = -1.26$, close to the true value since the camera is 1.65m high.

However since ransac is randomly choosing point, so there is a chance that it fits the point with inaccurate 3d location or that is false positively classified as a road. So we choose the one with less outliers(the vertical distance between the point and the plane is larger than 0.5) and apply a bigger iteration times.

Another way to do this is using MSE but since this is not robust to the points that has been classify as road but it is not road, and assume the calculation for 3d location is correct, then because of those point, the result plane will be higher than what we expected. So I choose to use RANSAC.

Question 5

**Plot each pixel in 3D (we call this a 3D point cloud).
On the same plot, show also the estimated ground plane.**

We use open3d to plot. With help of 3d location calculated by previous question, and the pixel color is it's own image, normalized by dividing 255.

For showing the plane, I generate 100000 value for X and Z(X is -150 to 150, Z is 0-300), and use the plane to fit Y. plot XYZ using red color on the same graph with the image.

Code:

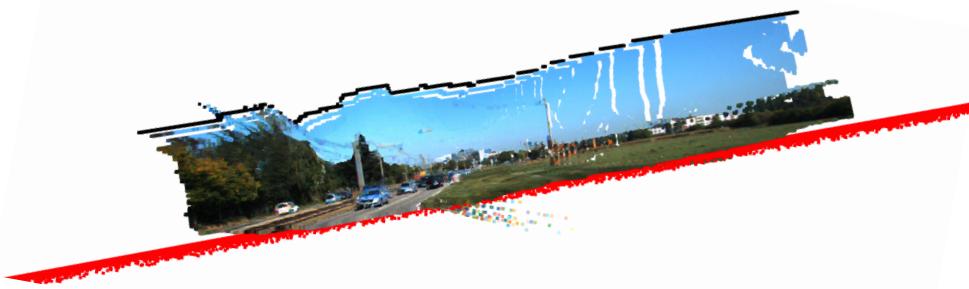
```
106 x = ((np.random.rand(100000))-0.5)*300
107 z = (np.random.rand(100000))*300
108 y = (optimal[3] - x*optimal[0] - z*optimal[2])/optimal[1]
109
110 pxyz = np.stack((x, y, z), axis=-1)
111 pxyzc1 = np.zeros((100000))
112 pxyzc2 = np.zeros((100000))
113 pxyzc3 = np.ones((100000))
114 pxyzc= np.stack((pxyzc1, pxyzc2, pxyzc3), axis=-1)
115
116 xyz = three_d.reshape((-1, 3))
117
118 xyz = np.vstack((xyz, pxyz))
119
120 pcd = o3d.geometry.PointCloud()
121 pcd.points = o3d.utility.Vector3dVector(xyz)
122 imgnl = './dis/road_038/um_00038.png'
123 left_img = cv.imread(imgnl).reshape((-1, 3))
124 C = left_img / 255
125 C = np.vstack((C, pxyzc))
126 C = C[:, ::-1]
127 pcd.colors = o3d.utility.Vector3dVector(C)
128 o3d.visualization.draw_geometries([pcd])
```

Result:

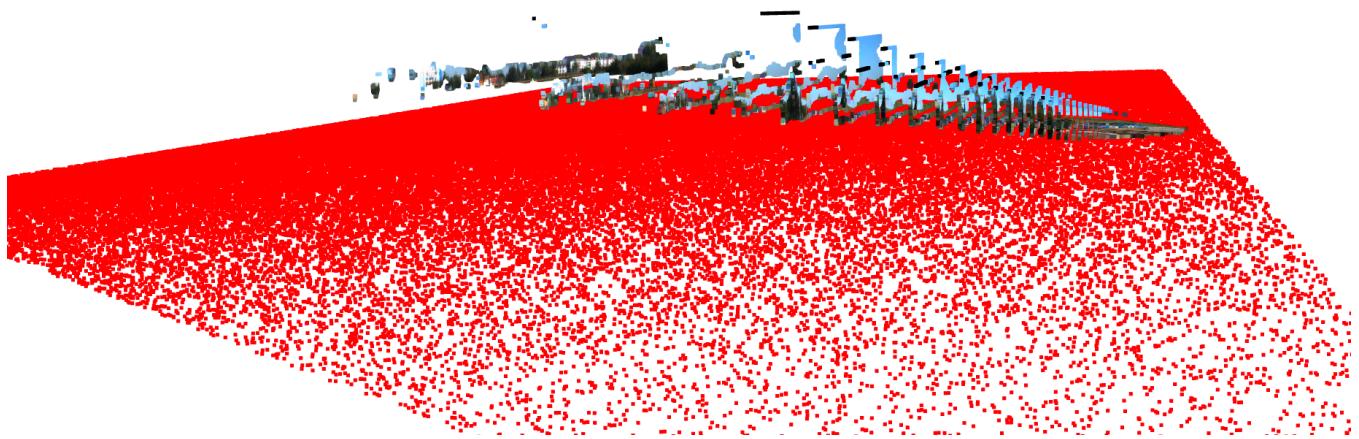
Input:



Capture from front:



Capture from side:



input:



Result:



Question 6

Detect cars in the image

To detect car in the image. I used YOLO v3. Which give the image as input then apply the model, then return the detected object's coordinate.

In my implementation, I will reshape all my input image to be size 416*416*3, Then apply multiple convolution.

