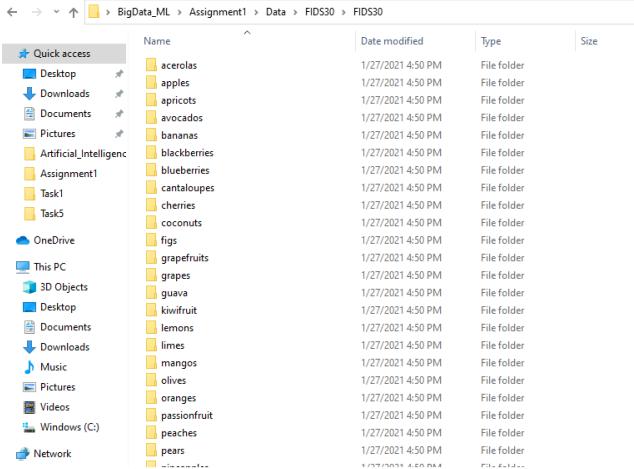


2 DOWNLOAD OR GENERATE A FRUITS AND VEGETABLES IMAGE DATA SET!

2.1 Download FIDS30 Data Set

Downloaded and unzipped the FIDS30 data from the <https://www.vicos.si/Downloads/FIDS30/>. [1] Figure 5 shows the screenshot of the data folder in my local system.



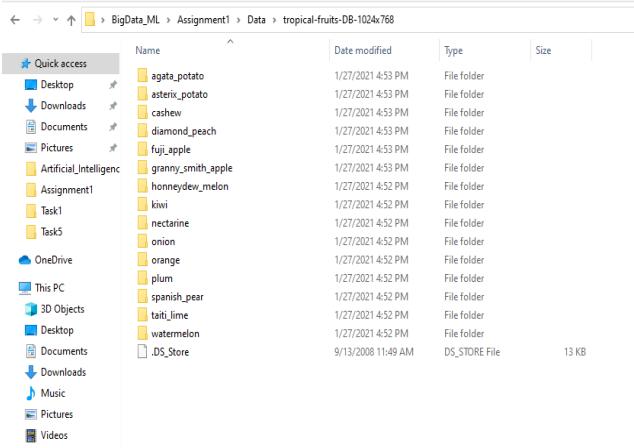
A screenshot of a Windows file explorer window. The path is 'BigData_ML > Assignment1 > Data > FIDS30 > FIDS30'. The table lists various fruit names as subfolders, all created on 1/27/2021 at 4:50 PM. The columns are Name, Date modified, Type, and Size.

Name	Date modified	Type	Size
acerolas	1/27/2021 4:50 PM	File folder	
apples	1/27/2021 4:50 PM	File folder	
apricots	1/27/2021 4:50 PM	File folder	
avocados	1/27/2021 4:50 PM	File folder	
bananas	1/27/2021 4:50 PM	File folder	
blackberries	1/27/2021 4:50 PM	File folder	
blueberries	1/27/2021 4:50 PM	File folder	
cantaloupes	1/27/2021 4:50 PM	File folder	
cherries	1/27/2021 4:50 PM	File folder	
coconuts	1/27/2021 4:50 PM	File folder	
figs	1/27/2021 4:50 PM	File folder	
grapefruits	1/27/2021 4:50 PM	File folder	
grapes	1/27/2021 4:50 PM	File folder	
guava	1/27/2021 4:50 PM	File folder	
kiwi_fruit	1/27/2021 4:50 PM	File folder	
lemons	1/27/2021 4:50 PM	File folder	
limes	1/27/2021 4:50 PM	File folder	
mangos	1/27/2021 4:50 PM	File folder	
olives	1/27/2021 4:50 PM	File folder	
oranges	1/27/2021 4:50 PM	File folder	
passionfruit	1/27/2021 4:50 PM	File folder	
peaches	1/27/2021 4:50 PM	File folder	
pears	1/27/2021 4:50 PM	File folder	

Figure 5: FIDS30 Data folder

2.2 Download Tropical Fruits Data set

Downloaded Tropical fruits data from <http://www.ic.unicamp.br/~rocha/pub/downloads/tropical-fruits-DB-1024x768.tar.gz>, unzipped the folder and stored in my local system. Figure 6 shows the screenshot the data folder I have created in my system.



A screenshot of a Windows file explorer window. The path is 'BigData_ML > Assignment1 > Data > tropical-fruits-DB-1024x768'. The table lists various fruit names as subfolders, all created on 1/27/2021 at 4:53 PM. The columns are Name, Date modified, Type, and Size.

Name	Date modified	Type	Size
agata_potato	1/27/2021 4:53 PM	File folder	
asterix_potato	1/27/2021 4:53 PM	File folder	
cashew	1/27/2021 4:53 PM	File folder	
diamond_peach	1/27/2021 4:53 PM	File folder	
fiji_apple	1/27/2021 4:53 PM	File folder	
granny_smith_apple	1/27/2021 4:53 PM	File folder	
honeydew_melon	1/27/2021 4:52 PM	File folder	
kiwi	1/27/2021 4:52 PM	File folder	
nectarine	1/27/2021 4:52 PM	File folder	
onion	1/27/2021 4:52 PM	File folder	
orange	1/27/2021 4:52 PM	File folder	
plum	1/27/2021 4:52 PM	File folder	
spanish_pear	1/27/2021 4:52 PM	File folder	
tai_t_lime	1/27/2021 4:52 PM	File folder	
watermelon	1/27/2021 4:52 PM	File folder	
DS_Store	9/13/2008 11:49 AM	DS_STORE File	13 KB

Figure 6: Tropical fruits Data folder

2.3 Choose three fruits.

I am choosing three fruits Banana, Guava and Pineapple. I am assigning labels image0 to banana, image1 to guava and image2 to pineapple to explain the rest of the assignment. Figure 7 to Figure 9 shows the images of the fruits I will be using for the feature extraction process.



Figure 7: image0



Figure 8: image1



Figure 9: image2

```

import cv2
import matplotlib.pyplot as plt

#read the image
banana_color = cv2.imread('banana.jpg')
guava_color = cv2.imread('guava.jpg')
pineapple_color = cv2.imread('pineapple.jpg')

#display the color channels
plt.axis('off')

#RGB for banana
plt.imshow(banana_color[:, :, 0])
plt.imshow(banana_color[:, :, 1])
plt.imshow(banana_color[:, :, 2])

#RGB for guava
plt.imshow(guava_color[:, :, 0])
plt.imshow(guava_color[:, :, 1])
plt.imshow(guava_color[:, :, 2])

#RGB for pineapple
plt.imshow(pineapple_color[:, :, 0])
plt.imshow(pineapple_color[:, :, 1])
plt.imshow(pineapple_color[:, :, 2])

```

Figure 10: Code for generating RGB scale images

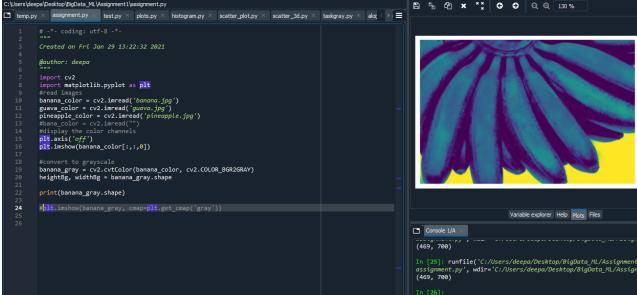


Figure 11: Red scale for image0

3 WRITE A SIMPLE CODE TO READ YOUR SELECTED IMAGES AND DISPLAY THEM ON THE ENVIRONMENT!

I am using opencv to read the images and matplotlib.pyplot to plot the images in python environment. I have generated RGB images for image0, image1 and image2 using the code shown in Figure 10. I have ran the code separately once passing R, then G and then B to get the respective scales for each image. Gray scale images for each fruit is generated separately by executing the code three times and corresponding dimensions are printed in console. Fig 11 to Figure 19 shows the execution of RGB code and Figure 20 to 22 corresponds to the Gray scale image code executed.

4 RESIZE THE IMAGES TO REDUCE THEIR DIMENSIONS!

4.1 Convert Image dimension divisible by 8

I have written a parametric function to calculate the divisibility by 8. Initially I have halved the dimensions in order to reduce the dimension. Then, I am passing these halved dimensions to the

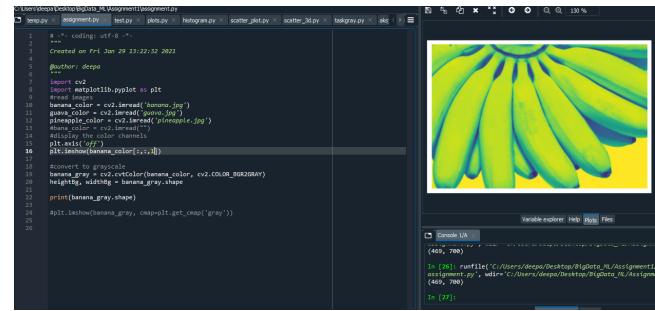


Figure 12: Green scale for image0

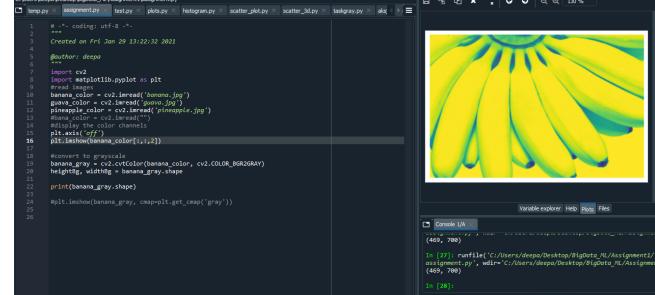


Figure 13: Blue scale for image0

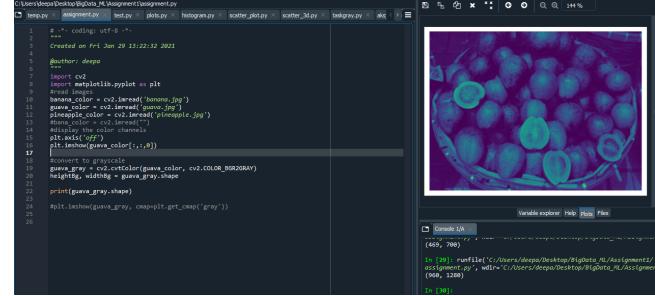


Figure 14: Red scale for image1

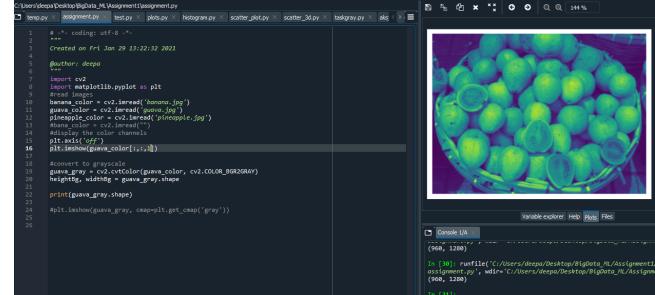


Figure 15: Green scale for image1

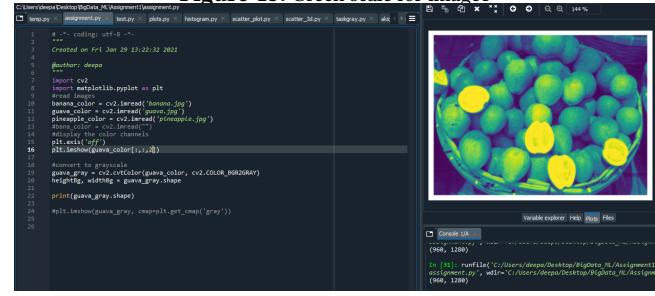


Figure 16: Blue scale for image1

```

1 # -*- coding: utf-8 -*-
2
3 # Created on Fri Jan 29 13:22:32 2021
4
5 #Author: deepa
6
7 import cv2
8 import matplotlib.pyplot as plt
9
10 read_image = cv2.imread('banana.jpg')
11 banana_color = cv2.cvtColor(banana_color, cv2.COLOR_BGR2GRAY)
12 pineapple_color = cv2.imread('pineapple.jpg')
13
14 plt.axis('off')
15 plt.imshow(np.concatenate([banana_color, pineapple_color], axis=0))
16
17 print(pineapple_color.shape)
18
19 #convert to grayscale
20 guava_gray = cv2.cvtColor(pineapple_color, cv2.COLOR_BGR2GRAY)
21 height, width = pineapple_gray.shape
22
23 print(pineapple_gray.shape)
24
25 plt.imshow(pineapple_gray, cmap=plt.get_cmap('gray'))
26
27 cv2.imwrite('pineapple_gray.jpg',pineapple_gray)
28
29

```

Figure 17: Red scale for image2

```

1 # -*- coding: utf-8 -*-
2
3 # Created on Fri Jan 29 13:22:32 2021
4
5 #Author: deepa
6
7 import cv2
8 import matplotlib.pyplot as plt
9
10 read_image = cv2.imread('banana.jpg')
11 banana_color = cv2.cvtColor(banana_color, cv2.COLOR_BGR2GRAY)
12 pineapple_color = cv2.imread('pineapple.jpg')
13
14 plt.imshow(np.concatenate([banana_color, pineapple_color], axis=0))
15
16 print(pineapple_color.shape)
17
18 #convert to grayscale
19 guava_gray = cv2.cvtColor(pineapple_color, cv2.COLOR_BGR2GRAY)
20 height, width = pineapple_gray.shape
21
22 print(pineapple_gray.shape)
23
24 plt.imshow(pineapple_gray, cmap=plt.get_cmap('gray'))
25
26 cv2.imwrite('pineapple_gray.jpg',pineapple_gray)
27
28

