

Big Data Assignment1: Feature Space Generation

Manish Shah

University of North Carolina at Greensboro

m_shah2@uncg.edu

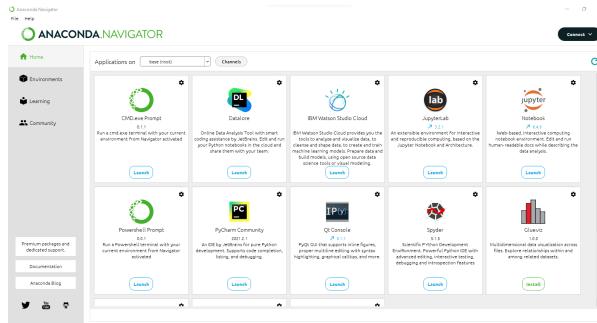


Fig. 1: Anaconda GUI

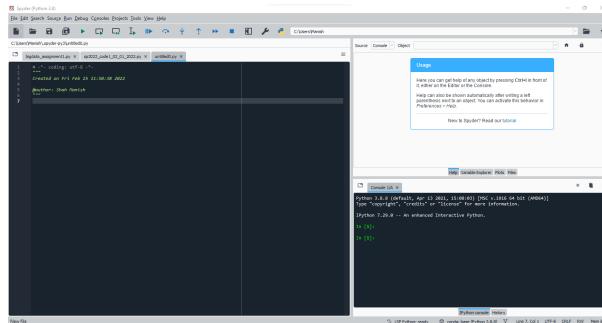


Fig. 2: SpyderIDE

I. TASK 1: BUILD YOUR PROGRAMMING ENVIRONMENT!

1.1 Download and install Anaconda software.

I have downloaded and install the latest version of Anaconda software on local machine. It has inbuilt python packages. Also JupyterNotebook and Spyder can be installed using Anaconda. The Figure 1 shows the anaconda installation.

1.2 Install Spyder IDE.

The latest version of Spyder can be downloaded from Anaconda software. It can be accessed from Anaconda bundle and also by directly going inside the Spyder directory in the local machine. The Figure 2 shows sydper IDE installed in the local machine.

1.3 Install Jupyter Notebook.

Similar to Spyder, the latest version of JupyterNotebook can be downloaded from Anaconda software. It can be accessed from Anaconda bundle and also by directly searching JupyterNotebook in the local machine.

1.4 Download and install OpenCV.

OpenCV can be dowloaded online or via anaconda prompt. In my case, I downloaded it online from here and added it to

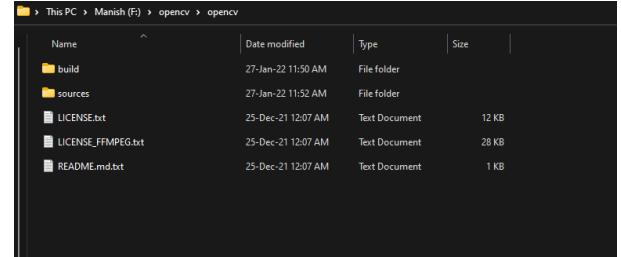


Fig. 3: OpenCV

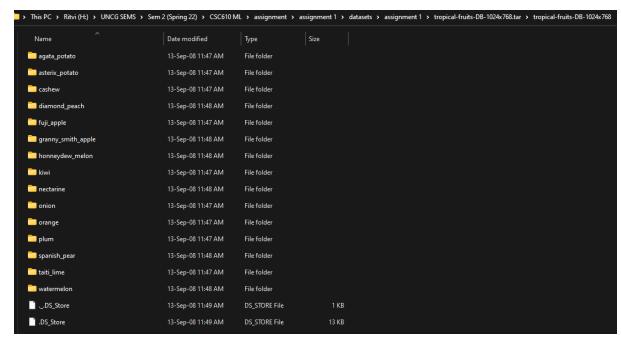


Fig. 4: Tropical Fruits Folder

the system variables. The Figure 3 shows OpenCV installed on the local machine.

II. TASK 2: DOWNLOAD OR GENERATE A FRUITS AND VEGETABLES IMAGE DATASET!

2.1 Tropical fruit dataset.

This dataset can be downloaded from here and contains images of tropical fruits. The folder is shown in Figure 4.

2.2 FIDS30 Dataset.

The FIDS dataset can be downloaded from here. The folder contains images of 32 kinds of fruits and is shown in Figure 5.

2.3 Image selection.

For this assignment, I am selecting images from FIDS30 dataset. I have selected cantaloupe, banana and tomato images. The images are selected such that there is no background except the fruit images. I am representing cantaloupe by 16.jpg, banana by 43.jpg and tomato by 5.jpg in the project. The selected images are displayed in Figure 6-Figure 8.

Name	Date modified	Type	Size
acerolas	03-Feb-22 9:36 AM	File folder	
apples	03-Feb-22 9:36 AM	File folder	
apricots	03-Feb-22 9:36 AM	File folder	
avocados	03-Feb-22 9:36 AM	File folder	
bananas	03-Feb-22 9:36 AM	File folder	
blackberries	03-Feb-22 9:36 AM	File folder	
blueberries	03-Feb-22 9:36 AM	File folder	
cantaloupes	03-Feb-22 9:36 AM	File folder	
cherries	03-Feb-22 9:36 AM	File folder	
coconuts	03-Feb-22 9:37 AM	File folder	
figs	03-Feb-22 9:37 AM	File folder	
grapefruits	03-Feb-22 9:37 AM	File folder	

Fig. 5: Tropical Fruits Folder



Fig. 6: Cantaloupe image

III. TASK 3:READ IMAGES AND DISPLAY IN THE PROGRAMMING ENVIRONMENT

3.1 Read image and display R G B channel.

For this assignment I am using Spyder IDE. OpenCV library is used to read color images of fruits and display the R,G,B color channels for each image and plot the images using matplotlib library. The code snippet is shown in Figure 9.

3.2 Convert images to gray scale and display dimensions.



Fig. 7: Banana image



Fig. 8: Tomato image

```
import pandas as pd
import cv2
import matplotlib.pyplot as plt
import numpy as np

# Read Images (e.g., .jpg, .png, and .tif)
cant_color = cv2.imread("C:/Users/Manish/Desktop/bigdata_assignment1/16.jpg")
bana_color = cv2.imread("C:/Users/Manish/Desktop/bigdata_assignment1/43.jpg")
toma_color = cv2.imread("C:/Users/Manish/Desktop/bigdata_assignment1/5.jpg")

# cant_color = cv2.imread("C:/Users/s_suthah/Desktop/Images/02_03.png")
# bana_color = cv2.imread("C:/Users/s_suthah/Desktop/Images/s_02_-04.tif")

# Display the color channels
plt.imshow(cant_color[:, :, 0])
plt.imshow(cant_color[:, :, 1])
plt.imshow(cant_color[:, :, 2])

plt.imshow(bana_color[:, :, 0])
plt.imshow(bana_color[:, :, 1])
plt.imshow(bana_color[:, :, 2])

plt.imshow(toma_color[:, :, 0])
plt.imshow(toma_color[:, :, 1])
plt.imshow(toma_color[:, :, 2])
```

Fig. 9: Code to display R,G,B values of fruits images

In this step, first colored images are converted to grayscale using cvtColor function of openCV library. To display dimensions of grayscale images, image.shape function is used. The code is displayed in Figure 10 and grayscale images are shown in Figure 11 - Figure 13.

```
@author: Shah Manish
"""
import pandas as pd
import cv2
import matplotlib.pyplot as plt
import numpy as np

# Read Images (e.g., .jpg, .png, and .tif)
cant_color = cv2.imread("C:/Users/Manish/Desktop/bigdata_assignment1/16.jpg")
bana_color = cv2.imread("C:/Users/Manish/Desktop/bigdata_assignment1/43.jpg")
toma_color = cv2.imread("C:/Users/Manish/Desktop/bigdata_assignment1/5.jpg")

# Convert to grayscale.
cantG = cv2.cvtColor(cant_color, cv2.COLOR_BGR2GRAY)
heightG, widthG = cantG.shape

banaG = cv2.cvtColor(bana_color, cv2.COLOR_BGR2GRAY)
heightbG, widthbG = banaG.shape

tomaG = cv2.cvtColor(toma_color, cv2.COLOR_BGR2GRAY)
heighttG, widthtG = tomaG.shape
```