Using the last layer's information, We can get the position of bounding box using logistic regression with non_max_suppression. Then I write the x,y,w,h in to a mat file for later reference and then draw a box around bounding box.

```
1  import cv2 as cv
2  import scipy.io as sio
3  import numpy as np
4
5  # Initialize the parameters
6  confThreshold = 0.5 #Confidence threshold
7  nmsThreshold = 0.4 #Non-maximum suppression threshold
8  inpWidth = 416 #Width of network's input image
9  inpHeight = 416 #Height of network's input image
10
11
12 # Load names of classes
13 classesFile = "coco.names"
14 classes = None
15 with open(classesFile, 'rt') as f:
16     classes = f.read().rstrip('\n').split('\n')
17
18 # Give the configuration and weight files for the model and load the network using them.
19 modelConfiguration = "yolov3.cfg"
20 modelWeights = "yolov3.weights"
21 frame_name = 'um_000087.png'
22 frame = cv.imread(frame_name)
```

```

26     def postprocess(frame, outs):
27         frameHeight = frame.shape[0]
28         frameWidth = frame.shape[1]
29
30         # Scan through all the bounding boxes output from the network and keep only the
31         # ones with high confidence scores. Assign the box's class label as the class with
32         # the highest score.
33         classIds = []
34         confidences = []
35         boxes = []
36         for out in outs:
37             for detection in out:
38                 scores = detection[5:]
39                 classId = np.argmax(scores)
40                 confidence = scores[classId]
41                 if confidence > confThreshold:
42                     center_x = int(detection[0] * frameWidth)
43                     center_y = int(detection[1] * frameHeight)
44                     width = int(detection[2] * frameWidth)
45                     height = int(detection[3] * frameHeight)
46                     left = int(center_x - width / 2)
47                     top = int(center_y - height / 2)
48                     classIds.append(classId)
49                     confidences.append(float(confidence))
50                     boxes.append([left, top, width, height])
51
52         # Perform non maximum suppression to eliminate redundant overlapping boxes with
53         # lower confidences.
54         indices = cv.dnn.NMSBoxes(boxes, confidences, confThreshold, nmsThreshold)
55         result_box = []
56         for i in indices:
57             i = i[0]
58             box = boxes[i]
59             left = box[0]
60             top = box[1]
61             width = box[2]
62             height = box[3]
63             cv.rectangle(frame, (left, top), (left + width, top + height), (0, 0, 255))
64             result_box.append(box)
65         cv.imwrite("bbbb.png", frame)
66         return result_box

```

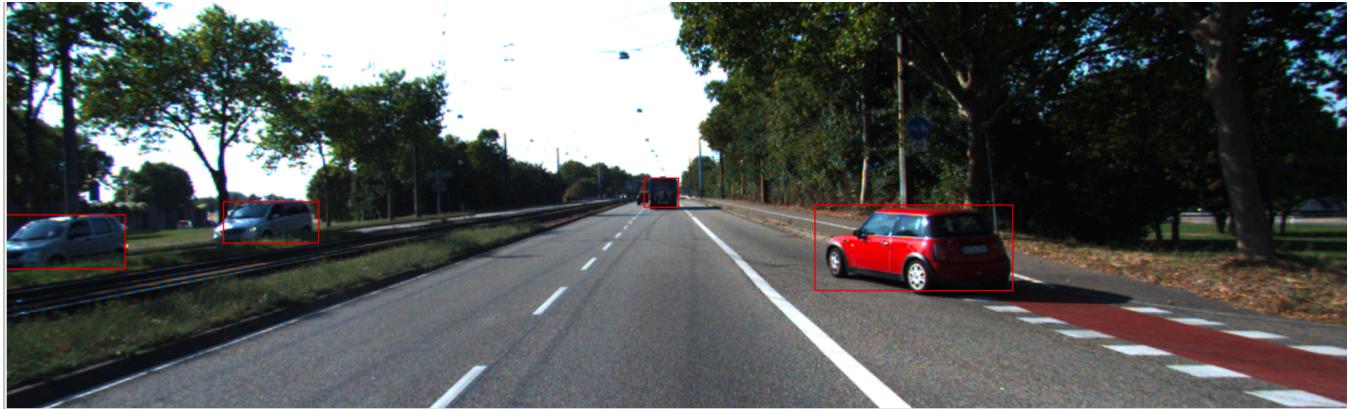
```

67 net = cv.dnn.readNetFromDarknet(modelConfiguration, modelWeights)
68 net.setPreferableBackend(cv.dnn.DNN_BACKEND_OPENCV)
69 net.setPreferableTarget(cv.dnn.DNN_TARGET_CPU)
70
71 blob = cv.dnn.blobFromImage(frame, 1 / 255, (inpWidth, inpHeight), [0, 0, 0], 1,
72                             crop=False)
73 net.setInput(blob)
74
75 # Runs the forward pass to get output of the output layers
76 layersNames = net.getLayerNames()
77 out_name = [layersNames[i[0] - 1] for i in net.getUnconnectedOutLayers()]
78 outs = net.forward(out_name)
79
80 # Remove the bounding boxes with low confidence
81 output = postprocess(frame, outs)
82 file_name = 'r{}.mat'.format(frame_name[:-4].zfill(6))
83 sio.savemat(file_name, {'ds': []})
84 for (left, top, width, height) in output:
85     right = left + width
86     bottom = top + height
87     cur = sio.loadmat(file_name)['ds']
88     if len(cur) == 0:
89         sio.savemat(file_name,
90                     {'ds': [left, right, top, bottom]})
91     else:
92         cur = np.append(cur, [[left, right, top, bottom]], axis=0)
93         sio.savemat(file_name, {'ds': cur})

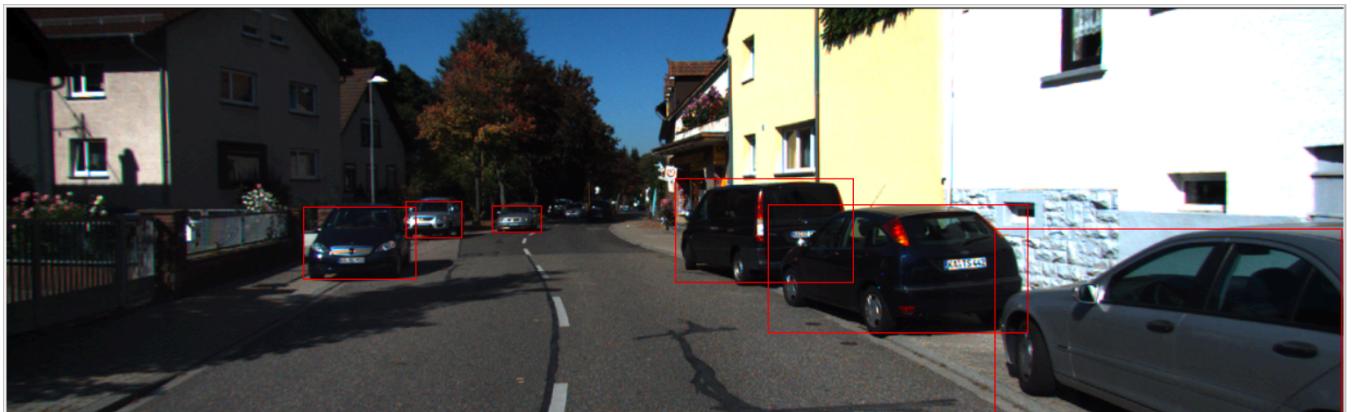
```

