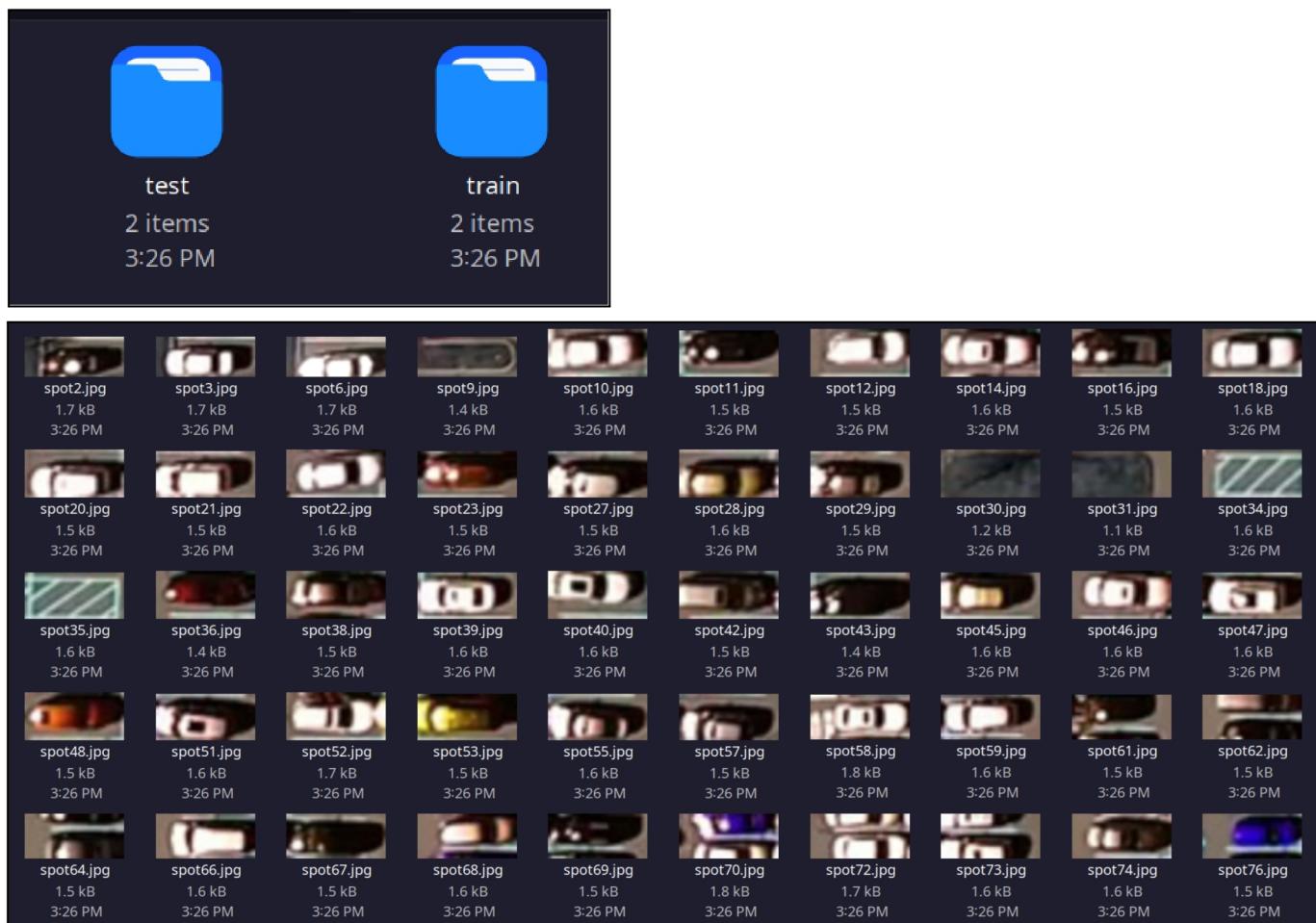


Task 1:Create a Car Parking Prediction Model using TinyML

1. Data Collection

Collect the required datasets for training and testing the model.



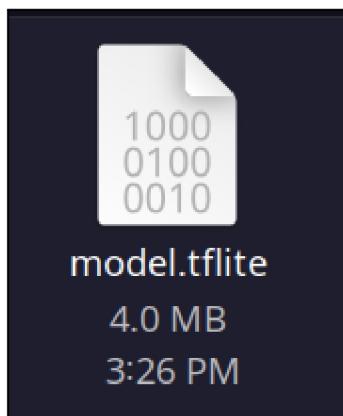
2. Train and prepare the model

Open and review the `train.ipynb` Jupyter notebook. Understand the data preprocessing, model selection, training process, and evaluation metrics used. Execute the code to train the model.

```
train.ipynb M ×  
TinyMLParking > train.ipynb > # Export to Tensorflow Lite model and label file in 'export_dir'.  
+ Code + Markdown | ▶ Run All ⌂ Restart ⌂ Clear All Outputs | ⌂ Variables ⌂ Outline ...  
tinyMLParking (Python 3.10.12)  
  
1 #from tflite_model_maker import image_classifier # type: ignore  
2 from tflite_model_maker.image_classifier import DataLoader # type: ignore  
[]  
  
1 # Load input data specific to an on-device ML app.  
2 train_data = DataLoader.from_folder('./train_images/train')  
3 test_data = DataLoader.from_folder('./train_images/test')  
[]  
  
1 # Customize the TensorFlow model.  
2 model = image_classifier.create(train_data)  
[]  
  
1 # Evaluate the model.  
2 loss, accuracy = model.evaluate(test_data)  
[]  
  
1 # Evaluate the model.  
2 loss, accuracy = model.evaluate(test_data)  
[]  
  
1 # Export to Tensorflow Lite model and label file in 'export_dir'.  
2 model.export(export_dir='./model')  
[]
```

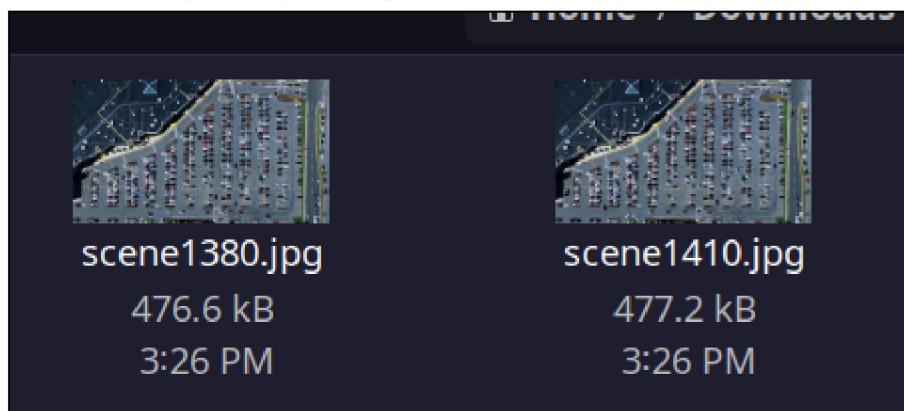
3. Model File

This directory contains the trained TinyML model file(s). These files are the result of training the model using the training data. Tflite file is generated successfully and are ready for deployment.



4. Checking the Accuracy by Testing Images

This directory contains images used for testing the parking spot detection algorithm. The model is tested for the given input images to check the accuracy



5. Identifying the Parking Spots:

Write the code in `identify_parking_spots.ipynb` Jupyter notebook. It contains code for identifying parking spots in the given images. Understand the algorithms and techniques used for parking spot detection. Execute the code and analyze the results.

```
1 from __future__ import division
2 import matplotlib.pyplot as plt
3 import cv2
4 import os, glob
5 import numpy as np
6 # from moviepy.editor import VideoFileClip
7 import tensorflow as tf
8
9 # cwd = os.getcwd()
10 # import tensorflow as tf

1 interpreter = tf.lite.Interpreter("./model/model.tflite")
2 interpreter.allocate_tensors()
3 input_details = interpreter.get_input_details()
4 output_details = interpreter.get_output_details()

INFO: Created TensorFlow Lite XNNPACK delegate for CPU.

1 def show_images(images, cmap=None):
2     cols = 2
3     rows = (len(images) + 1) // cols
4
5     plt.figure(figsize=(15, 12))
6     for i, image in enumerate(images):
7         plt.subplot(rows, cols, i + 1)
8         # use gray scale color map if there is only one channel
9         cmap = "gray" if len(image.shape) == 2 else cmap
10         plt.imshow(image, cmap=cmap)
11         plt.xticks([])
12         plt.yticks([])
13     plt.tight_layout(pad=0, h_pad=0, w_pad=0)
14     plt.show()
```

```
1 test_images = [plt.imread(path) for path in glob.glob("./test_images/*.jpg")]
2
3 show_images(test_images)
```

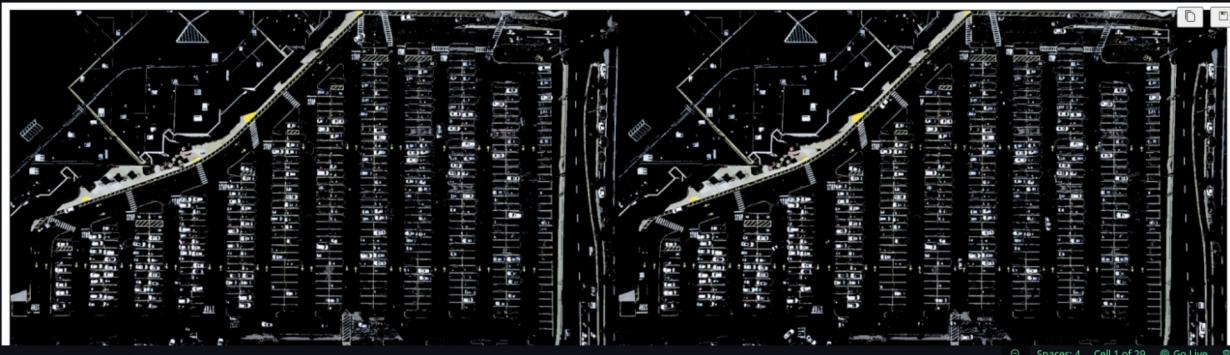
Python



Color Selection and Edge Detection

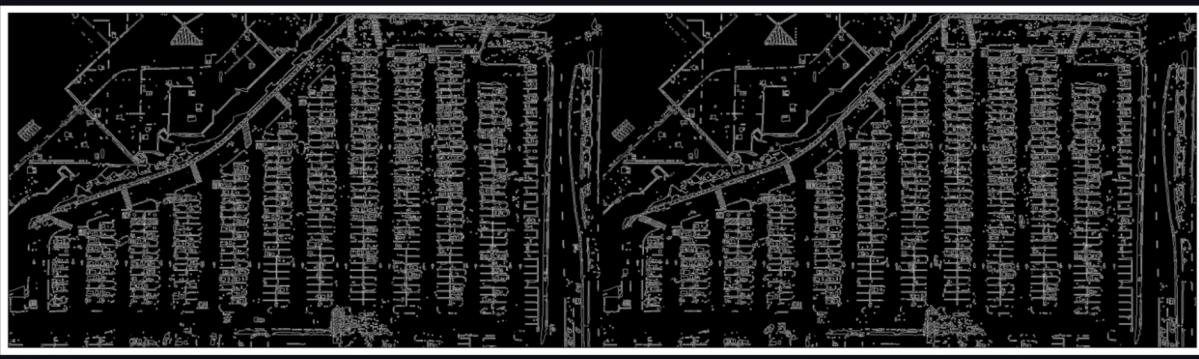
```
1 # image is expected be in RGB color space# image
2 def select_rgb_white_yellow(image):
3     ...# white color mask
4     ...lower = np.uint8([120, 120, 120])    Argument of type "list[int]" cannot be assigned to parameter "value" of type "_IntValue" in function "__init__" @ Type
5     ...upper = np.uint8([255, 255, 255])    Argument of type "list[int]" cannot be assigned to parameter "value" of type "_IntValue" in function "__init__" @ Type
6     ...white_mask = cv2.inRange(image, lower, upper)    No overloads for "inRange" match the provided arguments - (Pylance)
7     ...# yellow color mask
8     ...lower = np.uint8([190, 190, 0])    Argument of type "list[int]" cannot be assigned to parameter "value" of type "_IntValue" in function "__init__" @ Type
9     ...upper = np.uint8([255, 255, 255])    Argument of type "list[int]" cannot be assigned to parameter "value" of type "_IntValue" in function "__init__" @ Type
10    ...yellow_mask = cv2.inRange(image, lower, upper)    No overloads for "inRange" match the provided arguments - (Pylance)
11    ...# combine the mask
12    ...mask = cv2.bitwise_or(white_mask, yellow_mask)
13    ...masked = cv2.bitwise_and(image, image, mask=mask)
14    ...return masked
15
16
17 white_yellow_images = list(map(select_rgb_white_yellow, test_images))
18 show_images(white_yellow_images)
```

Python



```
1 def detect_edges(image, low_threshold=50, high_threshold=200):
2     ...return cv2.Canny(image, low_threshold, high_threshold)
3
4
5 edge_images = list(map(lambda image: detect_edges(image), gray_images))
6
7 show_images(edge_images)
```

Python





Identify area of interest

```
1 def filter_region(image, vertices):
2     """
3     Create the mask using the vertices and apply it to the input image
4     """
5     mask = np.zeros_like(image)
6     if len(mask.shape) == 2:
7         cv2.fillPoly(mask, vertices, 255)    No overloads for "fillPoly" match the provided arguments - (Pylance)
8     else:
9         cv2.fillPoly(
10             mask, vertices, (255,) * mask.shape[2]
11         ) # in case, the input image has a channel dimension
12     return cv2.bitwise_and(image, mask)
13
14
15 def select_region(image):
16     """
17     It keeps the region surrounded by the `vertices` (i.e. polygon). Other area is set to 0 (black).
18     """
19     # first, define the polygon by vertices
20     rows, cols = image.shape[:2]
21     pt_1 = [cols * 0.05, rows * 0.90]
22     pt_2 = [cols * 0.05, rows * 0.70]
23     pt_3 = [cols * 0.30, rows * 0.55]
24     pt_4 = [cols * 0.6, rows * 0.15]
25     pt_5 = [cols * 0.90, rows * 0.15]
26     pt_6 = [cols * 0.90, rows * 0.90]
27     # the vertices are an array of polygons (i.e. array of arrays) and the data type must be integer
28     vertices = np.array([[pt_1, pt_2, pt_3, pt_4, pt_5, pt_6]], dtype=np.int32)
29     return filter_region(image, vertices)
30
31
32 # images showing the region of interest only
33 roi_images = list(map(select_region, edge_images))
34
35 show_images(roi_images)
```