Fig. 10: Code to convert the images to grayscale

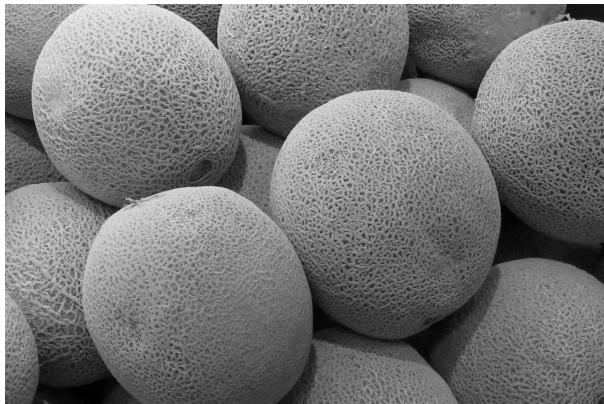


Fig. 11: Cantaloupe gray



Fig. 12: Banana gray



Fig. 13: Tomato gray

```

Author: Shah Manish
***
import cv2
#read the image
cantG = cv2.imread('C:/Users/Manish/Desktop/bigdata_assignment1/output_new/16_gray.jpg')
#print the original dimensions
print("Original Dimensions:",cantG.shape)
std_height = 261
height = cantG.shape[0]
width = cantG.shape[1]

#resize function
def newwidth(height , width):
    print("Original Height : ", height)
    print("Original Width : ", width)
    aspect_ratio = round(height/width)
    new_width = aspect_ratio * 261
    print("My width : ", new_width)
    return round(new_width)

#function to find divisible by 9
def div9(dim):
    remainder = dim % 9
    if remainder == 0:
        dim = dim
    else:
        dim = dim + (9-remainder)
    return dim

#values when both need to be divisible by 9
new_fruit_width = newwidth(height,width)
new_height = std_height

#values when both the height and width is to be divisible by 9
print("New Dimension of image divisible by 9: ", new_height, new_fruit_width)
image_resized = cv2.resize(cantG, dsize = (new_height , new_fruit_width))
cv2.imwrite('C:/Users/Manish/Desktop/bigdata_assignment1/output_new/16_grayresized.jpg' , image_resized)

```

Fig. 14: Function to resize grayscale images divisible by 9

```

Created Sun Feb 20 23:49:17 2022
Author: Shah Manish
***
import cv2
#read the image
cantG = cv2.imread('C:/Users/Manish/Desktop/bigdata_assignment1/output_new/16_gray.jpg')
#print the original dimensions
print("Original Dimensions:",cantG.shape)
std_height = 261
height = cantG.shape[0]
width = cantG.shape[1]

#resize function
def newwidth(height , width):
    print("Original Height : ", height)
    print("Original Width : ", width)
    aspect_ratio = round(height/width)
    new_width = aspect_ratio * 261
    print("My width : ", new_width)
    return round(new_width)

#function to find divisible by 9
def div9(dim):
    remainder = dim % 9
    if remainder == 0:
        dim = dim
    else:
        dim = dim + (9-remainder)
    return dim

#values when both need to be divisible by 9
new_fruit_width = newwidth(height,width)
new_height = std_height

#values when both the height and width is to be divisible by 9
print("New Dimension of image divisible by 9: ", new_height, new_fruit_width)
image_resized = cv2.resize(cantG, dsize = (new_height , new_fruit_width))
cv2.imwrite('C:/Users/Manish/Desktop/bigdata_assignment1/output_new/16_grayresized.jpg' , image_resized)

```

Fig. 15: Find new width of cantG image

IV. TASK 4: RESIZE THE IMAGES TO REDUCE THEIR DIMENSIONS!

4.1 Resize the grayscale images so that dimensions are divisible by 9.

In this I have written a function which reduces the grayscale images in pixel size, such that the output dimensions are divisible by 9. For this a parametric function is created in python, which checks divisibility of the dimensions (width/height) by 9. If not, reduce it, such that it is divisible by 9 by subtracting it from the remainder. The code snippet is shown in the Figure 14.

4.2 Resize the grayscale images to height of 261.

Here the height of grayscale images is fixed as 261. The code is written which maintains the aspect ratio and calculate the new width of the images such that it is divisible by 9. The code is applied on all the three fruit images individually and the new images are stored with a different names. Figure 15 shows code for finding the new width of cantaloupe.

V. TASK 5: GENERATE BLOCK-FEATURE VECTORS!

5.1 Divide image into block of 9*9 pixels and generate non overlapping feature space

Fig. 16: Generating non overlapping feature vectors

Fig. 17: Non overlapping Feature space for Cantaloupe having label 0

In this part, the resized grayscale images are used. The code for converting images to grayscale and resizing is added with the code of creating non overlapping blocks of 9×9 size to generate feature vectors of size 9×9 . The feature vectors are stored in a .csv file for each image. The labels for images of cantaloupe, banana and tomato in the feature space in the .csv file with labels of 0, 1, 2. All the images contains 841 blocks of 9×9 size. In the feature space there is a total of 81 features and 1 column for label and total of 841 observations excluding the top row. Figure 16 shows the code for creating non overlapping feature space and Figure 17 contains .csv file generated for cantaloupe fruit.

VI. TASK 6: GENERATE SLIDING BLOCK-FEATURE VECTORS!

6.1 Divide image into block of 9×9 pixels and generate overlapping feature vector space

In this part, an overlapping block of size 9×9 is used to generate feature vector space of 81 columns. Images of cantaloupe, banana and tomato are labelled with 0, 1 and 2 respectively. The feature vectors are stored in a .csv file which contains 64010 blocks each. The code contains converting colored image to grayscale and resizing the images to get width divisible by 9 and overlapping code. The code snippet is shown in Figure 18 and the generated csv in Figure 19.

VII. TASK 7: DERIVE STATISTICAL DESCRIPTORS!

7.1 Extract statistical information about number of observations and data dimensions

For non overlapping feature space there are 841 observations with 81 features for cantaloupe, banana and tomato images. In case of overlapping feature space, there are 64010 observations with 81 features for all images.

Fig. 18: Generating overlapping feature vectors

Fig. 19: Overlapping Feature space for Cantaloupe having label 0

7.2 Extract statistical information about mean of each feature

To calculate mean, I took the non overlapping .csv files and find the mean of each feature for each image and plot it which can be seen in Figure 21. The code snippet for finding the mean can be seen in Figure 20. From the two plots we can see that for the overlapping blocks, the difference between mean values of consecutive features is very less as compared to the non overlapping blocks. From mean we can see that banana has higher mean than both, cantaloupe has higher mean than tomato. This tells that the average values in banana is higher, which infers it has brighter image pixels than both. Secondly we can infer cantaloupe image is brighter than tomato image.

7.3 Extract statistical information about standard deviation of each feature

In this I have calculated standard deviation of each feature for every fruit in non overlapping blocks csv files. The code for calculating standard deviation is in the Figure 23 and the

```

import matplotlib.pyplot as plt
import pandas as pd

#read the csv file
cant = pd.read_csv('C:/Users/Manish/Desktop/bigdata_assignment1/16_nonoOverlap.csv')
bana = pd.read_csv('C:/Users/Manish/Desktop/bigdata_assignment1/43_nonoOverlap.csv')
toma = pd.read_csv('C:/Users/Manish/Desktop/bigdata_assignment1/5_nonoOverlap.csv')

#calculate the mean of the CSV files
mean_cant = round(cant.mean())
mean_bana = round(bana.mean())
mean_toma = round(toma.mean())

a =print(mean_cant)
b =print(mean_bana)
c =print(mean_toma)

#set x and y values
plt.title('Mean Plot Non-Overlapping Set')
plt.xlabel('Features')
plt.ylabel('Mean Values')
#Plot the mean values
plt.plot(mean_cant, 'r')
plt.plot(mean_bana, 'b')
plt.plot(mean_toma, 'g')

```

Fig. 20: Code to calculate mean of fruit data

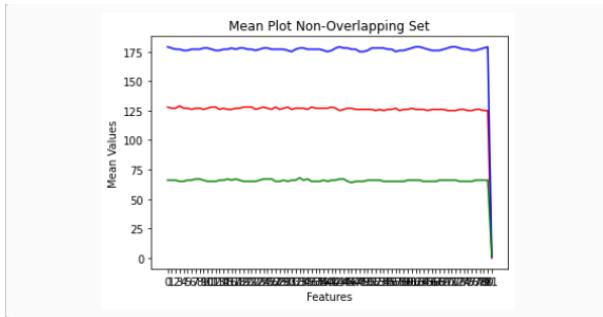


Fig. 21: Mean of non overlapping fruit data

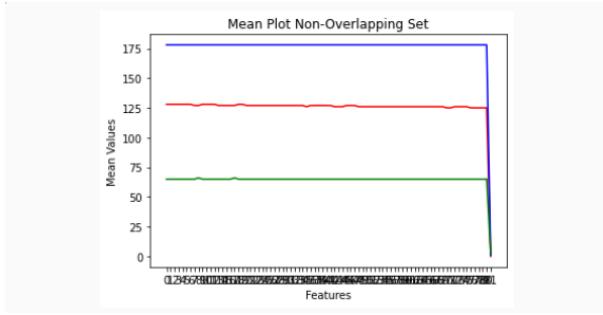


Fig. 22: Mean of overlapping fruit data

plot is shown in Figure 24. From the plot it can be seen the standard deviation is increasing among all fruit data. From the plot it can be seen that banana has higher standard deviation than others. So it is more spread out while other fruit values are more closely cluster around the mean.