Result image(good)





Bad result (only few car detected, but smaller car near vanishing point is not detected)



However with YOLO, the accuracy rate is not as good as model like faster RCNN, especially with smaller item, but the speed is very fast. But since we solving problem for automatic driving, smaller car(which in reality is usually far car) is not very important, but the speed for car detection is very important, so I chose to use YOLO.

Question 7

Train a classifier that predicts viewpoint for each car. The viewpoint labels are in 30° increments, thus train a 12-class classifier.

In this question, to train the classifier, I first used the bounding box data from training set to cut the car segment out and put them in their corresponding folder.

```
7 import numpy as np
8 import os
9 import pandas
10 from os import listdir
11 from os.path import isfile, join
12 from scipy import misc
```

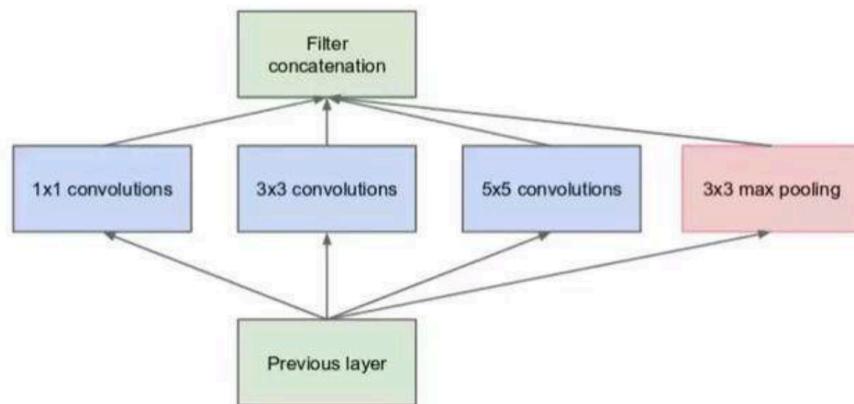
```

16 ▶ if __name__ == "__main__":
17     # PATHS
18     DETECTIONS_PATH = './training/label_2'
19     OUTPUT_PATH = './result'
20     IMAGES_PATH = './data_object/training/image_2'
21
22     # Get filenames
23     detection_files = [DETECTIONS_PATH + os.sep + f for f in
24         | listdir(DETECTIONS_PATH) if isfile(join(DETECTIONS_PATH, f))]
25
26     ▶ image_files = [IMAGES_PATH + os.sep + f for f in listdir(IMAGES_PATH) if
27         | isfile(join(IMAGES_PATH, f))]
28
29     image_files.sort()
30     detection_files.sort()
31
32     counter = 0
33     bins = np.arange(-180, 181, 30)
34     for i, f in zip(image_files, detection_files):
35         img = misc.imread(i)
36         datafram = pandas.read_csv(f, header=0, sep=' ')
37         dataset = datafram.values
38         cars = []
39         for d in dataset:
40             # If the object is a car.
41             if d[0] == 'Car':
42                 # Segment it.
43                 angle = float(d[14])
44                 bin = np.digitize(np.degrees(angle), bins)
45                 segmented = img[int(float(d[5])):int(float(d[7])),
46                     | int(float(d[4])):int(float(d[6])), :]
47                 segmented = misc.imresize(segmented, (128, 128), interp='bilinear',
48                     mode=None)
49                 os.makedirs(OUTPUT_PATH + os.sep + str(bin) + os.sep, exist_ok=True)
50                 misc.imsave(OUTPUT_PATH + os.sep + str(bin) + os.sep + str(
51                     | counter) + '.png', segmented)
52                 counter += 1

```

Then using the separated car segments, I used keras's InceptionV3 to train the model.