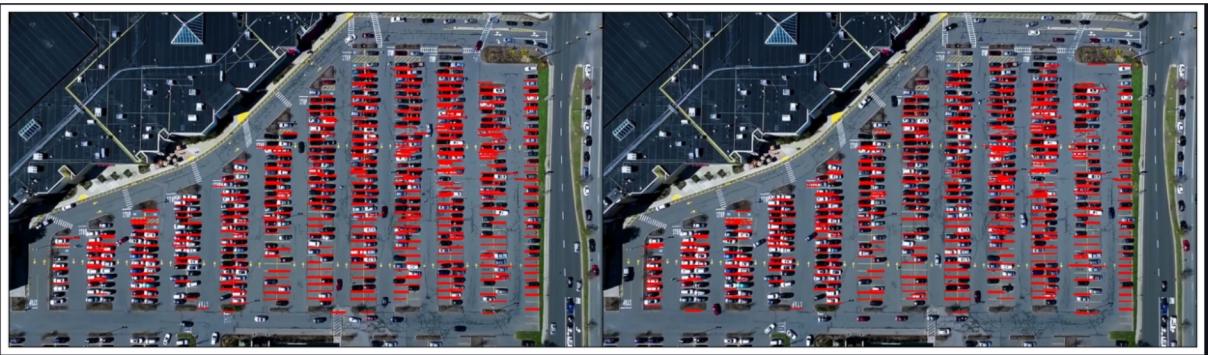
Hough line transform

```
1 def hough_lines(image):
2     """
3     'image' should be the output of a Canny transform.
4
5     Returns hough lines (not the image with lines)
6     """
7     return cv2.HoughLinesP(
8         image, rho=0.1, theta=np.pi / 10, threshold=15, minLineLength=9, maxLineGap=4
9     )
10
11
12 list_of_lines = list(map(hough_lines, roi_images))
```

Python

```
1 def draw_lines(image, lines, color=[255, 0, 0], thickness=2, make_copy=True):
2     # the lines returned by cv2.HoughLinesP has the shape (-1, 1, 4)
3     if make_copy:
4         image = np.copy(image) # don't want to modify the original
5     cleaned = []
6     for line in lines:
7         for x1, y1, x2, y2 in line:
8             if abs(y2 - y1) <= 1 and abs(x2 - x1) >= 25 and abs(x2 - x1) <= 55:
9                 cleaned.append((x1, y1, x2, y2))
10            cv2.line(image, (x1, y1), (x2, y2), color, thickness)
11    print(" No lines detected: ", len(cleaned))
12    return image
13
14
15 line_images = []
16 for image, lines in zip(test_images, list_of_lines):
17     line_images.append(draw_lines(image, lines))
18
19 show_images(line_images)
```

Python



Identify rectangular blocks of parking

```
1 def identify_blocks(image, lines, make_copy=True):
2     if make_copy:
3         new_image = np.copy(image)
4     # Step 1: Create a clean list of lines
5     cleaned = []
6     for line in lines:
7         for x1, y1, x2, y2 in line:
8             if abs(y2 - y1) <= 1 and abs(x2 - x1) >= 25 and abs(x2 - x1) <= 55:
9                 cleaned.append((x1, y1, x2, y2))
10
11 # Step 2: Sort cleaned by x1 position
12 import operator
13
14 list1 = sorted(cleaned, key=operator.itemgetter(0, 1))
15
16 # Step 3: Find clusters of x1 close together - clust_dist apart
17 clusters = {}
18 dIndex = 0
19 clus_dist = 10
20
21 for i in range(len(list1) - 1):
22     distance = abs(list1[i + 1][0] - list1[i][0])
23     # print(distance)
24     if distance <= clus_dist:
25         if not dIndex in clusters.keys():
26             clusters[dIndex] = []
27             clusters[dIndex].append(list1[i])
28             clusters[dIndex].append(list1[i + 1])
29         else:
30             dIndex += 1
31
32
```

Python

```

52
53     # Step 4: Identify coordinates of rectangle around this cluster
54     rects = {}
55     i = 0
56     for key in clusters:
57         all_list = clusters[key]
58         cleaned = list(set(all_list))
59         if len(cleaned) > 5:
60             cleaned = sorted(cleaned, key=lambda tup: tup[1])
61             avg_y1 = cleaned[0][1]
62             avg_y2 = cleaned[-1][1]
63             # print(avg_y1, avg_y2)
64             avg_x1 = 0
65             avg_x2 = 0
66             for tup in cleaned:
67                 avg_x1 += tup[0]
68                 avg_x2 += tup[2]
69                 avg_x1 = avg_x1 / len(cleaned)
70                 avg_x2 = avg_x2 / len(cleaned)
71             rects[i] = (avg_x1, avg_y1, avg_x2, avg_y2)
72             i += 1
73
74     print("Num Parking Lanes: ", len(rects))
75     # Step 5: Draw the rectangles on the image
76     buff = 7
77     for key in rects:
78         tup_topLeft = (int(rects[key][0] - buff), int(rects[key][1]))
79         tup_botRight = (int(rects[key][2] + buff), int(rects[key][3]))
80         # print(tup_topLeft, tup_botRight)
81         cv2.rectangle(new_image, tup_topLeft, tup_botRight, (0, 255, 0), 3)
82     return new_image, rects
83
84
85 # images showing the region of interest only
86 rect_images = []
87 rect_coords = []
88 for image, lines in zip(test_images, list_of_lines):
89     new_image, rects = identify_blocks(image, lines)
90     rect_images.append(new_image)
91     rect_coords.append(rects)
92
93 show_images(rect_images)

```

Num Parking Lanes: 12
Num Parking Lanes: 12

Python



Identify each spot and count num of parking spaces

Next step-

1. Based on width of each parking line segment into individual spots
2. draw a visualization of all parking spaces

```

1 def draw_parking(
2     image, rects, make_copy=True, color=[255, 0, 0], thickness=2, save=True
3 ):
4     if make_copy:
5         new_image = np.copy(image)
6         gap = 15.5
7         adj_y1 = {variable: dict[int, int] ID to its coords
8
9         adj_y1 = {
10            0: 20,
11            1: -10,
12            2: 0,
13            3: -11,
14            4: 28,
15            5: 5,
16            6: -15,
17            7: -15,
18            8: -10,
19            9: -30,
20            10: 9,
21            11: -32,
22        }

```

```

56     11: 0,
57   }
58   for key in rects:
59     # Horizontal lines
60     tup = rects[key]
61     x1 = int(tup[0] + adj_x1[key])
62     x2 = int(tup[2] + adj_x2[key])
63     y1 = int(tup[1] + adj_y1[key])
64     y2 = int(tup[3] + adj_y2[key])
65     cv2.rectangle(new_image, (x1, y1), (x2, y2), (0, 255, 0), 2)
66     num_splits = int(abs(y2 - y1) // gap)
67     for i in range(0, num_splits + 1):
68       y = int(y1 + i * gap)
69       cv2.line(new_image, (x1, y), (x2, y), color, thickness)
70     if key > 0 and key < len(rects) - 1:
71       # draw vertical lines
72       x = int((x1 + x2) / 2)
73       cv2.line(new_image, (x, y1), (x, y2), color, thickness)
74     # Add up spots in this lane
75     if key == 0 or key == (len(rects) - 1):
76       tot_spots += num_splits + 1
77     else:
78       tot_spots += 2 * (num_splits + 1)
79
80   # Dictionary of spot positions
81   if key == 0 or key == (len(rects) - 1):
82     for i in range(0, num_splits + 1):
83       cur_len = len(spot_dict)
84       y = int(y1 + i * gap)
85       spot_dict[(x1, y, x2, y + gap)] = cur_len + 1
86     else:
87       for i in range(0, num_splits + 1):
88         cur_len = len(spot_dict)
89         y = int(y1 + i * gap)
90         x = int((x1 + x2) / 2)
91         spot_dict[(x1, y, x, y + gap)] = cur_len + 1
92         spot_dict[(x, y, x2, y + gap)] = cur_len + 2
93
94   print("total parking spaces: ", tot_spots, cur_len)
95   if save:
96     filename = "with_parking.jpg"
97     cv2.imwrite(filename, new_image)
98   return new_image, spot_dict
99
100
101
102
103
104
105
106
107

```

```

108 delineated = []
109 spot_pos = []
110 for image, rects in zip(test_images, rect_coords):
111   new_image, spot_dict = draw_parking(image, rects)
112   delineated.append(new_image)
113   spot_pos.append(spot_dict)
114
115
116 show_images(delineated)

```

total parking spaces: 541 540
total parking spaces: 552 551



```
1 final_spot_dict = spot_pos[1]
```

Python

```
1 print(len(final_spot_dict))
```

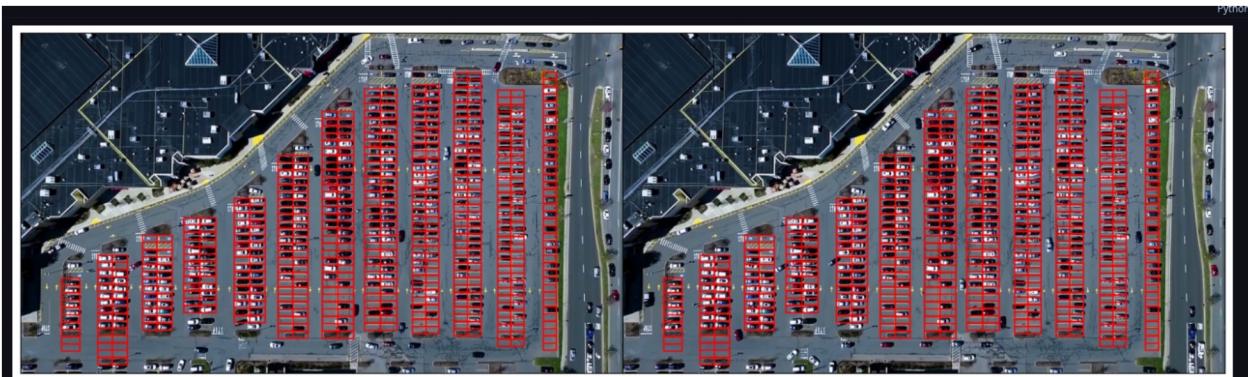
Python

552

```

1 def assign_spots_map(
2   image, spot_dict=final_spot_dict, make_copy=True, color=[255, 0, 0], thickness=2
3 ):
4   if make_copy:
5     new_image = np.copy(image)
6     for spot in spot_dict.keys():
7       (x1, y1, x2, y2) = spot
8       cv2.rectangle(
9         new_image, (int(x1), int(y1)), (int(x2), int(y2)), color, thickness
10        )
11   return new_image
12
13
14 marked_spot_images = list(map(assign_spots_map, test_images))
15 show_images(marked_spot_images)

```



```
1 # Save spot dictionary as pickle file
2 import pickle
3
4 with open("spot_dict.pickle", "wb") as handle:
5     pickle.dump(final_spot_dict, handle, protocol=pickle.HIGHEST_PROTOCOL)
```

Python

Save image for CNN model

```
1 def save_images_for_cnn(image, spot_dict=final_spot_dict, folder_name="for_cnn"):
2     for spot in spot_dict.keys():
3         (x1, y1, x2, y2) = spot
4         (x1, y1, x2, y2) = (int(x1), int(y1), int(x2), int(y2))
5         # crop this image
6         # print(image.shape)
7         spot_img = image[y1:y2, x1:x2]
8         spot_img = cv2.resize(spot_img, (0, 0), fx=2.0, fy=2.0)
9         spot_id = spot_dict[spot]
10
11         filename = "spot" + str(spot_id) + ".jpg"
12         print(spot_img.shape, filename, (x1, x2, y1, y2))
13
14         cv2.imwrite(os.path.join(folder_name, filename), spot_img)
15
16
17 # save_images_for_cnn(test_images[0])
```

[19]

Use trained CNN model to make predictions

```
1 from PIL import Image
2 # model = load_model(top_model_weights_path)
```

[20]

```
1 def make_prediction(img):
2     # Rescale image
3     # img = image/255.
4
5     # Convert to a 4D tensor
6     # image = np.expand_dims(img, axis=0)
7     # print(image.shape)
8
9     # make predictions on the preloaded model
10    # class_predicted = model.predict(image)
11    # inID = np.argmax(class_predicted[0])
12    # label = class_dictionary[inID]
13    # return label
14
15    new_img = cv2.resize(img, (224, 224))
16
17    # input_details[0]['index'] = the index which accepts the input
18    interpreter.set_tensor(input_details[0]["index"], [new_img])
19
20    # run the inference
21    interpreter.invoke()
22
23    # output_details[0]['index'] = the index which provides the input
24    output_data = interpreter.get_tensor(output_details[0]["index"])
25
26    return "empty" if output_data[0][0] > output_data[0][1] else "occupied"
```

[21]

```
1 from tqdm import tqdm
2
```

```

3 def predict_on_image(
4     image, spot_dict=final_spot_dict, make_copy=True, color=[0, 255, 0], alpha=0.5
5 ):
6     if make_copy:
7         new_image = np.copy(image)
8         overlay = np.copy(image)
9     else:
10        cnt_empty = 0
11        all_spots = 0
12        for spot in tqdm(spot_dict.keys()):
13            all_spots += 1
14            (x1, y1, x2, y2) = spot
15            (x1, y1, x2, y2) = (int(x1), int(y1), int(x2), int(y2))
16            # crop this image
17            spot_img = image[y1:y2, x1:x2]
18            spot_img = cv2.resize(spot_img, (48, 48))
19            label = make_prediction(spot_img)
20            # print(label)
21            if label == "empty":
22                cv2.rectangle(overlay, (int(x1), int(y1)), (int(x2), int(y2)), color, -1)
23                cnt_empty += 1
24
25    cv2.addWeighted(overlay, alpha, new_image, 1 - alpha, 0, new_image)
26
27    cv2.putText(
28        new_image,
29        "Available: %d spots" % cnt_empty,
30        (30, 95),
31        cv2.FONT_HERSHEY_SIMPLEX,
32        0.7,
33        (255, 255, 255),
34        2,
35    )
36
37    cv2.putText(
38        new_image,
39        "Total: %d spots" % all_spots,
40        (30, 125),
41        cv2.FONT_HERSHEY_SIMPLEX,
42        0.7,
43        (255, 255, 255),
44        2,
45    )
46
47    save = False
48
49    if save:
50        filename = "with_marking.jpg"
51        cv2.imwrite(filename, new_image)
52
53    return new_image
54
55
56 predicted_images = list(map(predict_on_image, test_images))
57 show_images(predicted_images)

```

[2] 1%|██████████| 3/552 [00:00<00:24, 22.30it/s]
100%|██████████| 552/552 [00:19<00:00, 28.27it/s]
100%|██████████| 552/552 [00:19<00:00, 27.84it/s]



6. Output Image - “with_parking.jpg” :

This is an output image showing parking spots identified or predicted by the algorithm. It has the identified empty slots where the car can be parked.