7.4 Histogram of fruit images

I have selected one feature for each image and plot histogram for both overlapping and non overlapping feature vectors. For cantaloupe, the features are skewed towards the middle for both overlapping and non overlapping. For banana, the features are skewed towards the right, for both overlapping

```
import matplotlib.pyplot as plt
import pandas as pd
#import numpy as np

#read the csv file
cant = pd.read_csv('C:/Users/Manish/Desktop/bigdata_assignment1/16_nonoverlap.csv')
bana = pd.read_csv('C:/Users/Manish/Desktop/bigdata_assignment1/43_nonoverlap.csv')
toma = pd.read_csv('C:/Users/Manish/Desktop/bigdata_assignment1/5_nonoverlap.csv')

std_cant = round(cant.std())
std_bana = round(bana.std())
std_toma = round(toma.std())

a =print(std_cant)
b =print(std_bana)
c =print(std_toma)

#Plot the mean values
plt.plot(std_cant, 'r')
plt.plot(std_bana, 'b')
plt.plot(std_toma, 'g')

#Setting the Co-ordinate values
plt.title('Standard Deviation for nonoverlapping data')
plt.xlabel('Features')
plt.ylabel('standard deviation')
```

Fig. 23: Code to find Standard deviation of non overlapping fruit data

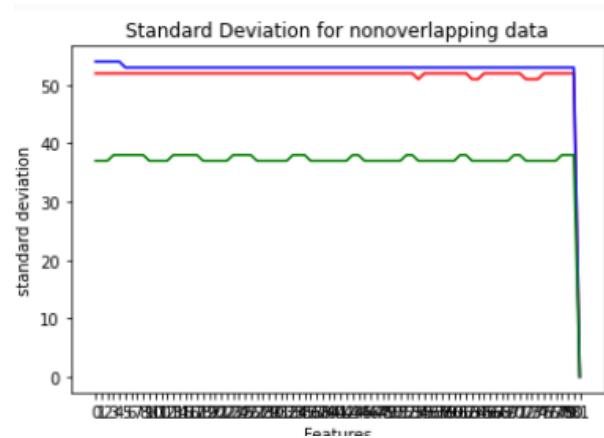


Fig. 24: plot of standard deviation of overlapping fruit data

```
import matplotlib.pyplot as plt
import pandas as pd

#read the fruit
cant = pd.read_csv('C:/Users/Manish/Desktop/bigdata_assignment1/overlap_16.csv')
bana = pd.read_csv('C:/Users/Manish/Desktop/bigdata_assignment1/overlap_43.csv')
toma = pd.read_csv('C:/Users/Manish/Desktop/bigdata_assignment1/overlap_5.csv')

feature = cant['16']

nbins = 30
plt.hist(feature, nbins )
plt.title("Cantaloupe Histogram for feature 16")
plt.xlabel('Feature 16')
plt.ylabel('value')
plt.show()
```

Fig. 25: Histogram generation code

and non overlapping. For tomato, the features are skewed towards the left for both overlapping and non overlapping. The code is shown in Figure 25 and histograms are shown below that.

7.5 Scatter Plot of fruit images

I have generate scatter plot of all the three fruits for both overlapping and non overlapping feature vectors. I have selected feature 16 and 27 for cantaloupe, features 45 and 64 for banana and features 23 and 68 for tomato. The code can be seen in Figure 32 and the scatter plots are mentioned below

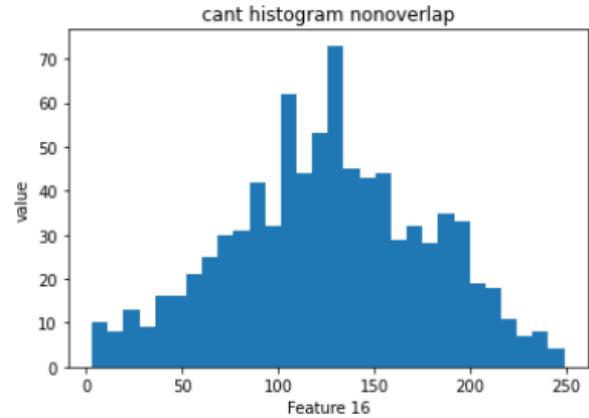


Fig. 26: Histogram of cant non overlap

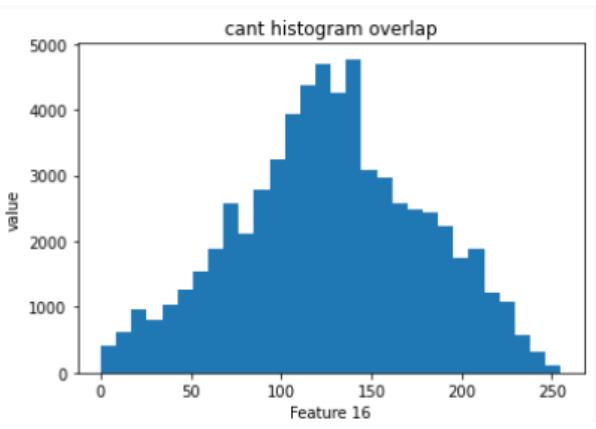


Fig. 27: Histogram of cant overlap

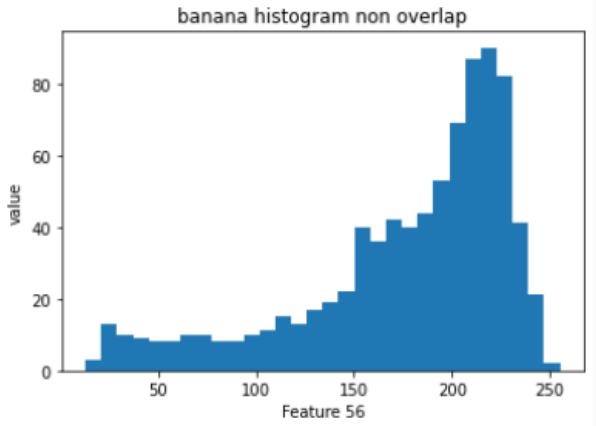


Fig. 28: Histogram of banana non overlap

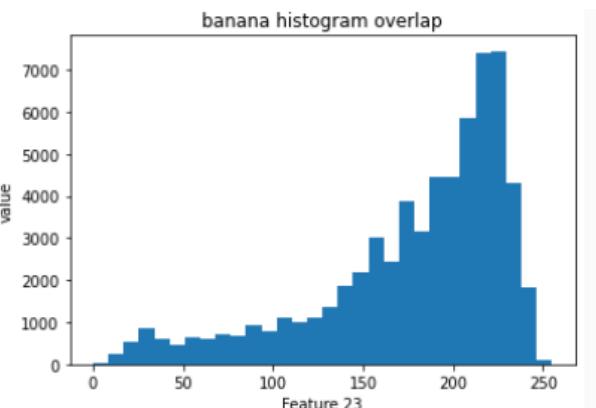


Fig. 29: Histogram of banana overlap

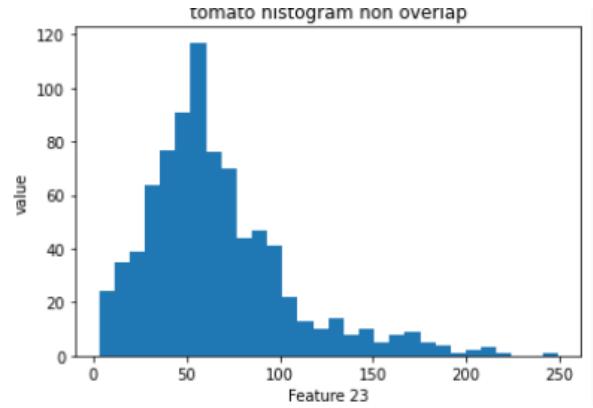


Fig. 30: Histogram of tomato non overlap

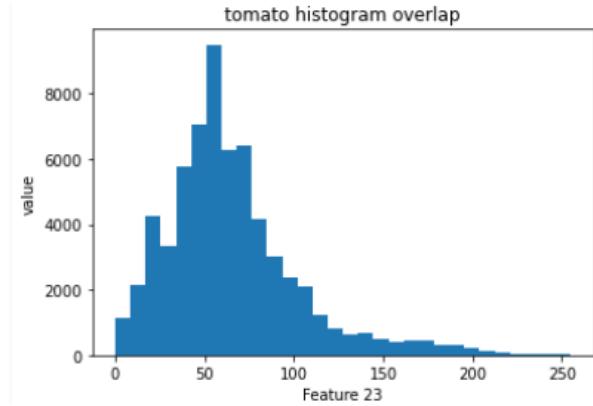


Fig. 31: Histogram of tomato overlap

it.

7.6 Answer the questions!

- 1 Data characteristics: Imbalanced, Incomplete, Inaccurate.

The data is balanced because the number of observations in all the three csv files are the same for both overlapping and non overlapping. As the number of class labels are the same for all so it is balanced. The data is not incomplete because it has no missing values in any csv files. The data is accurate as was it is generated properly and individually so it cannot be inaccurate.