InceptionV3 is still a model that using convolution. The basic idea for Inception is



And This is very computational costly. In order to decrease running time, we choose to add some 1*1 convolution layer to reduce dimension. Also separate 5*5 convolution matrix to two 3*3 convolution matrix to decrease running time, so intotal we have

type	patch size/stride or remarks	input size
conv	3x3/2	299x299x3
conv	3x3/1	149x149x32
conv padded	3x3/1	147x147x32
pool	3x3/2	147x147x64
conv	3x3/1	73x73x64
conv	3x3/2	71x71x80
conv	3x3/1	35x35x192
3xInception	As in figure 5	35x35x288
5xInception	As in figure 6	17x17x768
2xInception	As in figure 7	8x8x1280
pool	8 x 8	8 x 8 x 2048
linear	logits	1 x 1 x 2048
softmax	classifier	1 x 1 x 1000

Table 1. The outline of the proposed network architecture. The output size of each module is the input size of the next one. We are using variations of reduction technique depicted Figure 10 to reduce the grid sizes between the Inception blocks whenever applicable. We have marked the convolution with 0-padding, which is used to maintain the grid size. 0-padding is also used inside those Inception modules that do not reduce the grid size. All other layers do not use padding. The various filter bank sizes are chosen to observe principle 4 from Section 2.

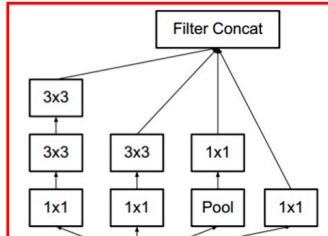


Figure 5. Inception modules where each 5×5 convolution is replaced by two 3×3 convolution, as suggested by principle 3 of Section 2.

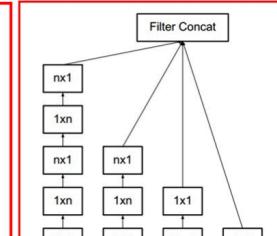


Figure 6. Inception modules after the factorization of the $n \times n$ convolutions. In our proposed architecture, we chose $n = 7$ for the 17×17 grid. (The filter sizes are picked using principle 3).

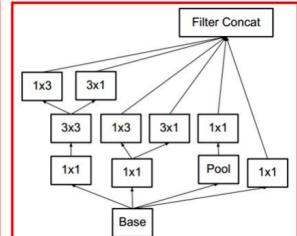


Figure 7. Inception modules with expanded the filter bank outputs. This architecture is used on the coarsest (8×8) grids to promote high dimensional representations, as suggested by principle 2 of Section 2. We are using this solution only on the coarsest grid, since that is the place where producing high dimensional sparse representation is the most critical as the ratio of local processing (by 1×1 convolutions) is increased compared to the spatial aggregation.

<http://blog.csdn.net/loveliuzz>

And also in InceptionV3, it factorized 7×7 convolution layer and add BatchNorm to support classification.

In my implementation, I put the entire segment in and resize it to $128 \times 128 \times 3$ and divide by 255 to normalize it to be within 0-1.

```

12 train_datagen = ImageDataGenerator(
13     rescale=1. / 255,
14     shear_range=0.2,
15     zoom_range=0.2,
16     horizontal_flip=False)
17
18 test_datagen = ImageDataGenerator(rescale=1. / 255)
19
20
21 train_generator = train_datagen.flow_from_directory(
22     './result', # this is the target directory
23     target_size=(128, 128), # all images will be resized to 128x128
24     batch_size=32,
25     class_mode='categorical')
26
27 input_tensor = Input(shape=(128, 128, 3))

```

Then add a HOG feature to the data set. Since the change direction for a patch is very important for viewpoint classification.

```

32     data_list = []
33     batch_index = 0
34     y = []
35
36     while batch_index <= train_generator.batch_index:
37         data = train_generator.next()
38         data_list.extend(data[0])
39         y.extend(data[1])
40         batch_index = batch_index + 1
41
42     data_array = np.asarray(data_list)
43     y_array = np.asarray(y)
44     his = np.empty((0, 128, 128, 1))
45
46     for i in range(data_array.shape[0]):
47         fd, hog_image = hog(data_array[i, :, :, :], orientations=8, pixels_per_cell=(16, 16),
48                               cells_per_block=(1, 1), visualize=True,
49                               multichannel=True)
50         hog_image = np.expand_dims(hog_image, axis=2)
51         hog_image = np.expand_dims(hog_image, axis=0)
52         his = np.vstack((his, hog_image))
53
54
55     data = np.append(data_array, his, 3)
56
57     # We use the InceptionV3 /GoogLeNet model but retrain it to classify our dataset.
58     base_model = InceptionV3(input_tensor=input_tensor, weights='imagenet', include_top=False)
59     # Add a global spatial average pooling layer.
60     x = base_model.output
61     x = GlobalAveragePooling2D()(x)
62     # Add a fully-connected layer.
63     x = Dense(1024, activation='relu')(x)
64     # Add a HEAVY dropout.
65     x = Dropout(0.7)(x)
66     predictions = Dense(12, activation='softmax')(x)
67
68     model = Model(input=base_model.input, output=predictions)