```

Figure 18: Green scale for image2

```

1 # -*- coding: utf-8 -*-
2
3 # Created on Fri Jan 29 13:22:32 2021
4
5 #Author: deepa
6
7 import cv2
8 import matplotlib.pyplot as plt
9
10 read_image = cv2.imread('banana.jpg')
11 banana_color = cv2.cvtColor(banana_color, cv2.COLOR_BGR2GRAY)
12 pineapple_color = cv2.imread('pineapple.jpg')
13
14 plt.imshow(np.concatenate([banana_color, pineapple_color], axis=0))
15
16 print(pineapple_color.shape)
17
18 #convert to grayscale
19 guava_gray = cv2.cvtColor(pineapple_color, cv2.COLOR_BGR2GRAY)
20 height, width = pineapple_gray.shape
21
22 print(pineapple_gray.shape)
23
24 plt.imshow(pineapple_gray, cmap=plt.get_cmap('gray'))
25
26 cv2.imwrite('pineapple_gray.jpg',pineapple_gray)
27
28

```

Figure 19: Blue scale for image2

```

1 # -*- coding: utf-8 -*-
2
3 # Created on Sat Feb 17 21:06:57 2021
4
5 #Author: deepa
6
7 import cv2
8 import matplotlib.pyplot as plt
9
10 read_image = cv2.imread('banana.jpg')
11 banana_color = cv2.cvtColor(banana_color, cv2.COLOR_BGR2GRAY)
12 guava_color = cv2.imread('guava.jpg')
13 pineapple_color = cv2.imread('pineapple.jpg')
14
15 plt.axis('off')
16
17 banana_gray = cv2.cvtColor(banana_color, cv2.COLOR_BGR2GRAY)
18
19 plt.imshow(np.concatenate([banana_gray, guava_color], axis=0))
20
21 print(pineapple_color.shape)
22
23 height, width = guava_color.shape[1]/2
24
25 print('Dimension of Banana/Imaged: ',height, width)
26 print('Dimension of Guava/Imaged: ',height, width)
27 print('Dimension of Pineapple/Imaged: ',height, width)
28
29

```

Figure 20: Gray image Banana with dimension printed

function to check their divisibility by 8. If they are not divisible, by using the remainder I am rounding them off them to nearest values that are divisible by 8. Figure 23 shows the code I have used to accomplish the task. Initial dimension of image0 was 700*469, image1 was 1280*960 and image2 was 1329*1000 and the resized image values are 352*232, 640*480 and 664*496 respectively. These images with new dimensions are uploaded in Data folder as proof of this task completion.

```

1 # -*- coding: utf-8 -*-
2
3 # Created on Sat Feb 17 21:06:57 2021
4
5 #Author: deepa
6
7 import cv2
8 import matplotlib.pyplot as plt
9
10 read_image = cv2.imread('banana.jpg')
11 banana_color = cv2.cvtColor(banana_color, cv2.COLOR_BGR2GRAY)
12
13 plt.imshow(np.concatenate([banana_color, guava_color], axis=0))
14
15 print(pineapple_color.shape)
16
17 banana_gray = cv2.cvtColor(banana_color, cv2.COLOR_BGR2GRAY)
18
19 plt.imshow(np.concatenate([banana_gray, guava_color], axis=0))
20
21 print(pineapple_color.shape)
22
23 height, width = guava_color.shape[1]/2
24
25 print('Dimension of Guava/Imaged: ',height, width)
26 print('Dimension of Pineapple/Imaged: ',height, width)
27
28

```

Figure 21: Gray image Guava with dimension printed

```

1 # -*- coding: utf-8 -*-
2
3 # Created on Sat Feb 17 21:06:57 2021
4
5 #Author: deepa
6
7 import cv2
8 import matplotlib.pyplot as plt
9
10 read_image = cv2.imread('banana.jpg')
11 banana_color = cv2.cvtColor(banana_color, cv2.COLOR_BGR2GRAY)
12
13 plt.imshow(np.concatenate([banana_color, guava_color], axis=0))
14
15 print(pineapple_color.shape)
16
17 banana_gray = cv2.cvtColor(banana_color, cv2.COLOR_BGR2GRAY)
18
19 plt.imshow(np.concatenate([banana_gray, guava_color], axis=0))
20
21 print(pineapple_color.shape)
22
23 height, width = guava_color.shape[1]/2
24
25 print('Dimension of Banana/Imaged: ',height, width)
26 print('Dimension of Guava/Imaged: ',height, width)
27 print('Dimension of Pineapple/Imaged: ',height, width)
28
29

```

Figure 22: Gray image Pineapple with dimension printed

```

1 import cv2
2
3 #read images
4 fruitGray = cv2.imread('guava_gray.jpg')
5 #Displaying image in gray scale
6 print('Dimension of Guava/Imaged: ',fruitGray.shape[1]/2)
7 height_fruit = int(fruitGray.shape[1]/2)
8 print('Original dimension: ',fruitGray.shape)
9 hw_divby8(dimension):
10     dimension_rem = dimension % 8
11     if dimension_rem == 0:
12         dimension = dimension
13
14     elif dimension_rem >=5:
15         dimension = dimension + (8-dimension_rem)
16     else:
17         dimension = dimension-dimension_rem
18
19     return dimension
20
21 new_height = round(hw_divby8(height_fruit))
22 new_width = round(hw_divby8(width_fruit))
23 print('Divisible by 8 and reduced dimension: ',new_height,new_width)
24
25 resized_fruit = cv2.resize(fruitGray, dsize=(new_width, new_height), interpolation=cv2.INTER_CUBIC)
26 cv2.imwrite('guava_div_8.png',resized_fruit)
27
28

```

Figure 23: Divisibility by 8 code

4.2 Convert the image to standard height 256

For this, I am reusing the code in the sub task 4.1. I am setting standard height to 256, and passing standard height and width to the function and calculating the new width as per the aspect ratio with the help of another function. The new width will be rounded off to the nearest value divisible by 8 according to the remainder generated. Then these values are passed through resize function and saved as new image. This process is repeated for image0, image1 and image2 to generate resized images. Resized image dimensions are 384*256 for image0, 344*256 for image1 and 336*256 for image2. Figure 25 shows the new images along with their dimensions saved in a folder.

5 GENERATE BLOCK FEATURE VECTORS

Each image is converted to gray scale image, converted in to standard height of 256 and resized width, then divided into non overlapping blocks of size 8*8 pixels to generate the block feature vectors. Each pixel is called Feature. image0 is split into 1536 Non

Figure 29: image1.csv with 1 as label in 64th column

Figure 30: image2.csv with 2 as label in 64th column

7 DERIVE STATISTICAL DESCRIPTORS!

7.1 Extract statistical information

In this task, Statistical information are obtained using Counting, calculating Mean of each feature, standard deviation, scatter plot and histogram.

7.2 Counting:

By looking at the csv files, below are the observations presented.

Non Overlapping data set. image0 has 64 features and 1536 Observations, 384×256 dimension and labelled as 0. image0 has 64 features and 1376 Observations, 344×256 dimension and labelled as 1. image2 has 64 features and 1346 observations, 336×256 dimension and labelled as 2. Cumulatively, there are 4258 non overlapping blocks of size 8×8 , 3 classes and 64 features and hence volume of the data is 4258.

Overlapping Data Set statistical information: image0 has 64 features and 93873 Observations, 384×256 dimension and labelled as 0. image0 has 64 features and 83913 Observations, 344×256 dimension and labelled as 1. image2 has 64 features and 81921 observations, 336×256 dimension and labelled as 2. Cumulatively, there

are 259707 overlapping blocks of size 8×8 , 3 classes and 64 features and hence the volume of the data is 259707

7.3 Mean and STD plot:

Figure 31 to 36 shows the Mean, STD plots and for image0, 1 and 2 for overlapping and Non overlapping data sets and the code used to generate the plots. Mean for image0 in Non overlapping Data set ranges between 198 to 200, image1 ranges between 103 to 104 and image2 ranges between 119 to 122. Mean for image0 in Overlapping Data set ranges between 198 to 200, image1 is 106, image2 ranges between 120 to 121. STD for image 0 in Non Overlapping Data set ranges between 43-44, image1 - 58, image2 - 65-66 STD for image 0 in Overlapping Data set ranges between 43, image1 - 58, image2 - 65-66

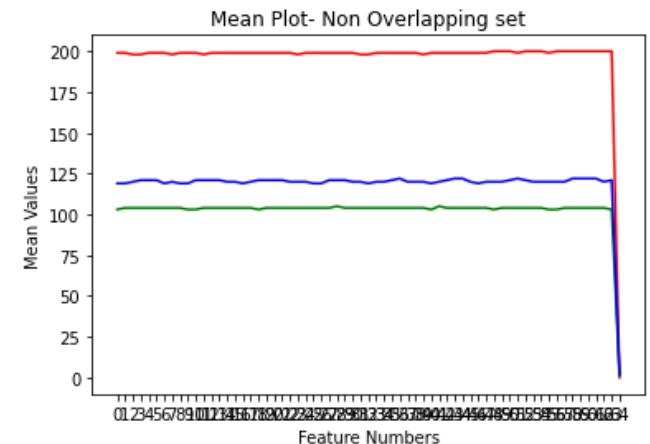


Figure 31: Mean Plot of all features for Non overlapping data set

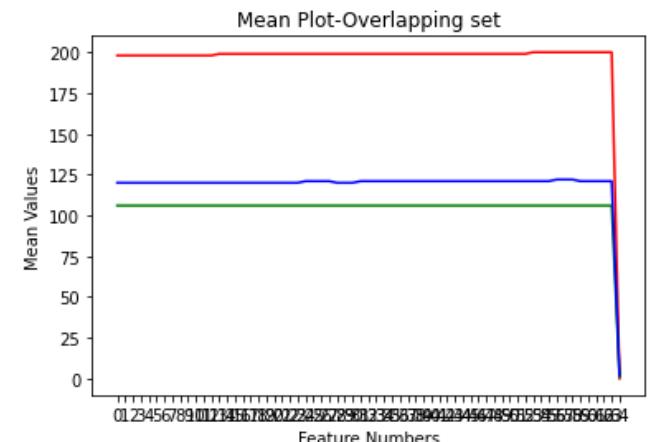


Figure 32: Mean Plot of all features for overlapping data set

7.4 Scatter Plot:

Scatter plots are mainly used to observe and determine the relationship between two features. I have selected feature10 and feature30 from image0 initially to analyze if there is any correlation. In Figure 37, I have plotted feature 10 on X-axis and feature 30 on Y-axis.

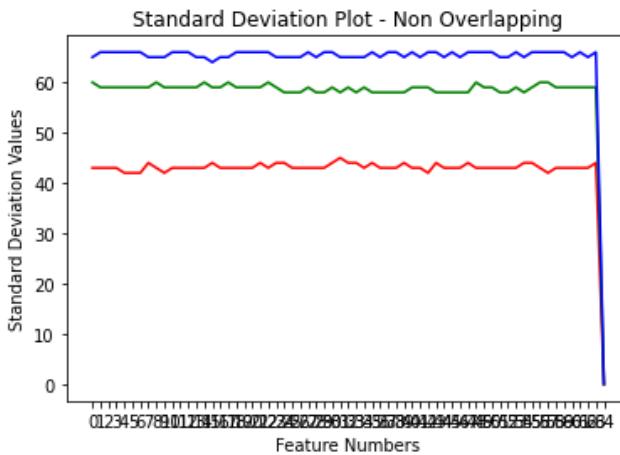


Figure 33: STD Plot of all features for Non overlapping data set

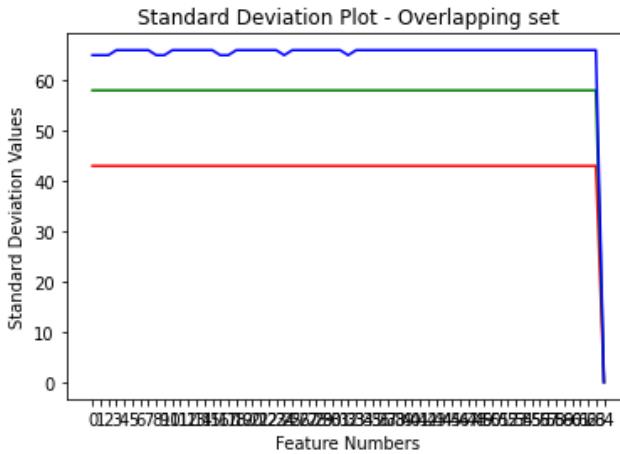


Figure 34: STD Plot of all features for Non overlapping data set

```
C:\Users\deepa\Desktop\BigData_ML\Assignment1\std_code.py
plots.py histogram.py scatter_plot.py scatter_3d.py assignment.py readmultipleimagesfromfolder.py no

1 # -*- coding: utf-8 -*-
2 """
3 Created on Mon Feb  8 18:18:35 2021
4
5 @author: deepa
6 """
7 import pandas as pd
8 import matplotlib.pyplot as plt
9 #reading csv files
10 data_banana = pd.read_csv('C:/Users/deepa/Desktop/BigData_ML/Assignment1/Task5/new/overlapping/banana.csv')
11 data_guava = pd.read_csv('C:/Users/deepa/Desktop/BigData_ML/Assignment1/Task5/new/overlapping/guava.csv')
12 data_pineapple = pd.read_csv('C:/Users/deepa/Desktop/BigData_ML/Assignment1/Task5/new/overlapping/pineapple.csv')
13
14 #calculating mean of all columns of all fruits
15 mean_banana = round(data_banana.mean())
16 mean_guava = round(data_guava.mean())
17 mean_pineapple = round(data_pineapple.mean())
18
19 #Plotting the mean values
20 plt.plot(mean_banana,color='r')
21 plt.plot(mean_guava,color='g')
22 plt.plot(mean_pineapple,color='b')
23
24 #setting the X and Y axis values
25 plt.title('Mean Plot-Overlapping set')
26 plt.xlabel('Feature Numbers')
27 plt.ylabel('Mean Values')
28 plt.axis('on')
29
```