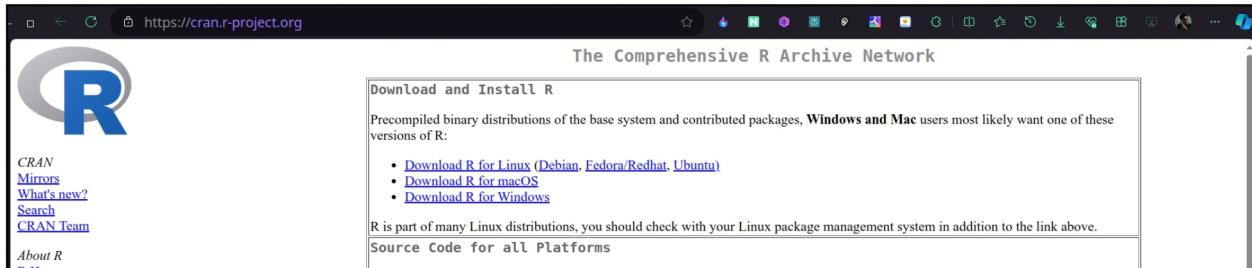


By following these steps, you can build a TinyML-based car parking prediction project. You'll be able to assess the effectiveness of the parking spot detection algorithm and the trained TinyML model.

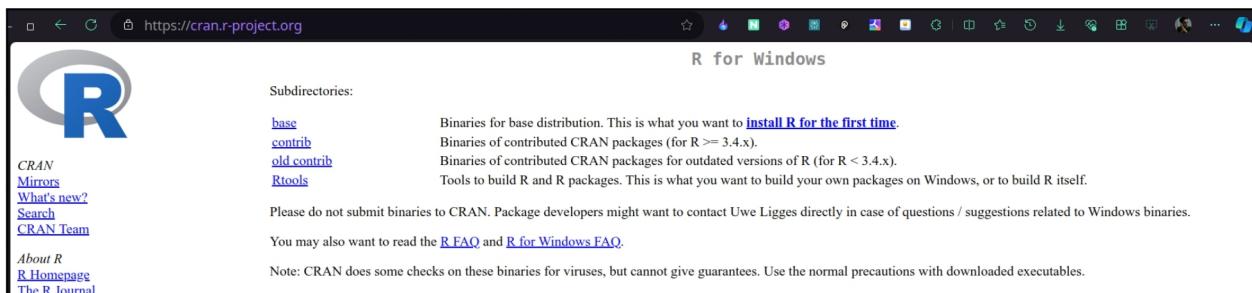
Task 2: Data Analytics using R programming

Setting up R and RStudio -

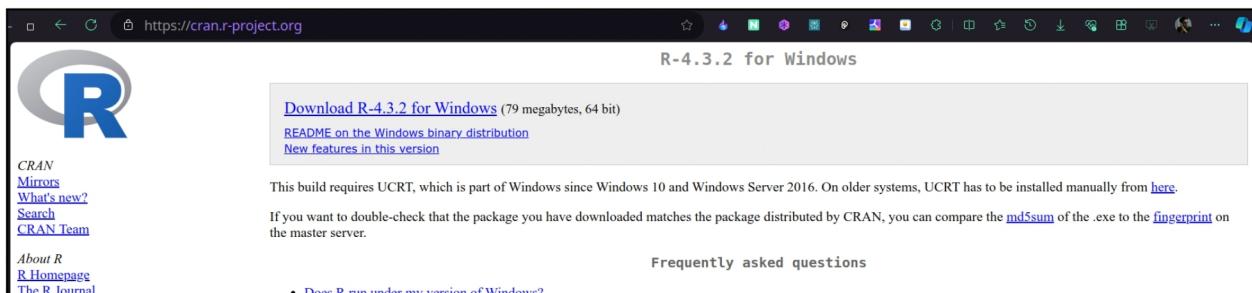
1. To install R, go to [CRAN](https://cran.r-project.org). Click Download R for Windows.



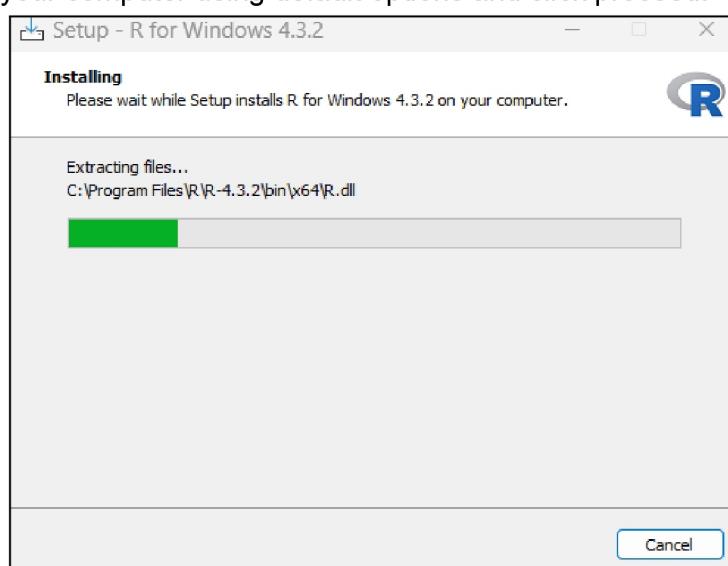
2. Install R Click on install R for the first time.



3. Click Download R for Windows. Open the downloaded file.



4. Setup R on your computer using default options and click proceed.



- After the installation of R , its desktop entry shown on the Desktop.



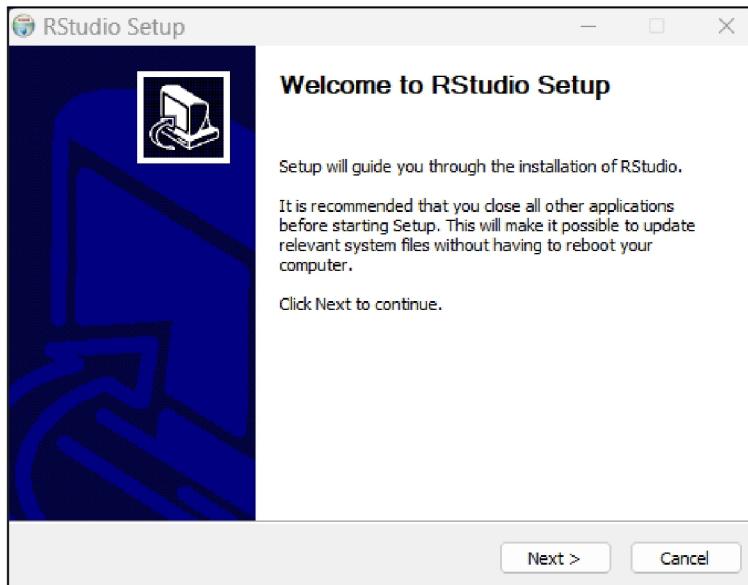
- Next, download RStudio. Go to [Download RStudio - Posit](#) .

The screenshot shows a web browser displaying the Posit website at <https://posit.co/downloads/>. The main heading is "RStudio IDE". Below it, a sub-headline reads "The most popular coding environment for R, built with love by Posit.". A paragraph describes RStudio as a set of tools for R and Python development, including a console, syntax-highlighting editor, plotting tools, and debugging features. A call-to-action button labeled "DOWNLOAD RSTUDIO" is visible. The browser's address bar shows the URL, and the top navigation bar includes links for PRODUCTS, SOLUTIONS, LEARN & SUPPORT, EXPLORE MORE, and PRICING.

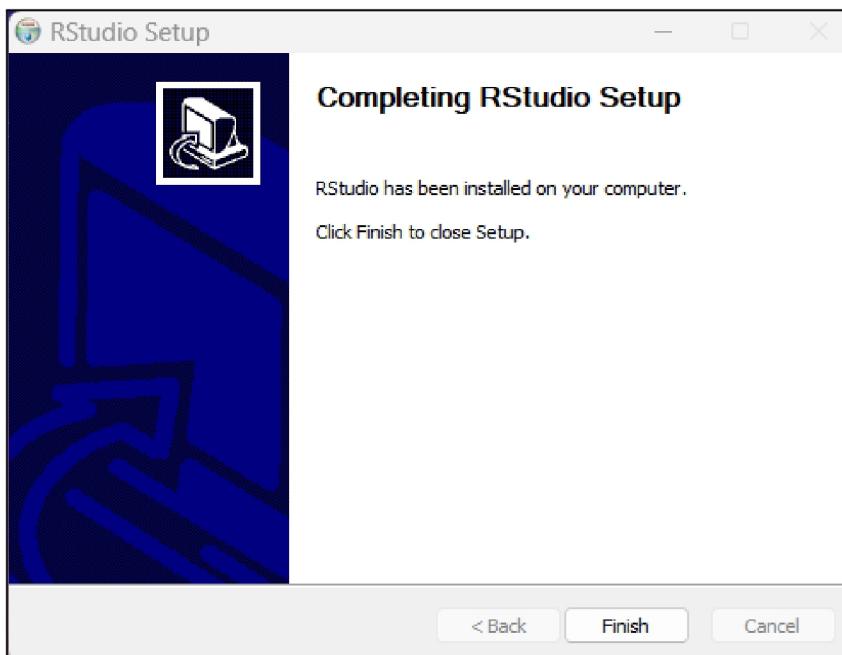
- Click Download RStudio.

The screenshot shows a web browser displaying the Posit website at <https://posit.co/downloads/>. The main heading is "1: Install R". Below it, a sub-headline reads "RStudio requires R 3.3.0+. Choose a version of R that matches your computer's operating system.". A call-to-action button labeled "DOWNLOAD AND INSTALL R" is visible. To the right, another section titled "2: Install RStudio" has a large "DOWNLOAD RSTUDIO DESKTOP FOR WINDOWS" button. Below this button, text provides file details: "Size: 215.66 MB | [SHA-256: D3C03C42](#) | Version: 2023.12.1+402 | Released: 2024-01-29". At the bottom, a dark bar contains the text "All Installers and Tarballs". The browser's address bar shows the URL, and the top navigation bar includes links for PRODUCTS, SOLUTIONS, LEARN & SUPPORT, EXPLORE MORE, and PRICING.

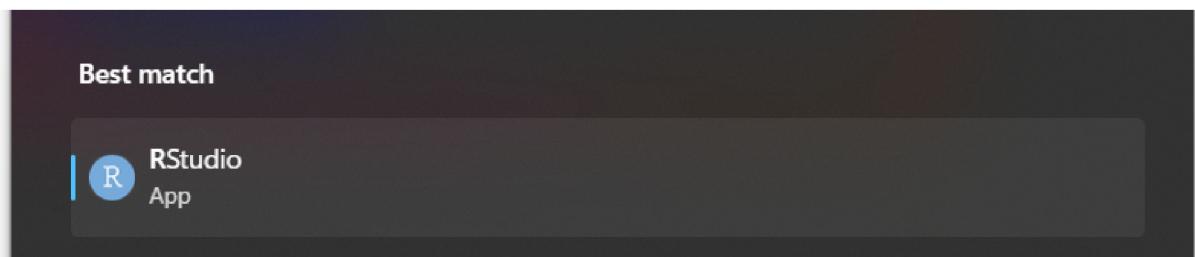
- Once the packet has downloaded, the Welcome to RStudio Setup Wizard will open. Click Next and go through the installation steps.



- Choose the default options and click finish.



- After the Setup Wizard finishing the installation, RStudio can be launched from the Search Menu.