- 2 Is it trivial data or possibly Big Data?

```
import matplotlib.pyplot as plt
import pandas as pd

#read the csv file
cant = pd.read_csv('C:/Users/Manish/Desktop/bigdata_assignment1/16_nonoverlap.csv')

f16 = cant['16']
f27 = cant['27']

plt.scatter(f16, f27, color =['green'], s=1)
plt.title("Feature Plot Cantaloupe")
plt.xlabel('Feature 16')
plt.ylabel('Feature 27')
plt.show()
```

Fig. 32: Scatter plot generation code

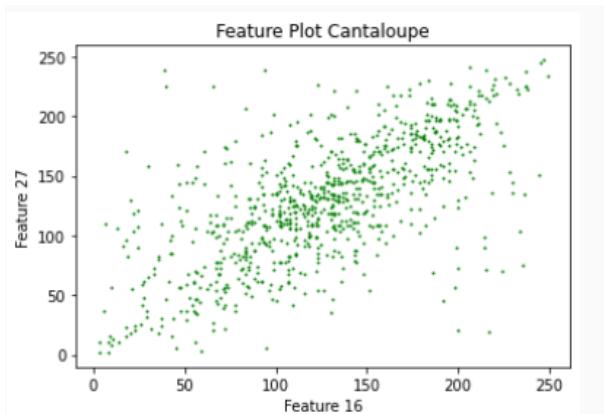


Fig. 33: Scatter plot for cantaloupe non overlap

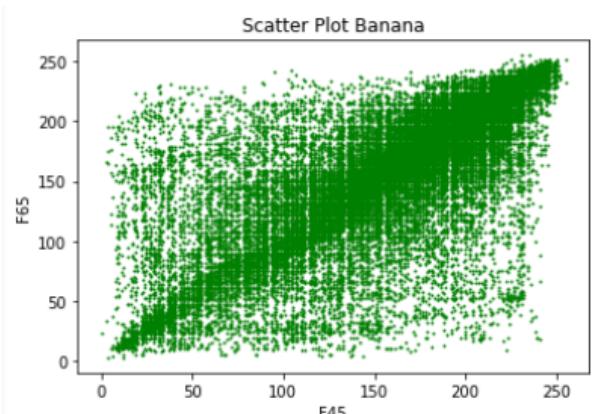


Fig. 36: Scatter plot for banana overlap

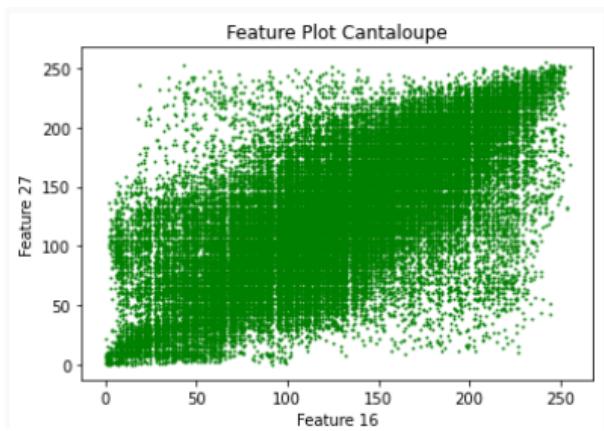


Fig. 34: Scatter plot for cantaloupe overlap

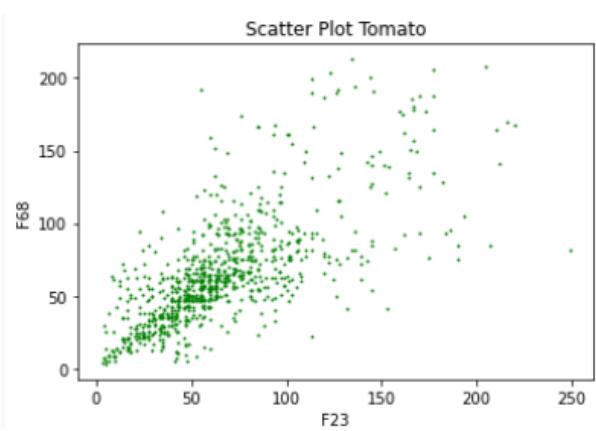


Fig. 37: Scatter plot for tomato non overlap

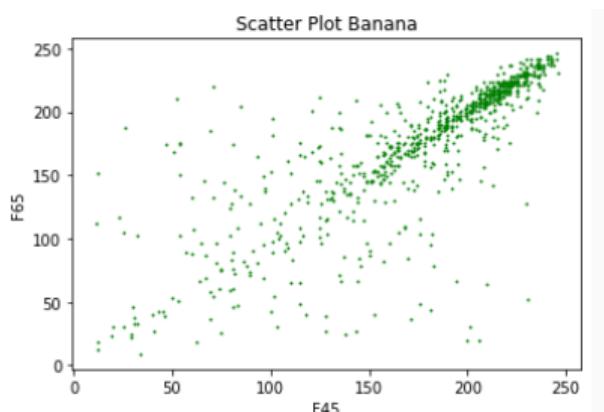


Fig. 35: Scatter plot for banana non overlap

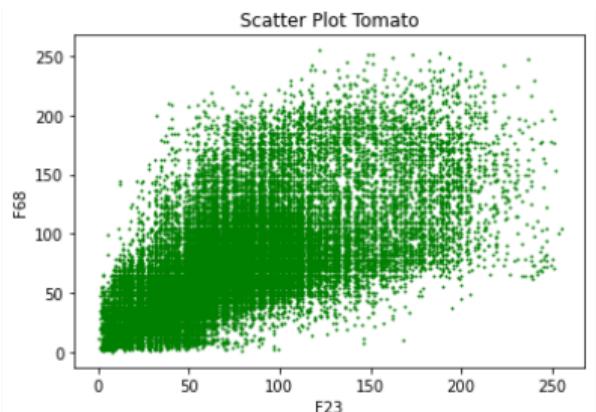


Fig. 38: Scatter plot for tomato overlap

```

import pandas as pd

cant = pd.read_csv('C:/Users/Manish/Desktop/bigdata_assignment1/16_nonoverlap.csv')
banana = pd.read_csv('C:/Users/Manish/Desktop/bigdata_assignment1/43_nonoverlap.csv')
toma = pd.read_csv('C:/Users/Manish/Desktop/bigdata_assignment1/5_nonoverlap.csv')

frames = [cant,banana]
merged = pd.concat(frames)

merged.to_csv('C:/Users/Manish/Desktop/bigdata_assignment1/image01.csv')

```

Fig. 39: Code to merge two files

The dataset generated is a trivial data and not big data as it is not exponentially increasing with time. It is trivial because for big data the number of classes should be many whereas in our case it's just three classes/labels. Secondly, the number of observations is more than the number of features so it is low dimensional data.

3 Does it have scalability problem? Are they high dimensional?

The data is not high dimensional because the number of observations are more than the number of rows. The data generated is not increasing with time so it doesn't have any scalability problem. As each feature vector has size 9*9 it has no scalability problem. If the number of feature increases then it becomes high dimensional and scalability problem might arise but right now it has no scalability problem.

4 Do you need to standardize? Do you need to normalize? How do they affect the data characteristics?

Yes, the data needs to be normalize and standardize. To have high accuracy all the values of the dataset needs to follow the same curve/data pattern. It is necessary for all the data points to fall in the same range and for that normalization and standardization is required. If the data is not standardize/normalized, then it would create issue while reading the data and will output wrong information. It has no effect on data characteristics.

VIII. TASK 8: CONSTRUCT A FEATURE SPACE!

In this part, I am making use of my files 16_nonoverlap.csv and 43_nonoverlap.csv which represents cantaloupe and banana files and use the concat function to merge the csv files and name the new created file as image01.csv. Secondly I use image01.csv to merge with 5_nonoverlap.csv and created a super merged file and name it as image012.csv. Lastly I use randomize function to randomize the placement of data in the image012.csv and generate a random dataset in it. The code to merge the datasets is shown in Figure 39, randomize code in Figure 42 and the generated dataset files are displayed after that.