Figure 35: Code to generate STD of all images

By looking at the plot, with the increase of x-axis values, there is increase in Y-axis values establishing a positive relation with each

```
C:\Users\deepa\Desktop\BigData_ML\Assignment1\plots.py
plots.py histogram.py scatter_plot.py scatter_3d.py assignment.py readmultipleimagesfromfolder.py no

1 # -*- coding: utf-8 -*-
2 """
3 Created on Mon Feb  8 18:18:35 2021
4
5 @author: deepa
6 """
7 import pandas as pd
8 import matplotlib.pyplot as plt
9 #reading csv files
10 data_banana = pd.read_csv('C:/Users/deepa/Desktop/BigData_ML/Assignment1/Task5/new/overlapping/banana.csv')
11 data_guava = pd.read_csv('C:/Users/deepa/Desktop/BigData_ML/Assignment1/Task5/new/overlapping/guava.csv')
12 data_pineapple = pd.read_csv('C:/Users/deepa/Desktop/BigData_ML/Assignment1/Task5/new/overlapping/pineapple.csv')
13
14 #calculating mean of all columns of all fruits
15 mean_banana = round(data_banana.mean())
16 mean_guava = round(data_guava.mean())
17 mean_pineapple = round(data_pineapple.mean())
18
19 #Plotting the mean values
20 plt.plot(mean_banana,color='r')
21 plt.plot(mean_guava,color='g')
22 plt.plot(mean_pineapple,color='b')
23
24 #setting the X and Y axis values
25 plt.title('Mean Plot-Overlapping set')
26 plt.xlabel('Feature Numbers')
27 plt.ylabel('Mean Values')
28 plt.axis('on')
29
```

Figure 36: Code to generate Mean of all images

other with few outliers and the data is more concentrated at the top right. Figure 38 shows the same plot with overlapping data. Similarly, scatter plots are plotted for image1 (Feature 10 and Feature 30) and image2 (Feature 15 and Feature 30). From the plots in Figure 39(image1 scatter plot), I can see that it is following a moderate, linear relationship with small clustering at the bottom left. Figure 40 shows the similar plot for overlapping data set. For image3 (Figure 41 and Figure 42), I have plotted scatter plot for feature 15 and 30 which also shows moderate linear relationship with few outliers.

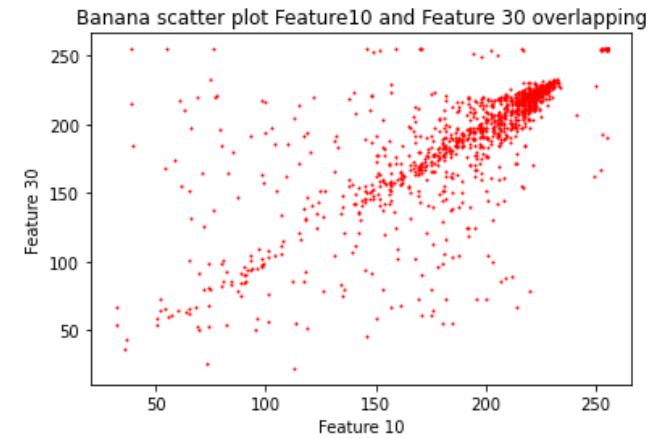


Figure 37: Scatter Plot of image0 feature10 and feature30 non sliding

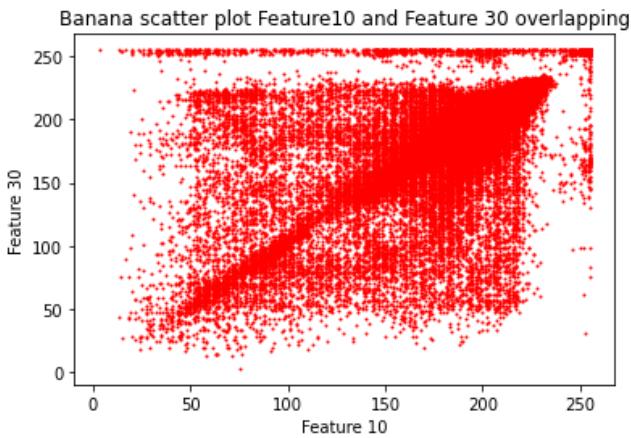


Figure 38: Scatter Plot of image0 feature10 and feature30 Overlapping

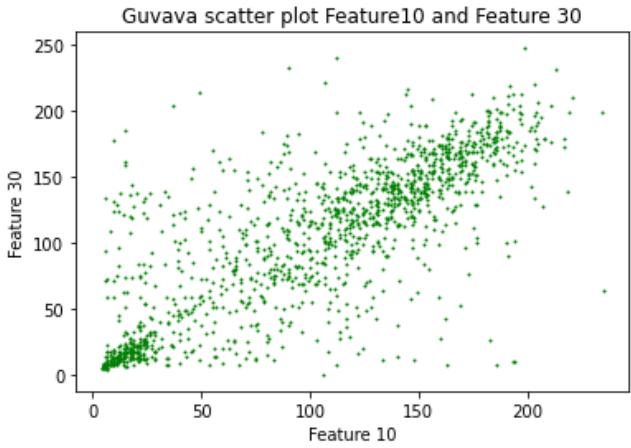


Figure 39: Scatter Plot of image1 feature10 and feature30 non sliding data

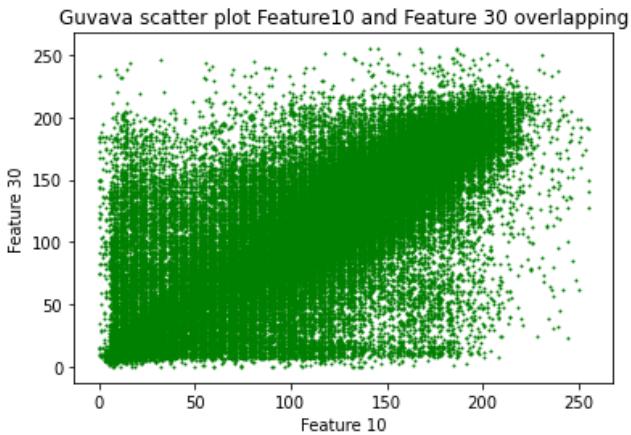


Figure 40: Scatter Plot of image1 feature10 and feature30 Overlapping data set

7.5 Histogram:

Histograms helps to visualize the distribution of the data numerically. I have plotted histograms to understand if the data follows Gaussian distribution or not. Figure 43 and Figure 44 are histogram

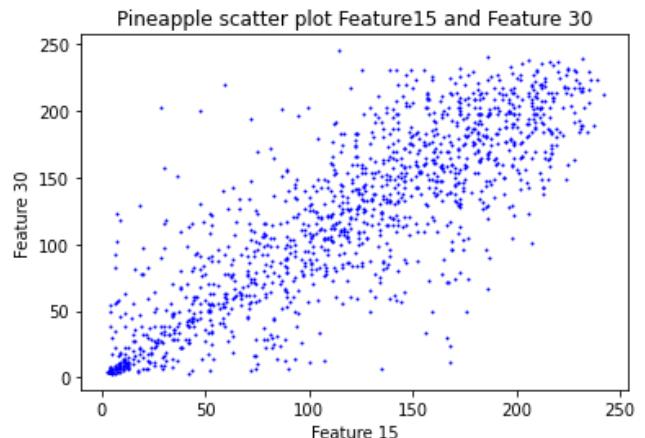


Figure 41: Scatter Plot of image2 feature15 and feature30 non sliding data

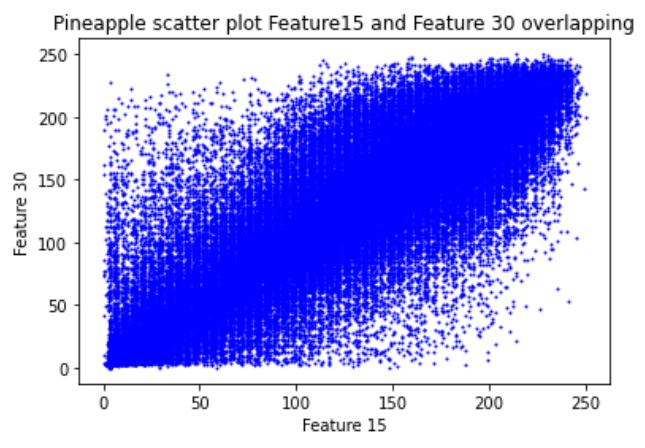


Figure 42: Scatter Plot of image2 feature15 and feature30 overlapping data set

plots for non sliding and sliding data set for image0 feature 10. The plot does not follow Gaussian distribution and is skewed on the left end and most of the data points are clustered on right end and the skewness could help differentiating classes. In histogram plot of Figure 45, feature 15, I can see that there are two peak values in non sliding data and data concentrated more in the middle. For image2, histogram plot for feature 15 in Figure 47 and 48, it is evident that the plot is almost equally distributed. So, none of the plots doesn't follow Gaussian distribution. Figure 49 has the code for generating histogram.

7.6 Is the data set imbalanced, inaccurate or incomplete?

The data set is imbalanced. For non sliding data set, the number of observations of image0 is greater than that of image1. Similarly, the observations of image1 is more than the observations in image2. For overlapping data set, image0 has maximum observations compared to image1 and image2. So it's an imbalanced data set. I have randomly selected feature 10 and 30 of image0 and feature15 and 25 from image0 and image2 to plot scatter plots. Figure 51 and 52 shows the plots obtained and they are imbalanced with some

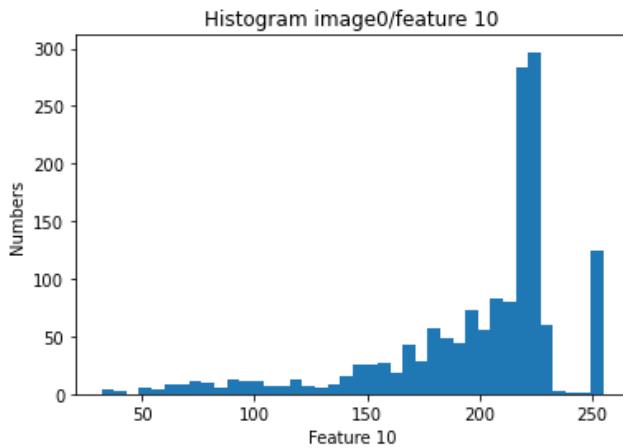


Figure 43: Histogram for feature10 of image0 non sliding data

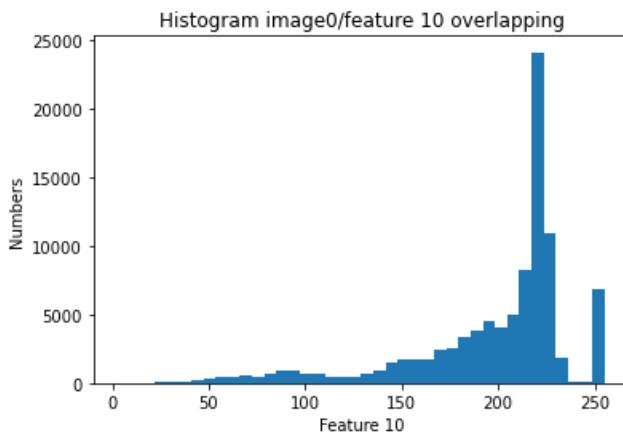


Figure 44: Histogram for feature10 of image0 non sliding data

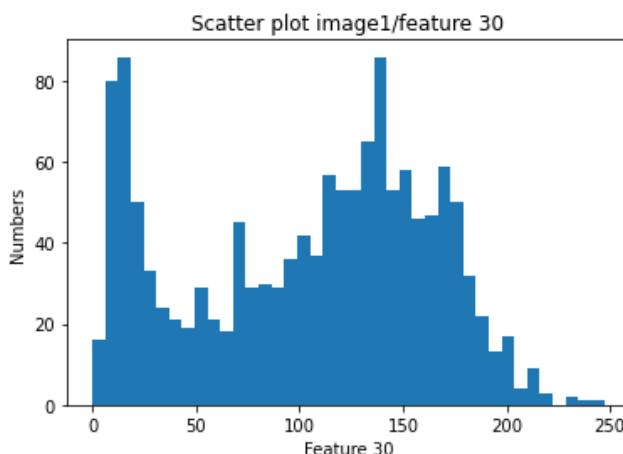


Figure 45: Histogram for feature 30 of image0 non sliding data

complex structures hidden. Also, none of the Histograms plotted follow Gaussian distribution, thus it makes them imbalanced data.