Analysis & Visualization of Bank Dataset-

11. Go to [Kaggle - Bank Marketing Dataset](#) and download the Bank Dataset.

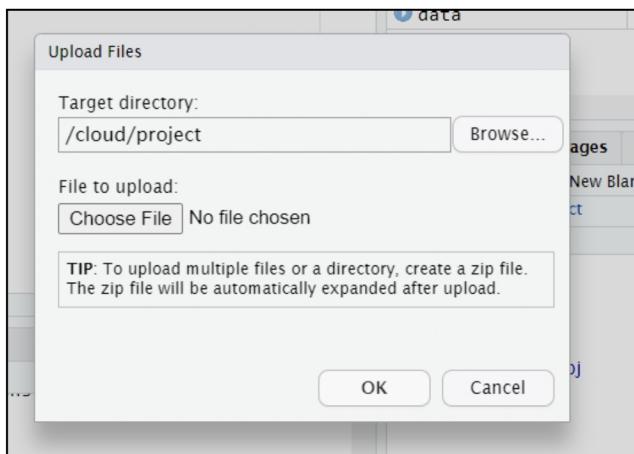
The screenshot shows the Kaggle interface with the following details:

- User: JANIO MARTINEZ BACHMANN - UPDATED 6 YEARS AGO
- Views: 459
- Actions: New Notebook, Download (146 kB)
- Title: Bank Marketing Dataset
- Description: Predicting Term Deposit Suscriptions
- Visual: A hand-drawn style graphic with the words "Marketing Strategy" in large letters, surrounded by various business-related terms like "social media", "management", "presenter", "recommendation", "innovation", etc.
- Links: Data Card, Code (183), Discussion (7)

12. Extract archive and the bank.csv will be extracted.



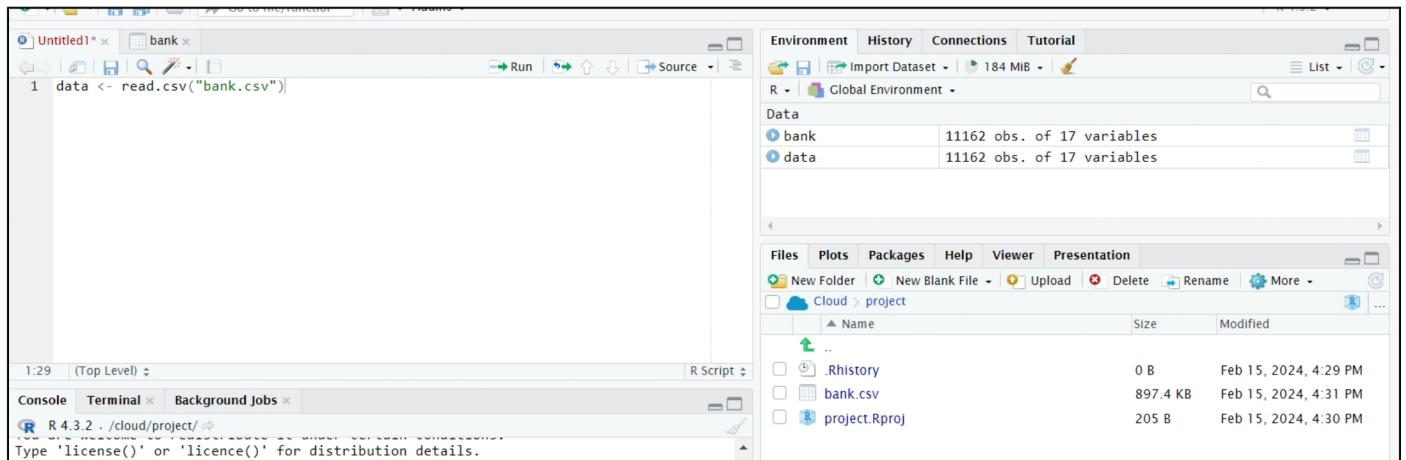
13. Open Rstudio and Import the bank.csv dataset.



The RStudio file browser interface. The left sidebar shows a tree view with "cloud" and "project" selected. The main area displays the contents of the "project" folder:

File	Size	Last Modified
.Rhistory	0 B	Feb 15, 2024, 4:29 PM
bank.csv	897.4 KB	Feb 15, 2024, 4:31 PM
project.Rproj	205 B	Feb 15, 2024, 4:30 PM

The "bank.csv" file is highlighted in blue, indicating it is selected.



14. Insert the following code in the RStudio Console line by line.

Code:

```

data <- read.csv("bank.csv")

#Visualization
install.packages("ggplot2")
# Create a pie chart for the 'poutcome' variable
ggplot(data, aes(x = "", fill = poutcome)) +
  geom_bar(width = 1, color = "white") +
  coord_polar("y", start = 0) +
  labs(title = "Pie Chart of Previous Campaign Outcome",
       fill = "Outcome",
       x = NULL,
       y = NULL) +
  theme_void() +
  theme(legend.position = "right")

#Histogram
# Create a histogram of the 'campaign' variable
ggplot(data, aes(x = campaign)) +
  geom_histogram(binwidth = 1, fill = "blue", color = "black") +
  labs(title = "Histogram of Campaign",
       x = "Number of Contacts in Current Campaign",
       y = "Frequency")

library(ggplot2)
library(dplyr)
# Bin the 'pdays' variable into categories
data <- mutate(data, pdays_category = cut(pdays, breaks = c(-1, 0, 100, 200, 300, 400, Inf)))

# Create a horizontal bar chart for the 'pdays' variable
ggplot(data, aes(x = pdays_category)) +

```

```

geom_bar(fill = "skyblue", color = "black") +
coord_flip() # This flips the axes to make it horizontal
labs(title = "Horizontal Bar Chart of pdays",
x = "Count",
y = "pdays Category") +
theme_minimal()

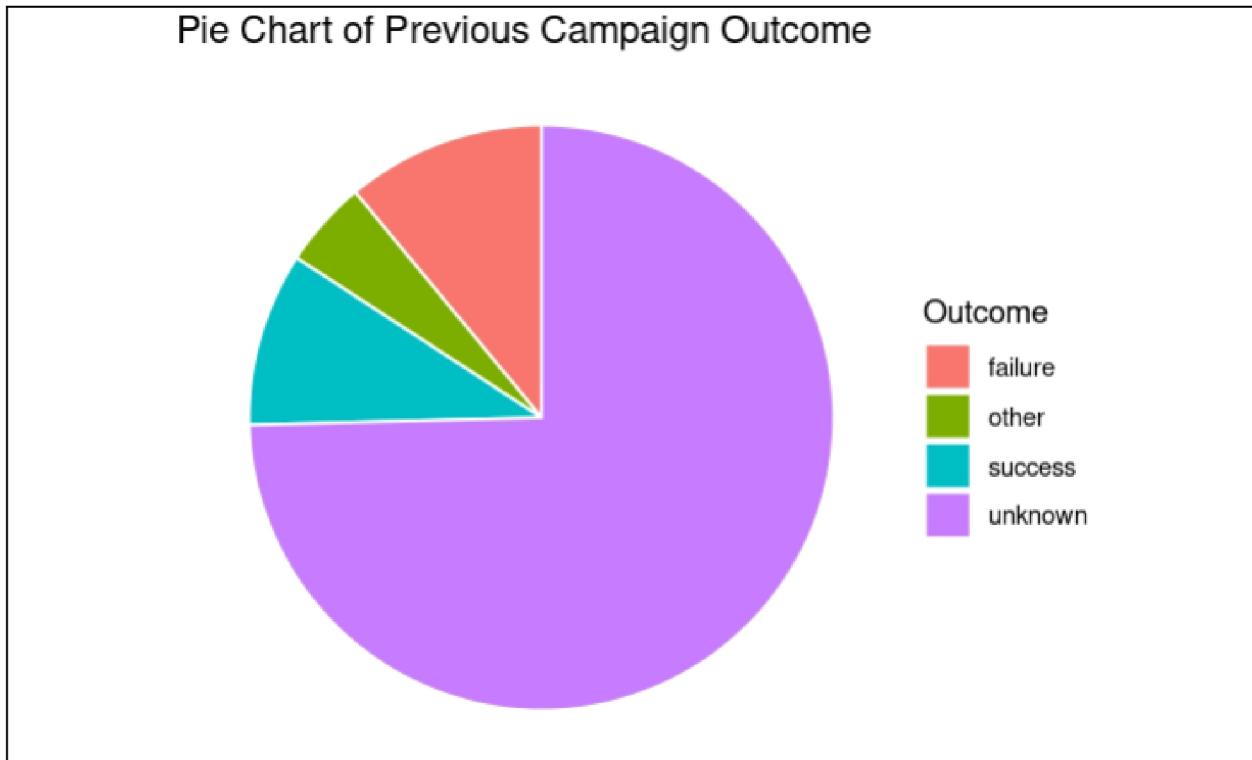
# Create a pie chart for the 'education' variable
install.packages("dplyr")
ggplot(data, aes(x = "", fill = education)) +
geom_bar(width = 1, color = "white") +
coord_polar("y", start = 0) +
labs(title = "Pie Chart of Education Level",
fill = "Education Level",
x = NULL,
y = NULL) +
theme_void() +
theme(legend.position = "right")

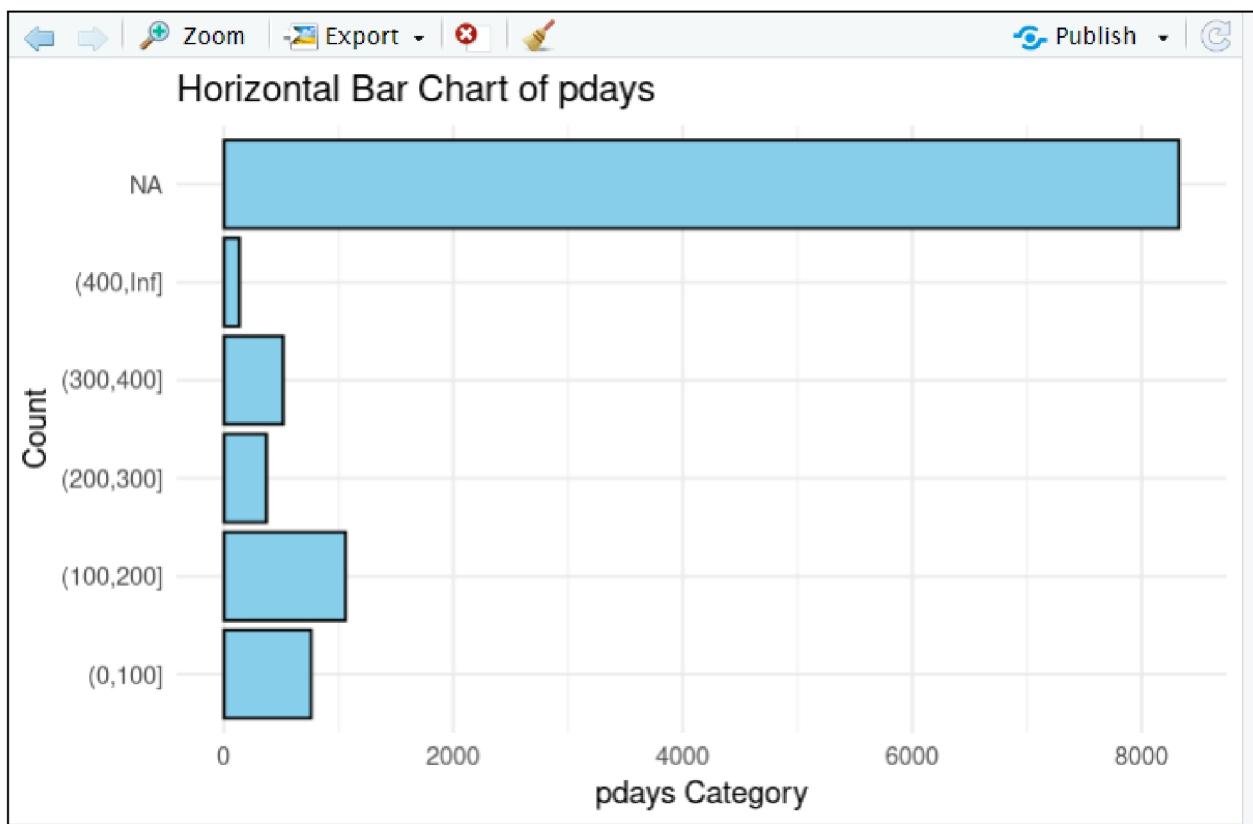
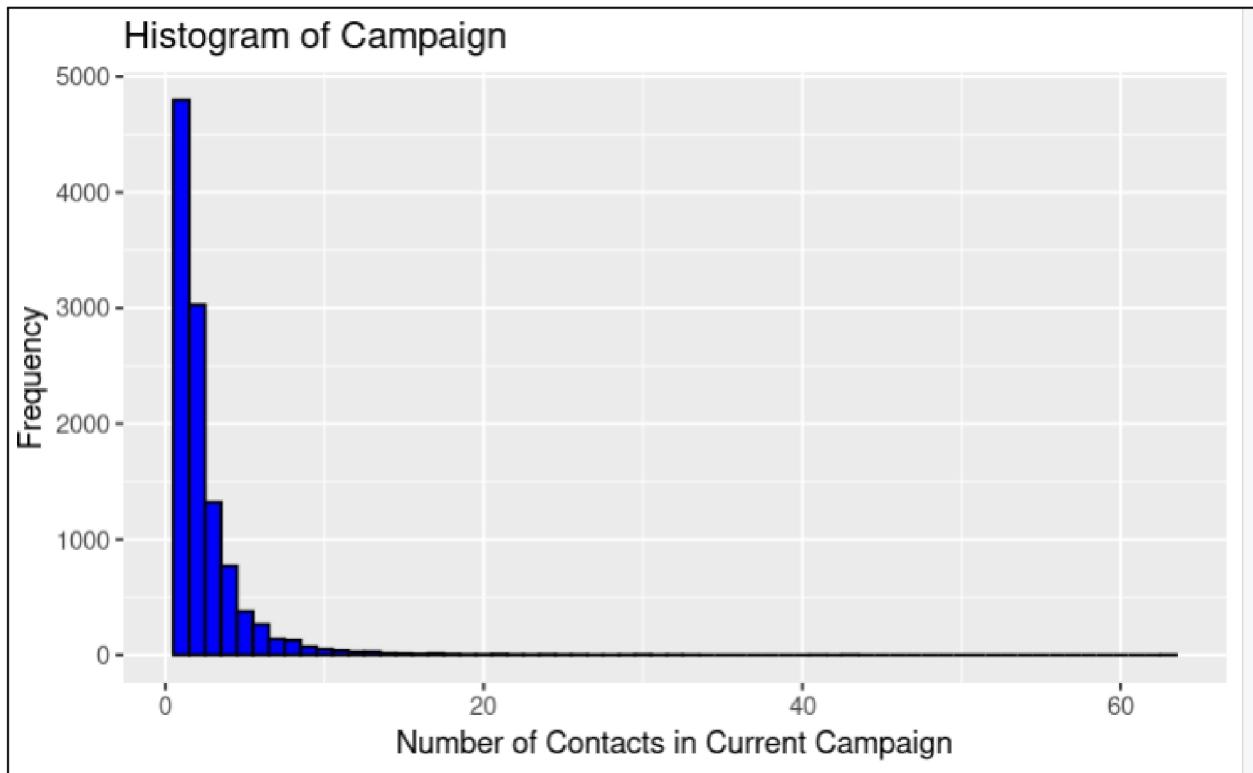
#Analysis
head(data) # Display first few rows of the data
str(data) # Get information about the data
summary(data) # statistics for each variable
names(data) # Get the names of the attributes

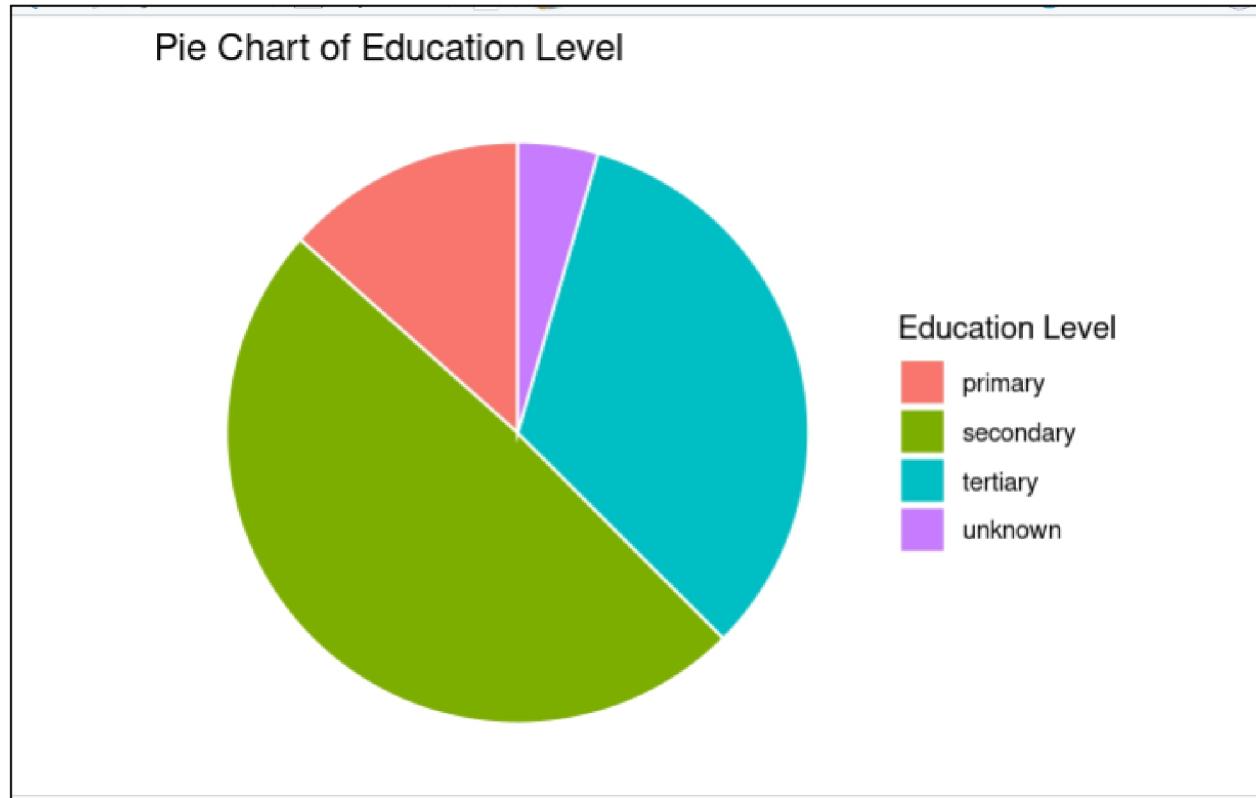
```