IX. TASK 9: DISPLAY SUBSPACES!

9.1 Plot 2D feature space

For creating 2D scatter plot, I am using files 16_nonoverlap.csv (cantaloupe file), 43_nonoverlap.csv (banana file), 5_nonoverlap.csv (tomato file). I have selected features 45 and 75 for when combining cantaloupe and banana images, features 23 and 57 when combining banana and tomato images and features 14 and 69 when combining

831	149	140	99	116	187	117	135	142	131	151	140	109	116	106	102	119	115	70	122	153	135	135	117	100	0
832	151	140	100	125	166	122	116	155	124	154	131	140	116	129	141	126	144	140	130	159	128	141	141	0	
833	152	141	101	126	167	123	117	156	125	157	132	141	117	130	142	127	145	141	131	159	129	142	142	0	
834	158	145	142	151	188	134	137	159	136	148	141	142	169	111	121	136	129	170	160	118	145	134	139	90	0
835	159	146	143	152	189	135	138	160	148	150	143	144	170	112	122	137	130	171	161	119	146	135	139	91	0
836	160	147	144	153	190	136	139	161	149	151	144	145	171	113	123	138	131	172	162	120	147	136	140	92	0
837	161	148	145	154	191	137	140	162	150	152	145	146	172	114	124	139	132	173	163	121	148	137	141	93	0
838	163	149	96	111	184	109	90	96	81	125	96	109	105	95	64	94	84	61	95	94	84	75	107	0	
839	164	150	97	112	185	110	91	97	82	126	97	109	106	96	75	95	94	85	95	95	85	108	0		
840	167	151	98	113	186	111	92	98	83	127	112	109	110	110	85	66	96	87	86	107	0	0	0		
841	167	152	99	114	187	112	93	99	84	128	113	110	111	111	86	67	97	88	87	108	0	0	0		
842	168	153	100	115	188	113	94	100	85	129	114	111	112	112	87	68	98	89	88	109	0	0	0		
843	169	154	101	116	189	114	95	101	86	130	115	112	113	113	88	69	99	89	88	110	0	0	0		
844	170	155	102	117	190	115	96	102	87	131	116	113	114	114	89	70	100	90	89	111	0	0	0		
845	171	156	103	118	191	116	97	103	88	132	117	114	115	115	90	71	101	91	90	112	0	0	0		
846	172	157	104	119	192	117	98	104	89	133	118	115	116	116	91	72	102	92	91	113	0	0	0		
847	173	158	105	120	193	118	99	105	90	134	119	116	117	117	92	73	103	93	92	114	0	0	0		
848	174	159	106	121	194	119	100	106	91	135	120	117	118	118	93	74	104	94	93	115	0	0	0		
849	175	160	107	122	195	120	101	107	92	136	121	118	119	119	94	75	105	95	94	116	0	0	0		
850	176	161	108	123	196	121	102	108	93	137	122	119	120	120	95	76	106	96	95	117	0	0	0		
851	177	162	109	124	197	122	103	109	94	138	123	120	121	121	96	77	107	97	96	118	0	0	0		
852	178	163	110	125	198	123	104	110	95	139	124	121	122	122	97	78	108	98	97	119	0	0	0		
853	179	164	111	126	199	124	105	111	96	140	125	122	123	123	98	79	109	99	98	120	0	0	0		
854	180	165	112	127	200	125	106	112	97	141	126	123	124	124	99	80	110	100	99	121	0	0	0		
855	181	166	113	128	201	126	107	113	98	142	127	124	125	125	100	81	111	101	100	122	0	0	0		
856	182	167	114	129	202	127	108	114	99	143	128	125	126	126	101	82	112	102	101	123	0	0	0		
857	183	168	115	130	203	128	109	115	100	144	129	126	127	127	102	83	113	103	102	124	0	0	0		
858	184	169	116	131	204	129	110	116	101	145	130	127	128	128	103	84	114	104	103	125	0	0	0		
859	185	170	117	132	205	130	111	117	102	146	131	128	129	129	104	85	115	105	104	126	0	0	0		
860	186	171	118	133	206	131	112	118	103	147	132	129	130	130	105	86	116	106	105	127	0	0	0		
861	187	172	119	134	207	132	113	119	104	148	133	130	131	131	106	87	117	107	106	128	0	0	0		
862	188	173	120	135	208	133	114	120	105	149	134	131	132	132	107	88	118	108	107	129	0	0	0		
863	189	174	121	136	209	134	115	121	106	150	135	132	133	133	108	89	119	109	108	130	0	0	0		
864	190	175	122	137	210	135	116	122	107	151	136	133	134	134	109	90	120	110	109	131	0	0	0		
865	191	176	123	138	211	136	117	123	108	152	137	134	135	135	110	91	121	111	110	132	0	0	0		
866	192	177	124	139	212	137	118	124	109	153	138	135	136	136	111	92	122	112	111	133	0	0	0		
867	193	178	125	140	213	138	119	125	110	154	139	136	137	137	112	93	123	113	112	134	0	0	0		
868	194	179	126	141	214	139	120	126	111	155	140	137	138	138	113	94	124	114	113	135	0	0	0		
869	195	180	127	142	215	140	121	127	112	156	141	138	139	139	114	95	125	115	114	136	0	0	0		
870	196	181	128	143	216	141	122	128	113	157	142	139	140	140	115	96	126	116	115	137	0	0	0		
871	197	182	129	144	217	142	123	129	114	158	143	140	141	141	116	97	127	117	116	138	0	0	0		
872	198	183	130	145	218	143	124	130	115	159	144	141	142	142	117	98	128	118	117	139	0	0	0		
873	199	184	131	146	219	144	125	131	116	160	145	142	143	143	118	99	129	119	118	140	0	0	0		
874	200	200	210	219	209	204	210	209	210	210	213	210	209	208	209	209	213	213	213	213	213	213	213	213	0
875	201	201	211	220	210	205	211	209	212	212	206	210	209	208	207	207	213	209	209	209	209	209	209	209	0
876	202	202	212	221	211	206	212	209	213	213	207	211	209	208	207	207	213	209	209	209	209	209	209	209	0
877	203	203	213	222	212	207	213	210	214	214	208	212	210	209	208	208	213	209	209	209	209	209	209	209	0
878	204	204	214	223	213	208	214	211	215	215	209	213	211	210	209	209	213	209	209	209	209	209	209	209	0
879	205	205	215	224	214	209	215	212	216	216	210	214	212	211	210	210	214	210	210	210	210	210	210	210	0
880	206	206	216	225	215	210	216	213	217	217	211	215	213	212	211	211	215	211	211	211	211	211	211	211	0
881	207	207	217	226	216	211	217	214	218	218	212	216	214	213	212	212	216	212	212	212	212	212	212	212	0
882	208	208	218	227	217	212	218	215	220	220	213	217	215	214	213	213	217	213	213	213	213	213	213	213	0
883	209	209	219	228	218	213	219	216	221	221	214	218	216	215	214	214	218	214	214	214	214	214	214	214	0
884	210	210	220	229																					

```

import matplotlib.pyplot as plt
import pandas as pd

#read the csv file
cant = pd.read_csv('C:/Users/Manish/Desktop/bigdata_assignment1/16_nonoverlap.csv')
bana = pd.read_csv('C:/Users/Manish/Desktop/bigdata_assignment1/43_nonoverlap.csv')

f45 = bana['45']
f75 = bana['75']

j45 = cant['45']
j75 = cant['75']

plt.scatter(f45, f75, color =[ 'yellow' ] , s=1 )
plt.scatter(j45, j75, color =[ 'blue' ] , s=1 )

plt.title("Scatter Plot Cantaloupe and Banana")
plt.xlabel('Feature 45')
plt.ylabel('Feature 75')
plt.show()