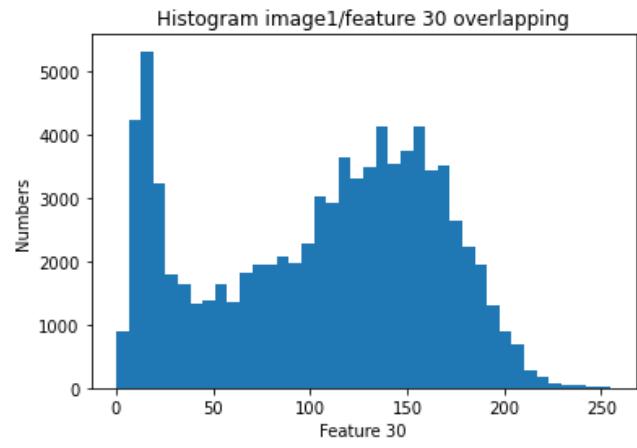


Figure 46: Histogram for feature 30 of image0 non sliding data

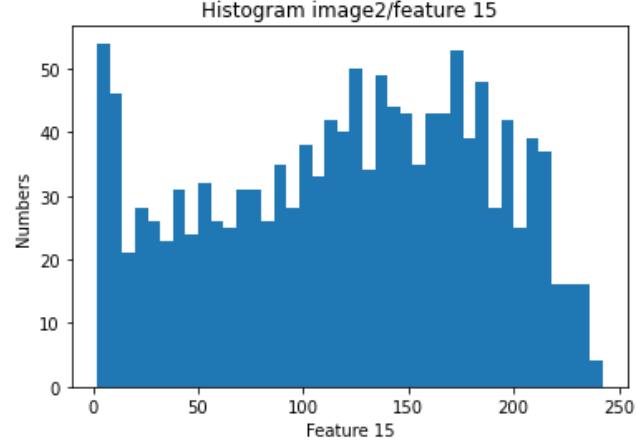


Figure 47: Histogram for feature10 of image0 non sliding data

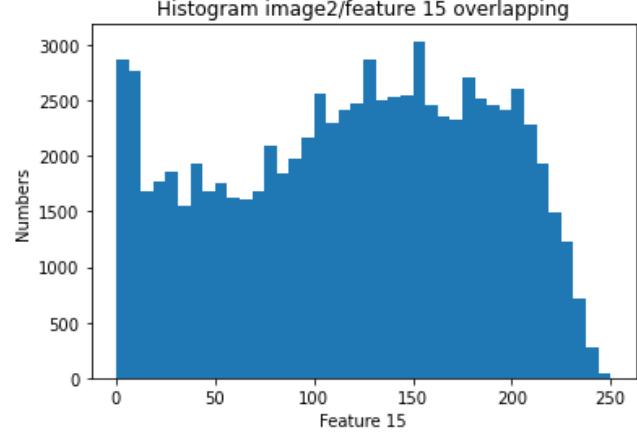


Figure 48: Histogram for feature10 of image0 non sliding data

7.7 Is it a trivial data or possibly a big data? Is it high dimensional?

This is trivial data and low dimensional because number of observations are less than number of features.

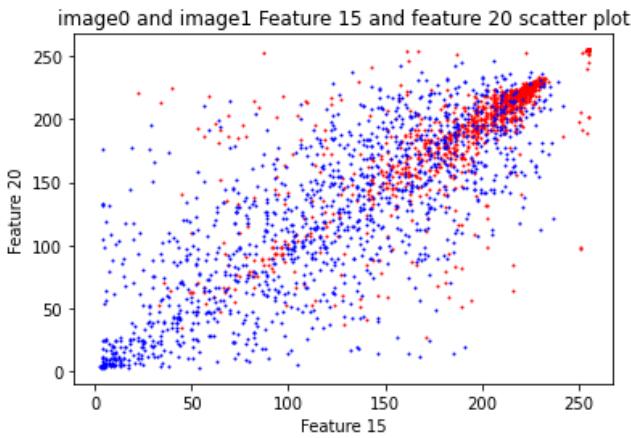


Figure 53: Two dimensional feature space image0 and image2 feature 15 and feature 25

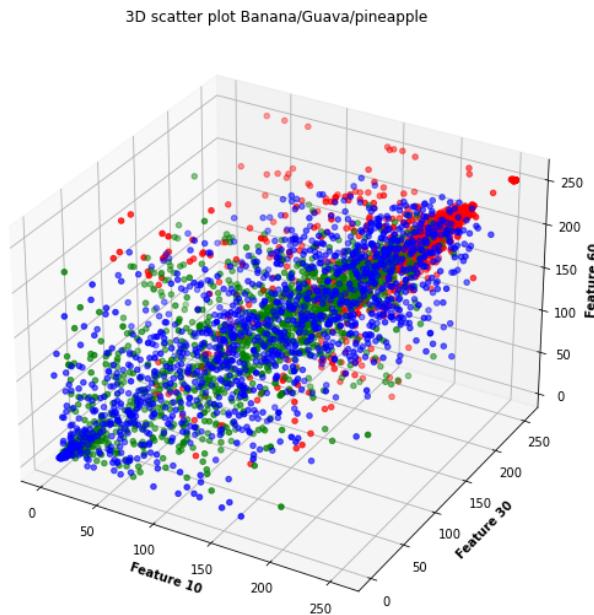


Figure 54: Three dimensional feature space

9.3 Discuss the figures:

By looking at the scatter plots in Figure 52, 53 and 54, the features are overlapping on each other making the class separation complex. Even though there are some clustering of image0 at the top right and image1 at bottom left, only that proof will not help in classification. The data needs to be randomized to make it balanced.

10 READ ANY NUMBERS OF IMAGES AND GENERATE FEATURE SPACE

Figure 55 shows the code use to read multiple images from folder and generate feature spaces. The code used in all the earlier tasks are combined and a new csv filename is generated and passed along with the image data to generate a feature space. For each image, csv file is generated with data appended with imagename. Figure 56

```
C:\Users\deepa\Desktop\BigData_ML_Assignment\readmultipleimagesfromfolder.py
scattered_3d.py x scattered_3d_stdarization.py x scatter_plot.py x standard deviation_code.py x readmultipleimagesfromfolder.py x

13
14     def resize_get_width(heightG,widthG):
15         aspect_ratio = widthG/heightG
16         new_fruit_width = aspect_ratio*std_height
17         return round(new_fruit_width)
18
19     def width_divisibilityBy8(height_fruit_width_fruit):
20         new_fruit_width = resize_get_width(height_fruit_width_fruit, width_fruit)
21         dimension_rem = new_fruit_width % 8
22         if dimension_rem == 0:
23             return new_fruit_width
24
25         elif dimension_rem < 8:
26             new_fruit_width = new_fruit_width + (8-dimension_rem)
27         else:
28             new_fruit_width = new_fruit_width-dimension_rem
29
30     return new_fruit_width
31
32     def generate_feature(img,filename):
33         fruit_gray = cv2.cvtColor(img, cv2.COLOR_BGR2GRAY)
34         height, width, fruit_gray.shape
35         new_fruit_width = width_divisibilityBy8(heightG,widthG)
36         resized_fruit = cv2.resize(fruit_gray, dsize=(new_fruit_width, std_height), interpolation=cv2.INTER_CUBIC)
37         heightC, widthC = resized_fruit.shape
38         cc = round((heightC)/(widthC)/8)
39         flatc = np.ones((cc, 64), np.uint8)
40         k = 0
41         for i in range(0,heightC,8):
42             for j in range(0,widthC,8):
43                 crop_tmp1 = resized_fruit[i:i+8,j:j+8]
44                 flatc[k] = crop_tmp1.flatten()
45                 k = k + 1
46         fspaced = pd.DataFrame(flatc)
47         #panda object
48         fspaced.to_csv(filename, index=False)
```

Figure 55: Code to read multiple images from folder

	Name	Date modified	Type	Size
Quick access				
Desktop	dataset1	2/18/2021 4:18 PM	Microsoft Excel C...	21,715 KB
Downloads	dataset2	2/18/2021 4:18 PM	Microsoft Excel C...	21,591 KB
Documents	dataset3	2/18/2021 4:19 PM	Microsoft Excel C...	20,717 KB
Pictures	dataset4	2/18/2021 4:19 PM	Microsoft Excel C...	21,168 KB
Assignment1	dataset5	2/18/2021 4:19 PM	Microsoft Excel C...	19,666 KB
gray	dataset6	2/18/2021 4:19 PM	Microsoft Excel C...	20,732 KB
Screenshots	dataset7	2/18/2021 4:19 PM	Microsoft Excel C...	20,713 KB
Task7	dataset8	2/18/2021 4:18 PM	Microsoft Excel C...	20,702 KB
OneDrive	dataset9	2/18/2021 4:18 PM	Microsoft Excel C...	15,558 KB
This PC	dataset10	2/18/2021 4:18 PM	Microsoft Excel C...	21,097 KB
3D Objects	dataset11	2/18/2021 4:18 PM	Microsoft Excel C...	19,973 KB
Desktop	dataset12	2/18/2021 4:18 PM	Microsoft Excel C...	20,049 KB
Downloads	dataset13	2/18/2021 4:18 PM	Microsoft Excel C...	14,907 KB
Music	dataset14	2/18/2021 4:18 PM	Microsoft Excel C...	21,362 KB
Pictures	dataset15	2/18/2021 4:18 PM	Microsoft Excel C...	19,969 KB
Videos	dataset16	2/18/2021 4:18 PM	Microsoft Excel C...	21,732 KB
Windows (C)	dataset17	2/18/2021 4:18 PM	Microsoft Excel C...	11,205 KB
Network	dataset18	2/18/2021 4:18 PM	Microsoft Excel C...	10,058 KB
	dataset19	2/18/2021 4:18 PM	Microsoft Excel C...	21,195 KB
	dataset20	2/18/2021 4:18 PM	Microsoft Excel C...	20,479 KB
	dataset21	2/18/2021 4:19 PM	Microsoft Excel C...	21,144 KB
	dataset22	2/18/2021 4:19 PM	Microsoft Excel C...	23,414 KB

Figure 56: multiple csv files generated

shows the resulted overlapping data in csv files. Non overlapping data sets are generated and csv files are uploaded for reference.

11 DESCRIBE THE EFFECTS OF BLOCK SIZE ON DIMENSIONALITY

I have considered blocks of size 8×8 which generates feature vector of size 64. But, if I were to chose 4×4 blocks that would result in a 16 vector size and the number of observations would increase adding more volume and complexity to the data. More observations will make it difficult to understand and predict the data which will finally pose problems to classification of data.

REFERENCES

- [1] Š. Marko, "Automatic fruit recognition using computer vision," Mentor: Matej Kristan), Fakulteta za racunalništvo in informatiko, Univerza v Ljubljani, 2013.

12 COMPLETE AND EXTEND THE TASKS OF ASSIGNMENT 1

All the tasks in Assignment 1 are completed. Overlapping feature vectors for two class, non-overlapping feature vectors for two class, overlapping feature vectors for three class and non-overlapping feature vectors for three class are prepared for carrying out Machine learning process.

12.1 Dividing the Data set to 80:20

For machine learning purpose, complete Data has to be divided into 80 percent as Training set and 20 percent has Testing set. This 80 percent data helps to Train the model so that we can Test the Model using the other 20 percent data. I have written Python code and saved the data sets as X_Train and X_Test csv files separately for image01 and image012 merged data sets. Figure 57 shows the code used to achieve the task and Fig 58 shows the generated files stored in a folder.

```
import pandas as pd
import numpy as np

image01_N = pd.read_csv('C:/Users/deepa/Desktop/BigData_ML/Assignment2/files/nonoverlap/image01.csv')
image012_N = pd.read_csv('C:/Users/deepa/Desktop/BigData_ML/Assignment2/files/nonoverlap/image012.csv')
image01_O = pd.read_csv('C:/Users/deepa/Desktop/BigData_ML/Assignment2/files/overlap/image01.csv')
image012_O = pd.read_csv('C:/Users/deepa/Desktop/BigData_ML/Assignment2/files/overlap/image012.csv')

X = image012_N
row, col = X.shape
TR = round(row*0.8)
X1 = np.array(X)
TT = row-TR
X_train = X1[0:TR,:]
X_test = X1[TR:row,:]

data_train = pd.DataFrame(X_train)
data_train.to_csv('C:/Users/deepa/Desktop/BigData_ML/Assignment2/files/nonoverlap/X_Train_012.csv')
data_test = pd.DataFrame(X_test)
data_test.to_csv('C:/Users/deepa/Desktop/BigData_ML/Assignment2/files/nonoverlap/X_Test_012.csv',
```

Figure 57: Code executed to generate Train and Test Data set

Name	Date modified	Type	Size
image0	3/10/2021 11:09 AM	Microsoft Excel C...	339 KB
image01	3/10/2021 12:09 PM	Microsoft Excel C...	669 KB
image1	3/10/2021 12:07 PM	Microsoft Excel C...	330 KB
image2	2/12/2021 9:56 PM	Microsoft Excel C...	305 KB
image012	3/12/2021 9:31 AM	Microsoft Excel C...	973 KB
X_Test_01	3/19/2021 12:24 PM	Microsoft Excel C...	134 KB
X_Test_012	3/19/2021 10:13 AM	Microsoft Excel C...	196 KB
X_Train_01	3/19/2021 12:18 PM	Microsoft Excel C...	535 KB
X_Train_012	3/19/2021 10:30 AM	Microsoft Excel C...	778 KB

Figure 58: Folder where Train and Test Data sets are saved locally

12.2 Draw histograms of Two features from Training and Test Data set

I have selected Feature 30 and Feature 60 as there was no huge variation in the mean and Standard deviation of all features. From Figure 59 to 62, we can observe that histogram of Train data set of overlapping/non-overlapping of two class and overlapping and non-overlapping data set of three class follow the same distribution as that of their respective Test data sets. The Mean and Standard deviation of the features are calculated and a comparison table is drawn in Figure 63 and 64.