```

The base mode doesn't have top layer when build it, but we add a Rectified Linear Unit on top. Also I add a average pooling layer before that and I use Dropout to avoid overfitting.

```

69
70     for layer in base_model.layers:
71         layer.trainable = False
72
73     # compile the model (should be done *after* setting layers to non-trainable)
74     model.compile(optimizer='rmsprop', loss='categorical_crossentropy', metrics=['accuracy'])
75
76     model.fit(data, y_array, epochs=10)
77

```

First we froze all layers for base model to capture correct bottleneck,

```

82
83     for layer in model.layers[:172]:
84         layer.trainable = False
85     for layer in model.layers[172:]:
86         layer.trainable = True
87
88     model.compile(optimizer=SGD(lr=0.01, decay=1e-6, momentum=0.9))
89     model.fit(data, y_array, epochs=15)
90
91     model.save('object_classifierfit2.h5') # Save the final model.
92

```

Then we turn on trainable for some inception blocks, refit it and save the model.

The accuracy rate is around 0.88 after training.

Then we can use the trained classification model to classify viewpoint.

```

5   import cv2
6   import os
7   import scipy.io as sio
8   from keras.models import load_model
9   import numpy as np
10  from scipy.misc import imsave, imread
11  import csv
12
13  # Set paths
14
15  # Detections in .mat format from DPM / matlab
16  DETECTIONS_DIR = './mat_result'
17  # Image Files for the detections - must have the same basename
18  IMAGES_DIR = './data_object/testing/image_2'
19  # Path to output everything to
20  OUTPUT_DIR = './angle_re'
21
22  # Get filepaths for images
23  image_files = [os.path.join(IMAGES_DIR, f) for f in os.listdir(IMAGES_DIR) if f.endswith('.png')]
24
25  # Get filepaths for detections
26  detection_files = [os.path.join(DETECTIONS_DIR, f) for f in
27                      os.listdir(DETECTIONS_DIR) if
28                      f.endswith('.mat')]
29
30  # Sort them so that they're in the same order
31  image_files.sort()
32  detection_files.sort()

```

```

40     model = load_model('object_classifierfit2.h5')
41     # Iterate through images and dpm detections
42     for i, d in zip(image_files, detection_files):
43         angles = []
44         # Filter detections
45         filtered = []
46         for m in sio.loadmat(d) ['ds']:
47             filtered.append(m)
48         # Load image
49         image = imread(i)
50         # Iterate through all detected objects
51         for d in filtered:
52             segment = image[int(d[2]):int(d[3]), int(d[0]):int(d[1]),:]
53             if (segment.shape[0] == 0 or segment.shape[1] == 0):
54                 continue
55             segment = cv2.resize(segment, (128, 128))
56             segment = np.array([segment])
57             segment = segment / 255.
58             # Generate a prediction
59             prediction = model.predict(segment)
60             angle = (np.argmax(prediction, axis=1)) * 30 - 180
61             x = d.tolist()
62             x.append(angle[0])
63             angles.append(x)
64         # Output to a csv
65         with open(os.path.basename(i).split('.')[0] + '.csv', 'w') as f:
66             w = csv.writer(f)
67             w.writerows(angles)

```

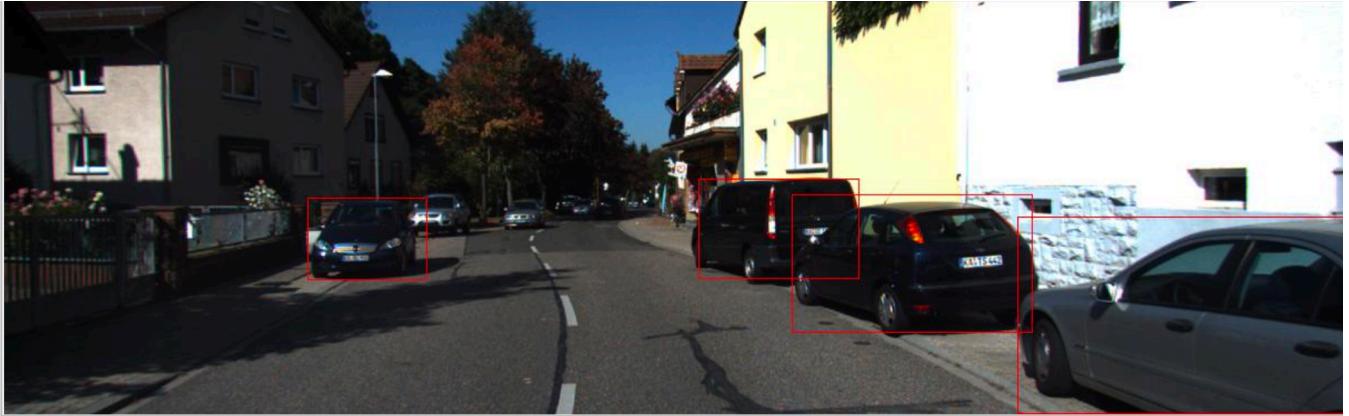
For input we need the image's name and .mat file that contains the detection to be in the same sorted order. In .mat file we have the order per each row to be left, right, up, down of bounding box, same as the output in csv. And the last element is the angle.

Result

Input image and the corresponding detect element is



The returned csv is 754,927,189,268,-30 Which is accurate.