```

Fig. 44: Scatter plot 2D code

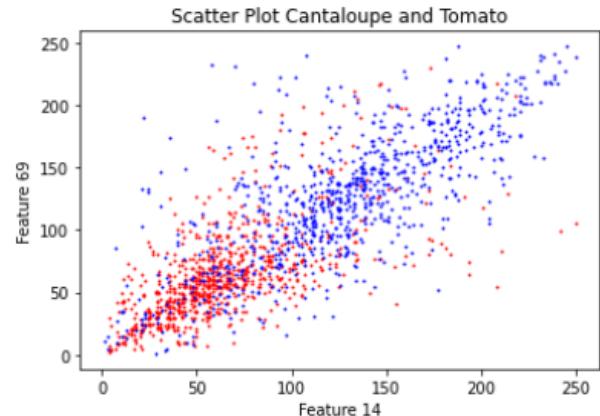


Fig. 47: 2d plot cant/toma

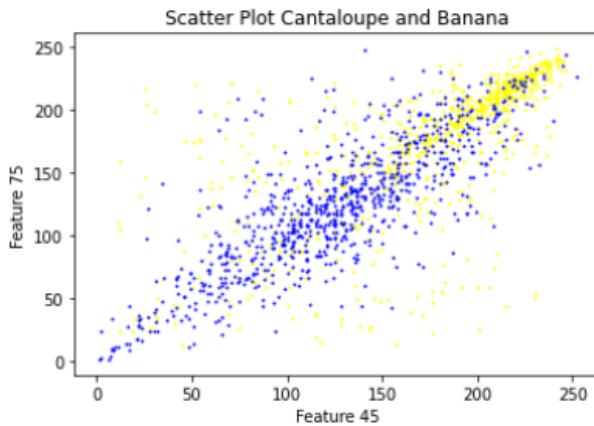


Fig. 45: 2d plot cant/bana

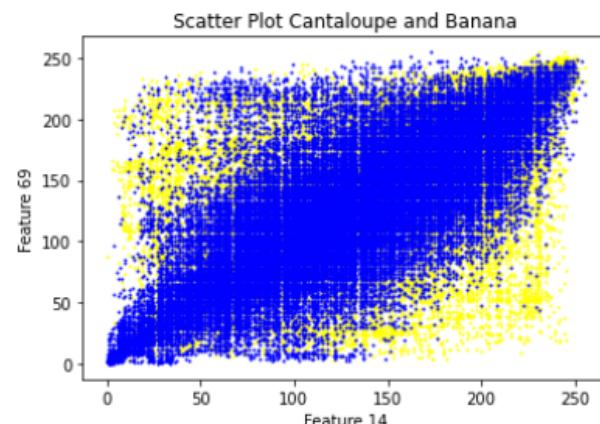


Fig. 48: 2d plot cant/bana overlapping

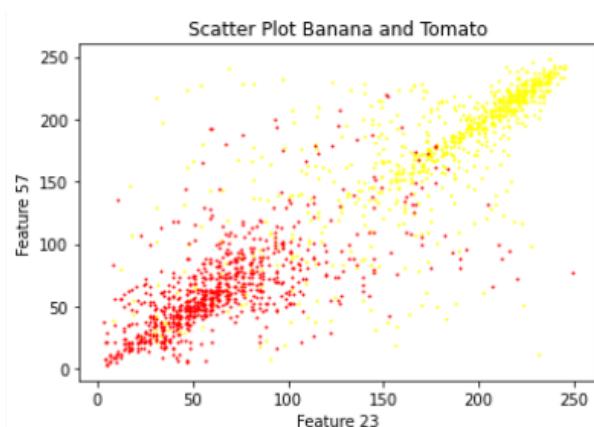


Fig. 46: 2d plot bana/toma

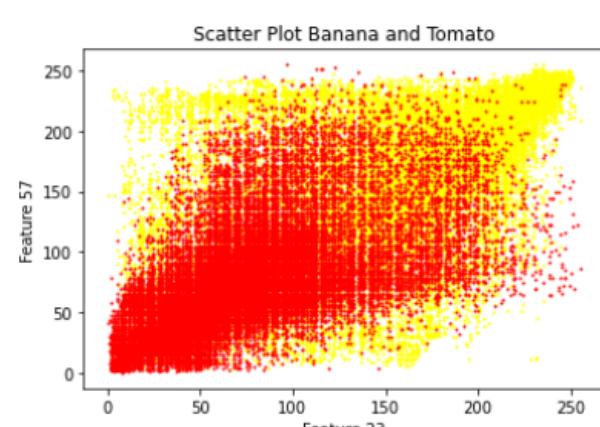


Fig. 49: 2d plot bana/toma overlapping

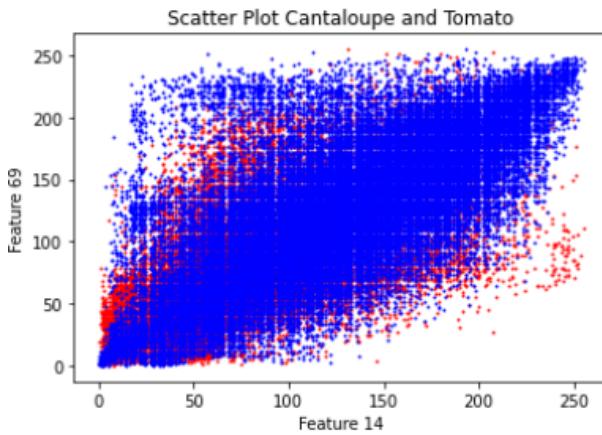


Fig. 50: 2d plot cant/toma overlapping

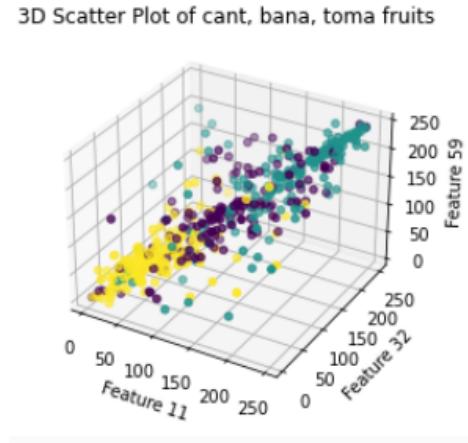


Fig. 52: 3d scatter plot for fruits

```
import pandas as pd
import cv2
import matplotlib.pyplot as plt
import numpy as np
from mpl_toolkits.mplot3d import Axes3D
ax = plt.axes(projection='3d')

II = pd.read_csv("C:/Users/Manish/Desktop/bigdata_assignment1/Image012_random.csv", header=None)
II = np.array(II)
NN = 500
ax.scatter(3D[II[1:NN,11], II[1:NN,32], II[1:NN,59], c=II[1:NN,81])
plt.show()
```

Fig. 51: Code for 3d scatter plot

it seems the points are overlapping and it becomes difficult to classify them. But from some other angle it can be seen more clearly and maybe points are not overlapping and a plane can pass through them which can classify them.

9.3 Discuss these figures and describe your observations in terms of their separable features

From 2D scatter plot we can see that the bananas are skewed towards right, cantaloupes are spread at the center and tomatoes are skewed towards left. This can also be verify from the histogram. So it can become easy to classify based on their skewness. From 3D scatter plot we can see that points are overlapping over each other from the given angle. But maybe from other angle it can happen that the points are not overlapping over each other but are at different distances from each other, then it will become easy to classify them.

X. TASK 10: READ MULTIPLE FILES AND GENERATE FEATURE SPACE!

For this task I have written a for loop which iterates through the folder files and apply feature space creation code to it and outputs feature space for all the images in individual csv files. The code can be seen in Figure 53 and the screenshot of output is after that.

XI. TASK 11: EFFECTS OF BLOCK SIZE ON FEATURE SPACE DIMENSIONS AND VECTOR SPACE

The block size that we are selecting in this assignment is 9*9. This results in creating a feature space of 81 dimensions and hence 81 features. If we reduce the block size to 8*8 then the number of vectors/observations increases but the dimension decreases. If we increase the block size to 10*10,

```

import cv2
import json
import numpy as np
import pandas as pd

std_height = 261

def new_width(heightG,widthG):
    aspect_ratio=heightG/widthG
    new_width = aspect_ratio*std_height
    return round(new_width)

def divide(height_fruit,width_fruit):
    new_fruit_width = height_fruit * width_fruit / (height_fruit + width_fruit)
    if rem == 0:
        new_fruit_width = int(new_fruit_width)
    elif rem > 0:
        new_fruit_width = int(new_fruit_width) + 1
    else:
        new_fruit_width = int(new_fruit_width) - 1
    new_fruit_width = new_fruit_width + (rem)
    return new_fruit_width

def featurespace(img,filename):
    convertToShape = cv2.cvtColor(img, cv2.COLOR_BGR2GRAY)
    heightsG, widthsG = convertToShape.shape
    new_width = new_width(heightsG, widthsG)
    resized_fruit = cv2.resize(convertToShape, dsize=(new_fruit_width, std_height), interpolation=cv2.INTER_CUBIC)
    height, width = resized_fruit.shape
    for i in range(0, height):
        for j in range(0, width):
            flatc[i][j] = crop_top_left[i][j]*255
    flatc[i][j] = crop_top_left[i][j].flatten()
    k = 0
    for i in range(0, heightsG):
        for j in range(0, widthsG):
            crop_top_left[i][j] = flatc[k][j]
            k = k + 1
    fspaced = pd.DataFrame(flatc)
    #spaced object
    #fspaced.to_csv(filename, index=False)

location = '/Users/SEAN/Downloads/Sem 2 (Spring 22)/CSCI610 HW/assignment/assignment 1/datasets/assignment 1/FID038/apricots'
images = []
for file in os.listdir(location):
    img = cv2.imread(os.path.join(location,file))
    if img != None:
        images.append(file)
        img_name = file
        img_name = img_name.replace('.','')
        img_name = img_name + '_resized'
        img_name = img_name + '.png'
        split_file_name = img_name.split('_')
        name_file = split_file_name[0] + '_resized' + split_file_name[1]
        feature_space(img_name,rename_csvfile)
        fspaced.to_csv(name_file+'.csv', index=False)

```

Fig. 53: Code to read multiple files and generate feature space

then the number of feature increases which can lead to high dimensionality issues.

XII. ASSIGNMENT 2: FEATURE SPACE TO A CLASSIFIER

XIII. TASK 12:COMPLETE AND EXTEND THE TASKS OF ASSIGNMENT 1

12.1 Divide the data domain of the datasets into 78:22 and save corresponding files

For this task, I have divided the overlapping and non overlapping files for two class and three class generated in assignment1 into 78% of training and 22% of testing data using the shape function of python. The code for the task is shown in Figure 55.