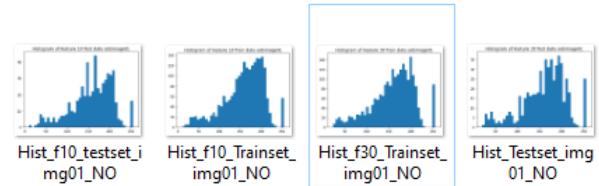


Figure 59: Histogram of Feature 10 and Feature 30 01 Non oVerlapping Data set



Figure 60: Histogram of Feature 10 and Feature 30 01 oVerlapping Data set

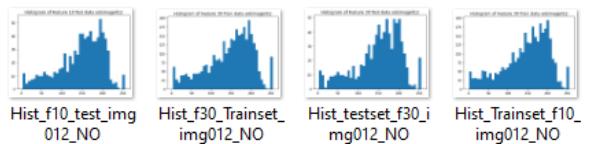


Figure 61: Histogram of Feature 10 and Feature 30 012 Non oVerlapping Data set

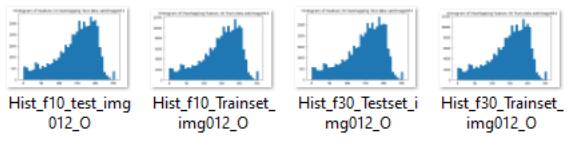


Figure 62: Histogram of Feature 10 and Feature 30 012 oVerlapping Data set

12.3 Draw Scatter plot of Two features from Training and Test Data set

Figure 65 to 72 to shows scatter plot of Feature 10 and 30 for non-overlapping and overlapping data sets.

13 IMPLEMENTING REGRESSION BASED MODEL

13.1 Implementing Lasso Regression Model

I have implemented Lasso regression model. Trained the model using X_Train_01 sets of overlapping and non overlapping data sets generated in earlier task. Training generates the coefficients for each feature/column in the data set. These coefficients are multiplied with feature matrix to get the predicted labels. Initially lamda is set to 0 during execution and increased up to 17. But the accuracy of the model started reducing after 17, so I have kept the lamda

In [8]: feature_10.mean() Out[8]: 158.75724137931036	In [13]: feature_10.mean() Out[13]: 158.75690607734808
In [9]: feature_30.mean() Out[9]: 159.72919540229884	In [14]: feature_10.std() Out[14]: 47.26185648551175
In [10]: feature_10.std() Out[10]: 47.20729450475824	In [15]: feature_30.mean() Out[15]: 160.88397790055248
In [11]: feature_30.std() Out[11]: 47.544938191446676	In [16]: feature_30.std() Out[16]: 48.837756930883174
Nonoverlap_01_Train	Nonoverlap_01_Test
In [39]: feature_10.mean() Out[39]: 157.1648073328107	In [44]: feature_10.mean() Out[44]: 157.30014999962313
In [40]: feature_30.mean() Out[40]: 157.31158415244528	In [45]: feature_30.mean() Out[45]: 157.21767281991754
In [41]: feature_10.std() Out[41]: 45.66635526544437	In [46]: feature_10.std() Out[46]: 45.55917374876527
In [42]: feature_30.std() Out[42]: 45.29519121833504	In [47]: feature_30.std() Out[47]: 45.531036602680544
Overlap_01_Train	Overlap_01_Test

Figure 63: Comparison Table of Mean and STD for image01 data set

In [18]: feature_10.mean() Out[18]: 146.04983082128575	In [23]: feature_10.mean() Out[23]: 147.1168511685117
In [19]: feature_10.std() Out[19]: 57.53335845756468	In [24]: feature_10.std() Out[24]: 54.58987082308333
In [20]: feature_30.mean() Out[20]: 146.2177791448785	In [25]: feature_30.mean() Out[25]: 149.69987699876998
In [21]: feature_30.std() Out[21]: 57.89490742068153	In [26]: feature_30.std() Out[26]: 55.87169859196968
Nonoverlap_012_Train	Nonoverlap_012_Test
In [28]: feature_10.mean() Out[28]: 145.0121591895219	In [33]: feature_10.mean() Out[33]: 144.85721781598758
In [29]: feature_30.mean() Out[29]: 145.1390486569393	In [34]: feature_10.std() Out[34]: 55.625432231536415
In [30]: feature_10.std() Out[30]: 55.86148253001543	In [35]: feature_10.std() Out[35]: 55.625432231536415
In [31]: feature_30.std() Out[31]: 55.79853536160791	In [36]: feature_30.mean() Out[36]: 144.83746039434118
Overlap_012_Train	Overlap_012_Test

Figure 64: Comparison Table of Mean and STD for image012 data set

value as 17 for analysing accuracy measures. Figure 73 shows the screenshot of code and execution in the local environment.

13.2 Testing and adding predicted labels

In the above step model is trained using Train data set and now I am testing the accuracy of the model for Test data set. Testing is done by matrix multiplication of parameters or coefficients obtained through Training and X_Test_01 feature vectors. This process gives us the predicted labels for Test data. Predicted labels are saved in a lasso_predicted.csv file and later placing them as 66th column in X_Test_01 files of both overlapping and non-overlapping data sets. I started off the testing by initially setting the lamda value to 0 and

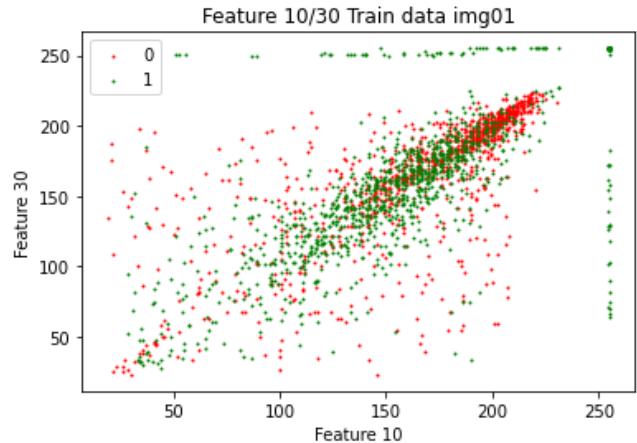


Figure 65: Scatter Plot 01 Non overlapping Train Dataset

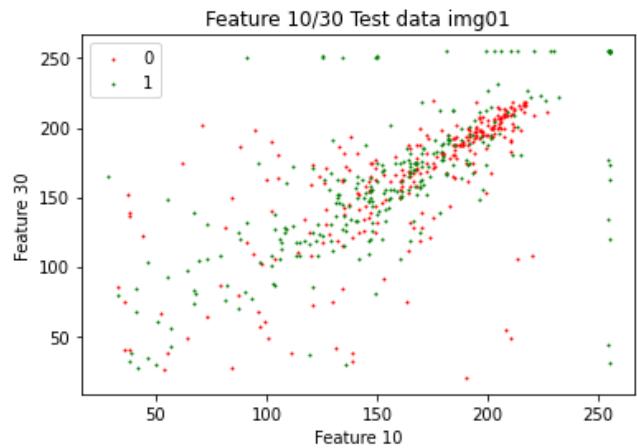


Figure 66: Scatter Plot 01 overlapping Test Dataset

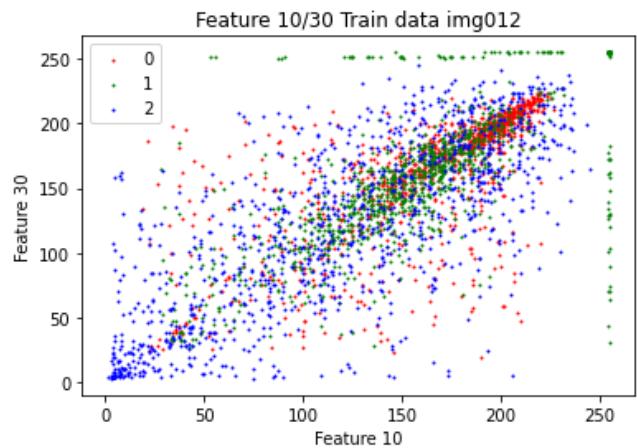


Figure 67: Scatter Plot 012 Non overlapping Train Dataset

gradually increasing up to 17. But, there was only slight increase in the accuracy level and it started decreasing after value 17.

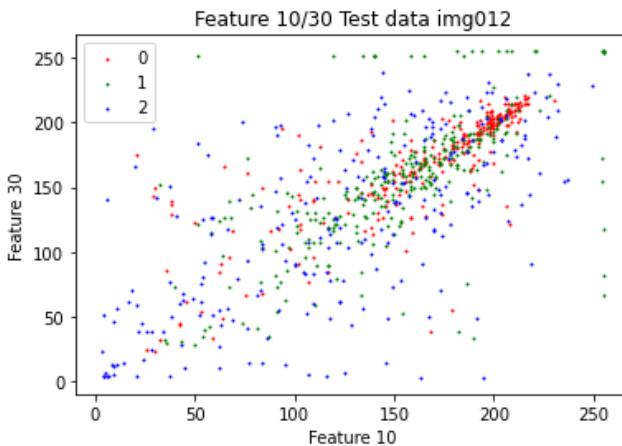


Figure 68: Scatter Plot 012 Non overlapping Train Dataset

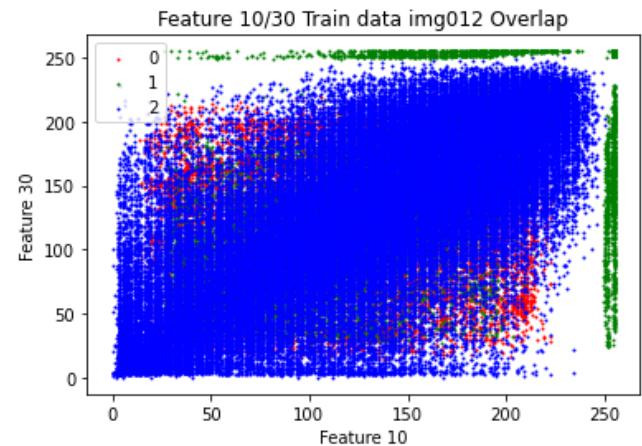


Figure 71: Scatter Plot of img012 Overlapping Train Data set

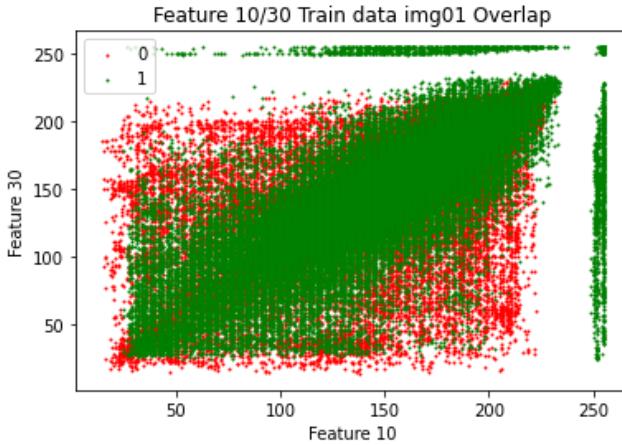


Figure 69: Scatter Plot of img01 Overlapping Train Data set

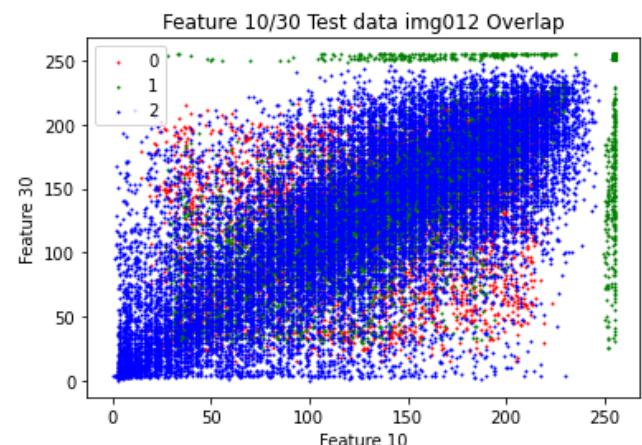


Figure 72: Scatter Plot of img012 Overlapping Test Data set

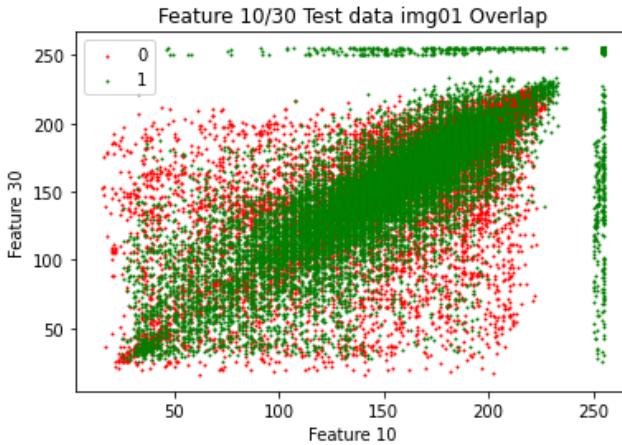


Figure 70: Scatter Plot of img01 Overlapping Test Data set

13.3 Create Confusion Matrices

Confusion matrix is created by passing 65th column and 66th column generated in previous task to the confusion matrix function of `sklearn.metrics` library. The matrix is printed on the command