Output:

Visualization -







Analysis-

```
> head(data) # Show the first few rows of the data
   age      job marital education default balance housing loan contact day month
1 59 admin. married secondary no 2343 yes no unknown 5 may
2 56 admin. married secondary no 45 no no unknown 5 may
3 41 technician married secondary no 1270 yes no unknown 5 may
4 55 services married secondary no 2476 yes no unknown 5 may
5 54 admin. married tertiary no 184 no no unknown 5 may
6 42 management single tertiary no 0 yes yes unknown 5 may
  duration campaign pdays previous poutcome deposit balance_bins pdays_category
1 1042 1 -1 0 unknown yes (2e+03,3e+03] <NA>
2 1467 1 -1 0 unknown yes (0,1e+03] <NA>
3 1389 1 -1 0 unknown yes (1e+03,2e+03] <NA>
4 579 1 -1 0 unknown yes (2e+03,3e+03] <NA>
5 673 2 -1 0 unknown yes (0,1e+03] <NA>
```

```
> str(data) # Get information about the data
'data.frame': 11162 obs. of 19 variables:
 $ age : int 59 56 41 55 54 42 56 60 37 28 ...
 $ job : chr "admin." "admin." "technician" "services" ...
 $ marital : chr "married" "married" "married" "married" ...
 $ education : chr "secondary" "secondary" "secondary" "secondary" ...
 $ default : chr "no" "no" "no" "no" ...
 $ balance : int 2343 45 1270 2476 184 0 830 545 1 5090 ...
 $ housing : chr "yes" "no" "yes" "yes" ...
 $ loan : chr "no" "no" "no" "no" ...
 $ contact : chr "unknown" "unknown" "unknown" "unknown" ...
 $ day : int 5 5 5 5 5 5 6 6 6 6 ...
```

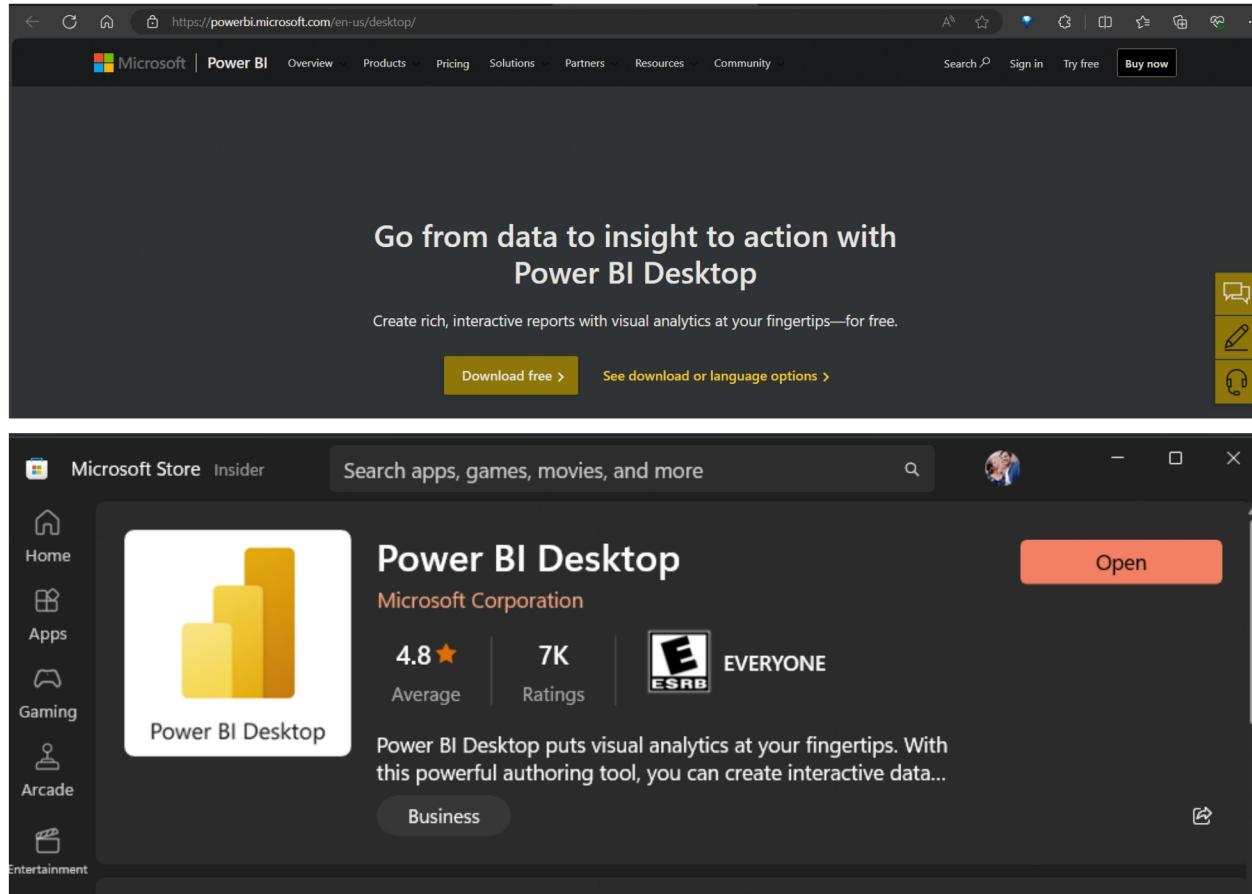
```
> summary(data) # Descriptive statistics for each variable
      age          job         marital        education
Min. :18.00  Length:11162    Length:11162    Length:11162
1st Qu.:32.00 Class :character  Class :character  Class :character
Median :39.00 Mode  :character  Mode  :character  Mode  :character
Mean   :41.23
3rd Qu.:49.00
Max.   :95.00

      default        balance        housing        loan
Length:11162    Min.   :-6847    Length:11162    Length:11162
Class :character 1st Qu.: 122     Class :character  Class :character
Mode  :character  Median : 550     Mode  :character  Mode  :character
                  Mean   : 1529
                  3rd Qu.: 1708
                  Max.   :81204
```

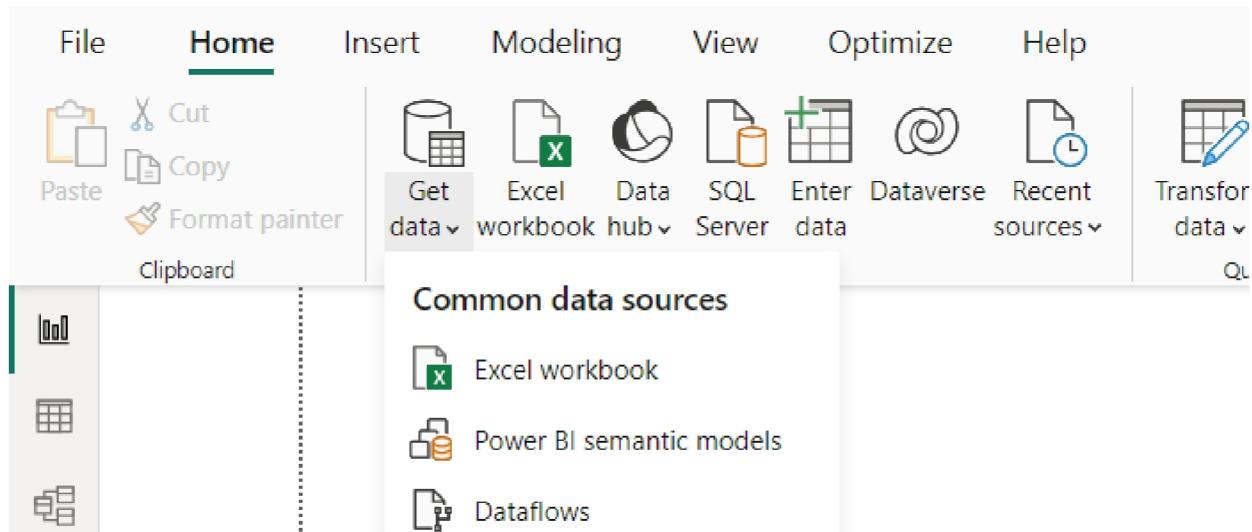
```
> names(data) # Get the names of the attributes
[1] "age"          "job"          "marital"       "education"
[5] "default"      "balance"      "housing"      "loan"
[9] "contact"      "day"          "month"        "duration"
[13] "campaign"     "pdays"        "previous"     "poutcome"
[17] "deposit"      "balance_bins" "pdays_category"
```