Fig. 54: Creation of multiple feature spaces

```

import pandas as pd
import numpy as np

# Read a feature space
input_data = pd.read_csv("C:/Users/Manish/Desktop/data/overlapping/merge/random/image012_overlap_random.csv", header=0)

X = input_data
row, col = X.shape

#Split the data into 78% - Training Dataset
TR = round(col*0.78)
X1 = np.array(X)
TT = row-TR

#remaining is the Testing Dataset
X_train = X1[0:TR,:]
X_test = X1[TR:row,:]

#SaveTraining dataset
data_train = pd.DataFrame(X_train)
data_train.to_csv('C:/Users/Manish/Desktop/splitted data/train/image012_train_overlap.csv', index=False)

#Save testing Dataset
data_test = pd.DataFrame(X_test)
data_test.to_csv('C:/Users/Manish/Desktop/splitted data/test/image012_test_overlap.csv', index=False)

```

Fig. 55: Code to split data into 78:22

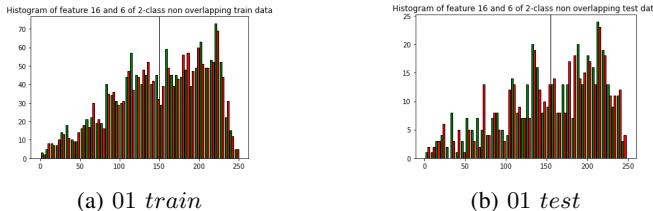


Fig. 56: Histogram of feature 16 & 6 of 2-class non overlap train/test data

12.2 Select two features in each category of training and testing datasets and plot histograms to see if they follow same distribution.

I have selected two features: feature16 and feature6 and plot histograms for training and test dataset generated in the previous steps. From the histograms in Figure 56 to Figure 59, it is clear that both train and test data follow the same distribution for overlapping feature space. For non overlapping data, train data don't have much difference among features 6 and 16 but for test data there happens to be a considerable difference between the features because of small data size. The information about the mean can also be seen in the histograms shown by the solid dark line.

The mean for 2 class non overlapping train data is around 150.53, but the mean of 2 class non overlapping test data is around 154.84. Similarly for the 3 class, non overlapping train data is around 122.93, but the mean of 2 class non overlapping test data is around 123.72. So the test data is bit skew to the left side than the train data.

12.3 Use the same two features to generate scatter plots for each category of training and testing datasets.

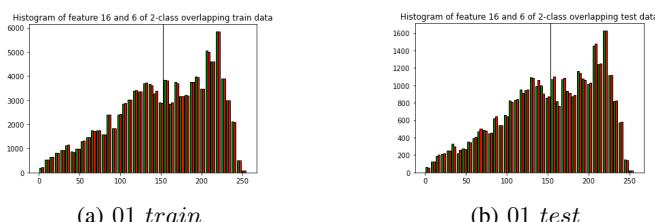
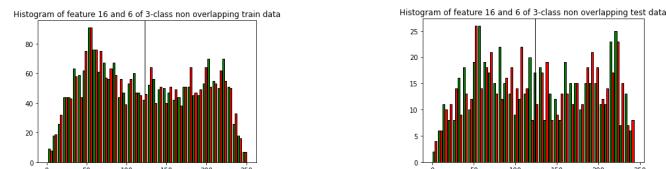


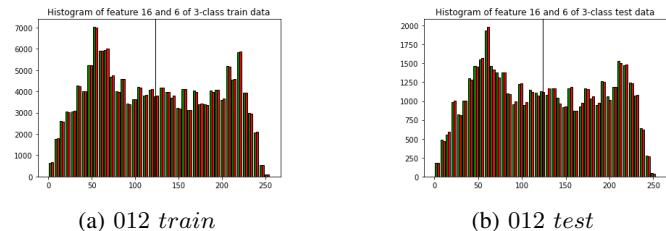
Fig. 57: Histogram of feature 16 & 6 of 2-class overlap train/test data



(a) 012 train

(b) 012 test

Fig. 58: Histogram of feature 16 & 6 of 3-class non overlap train/test data



(a) 012 train

(b) 012 test

Fig. 59: Histogram of feature 16 & 6 of 3-class overlapping train/test data

The same features: feature16 and feature6 is selected for scatter plots. Figure 60 shows the code and Figure 61 to Figure 68 shows the scatter plots for the same.

XIV. TASK 13: IMPLEMENTING A REGRESSION-BASED MODEL

13.1 Implement and train lasso regression or elastic-net regression as a two-class classifier using the training sets of the overlapping and non-overlapping feature vectors.

I have implemented lasso regression model for this task. I have calculated the parameter 'a' and the model 'y' of the lasso regression using the formula provided in the textbook. I have used lambda as 0.1 for calculation. I trained the model with the training data files for 2 class classification. Then predict the test data on the model. The code is shown in the Figure 69.

13.2 Add the predicted labels next to their actual labels in the corresponding spreadsheets.

In this step, the predicted values are merged with the test data actual label values as 83rd column in a new csv file. Figure 70 shows the added column.

```

import seaborn as sns
import matplotlib.pyplot as plt
import pandas as pd
df = pd.read_csv('C:/Users/Manish/Desktop/splitted data/train/image012_train_overlap.csv')
df = df.head(2000)
f6 = df['6']
f16 = df['16']
label = df['81']

fig = plt.figure(figsize = (20 , 20))
ax1 = fig.add_subplot(221)
ax1.set_title("image012_train dataset overlapping")
ax1.set_xlabel("Feature 6")
ax1.set_ylabel("Feature 16")
sns.scatterplot(x = f6 , y= f16 , hue= label)
plt.show()