A screenshot of the Spyder Python IDE. The code editor shows a script named `lasso_correct.py` with the following content:

```

1 #Author: deepo
2 import pandas as pd
3 import numpy as np
4 from sklearn import linear_model
5 from sklearn.metrics import confusion_matrix
6
7 X_train = pd.read_csv('C:/Users/deepo/Desktop/BIGData_ML/Assignment2/files/nonoverlap/X_Train_01.csv')
8 X_test = pd.read_csv('C:/Users/deepo/Desktop/BIGData_ML/Assignment2/files/nonoverlap/X_Test_01.csv')
9 Y = X_train['64']
10 X_train.drop(['64'], axis=1, inplace=True)
11 X_train = np.array(X_train)
12 X_test = np.array(X_test)
13 lasso = Lasso()
14 lasso.fit(X_train, Y)
15 X2 = X1.transpose()
16 X2dash = np.matmul(X2, X1)
17 IXdash = inv(X2dash)
18 ymatX = np.matmul(X2, Y)
19 # inverse of the square of the feature matrix
20 first = np.matmul(ymatX, IX)
21 Z1 = np.matmul(first, X2)
22 one = (X1.T).dot(Z1)
23 three = np.matmul(one, X1)
24 Z2 = np.matmul(X1, three)
25 Z2dash = inv(Z2)
26 yhatTrain = Z2dash.astype(int)
27 yhatTest = Z1.astype(int)
28 yhatTrain_secd = pd.DataFrame(yhatTrain)
29 yhatTest_secd = pd.DataFrame(yhatTest)
30 CC = confusion_matrix(Y, yhatTrain)
31 CC = confusion_matrix(Y, yhatTest)
32

```

The IPython console shows the execution results:

```

In [9]: runfile('C:/Users/deepo/Desktop/BIGData_ML/Assignment2/files/nonoverlap/lasso_correct.py', wdir='C:/Users/deepo/Desktop/BIGData_ML/Assignment2/files/nonoverlap')
Train_Accuracy_Score: 0.472880294117646
Train_Precision_Score: 0.473512134119782
Train_Negative_Precision_Score: 0.472880294117646
Train_Specificity_Score: 0.4668463197935516
Test_Accuracy_Score: 0.420955882329411
Test_Precision_Score: 0.41981488458704225
Test_Negative_Precision_Score: 0.420955882329411
Test_Specificity_Score: 0.4

```

Figure 73: Lasso regression execution in local environment

prompt. Figure 76 shows the code of confusion matrix. I am considering class 0 as the True positive class and qualitative analysis is performed based on it.

13.4 Comparing performance of Data sets

By using Confusion matrix, accuracy, precision, sensitivity and specificity is calculated for nonoverlapping and overlapping sets

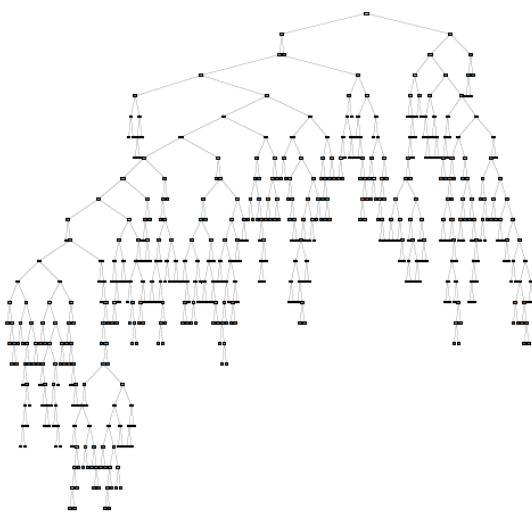


Figure 80: Random Forest tree for non overlapping data as two class classifier

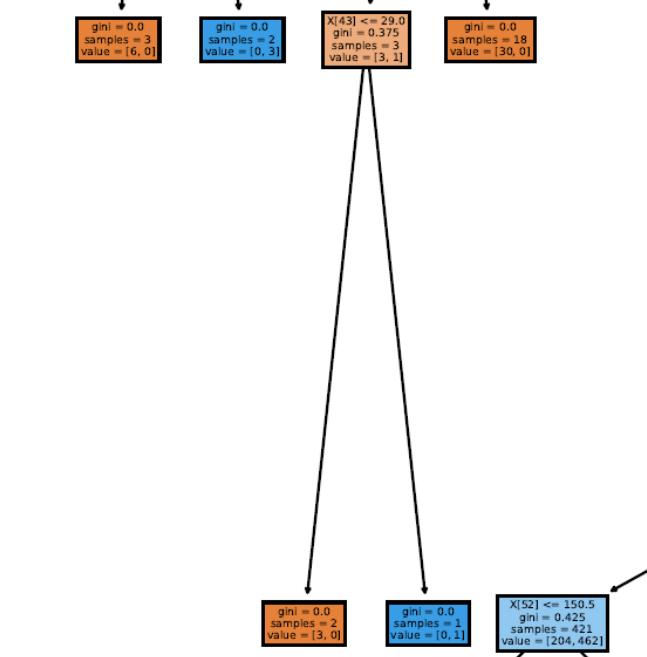


Figure 81: Zoomed in tree

14.3 Applying trained model to Test data

Trained data of both two class and three class classifier is applied to Test data to evaluate the performance of the model. Figure 83 has the code used for testing. Model is trained and tested by changing the n_estimators like 50, 100, 500 and 1000 to see the variations in performance. The best model's predicted labels are saved in 66th column of Test set.

```

import pandas as pd
import numpy as np
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score,precision_score,recall_score
from sklearn.metrics import classification_report
from sklearn.ensemble import RandomForestClassifier

x_train = pd.read_csv('C:/Users/deepa/Desktop/BigData_ML/Assignment2/files/RandomForest/overlap.csv')
x_test = pd.read_csv('C:/Users/deepa/Desktop/BigData_ML/Assignment2/files/RandomForest/overlap.csv')

y = x_train[64]
X = np.array(x_train)
X = X.drop(64, axis=1, inplace=True)
X1 = np.array(X)

#training the model
randomClassifier = RandomForestClassifier(random_state=0,n_estimators=1000,oob_score=True, n_jobs=-1)
model = randomClassifier.fit(X1, y)

importance = model.feature_importances_
indices = importance.argsort()[::-1]
oob_error = 1 - randomClassifier.oob_score_
  
```

Figure 82: Training the Random Classifier Code snippet for three class classifier

```

#Testing the model using Trained results
test_y = x_test[64]
Y_test = np.array(test_y)
X_test = np.array(x_test)
X_test.drop(64, axis=1, inplace=True)
X_test_np = np.array(X_test)
y_pred = randomClassifier.predict(X_test_np)
y_hat_saved = pd.DataFrame(y_pred)
y_hat_saved.to_csv('C:/Users/deepa/Desktop/BigData_ML/Assignment2/files/Randomforest/nonoverlap.csv')
  
```

Figure 83: Testing the model code snippet

14.4 Constructing confusion matrices

After testing the model, the predicted results in 66th column and original values in 65th column are fed to Confusion matrix. I am taking class 0 as the true positive class and the evaluation is done based on that. Figure 84 and 85 shows the code snippet for both two class and three class Confusion matrices respectively.

```

CC_test = confusion_matrix(test_y, y_pred)
# 0 as my positive class
TN = CC_test[1,1]
FP = CC_test[1,0]
FN = CC_test[0,1]
TP = CC_test[0,0]
  
```

Figure 84: Constructing Confusion matrices code snippet Two class

```

TP = CC_test[0,0]
TN = CC_test[1,1]+CC_test[1,2]+CC_test[2,1]+CC_test[2,2]
FP = CC_test[0,1]+CC_test[0,2]
FN = CC_test[1,0] + CC_test[2,0]
  
```

Figure 85: Constructing Confusion matrices code snippet three class

14.5 Performance evaluation of the random forest model

After constructing the confusion matrices, the important task is evaluating the model for it's effectiveness in terms on accuracy, precision, specificity and sensitivity. All the classifiers i.e two class and three class are trained and tested by setting number of trees option to 50, 100, 500 and 1000. The results obtained are shown in Figure 86 to Figure 88. Figure shows the two class classifier for non overlapping data set. Here we can see that as the number of estimators increases, performance was decreasing slightly. In contrast overlapping data set performed much better with increasing estimators. In three class non overlapping data set in Figure 86 demonstrates when estimator value is 100 it is giving good results and as the estimator value increased performance got slightly down. For three class overlapping classifier, as the estimators increased,

15.3 Comparing the model to find the best one

From Figure 92, we can see that random forest model is performing better than Lasso regression model. There is only slight difference in the calculations of in-built Confusion based metrics and imported libraries. In most cases, both of them gave the same results. Random forest gave best results for overlapping data set compared to non-overlapping data set. Random forest model trained using 500 estimators for overlapping data for three class classifier was the best model among all others with 94.2% accuracy.

16 ASSIGNMENT 3: A CLASSIFIER ON A BIG DATA SYSTEM

17 MAKING SURE ALL TASKS OF ASSIGNMENT 1 AND ASSIGNMENT 2 ARE COMPLETED

All the below tasks related to assignment 1 and assignment 2 are completed.

- In assignment 1, all four feature vectors/ data frames of two class non-overlapping, two class overlapping, three class non-overlapping and three class overlapping are generated and stored in respective csv files.
- In assignment 2, implementation of Lasso regression as two class classifier and random forest as two and three class classifier is performed in local environment and the best model and the best feature set is selected.
- In addition to above tasks, I have generated single csv file with multiple feature features generated from folder of image0 (banana) and image1 (guava). Then, I have merged all these multiple feature vectors to generate a large volume data set to experiment on Big data environment. So, this data frame is prepared with around 40 banana and 40 guava images I took from FIDS30 folder. I will be trying to run this data around 1 GB in local as well as big data environment to analyze the performance.

18 UNDERSTANDING AND USING BIGDATA SYSTEM USING DATABRICKS SYSTEM

18.1 Signing up Databricks community version

Created my account by signing up to Databricks community using UNCG email id. I chose Google Cloud platform as the backend.

18.2 Exploring Quickstart tutorial and setting up Big data environment in Databricks community

I created my community Databricks account by choosing gcp as backend platform which gives around 15GB of free memory. Then, created cluster CSC610 using "Create Cluster" option. I chose 7.4 ML version apache spark 3.0.1, scala 2.12 as there was no 7.3 version present in the list. It took around 5 minutes for the cluster to be up and running. There were two options available to create tables, one was using UI , other was using notebook. I uploaded all the data set generated in assignment 1 and created the tables using notebook option since there was an internal error thrown while trying to create table using UI. Figure 95 and 97 shows the creation

A screenshot of the Databricks workspace interface. On the left is a sidebar with icons for Home, Workspace, Recents, Data, Clusters, Jobs, and Search. The main area is titled 'Workspace' and shows a list of notebooks. The list includes: 'd_jayanna@uncg.edu' (selected), '2021-04-09 - DBFS Example (1)', '2021-04-09 - DBFS Example (2)', '2021-04-09 - DBFS Example (3)', '2021-04-09 - DBFS Example (4)', '2021-04-09 - DBFS Example (5)', '2021-04-14 - DBFS Example', 'assignment3', 'nonoverlap01', 'nonoverlap012', 'overlap01', and 'overlap012'. At the bottom right of the workspace, it says 'Driver' and 'Apache Spark 2.4.5... Comm...'.

Figure 93: Account created in Databricks

of tables using notebook. I first selected the required parameters like selecting the cluster, table name, first row header checkbox and infer schema checked using UI option itself. After selecting all these parameters when I hit create table, an Internal error was thrown and it persisted for other table creations as well. So, I chose create table with notebook and ran all the queries to generate tables. I followed the same procedure to create four tables by uploading 4 csv file, overlapping two class and three class and non overlapping two class and three class. Next step was to create notebook in order to write the code and run the models. I created 4 different notebooks to execute models all 4 types of dataset.

18.3 Cluster creation:

Cluster is created with 7.4 ML version apache spark 3.0.1, scala 2.12 since version 7.2 was not available in the list. Figure 94 shows the screenshot of cluster created in Databricks community version.