The result csv file is

```
709,925,174,297,-30  
911,1220,193,371,0  
271,382,175,252,150  
621,773,156,253,-30
```

We can see that the right most car has should have a negative viewpoint angle but since the image didn't capture the entire car, it is classified as 0. So this is a limitation for this method. It depend on the detection box. When the car is not fully showed, there might be classification error.

The reason why I choose InceptionV3 with GAP(global average pooling) is because First CNN is good for segment classification and InceptionV3 is very good model under CNN. It has increased accuracy and speed than regular CNN. The reason why I choose GAP but not GMP(global maximum pooling) and FC(fully connected) is because some data shows that FS's accuracy decreased after several epochs, might caused by over fitting, and GMP has a lower accuracy than GAP.

(code reference : <https://github.com/Colliflower/Computer-Vision-Project>)

Question 8

Show a test image with the detected car bounding boxes and show the estimated viewpoints by plotting an arrow in the appropriate direction.

The arrow can be generate by the viewpoint and detect box but need to calibrate.

When angle = 0, the arrow is vertical. Thus change in vertical is $-\sin(\text{angle}+\pi/2)$, change in horizontal is $\cos(\text{angle}+\pi/2)$,

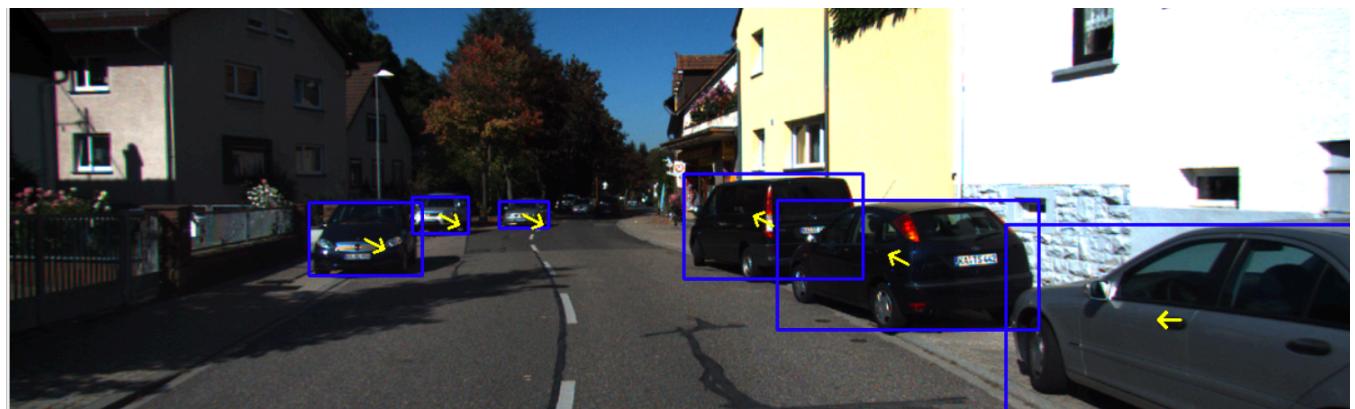
```

1 import csv
2 import cv2
3 import math
4
5 img = './data_object/testing/image_2/um_000087.png'
6 image = cv2.imread(img)
7 angle = 'um_000087.csv'

9 detect = []
10 with open(angle) as csv_file:
11     csv_reader = csv.reader(csv_file, delimiter=',')
12     for row in csv_reader:
13         detect.append(row)
14
15 for data in detect:
16     data = [int(x) for x in data]
17     color = (255, 0, 0)
18     thickness = 2
19     sp = (data[0], data[2])
20     ep = (data[1], data[3])
21     cv2.rectangle(image, sp, ep, color, thickness)
22     start_point = (int((data[2]+data[3])/2), int((data[0]+data[1])/2))
23
24 # End coordinate
25 data[4] = data[4]+math.pi/2
26 end_point = (start_point[0]-int(20*math.sin(data[4])), start_point[1]+int(20*math.cos(data[4])))
27
28 print(end_point)
29 # Yellow color in BGR
30 color = (0, 255, 255)
31 # Line thickness of 9 px
32 thickness = 2
33 # Using cv2.arrowedLine() method
34 # Draw a red arrow line
35 # with thickness of 9 px and tipLength = 0.5
36 image = cv2.arrowedLine(image, start_point[::-1], end_point[::-1],
37                         color, thickness, tipLength=0.5)
38
39 cv2.imwrite("arrow87.png", image)
40

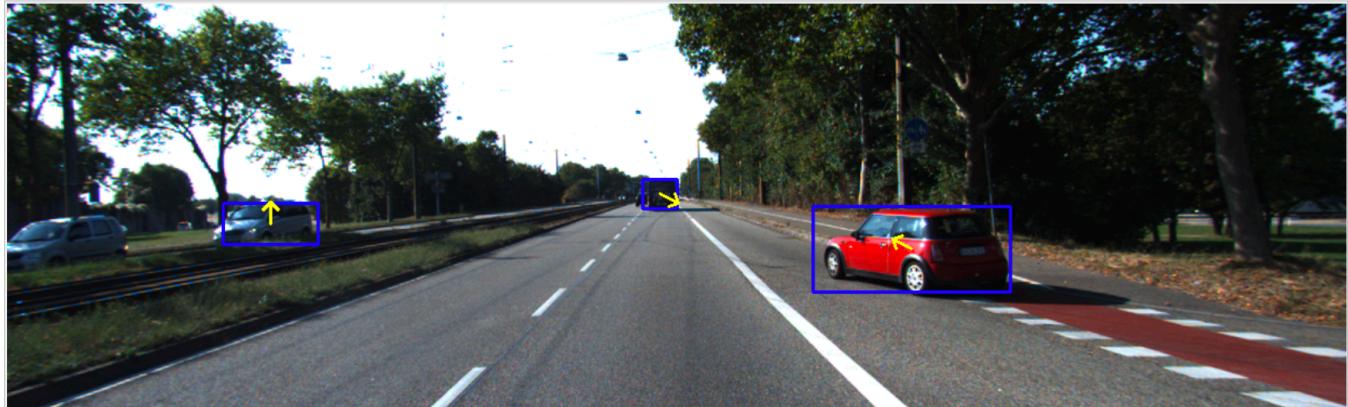
```

Result



Since we only have 12 classifier, and thus when we have a angle in between it will be less accurate.