Task 3: Data Visualization using PowerBI

1. Install Powerbi by the link provided - <https://powerbi.microsoft.com/en-us/desktop/>



2. Install the required dataset and import it in Powerbi

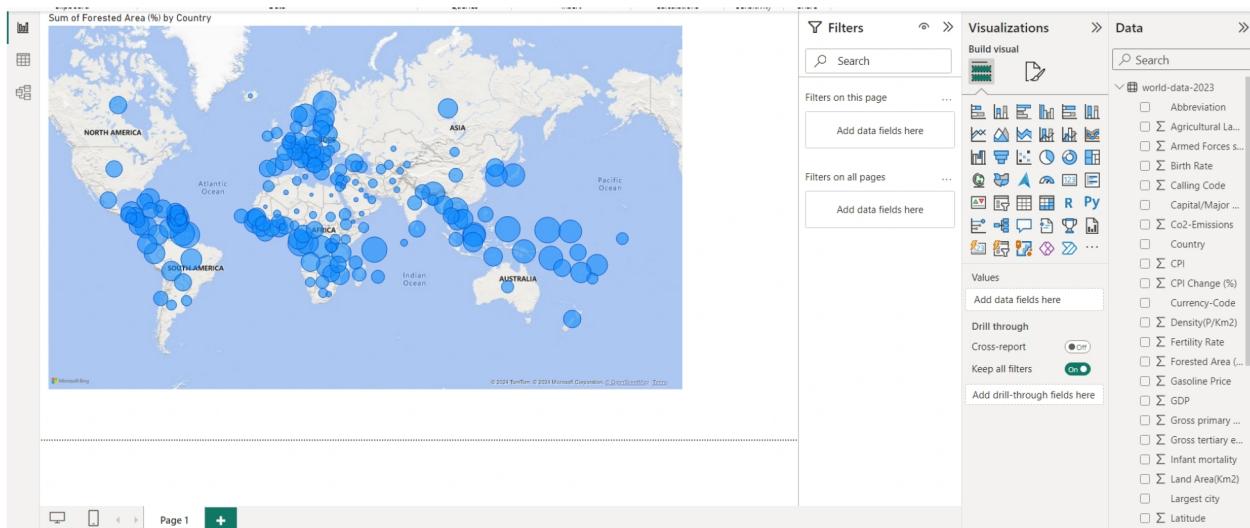


3. Click on Load Dataset.

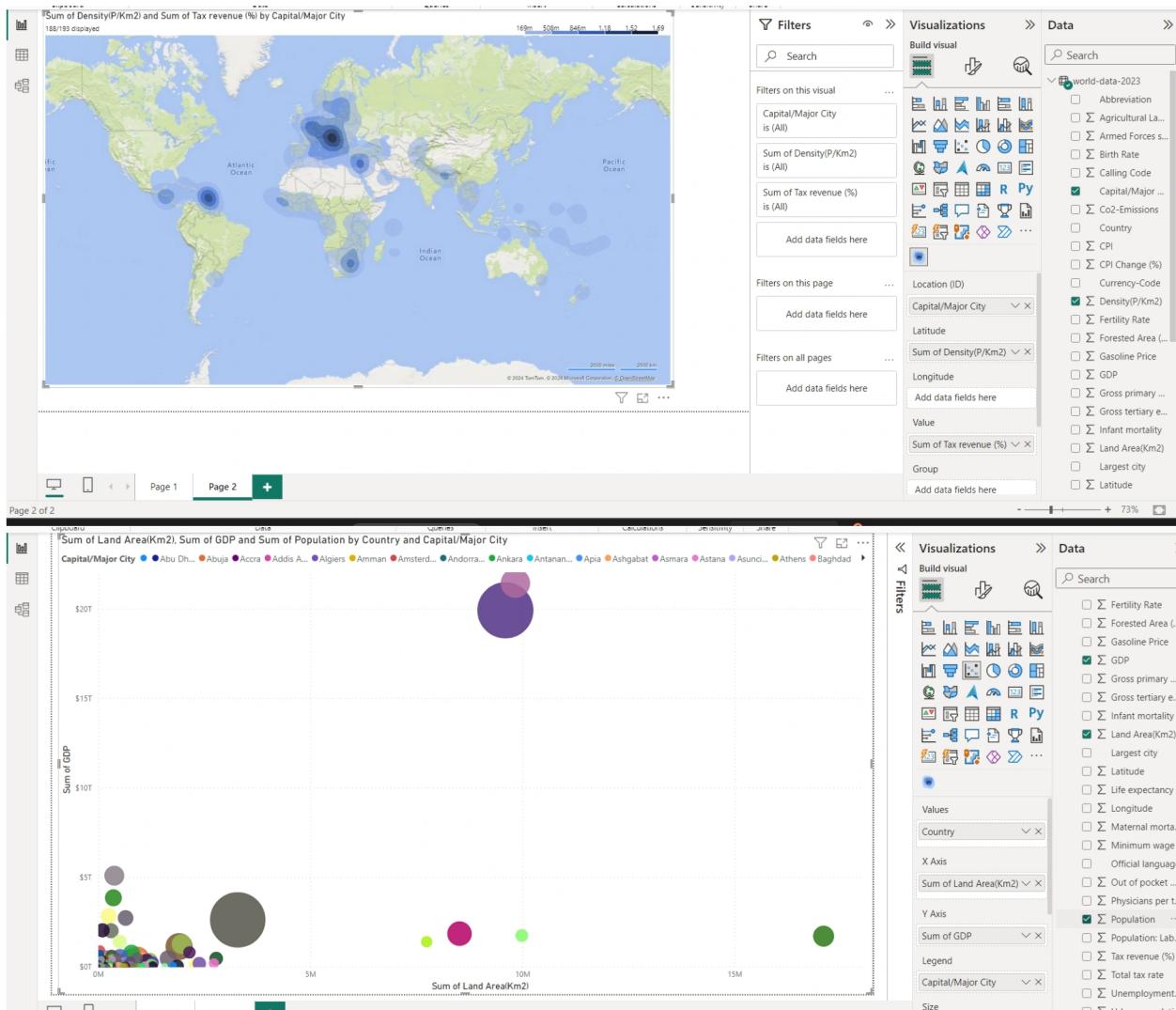
The screenshot shows the 'world-data-2023.csv' file being loaded into Power BI. The 'Data Type Detection' panel indicates the data is based on the first 200 rows. The table contains columns for Country, Density (P/Km2), Abbreviation, Agricultural Land (%), Land Area(Km2), Armed Forces size, Birth Rate, and Calling Code. The preview shows data for various countries like Afghanistan, Albania, Algeria, Andorra, Angola, Antigua and Barbuda, Argentina, Armenia, Australia, Austria, Azerbaijan, The Bahamas, Bahrain, Bangladesh, Barbados, Belarus, Belgium, Belize, Benin, and Bhutan. At the bottom, there are buttons for 'Extract Table Using Examples', 'Load', 'Transform Data', and 'Cancel'.

Country	Density (P/Km2)	Abbreviation	Agricultural Land(%)	Land Area(Km2)	Armed Forces size	Birth Rate	Calling Code
Afghanistan	60	AF	58.10%	652230	323000	32.49	
Albania	105	AL	43.10%	28748	9000	11.78	3
Algeria	18	DZ	17.40%	2381741	317000	24.28	2
Andorra	164	AD	40.00%	468	null	7.2	3
Angola	26	AO	47.50%	1246700	117000	40.73	2
Antigua and Barbuda	223	AG	20.50%	443	0	15.33	
Argentina	17	AR	54.30%	2780400	105000	17.02	
Armenia	104	AM	58.90%	29743	49000	13.99	3
Australia	3	AU	48.20%	7741220	58000	12.6	
Austria	109	AT	32.40%	83871	21000	9.7	
Azerbaijan	123	AZ	57.70%	86600	82000	14	9
The Bahamas	39	BS	1.40%	13880	1000	13.97	
Bahrain	2239	BH	11.10%	765	19000	13.99	9
Bangladesh	1265	BD	70.60%	148460	221000	18.18	8
Barbados	668	BB	23.30%	430	1000	10.65	
Belarus	47	BY	42.00%	207600	155000	9.9	3
Belgium	383	BE	44.60%	30528	32000	10.3	
Belize	17	BZ	7.00%	22966	2000	20.79	5
Benin	108	BJ	33.30%	112622	12000	36.22	2
Bhutan	20	BT	13.60%	38394	6000	17.26	9

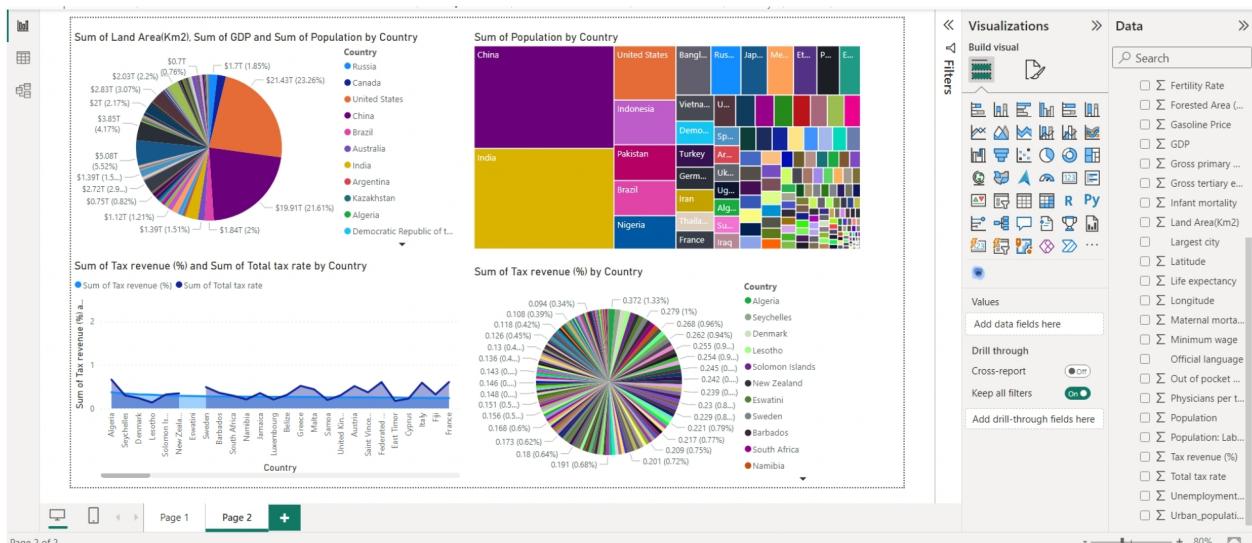
3. Generate the map according to country density.



4. Generate the Heatmap according to country density/poulation or land area.



6. Do visualization of the dataset by various graphs like pie chart, line graph, etc.



TASK-4 CREATE and manipulate database in SPARK

Steps:

1. Install Spark in Google Collaboratory.

```
1 !pip install gspread-pandas
2 !pip install pyspark

[ ] 1 from pyspark.sql import SparkSession
2
3 spark = SparkSession.builder\
4     .master("local")\
5     .appName("Colab")\
6     .config('spark.ui.port', '4050')\
7     .getOrCreate()
8 spark

SparkSession - in-memory
SparkContext
Spark UI
Version
v3.5.1
Master
local
AppName
Colab
```

2. Use RDD & Dataframes.

```
1 # Import necessary libraries
2 import gspread
3 from google.auth import default
4 from pyspark.sql.types import StructType, StructField, StringType
5
6 from google.colab import auth
7 auth.authenticate_user()
8
9 # from google.colab import drive
10 # drive.mount('/content/drive')
11
12 # Get Google Sheets credentials
13 creds, _ = default()
14
15 # Authorize gspread
16 gc = gspread.authorize(creds)
17
18 # Open the Google Sheets document by key
19 sheets_key = '1EuNWRKUOj0waKIgxf_hwstrFQezmRpzgtIYpkNnJj5M'
20 worksheet = gc.open_by_key(sheets_key).sheet1 # Adjust the sheet index as needed
21
22 # Get all values from the worksheet
23 data = worksheet.get_all_values()
24
25 # Define the Spark session
26 spark = SparkSession.builder.appName("GoogleSheetsToPySpark").getOrCreate()
27
28 # Define the schema based on the header
29 header = data[0]
30 schema = StructType([StructField(col, StringType(), True) for col in header])
31
32 # Create a PySpark DataFrame
33 rdd = spark.sparkContext.parallelize(data[1:])
34 df = spark.createDataFrame(rdd, schema)
35
36 # Show the PySpark DataFrame
37 df.show()
```

Mounted at /content/drive

	Serial No.	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
1	337	118	4	4.5	4.5	9.65	1	1	1
2	324	107	4	4	4.5	8.87	1	1	0.76
3	316	104	3	3	3.5	8	1	1	0.72
4	322	110	3	3.5	2.5	8.67	1	1	0.8
5	314	103	2	2	3	8.21	0	0	0.65
6	338	115	5	4.5	3	9.34	1	1	0.9
7	321	109	3	3	4	8.2	1	1	0.75
8	308	101	2	3	4	7.9	0	0	0.68
9	302	102	1	2	1.5	8	0	0	0.5
10	323	108	3	3.5	3	8.6	0	0	0.45
11	325	106	3	3.5	4	8.4	1	1	0.52
12	327	111	4	4	4.5	9	1	1	0.84
13	328	112	4	4	4.5	9.1	1	1	0.78
14	307	109	3	4	3	8	1	1	0.62
15	311	104	3	3.5	2	8.2	1	1	0.61
16	314	105	3	3.5	2.5	8.3	0	0	0.54
17	317	107	3	4	3	8.7	0	0	0.66
18	319	106	3	4	3	8	1	1	0.65
19	318	110	3	4	3	8.8	0	0	0.63
20	303	102	3	3.5	3	8.5	0	0	0.62

only showing top 20 rows

3. Implement linear regression using sparkMlib.

```

1  from pyspark.ml.feature import VectorAssembler
2  from pyspark.ml.regression import LinearRegression
3  from pyspark.ml import Pipeline
4
5  # Select features and target column
6  feature_columns = ['GRE Score', 'TOEFL Score', 'University Rating', 'SOP', 'LOR', 'CGPA', 'Research']
7  target_column = 'Chance of Admit'
8
9  from pyspark.sql.functions import col
10 from pyspark.sql.types import FloatType
11
12 # Convert columns to appropriate numeric types
13 for col_name in feature_columns + [target_column]:
14     df = df.withColumn(col_name, col(col_name).cast(FloatType()))
15
16 # Verify the schema after conversion
17 df.printSchema()
18
19
20 # Create a VectorAssembler to assemble the features into a single vector column
21 assembler = VectorAssembler(inputCols=feature_columns, outputCol="features")
22
23 # Create a Linear Regression model
24 lr = LinearRegression(featuresCol="features", labelCol=target_column)
25
26 # Create a pipeline with the VectorAssembler and Linear Regression model
27 pipeline = Pipeline(stages=[assembler, lr])
28
29 # Fit the pipeline to the data
30 model = pipeline.fit(df)
31
32 # Make predictions
33 predictions = model.transform(df)
34
35 # Show the predictions
36 predictions.select("features", target_column, "prediction").show()
37

```

```

root
|-- Serial No.: string (nullable = true)
|-- GRE Score: float (nullable = true)
|-- TOEFL Score: float (nullable = true)
|-- University Rating: float (nullable = true)
|-- SOP: float (nullable = true)
|-- LOR: float (nullable = true)
|-- CGPA: float (nullable = true)
|-- Research: float (nullable = true)
|-- Chance of Admit: float (nullable = true)

+-----+-----+
|      features|Chance of Admit|      prediction|
+-----+-----+
|[337.0,118.0,4.0, ...|      1.0|0.9527419787810929|
|[324.0,107.0,4.0, ...|      0.76|0.8042442709169288|
|[316.0,104.0,3.0, ...|      0.72| 0.653362064831734|
|[322.0,110.0,3.0, ...|      0.8| 0.744800849900342|
|[314.0,103.0,2.0, ...|      0.65| 0.631671909170074|
|[330.0,115.0,5.0, ...|      0.9|0.8748625451723717|
|[321.0,109.0,3.0, ...|      0.75|0.7089061133853698|
|[308.0,101.0,2.0, ...|      0.68|0.5965489750932171|
|[302.0,102.0,1.0, ...|      0.5|0.5503679652106606|
|[323.0,108.0,3.0, ...|      0.45|0.7169636074292443|
|[325.0,106.0,3.0, ...|      0.52|0.7325297349799527|
|[327.0,111.0,4.0, ...|      0.84|0.836486918374747|
|[328.0,112.0,4.0, ...|      0.78|0.8530258541953695|
|[307.0,109.0,3.0, ...|      0.62|0.6436557091259056|
|[311.0,104.0,3.0, ...|      0.61|0.6431819652204682|
|[314.0,105.0,3.0, ...|      0.54|0.6476620556847543|
|[317.0,107.0,3.0, ...|      0.66| 0.715595016570292|
|[319.0,106.0,3.0, ...|      0.65|0.6577263403901585|
|[318.0,110.0,3.0, ...|      0.63|0.7377331438983201|
|[303.0,102.0,3.0, ...|      0.62|0.6508501135082854|
+-----+-----+
only showing top 20 rows