A screenshot of the Databricks Cluster creation interface. The top bar shows 'Clusters / CSC610' with buttons for Edit, Clone, Restart, Terminate, and Delete. Below the bar, it says 'Databricks Runtime Version: 7.4 ML (includes Apache Spark 3.0.1, Scala 2.12)'. Under 'Driver Type', it says 'Community Optimized' with '15.3 GB Memory, 2 Cores, 1 DBU'. The 'Instance' section notes 'Free 15GB Memory. As a Community Edition user, your cluster will automatically terminate after an idle period of two hours. For more configuration options, please upgrade your Databricks subscription.' At the bottom, tabs for Instances, Spark, JDBC/ODBC, and Permissions are shown, along with an 'Availability Zone' dropdown set to 'us-west-2c'.

Figure 94: Cluster creation

18.4 Create Dataframes

Four data frames are created by uploading two class overlapping and non-overlapping, three class overlap and non-overlap feature vectors into big data system. These are the same feature vectors

used in previous assignments. Overlap01, overlap012, nonoverlap01 and nonoverlap012 tables are created. Similarly, large_dataset folder which is around 1.5 GB is also uploaded to see the performance of big data system. An end to end process for generating a table for one of the dataset is shown from figure 95 to 97. Same code and procedure is used to upload other data sets as well.

```

2021-04-09 - DBFS Example (5) (Python)
db  CSC610 Overview

1 # File location and type
2 file_location = "/FileStore/tables/nonoverlap01.csv"
3 file_type = "csv"
4
5 # CSV options
6 infer_schema = "true"
7 first_row_is_header = "true"
8 delimiter = ","
9
10 # The option below are ignored for CSV files. For other file types, these will be ignored.
11 # header = "true", inferSchema = true, firstRowIsHeader = true
12 .option("inferSchema", "true") \
13 .option("header", "true") \
14 .option("sep", ",") \
15 .load(file_location)
16
17 display(df)

> (1) Spark Jobs
> df: pyspark.sql.dataframe.DataFrame = [0: integer, 1: integer ... 63 more fields]

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 132 | 124 | 144 | 157 | 175 | 186 | 184 | 187 | 149 | 121 |
| 135 | 134 | 132 | 132 | 129 | 121 | 119 | 118 | 131 | 132 |
| 136 | 140 | 145 | 150 | 137 | 145 | 146 | 147 | 135 | 139 |

```

Figure 95: Dataframe creation part1

```

Cmd 3
1 # Create a view or table
2
3 temp_table_name = "overlap01"
4
5 df.createOrReplaceTempView(temp_table_name)

Command took 0.03 seconds -- by d_jayanna@uncg.edu at 4/9/2021, 4:40:42 PM on csc610

Cmd 4
1 %sql
2
3 /* Query the created temp table in a SQL cell */
4
5 select * from `overlap01`;

> (1) Spark Jobs

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 132 | 124 | 144 | 157 | 175 | 186 | 184 | 187 | 149 | 121 |
| 135 | 134 | 132 | 132 | 129 | 121 | 119 | 118 | 131 | 132 |
| 136 | 140 | 145 | 150 | 137 | 145 | 146 | 147 | 135 | 139 |

```

Figure 96: Dataframe creation part2

```

2021-04-09 - DBFS Example (5) (Python)
db  CSC610 Overview

1 # With this registered as a temp view, it will only be available to this particular notebook. If you'd like other users to be able to query this table, you can also create a table from the DataFrame.
2 # Create a permanent table with persist across cluster restarts as well as allow various users across different notebooks to query this data.
3 # To do so, choose your table name and uncomment the bottom line.
4
5 permanent_table_name = "overlap01"
6 df.write.format("parquet").saveAsTable(permanent_table_name)

> (1) Spark Jobs

Command took 12.47 seconds -- by d_jayanna@uncg.edu at 4/9/2021, 4:41:04 PM on csc610

Shift+Enter to run

```

Figure 97: Dataframe creation part3

19 CONVERTING THE CODE TO RUN YOUR DATA SETS

In assignment 2, I had ran multiple combinations by changing the n_estimators values to analyse its impact on accuracy and to find out the best model among each. The best model in each category is predicted as below.

- Non overlap two class - 50 estimators with 80% accuracy.
- Non overlap three class - 100 estimators with 82% accuracy.
- Overlap two class - 1000 estimators with 92.9% accuracy.
- Overlap three class - 500 estimators with 94.2% accuracy.

For this task, I am executing aforementioned models, in the databricks platform as well as in local environment. General steps followed to execute the model are as follows:

19.1 pyspark

Imported pyspark as shown in Figure 98 for computational purpose and fetching the data from pyspark.sql.dataframe.DataFrame uploaded earlier in Task2 using the select query.

```

nonoverlap01 (Python)
db  CSC610 Overview

Cmd 1
1 import pyspark

Command took 0.04 seconds -- by d_jayanna@uncg.edu at 4/24/2021, 6:25:10 PM on CSC610

Cmd 2
1 df = sqlContext.sql("SELECT * FROM nonoverlap01")

> df: pyspark.sql.dataframe.DataFrame = [0: integer, 1: integer ... 63 more fields]
Command took 0.87 seconds -- by d_jayanna@uncg.edu at 4/24/2021, 6:25:26 PM on CSC610

Cmd 3
1 input_data = df.select("+").toPandas()

> (1) Spark Jobs
Command took 1.83 seconds -- by d_jayanna@uncg.edu at 4/24/2021, 6:25:34 PM on CSC610

```

Figure 98: importing pyspark and fetching Dataframe

19.2 Generating histograms and boxplots

Figure 99 and 100 shows the histogram and boxplots respectively for Feature 30 using seaborn package.

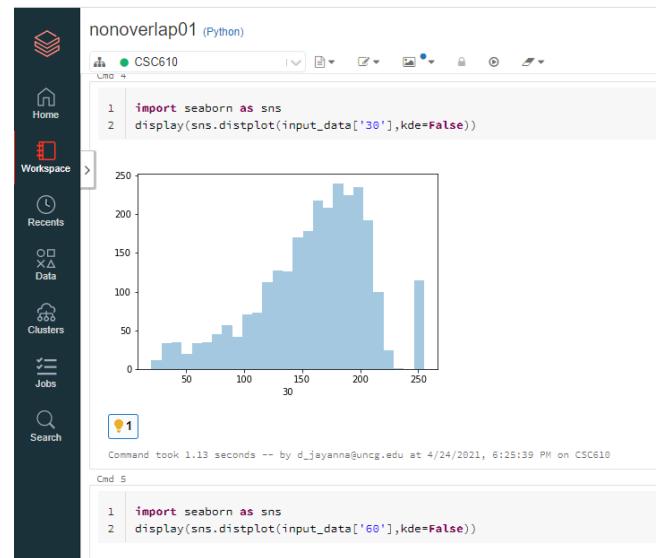


Figure 99: Histogram for Feature 30

19.3 Splitting the data

Checked the input data for any missing values before applying the model using isna() function. Then, splitting the data into 80% Train and 20% Test set into by assigning test_split parameter to 0.2. Splitting is done using train_test_split imported from sklearn.model_selection module. Figure 101 shows the execution of this task.

19.4 Training RandomForest Classifier

RandomForestClassifier is trained using 80% Train data set for 50 estimators(as this was the best predicted model in this scenario).

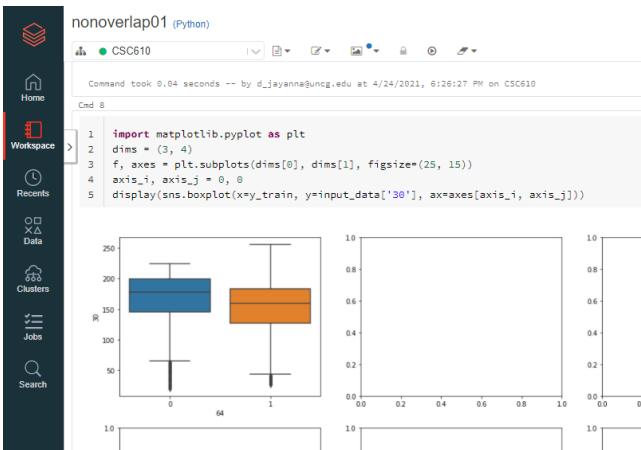


Figure 100: Boxplot for feature 30

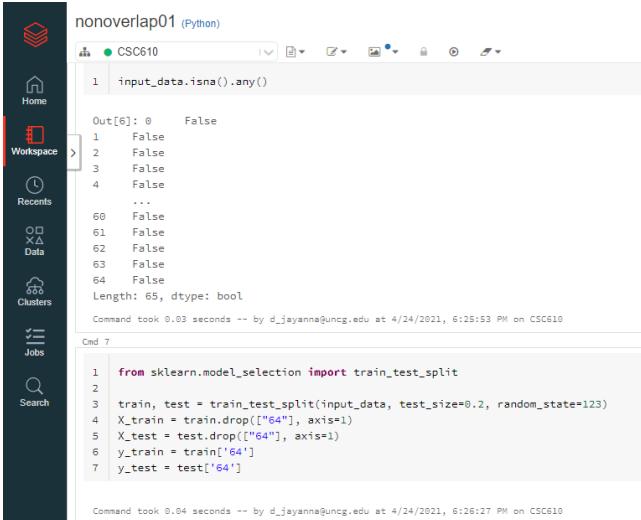


Figure 101: Splitting the data to Train and Test data

This process took around 0.64 seconds in Databricks environment for training the model. Then the model is tested for X_Test values to predict the output values. The proof of execution is as shown in Figure 102.

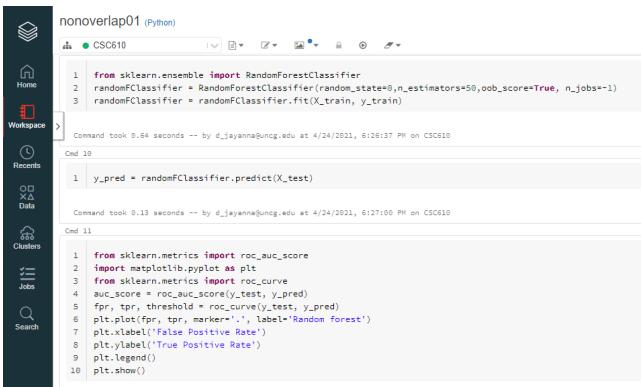


Figure 102: Random Forest classifier

19.5 Generating ROC Curve

ROC curve is drawn using roc_auc_score by passing test and predicted outputs. Also, predict_probabilities with multiple threshold values among tpr along y-axis and fpr along x-axis. Figure 103 and 104 shows the generated ROC curve.

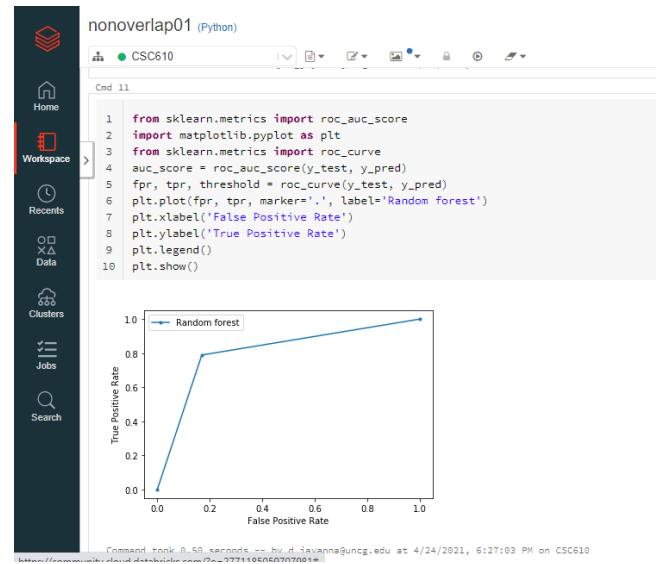


Figure 103: ROC curve

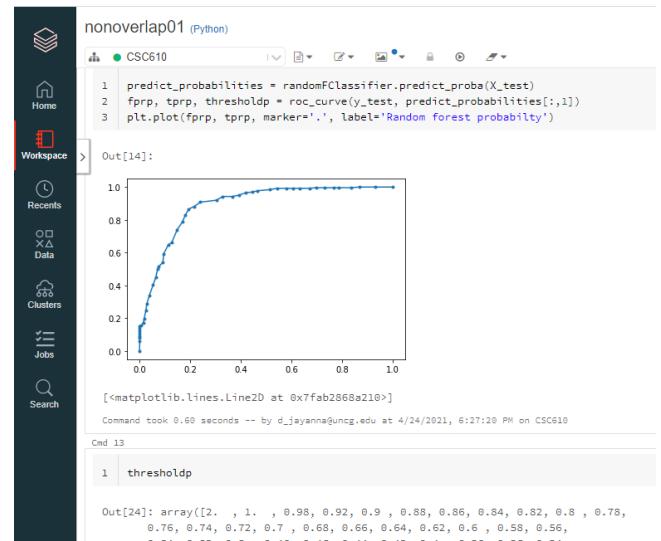


Figure 104: ROC curve with multiple threshold values

19.6 Calculating the feature importance

Importance of every feature is calculated to find the feature that is maximum impacting in classification. Looking in to figure 105, feature 63 seems to have highest importance compared to other features.

```

nonoverlap01 (Python)
In [1]: thresholdp

Out[24]: array([0.98, 0.97, 0.9, 0.88, 0.86, 0.84, 0.82, 0.8, 0.78,
   0.76, 0.74, 0.72, 0.7, 0.68, 0.66, 0.64, 0.62, 0.6, 0.58, 0.56,
   0.54, 0.52, 0.5, 0.48, 0.46, 0.44, 0.42, 0.4, 0.38, 0.36, 0.34,
   0.32, 0.3, 0.28, 0.26, 0.24, 0.22, 0.2, 0.18, 0.16, 0.14, 0.12,
   0.1, 0.08, 0.06, 0.04, 0.02, 0. ])
```