Bad result



When the segment is very small, since we will resize it to same size with others, it will lost a lot of feature. But how ever in autonomous driving, car with small segment are very far away and we don't need to analysis it yet.

Question 9

Given the ground plane, estimated depth, and the location of the car's bounding box,

how would you compute a 3D bounding box around each detected car? Add the 3D

bounding boxes to your plot from 5.

To solve this question, first we need to get some data from previous parts. Including 2D box location, orientation of the car (the more accurate, the better result. So, I fixed some values manually) and the 3D location of the car. Besides, we also need the size of the car (I assume each car have similar size) and the camera rotation matrices provided to compute 3D value to 2D.

Assume we are given the size of each car, then we can easily compute 8 points to generate a 3D box, for example, the front left upper corner can be $-1/2 \times \text{width}$, height, $-1/2 \times \text{length}$ to the center of the car. with the 8 points are calculated then rotate it by the given rotation angle, we next can add the car

center position to it and get the relative location of each point in real world to the camera. The last step is applying camera rotation matrices to each of those points and get the 2D positions in the image. Connect pair of points to draw the 3D bounding boxes.

As for the 3d Box in the cloud points, it will follow the same steps above instead of transferring the points back to 2D image, this time we plot the points in the cloud with the calculated 3D location, then draw the box.

Code for 3D box in 2d image:

```
def project_to_image(P, pts):
    n_pts = pts.shape[1]
    # print('pts:', pts)
    #print('ones', np.ones((1, n_pts)))
    X = np.concatenate((pts, np.ones((1, n_pts))), axis=0)
    X = np.dot(P, X)
    X = X[:2, :] / X[2:3, :]
    return X

def kitti_plot_3d_bbox(img, label:KittiLabel, cam_to_img):
    colors = {
        'green': (0, 255, 0),
        'pink': (255, 0, 255),
        'blue': (0, 0, 255)
    }
    color = colors['pink']

    # Object dimensions
    h, w, l = label.dimensions

    # Coordinate of object centers
    # Note that label.location is location of the point at the bottom of the object
    # Need to add half object height to get object center
    tx, ty, tz = label.location

    theta_ray = np.arctan2(tx, tz)
    theta_local = label.alpha
    theta = theta_ray + theta_local

    ry = theta

    obj2body_R = R.from_euler('y', ry).as_dcm()
    obj2body_t = np.array([tx, ty, tz], dtype=float).reshape(-1, 1)
```

```

# Box corners
# x-axis on object (car) points from back to front of car
x_corners = [1/2, 1/2, -1/2, -1/2, 1/2, 1/2, -1/2, -1/2]
y_corners = [0,0,0,0,-h,-h,-h,-h]
z_corners = [w/2, -w/2, -w/2, w/2, w/2, -w/2, -w/2, w/2]

corners_in_obj = np.asarray([x_corners, y_corners, z_corners], dtype=float)
corners_in_body = np.dot(obj2body_R, corners_in_obj) + obj2body_t

# Project to image
corners_in_img = project_to_image(cam_to_img, corners_in_body)
face_idxes = [[0, 1, 5, 4], # front face
              [1, 2, 6, 5], #left face
              [2, 3, 7, 6], # back face
              [3, 0, 4, 7], #right face
              ]

for idx in face_idxes:
    p1 = corners_in_img[:, idx].astype(int)
    p2 = np.roll(p1, 1, axis=1)
    for i in range(p1.shape[1]):
        cv2.line(img, (p1[0, i], p1[1, i]), (p2[0, i], p2[1, i]), color, 2)
return img

def main():
    img_file = './kitti_2d/training/image_2/000001.png'
    label_file = './kitti_2d/training/label_2/000001.txt'
    calib_file = './kitti_2d/training/calib/000001.txt'
    image = cv2.imread(img_file, cv2.COLOR_BGR2RGB)
    image = cv2.cvtColor(image, cv2.COLOR_BGR2RGB)

    labels = kitti_read_label_file(label_file)
    labels = [l for l in labels if l.type != 'DontCare']

    cam_to_img = kitti_get_calibration_cam_to_image(calib_file)
    print(cam_to_img)

    for label_info in labels:
        if label_info.truncated > 0:
            continue

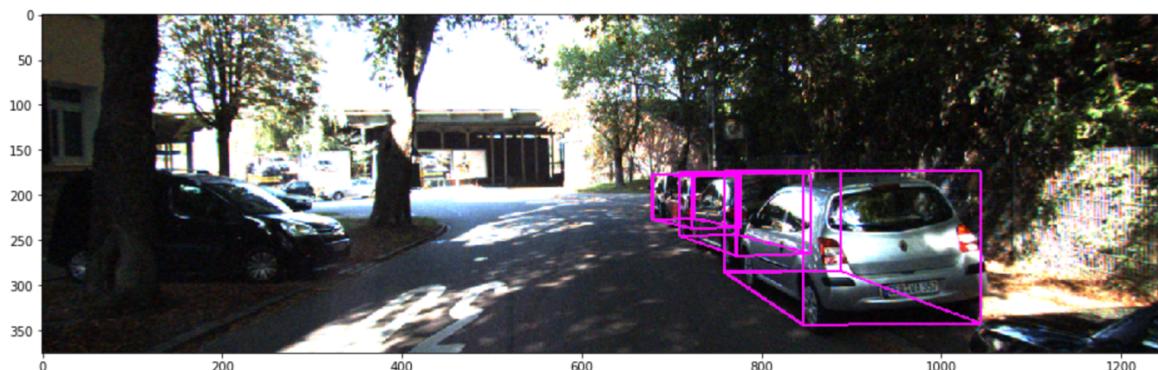
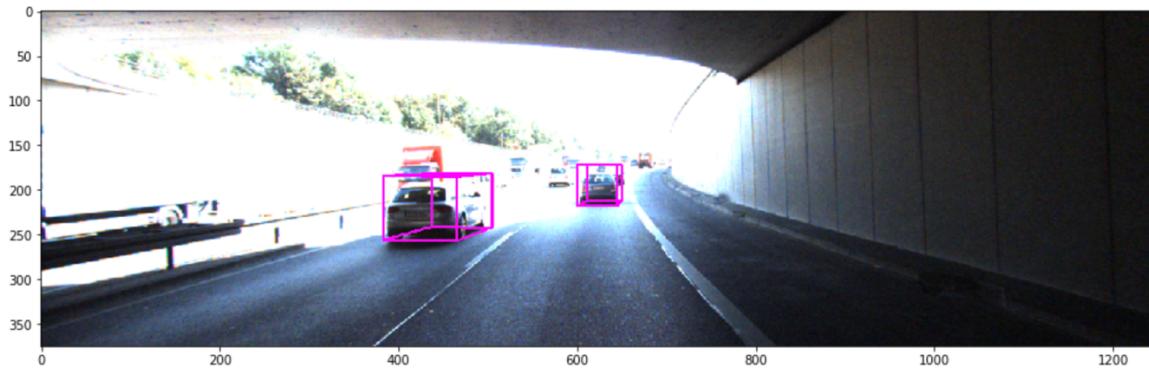
        # Overlay 3D box using computed location
        image = kitti_plot_3d_bbox(image, label_info, cam_to_img)
    return image

img = main()
plt.figure(figsize=(16, 16))
plt.imshow(img)
plt.savefig('img_with_box3d.png')