```

```

1  from pyspark.ml.evaluation import RegressionEvaluator
2
3  # Evaluate the model using Mean Squared Error (MSE)
4  evaluator = RegressionEvaluator(labelCol=target_column, predictionCol="prediction", metricName="mse")
5  mse = evaluator.evaluate(predictions)
6  print(f"Mean Squared Error (MSE): {mse}")
7
8  # Evaluate the model using R-squared
9  evaluator = RegressionEvaluator(labelCol=target_column, predictionCol="prediction", metricName="r2")
10 r2 = evaluator.evaluate(predictions)
11
12 # Convert R-squared to percentage
13 r2_percentage = r2 * 100
14
15 print(f"R-squared: {r2_percentage:.2f}%")

```

Mean Squared Error (MSE): 0.0035432110127672583
R-squared: 82.25%

Full Google Colab

TASK-5 MLops using Heroku and Github actions

- Submitted by Shobhit (2021UCS1618)

1. Make a new project - DEMO using ‘Anaconda’.

```
Select Anaconda Prompt (anaconda3)

(base) E:\>conda activate flask_demo

(flask_demo) E:\>cd __Demo
```

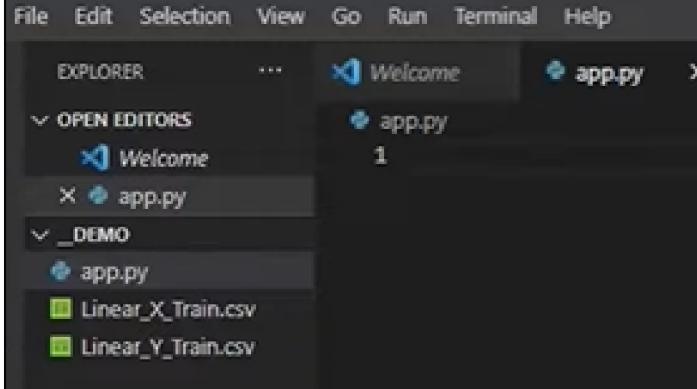
2. Install flask using ‘Anaconda pip’.

```
(flask_demo) E:\__Demo>pip install flask
Requirement already satisfied: flask in c:
```

Dataset and Preprocessing

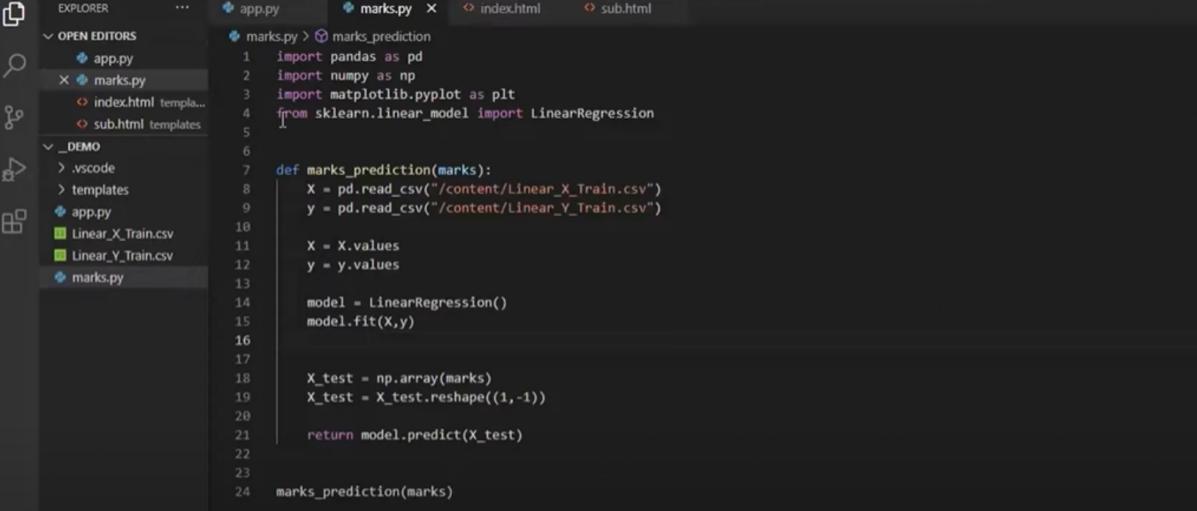
3. Import Class Score dataset (further used in linear regression) in this folder.

	A	B	
1	Hours	Scores	
2	2.5	21	
3	5.1	47	
4	3.2	27	
5	8.5	75	
6	3.5	30	
7	1.5	20	
8	9.2	88	



Encoding and Machine Learning Model

4. Create a new py file marks.py using that we will create and run the regression model.



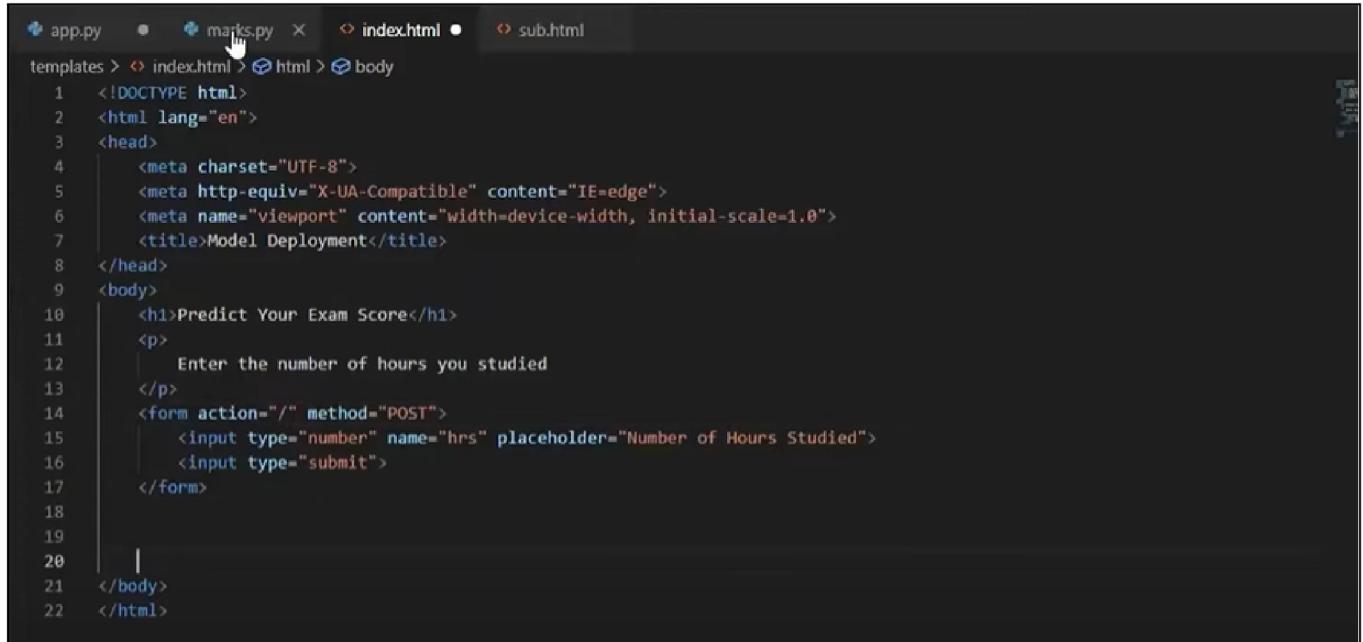
```
EXPLORER ... app.py marks.py index.html templates sub.html

marks.py > marks_prediction
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.linear_model import LinearRegression
5
6 def marks_prediction(marks):
7     X = pd.read_csv("/content/Linear_X_Train.csv")
8     y = pd.read_csv("/content/Linear_Y_Train.csv")
9
10    X = X.values
11    y = y.values
12
13    model = LinearRegression()
14    model.fit(X,y)
15
16
17    X_test = np.array(marks)
18    X_test = X_test.reshape((1,-1))
19
20
21    return model.predict(X_test)
22
23
24 marks_prediction(marks)
```

Deployment

5. Create 3 files - index.html, sub.html & app.py.

These files will be used to deploy our model.



```
app.py  ●  marks.py  ✘  index.html ●  sub.html
templates > index.html > html > body
1  <!DOCTYPE html>
2  <html lang="en">
3  <head>
4      <meta charset="UTF-8">
5      <meta http-equiv="X-UA-Compatible" content="IE=edge">
6      <meta name="viewport" content="width=device-width, initial-scale=1.0">
7      <title>Model Deployment</title>
8  </head>
9  <body>
10     <h1>Predict Your Exam Score</h1>
11     <p>
12         Enter the number of hours you studied
13     </p>
14     <form action="/" method="POST">
15         <input type="number" name="hrs" placeholder="Number of Hours Studied">
16         <input type="submit">
17     </form>
18
19
20
21 </body>
22 </html>
```

Index.html

```
from flask import Flask, render_template, request
import marks

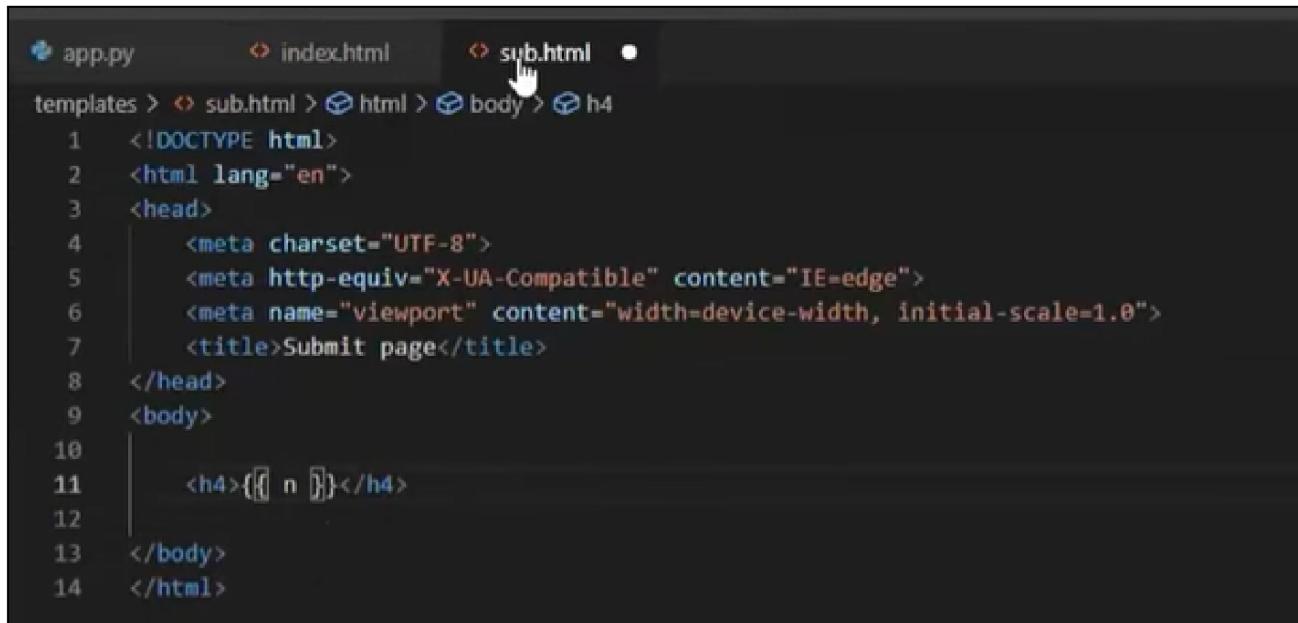
app = Flask(__name__)

@app.route("/", methods=['POST'])
def hello():
    if request.method == "POST":
        hrs = request.form['hrs']
        marks_pred = marks.marks_prediction(hrs)
        print(marks_pred)

    return render_template("index.html")

if __name__ == "__main__":
    app.run(debug=True)
```

app.py



```

app.py          index.html      sub.html •
templates > sub.html > html > body > h4
1  <!DOCTYPE html>
2  <html lang="en">
3  <head>
4      <meta charset="UTF-8">
5      <meta http-equiv="X-UA-Compatible" content="IE=edge">
6      <meta name="viewport" content="width=device-width, initial-scale=1.0">
7      <title>Submit page</title>
8  </head>
9  <body>
10
11      <h4>{{ n }}</h4>
12
13  </body>
14  </html>

```

sub.html

- On running app.py, our server will get started on localhost:5000.

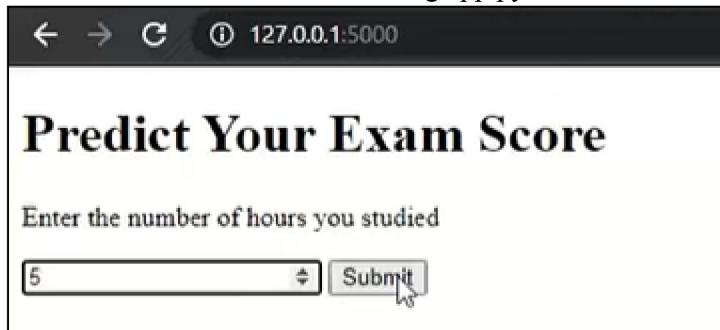


```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL
NameError: name 'marks' is not defined
PS E:\_Demo> & "C:/Users/Syed Junaid Iqbal/anaconda3/envs/Flask_Demo/python.exe" e:/_Demo/app.py
* Serving Flask app "app" (lazy loading)
* Environment: production
WARNING: This is a development server. Do not use it in a production deployment.
Use a production WSGI server instead.
* Debug mode: on
* Restarting with stat
* Debugger is active!
* Debugger PIN: 111-987-899
* Running on http://127.0.0.1:5000/ (Press CTRL+C to quit)

```

- Go to localhost:5000 after re-running app.py and enter number of hours you studied.