Command took 0.02 seconds -- by d_jayanna@uncg.edu at 4/24/2021, 8:42:34 AM on CSC610

```

1: importance = randomClassifier.feature_importances_
2: indices = importance.argsort()[:-1]
3: print(importance)

[0.01406511 0.01427903 0.01516428 0.0199635 0.01243488 0.01408556
 0.01918101 0.01596193 0.0187415 0.01480272 0.00848557 0.01798811
 0.01251177 0.01263585 0.01801622 0.01975255 0.01655655 0.01030317
 0.01246253 0.01532984 0.01473388 0.01653953 0.01715871 0.02291613
 0.01366746 0.00966782 0.01185785 0.02789232 0.02396921 0.01184294
 0.01434776 0.01737822 0.01386685 0.01216219 0.01020599 0.01561882
 0.01809729 0.01514133 0.01749658 0.01896083 0.01714388 0.01615925
 0.01190612 0.01337729 0.01259179 0.01101501 0.01228757 0.01678201
 0.03219542 0.01636825 0.01199179 0.01581682 0.01175320 0.01383667
 0.01269465 0.01906779 0.01675467 0.01690286 0.01452614 0.01447623
 0.0164046 0.01251455 0.0172898 0.02779535]
```

Command took 0.13 seconds -- by d_jayanna@uncg.edu at 4/24/2021, 8:28:58 PM on CSC610

Figure 105: Feature importance evaluation

19.7 Calculating Accuracy scores, AUC scores and Confusion matrix

Accuracy of the model is calculated using `sklearn.metrics.accuracy_score`, AUC score is calculated and confusion matrix is printed. Figure 106 shows the execution of accuracy and confusion matrix. Similar model (100 estimators and merged files data set) ran in both

```

nonoverlap01 (Python)
In [15]: from sklearn.metrics import accuracy_score
2: accuracy = accuracy_score(y_test, y_pred)
3: print(accuracy)

0.8106617647058824
```

Command took 0.13 seconds -- by d_jayanna@uncg.edu at 4/24/2021, 8:30:58 PM on CSC610

```

1: from sklearn.metrics import confusion_matrix
2: CC_test = confusion_matrix(y_test, y_pred)
3: print(CC_test)

[[235 48]
 [ 55 206]]
```

Command took 0.02 seconds -- by d_jayanna@uncg.edu at 4/24/2021, 8:30:13 PM on CSC610

Shift+Enter to run

Figure 106: Accuracy and Confusion matrix

big data system and local system which has good results i.e where the big data system is performing well is shown from Figure 107 to Figure 116.

```

overlap012 (Python)
In [1]: import pyspark
```

Command took 0.01 seconds -- by d_jayanna@uncg.edu at 4/27/2021, 1:18:14 PM on CSC610

```

1: df = sqlContext.sql("SELECT * FROM overlap012_3")
2: input_data = df.select("*").toPandas()
```

(1) Spark Jobs

```

3: df_pyspark = pandas.DataFrame.DataFrame(df, 0, integer, 1, integer, 63 more fields)
```

Connect took 12.05 seconds -- by d_jayanna@uncg.edu at 4/27/2021, 1:18:14 PM on CSC610

```

1: import seaborn as sns
2: display(sns.distplot(input_data['30'], kde=False))
```

Figure 107: Reading the tables

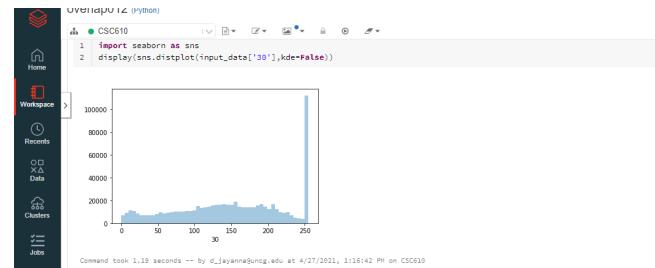


Figure 108: Histogram

```

overlap012 (Python)
In [5]: input_data.isna().any()

Out[10]: 0    False
1    False
2    False
3    False
4    False
...
60   False
61   False
62   False
63   False
64   False
Length: 65, dtype: bool
```

Command took 0.08 seconds -- by d_jayanna@uncg.edu at 4/27/2021, 1:38:46 PM on CSC610

```

1: from sklearn.model_selection import train_test_split
2:
3: train, test = train_test_split(input_data, test_size=0.2, random_state=123)
4: X_train = train.drop(['64'], axis=1)
5: X_test = test.drop(['64'], axis=1)
6: y_train = train['64']
7: y_test = test['64']
```

Command took 0.42 seconds -- by d_jayanna@uncg.edu at 4/27/2021, 1:16:50 PM on CSC610

Figure 109: Splitting the data set

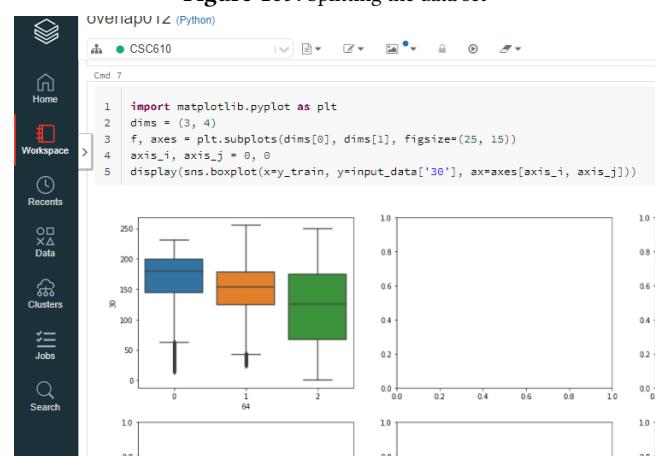


Figure 110: Box plot

20 COMPARING THE RESULTS ON A LOCAL SYSTEM VERSUS A BIG DATA SYSTEM

20.1 Compiling the results from Local system

Two class Nonoverlap data set executed the model in just 0.169 seconds. Three class Nonoverlap data set just took 0.44 seconds to execute or train the model. Two class overlap dataset took 5.4 minutes and three class overlap dataset took 11.50 minutes and three class overlap took 27.26 minutes. I increased the data set by merging 9 files i.e 3 data set from all three categories to increase the overall data set size. This particular model with above said data set

The screenshot shows a Jupyter Notebook interface with the title "OverlappedPU 1.2 (Python)". It contains two code cells. The first cell imports `RandomForestClassifier` from `sklearn.ensemble` and defines a classifier with random_state=0, n_estimators=100, oob_score=True, and n_jobs=1. The second cell prints the prediction of the classifier on the test set. Both commands took approximately 8.74 minutes.

Figure 111: Training and predicting the model

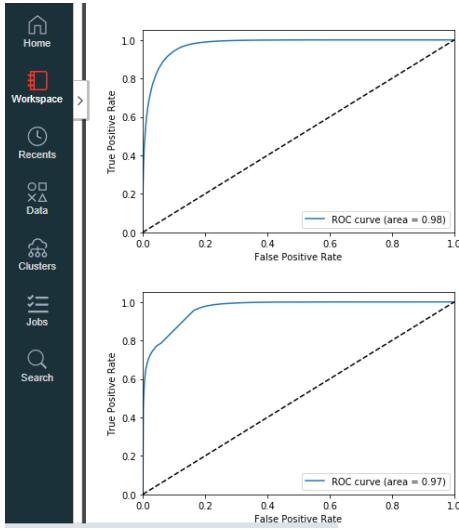


Figure 112: ROC and AUC curve class1 and class2

The screenshot shows a Jupyter Notebook interface with the title "CSC610". It contains two code cells. The first cell prints the prediction of the classifier on the test set. The second cell prints the feature importances and their indices. Both commands took approximately 8.05 seconds.

Figure 113: Feature importance

took 31.15 minutes, 20.65 minutes, 10.21 minutes for 300, 200 and 100 estimator values respectively. But, it failed to allocate memory for 500 estimator value. All the executed models screenshots are attached in the uploaded screenshot folder.

The screenshot shows a Jupyter Notebook interface with the title "overlapped012 (Python)". It contains two code cells. The first cell imports `roc_curve` and `auc` from `sklearn.metrics` and defines variables fpr, tpr, and roc_auc. The second cell iterates over n_classes and prints the ROC curve and AUC for each class. Both commands took 0.06 seconds.

Figure 114: ROC and AUC class0

The screenshot shows a Jupyter Notebook interface with the title "overlap012 (Python)". It contains two code cells. The first cell imports `accuracy_score` from `sklearn.metrics` and defines variables accuracy and CC_test. The second cell prints the confusion matrix CC_test. Both commands took approximately 0.03 seconds.

Figure 115: Accuracy and Confusion matrix

The screenshot shows a Jupyter Notebook interface with the title "C:\Users\Deepa\BigData\P4\Assignment2\RandomForestTreeClassif.py". It displays a command line interface showing various system metrics and command outputs. It also shows a terminal window with some Python code related to a classifier and its performance metrics.

Figure 116: Local environment results

20.2 Compiling the results from Big data system

A similar experiments are conducted in big data system as local system. Two class Nonoverlap executed the model in just 0.169 seconds. Three class Nonoverlap - just took 0.44 seconds to execute or train the model. Two class overlap - took 5.4 minutes and three class overlap took 11.50 minutes. Three class overlap took 27.26 minutes. I increased the data set by merging 9 files i.e 3 data set from all three categories to increase the overall data set size. This particular model with above said data set took 31.15 minutes, 20.65 minutes, 10.21 minutes for 300, 200 and 100 estimator values respectively. But, it failed to allocate memory for 500 estimator value. All the executed models screenshots are attached in the uploaded screenshot folder.

20.3 Comparing the results

I started off with the best models I summarized in assignment 2 shown in Figure 87 and Figure 88. For Nonoverlap two class, local system took 0.169 second to complete training the model for 50 estimators. And the big data system took 0.64 seconds. For Nonoverlap three class, model was trained with 100 estimators. Local system took 0.44 seconds and bigdata system took 1.63 seconds. Overlap

two class was trained for 1000 estimators. Local system took 5.4 minutes and big data system took 12.78 minutes. Overlap three class trained for 1000 estimators. Local system took around 28 minutes and big data system took 25 minutes for execution. The table showing the summary of these results is shown in figure 117. I

	Nonoverlap01	Nonoverlap012	Overlap01	Overlap012
Estimators	50	100	1000	1000
Local env	0.169s	0.44s	5.4m	27.26m
Bigdata env	0.64s	1.63s	12.78m	25m
Accuracy	80.33	82.1	92.95	93.8

Figure 117: Comparing local and Bigdata environment

observed local system was performing better than big data system for small data set like nonoverlap two, three class and two class overlap data set. Overlap three class data set has way more observations compared to others. This set started performing better than local system. So I chose to increase this specific dataset with more observations. This got me to an understanding that big data system really needs to have large data set to perform better than local systems. I created a new merged data set of 9 files, three from each category of fruits. Then, I carried few experiments to have a better understanding on Big data environment performance. Four

Overlap three class multiple files merged data				
Estimator values	Bigdata	Local	Accuracy Local	Accuracy Bigdata
500	44.7	Memory error	NA	87.5
300	25.48	31.15	87	87
200	17.38	20.65	89.8	86.94
100	8.74	10.21	89.8	86.87

All numbers are in minutes

Figure 118: Comparing performance of local and big data system with overlap data

experiments are conducted on the above said dataset by changing the estimator values and setting n_jobs to 1.

- Experiment 1:- 100 estimators - In this trial, both local and big data system are trained with 100 estimators. Local system ended up taking 10.21 minutes and big data system took 8.74.
- Experiment 2: 200 estimators - In this trial, both local and big data system are trained with 200 estimators. Local system ended up taking 20.65 minutes and big data system took 17.38 min.
- Experiment 3: 300 estimators - In this trial, both local and big data system are trained with 300 estimators. Local system ended up taking 31.15 minutes and big data system took 25.48 min.
- Experiment 4: 500 estimators - In this trial, both local and big data system are trained with 400 estimators. Local system failed in the middle throwing "Could not allocate memory". But, the big data environment successfully executed the model in 44.7 minutes.

A cumulative table with all the execution time and the accuracy of specific models of both big data and local system is summarized in Figure 118. In conclusion, big data environment started performing better after increasing the data set, as well as increasing the estimator values. From figure 117 and 118, we can see that bigdata environment time increased by a lesser number compared to local system while increasing the estimator values as well.