```

Reference: <https://github.com/nhonth/visualize-3D-BoundingBox.git>

Examples:



Code for 3D box in cloud:

Generate 8 points by given car size

```

def get_box_points(location, angle):
    tx, ty, tz = location
    theta_ray = np.arctan2(tx, tz)
    theta_local = angle
    ry = theta_ray + theta_local

    obj2body_R = R.from_euler('y', ry).as_dcm()
    obj2body_t = np.array([tx, ty, tz], dtype=float).reshape(-1, 1)
    l, h, w = car_size
    x_corners = [l / 2, -l / 2, -l / 2, l / 2, l / 2, -l / 2, -l / 2]
    y_corners = [0, 0, 0, 0, -h, -h, -h, -h]
    z_corners = [w / 2, -w / 2, -w / 2, w / 2, w / 2, -w / 2, -w / 2]
    corners_in_obj = np.asarray([x_corners, y_corners, z_corners], dtype=float)
    corners_in_body = np.dot(obj2body_R, corners_in_obj) + obj2body_t
    return corners_in_body.T

```

```

def draw_box(points):
    lines = [[0, 1], [1, 5], [2, 3], [4, 5], [1, 2], [3, 0], [6, 7],
             [0, 4], [1, 5], [2, 6], [3, 7], [4, 7], [5, 6]]
    colors = [[0, 255, 255] for _ in range(len(lines))]
    line_set = o3d.geometry.LineSet()
    line_set.lines = o3d.utility.Vector2iVector(lines)
    line_set.colors = o3d.utility.Vector3dVector(colors)
    print(points)
    line_set.points = o3d.utility.Vector3dVector(points)
    return line_set

```

To get the 3D location of the car, first I get the center part of the car patch and get its 3D location. Then apply sin and cos function to calculate how far from that point to the car center point. Then add them up to get the final location for the center point of the car. However, it will not accurate enough sometimes.

```

def get_location(center, ori):
    angle = -ori / 180 * np.pi
    # new3d_img = calculate_3d(imgnl, depth)
    location = three_d[center[0], center[1], :]
    location[0] = location[0] + np.sin(angle) * 1 / 2 * car_size[0]
    location[1] = 0
    location[2] = location[2] + np.cos(angle) * 1 / 2 * car_size[2]
    return location, angle

```

```

location, angle = get_location(center, angle)
line_set = draw_box(get_box_points(location, angle))
print('detail', location, angle)

```

Example:

By using the data from previous parts, I draw the boxes for the detected cars. However, the data is definitely not accurate enough, thus the error in the result 3D cloud is huge.



After I fixed the data manually by giving a more accurate center point for the red car on the right side, the result performs better for the red car.