Predict Your Exam Score

Enter the number of hours you studied

- It will show predicted marks.

Predict Your Exam Score

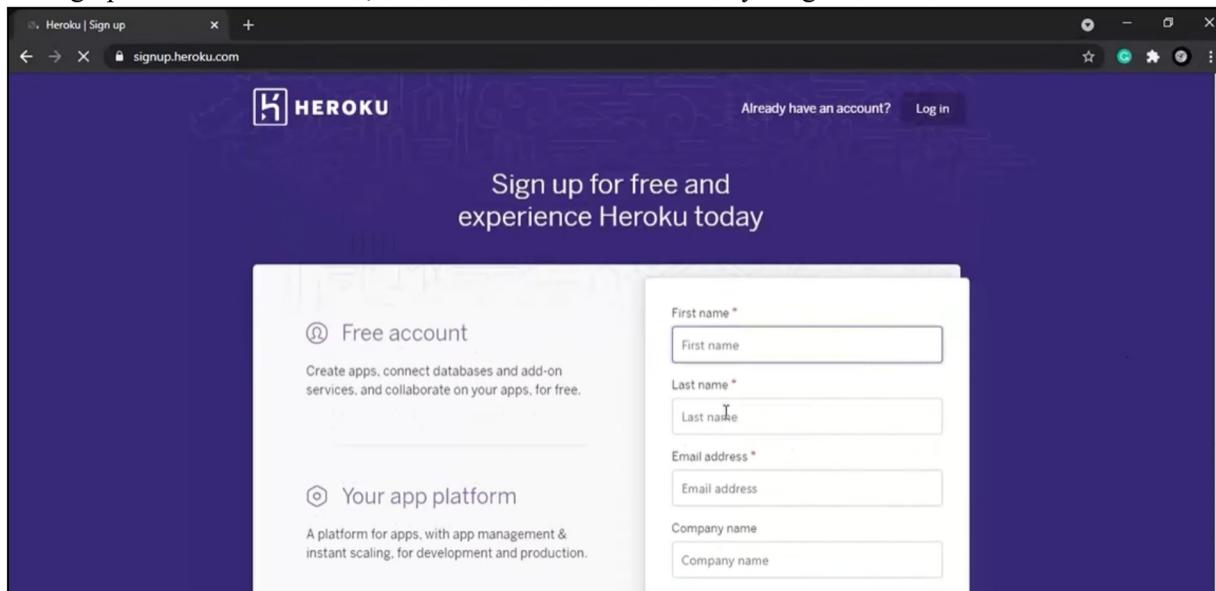
Enter the number of hours you studied

Your predicted marks are

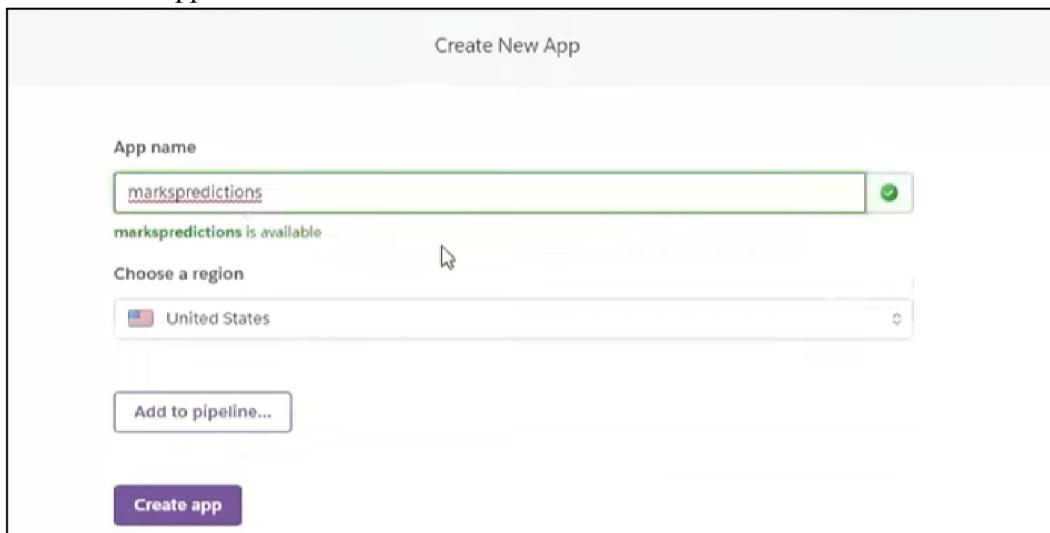
[3.72798828] 

Deployment on Heroku

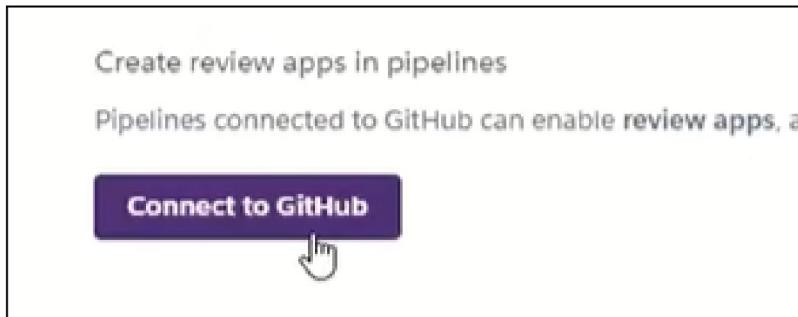
9. Setting up heroku - On Heroku, create new account & connect your github to it.



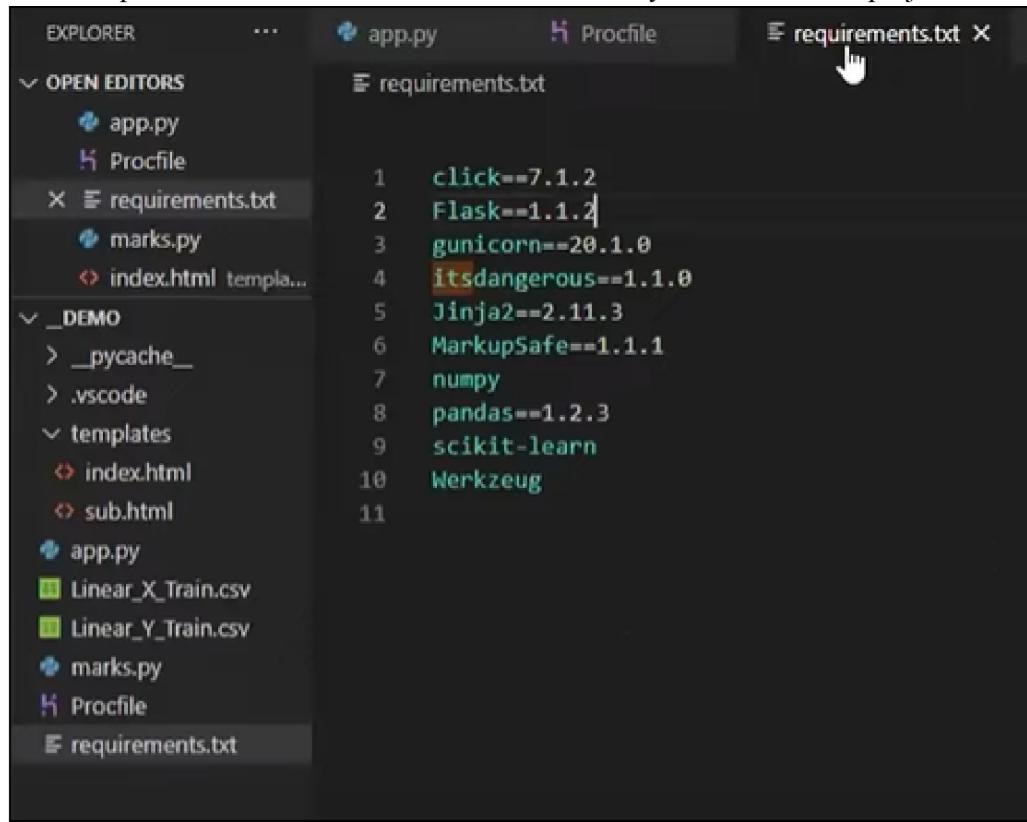
10. Create a new app on heroku



11. Connect to Github.



12. Create requirements.txt which will contain all necessary modules for this project.



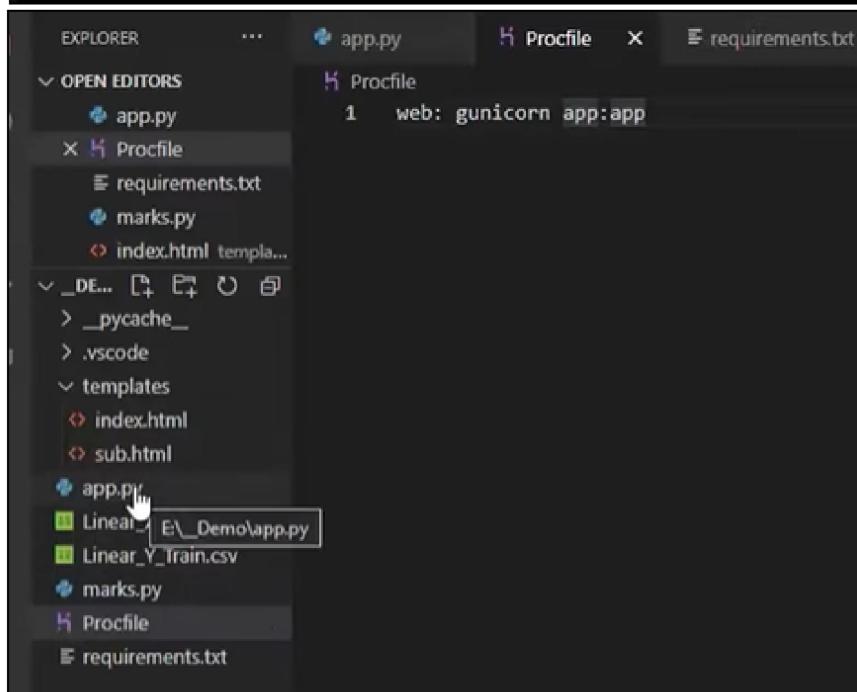
The screenshot shows the VS Code interface with the 'OPEN EDITORS' sidebar on the left. The 'requirements.txt' file is open in the editor tab, displaying the following dependencies:

```
click==7.1.2
Flask==1.1.2
gunicorn==20.1.0
itsdangerous==1.1.0
Jinja2==2.11.3
MarkupSafe==1.1.1
numpy
pandas==1.2.3
scikit-learn
Werkzeug
```

13. Install gunicorn in the environment and create Procfile.

```
(flask_demo) E:\__Demo>pip freeze > requirements.txt

(flask_demo) E:\__Demo>pip install gunicorn
Requirement already satisfied: gunicorn in c:\users\syed junaid iqbal\anaconda3\envs\flask_demo\lib
e-packages (20.1.0)
Requirement already satisfied: setuptools>=3.0 in c:\users\syed junaid iqbal\anaconda3\envs\flask_d
lib\site-packages (from gunicorn) (52.0.0.post20210125)
```



The screenshot shows the VS Code interface with the 'OPEN EDITORS' sidebar on the left. The 'Procfile' file is open in the editor tab, containing the following configuration:

```
web: gunicorn app:app
```

Incorporating GitHub actions

14. Create a new repository on github in which we will upload our project.

Create a new repository

A repository contains all project files, including the revision history. Already have a project repository elsewhere? [Import a repository.](#)

Required fields are marked with an asterisk (*).

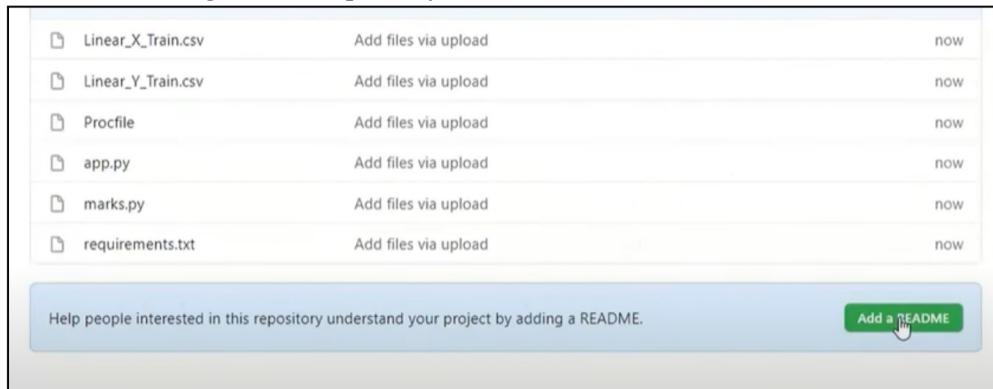
Owner * Repository name *

 prachisah07 / Model

Model is available.

Great repository names are short and memorable. Need inspiration? How about [psychic-fiesta](#) ?

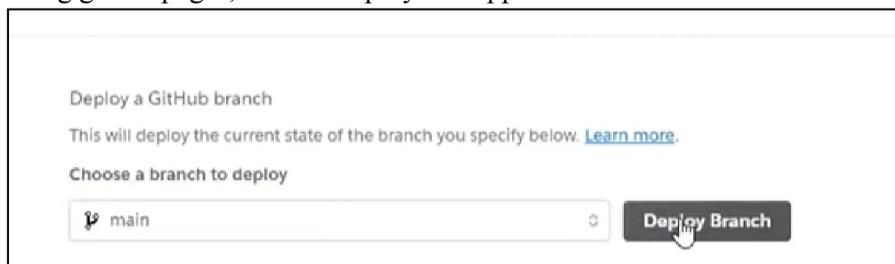
15. Commit the changes in the repository.



16. Search for your repository in github section of heroku and connect it.



17. Using github pages, we will deploy this app on heroku.



18. The app has been deployed on heroku and the link has been generated for the web page.