```

Fig. 60: Code for scatter plot for train/test data

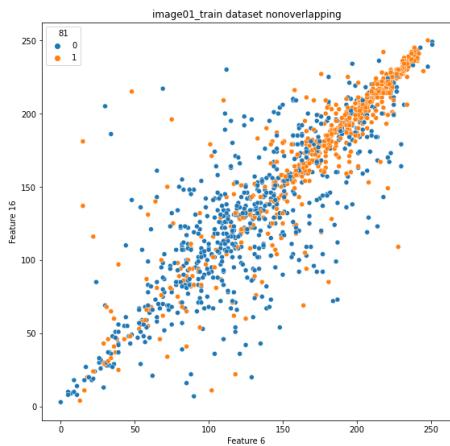


Fig. 61: Scatter plot 2-class train data non overlap

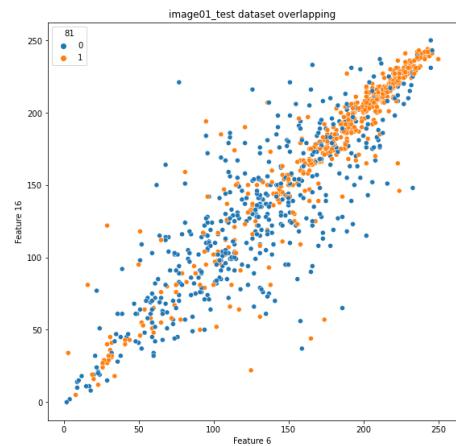


Fig. 64: Scatter plot 2-class test data overlap

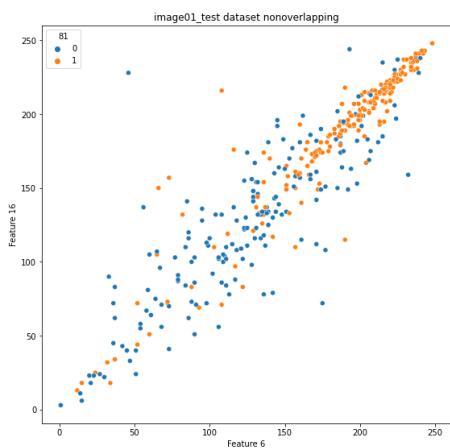


Fig. 62: Scatter plot 2-class test data non overlap

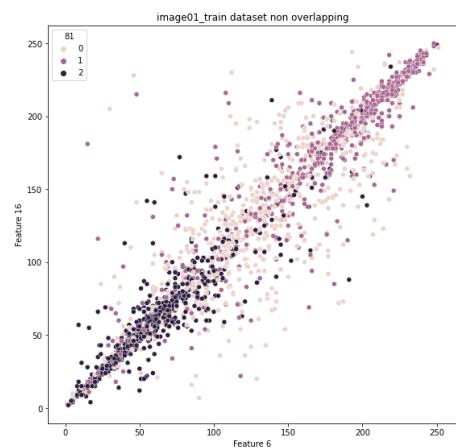


Fig. 65: Scatter plot 3-class train data non overlap

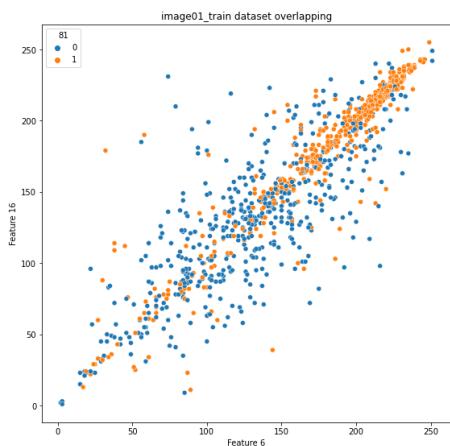


Fig. 63: Scatter plot 2-class train data overlap

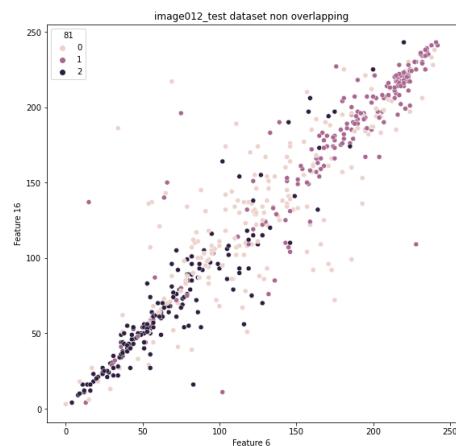


Fig. 66: Scatter plot 3-class test data non overlap

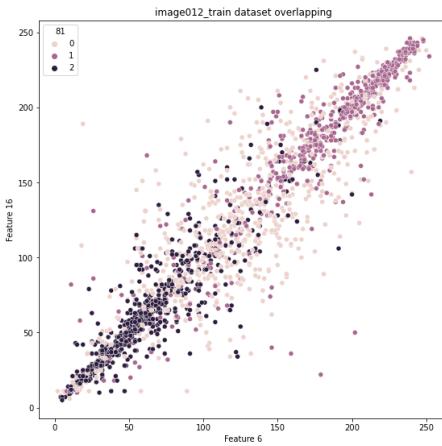


Fig. 67: Scatter plot 3-class train data overlap

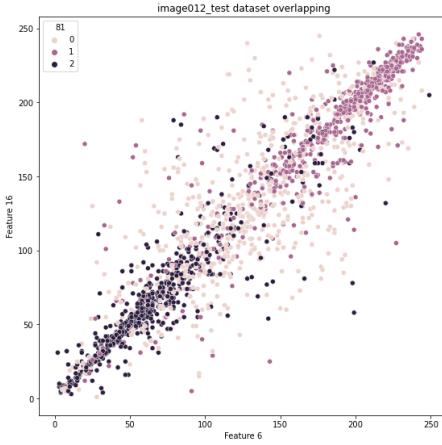


Fig. 68: Scatter plot 3-class test data overlap

```
import pandas as pd
import numpy as np
from numpy.linalg import inv
from sklearn.metrics import confusion_matrix

#read the training and the testing dataset
input_data_train = pd.read_csv('C:/Users/Hanish/Desktop/splitted data/train/image01_train_overlap.csv', header=None)
input_data_test = pd.read_csv('C:/Users/Hanish/Desktop/splitted data/test/image01_test_overlap.csv', header=None)

#extract the labels column
y = input_data_train[81]
Y = np.array(y)

#extract features data
input_data_train.drop(81, axis=1, inplace=True)
X = input_data_train
X1 = np.array(X)

#set lambda as 0.1
lambda_ = 0.1

# Calculate X transpose
X2 = X1.T

#Perform xx
XX = np.matmul(X2, X1)

#Calculating inverse
IX = Inv(XX)

#Calculate yx'
yx = np.matmul(Y,X1)
```

Fig. 69: Code for Lasso Regression

81	64	85	66	67	66	69	70	71	72	75	74	76	76	77	78	79	80	81	82
108	182	156	162	184	149	188	189	182	133	188	184	186	187	169	162	176	158	1	1
143	143	137	163	127	123	120	120	120	120	120	120	120	120	120	120	120	120	120	1
12	17	30	36	40	48	38	67	71	15	16	26	54	41	51	55	64	66	1	1
136	200	212	208	202	198	187	213	172	214	162	202	226	194	226	211	217	193	8	1
225	220	220	220	220	220	220	220	220	220	220	220	220	220	220	220	220	220	220	1
117	144	142	139	136	135	139	139	135	142	137	135	146	121	133	139	128	135	8	1
121	124	124	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	1
147	134	66	137	70	105	145	150	92	146	145	146	136	72	118	117	148	172	8	1
129	197	159	159	157	157	159	158	162	172	154	154	154	156	155	152	155	178	1	1
110	110	110	110	110	110	110	110	110	110	110	110	110	110	110	110	110	110	110	1
215	216	222	223	225	223	213	231	234	213	214	214	226	222	223	234	231	235	235	1
229	229	123	148	149	149	149	149	149	149	149	149	149	149	149	149	149	149	149	1
111	109	92	66	81	103	108	116	99	112	97	104	110	107	114	84	103	81	8	1
207	220	187	218	211	222	201	210	221	221	228	206	208	225	224	226	219	217	212	8
213	209	214	214	214	214	214	214	214	214	214	214	214	214	214	214	214	214	214	1
153	143	143	143	143	143	143	143	143	143	143	143	143	143	143	143	143	143	143	1
36	56	56	56	56	56	56	56	56	56	56	56	56	56	56	56	56	56	56	1
143	140	108	91	104	147	156	106	122	169	141	79	145	172	172	164	146	79	8	1
175	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	1
125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	1
243	239	236	217	201	212	218	244	238	239	245	245	251	258	264	214	252	217	8	1
225	225	160	160	160	160	160	160	160	160	160	160	160	160	160	160	160	160	160	1
209	206	204	204	203	203	199	199	199	199	199	199	199	199	199	199	199	199	199	1
210	199	199	199	199	199	199	199	199	199	199	199	199	199	199	199	199	199	199	1
83	131	131	131	131	130	68	80	83	100	97	108	121	108	144	106	85	113	8	1
199	130	130	130	166	214	133	185	120	113	188	146	130	152	143	186	161	189	8	1
124	124	124	124	124	124	124	124	124	124	124	124	124	124	124	124	124	124	124	1
218	222	228	232	229	227	154	209	228	221	221	227	229	228	226	226	226	226	226	1
106	80	96	90	77	60	28	29	29	77	93	85	82	71	66	24	29	28	8	1

Fig. 70: Added predicted label next to actual label in test data

```
#confusion matrix building
CC_test = confusion_matrix(test_y, y2)

TN = CC_test[0,0]
FP = CC_test[0,1]
FN = CC_test[1,0]
TP = CC_test[1,1]
PPFN = FP+FN
TPTN = TP+TN

Accuracy = 1/(1+(PPFN/TPTN))
print("Our_Accuracy_Score:",Accuracy)

Precision = 1/(1+(FP/TP))
print("Our_Precision_Score:",Precision)

Sensitivity = 1/(1+(FN/TP))
print("Our_Sensitivity_Score:",Sensitivity)

Specificity = 1/(1+(FP/TN))
print("Our_Specificity_Score:",Specificity)
```

Fig. 71: Confusion matrix creation

13.2 Construct confusion matrices using the actual and predicted labels in the 82nd and 83rd columns of the testing datasets.

The code for constructing the confusion matrix using the actual and predicted labels can be seen in the Figure 71. The result of confusion matrix can be seen in Figure 72 - Figure 73.

13.3 Select two qualitative measures that use the concept of confusion matrix to quantify the performance of a machine learning model.

The qualitative measures accuracy, precision, specificity,

```
In [92]: CC_test
Out[92]:
array([[ 9140,  4984],
       [ 2744, 11296]], dtype=int64)
```

Fig. 72: Output of confusion matrix for overlapping data

```
In [100]: CC_test
Out[100]:
array([[116,  64],
       [ 56, 134]], dtype=int64)
```

Fig. 73: Output of confusion matrix for non overlapping data

```
.... print( Our_Specificity_Score: ,spec
Our_Accuracy_Score: 0.7256071580741372
Our_Precision_Score: 0.6938574938574938
Our_Sensitivity_Score: 0.8045584045584045
Our_Specificity_Score: 0.6471254602095724
```

Fig. 74: Output of confusion matrix for image01 overlapping dataset

```

import pandas as pd
import numpy as np
from sklearn.metrics import confusion_matrix,plot_confusion_matrix
from sklearn.metrics import accuracy_score,precision_score,recall_score
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt

input_data_train = pd.read_csv('C:/Users/Manish/Desktop/splitted data/train/image01_train.csv', header = None)
input_data_test = pd.read_csv('C:/Users/Manish/Desktop/splitted data/test/image01_test.csv', header = None)

splitting the labels
y = input_data_train[81]
Y = np.array(y)
splitting the features
input_data_train.drop(81,axis=1,inplace=True)
input_data_train
X1 = np.array(X)

# Train the model
rf = RandomForestClassifier(random_state=0,n_estimators=100,oob_score=True, n_jobs=1)
rf = rf.fit(X1, Y)

saving the test dataset to peek
peek = pd.DataFrame(input_data_test)

#Testing the model using Trained results
test_y = input_data_test[81]
Y_test = np.array(test_y)

splitting the features from the test dataset to test the model
input_data_test.drop(81,axis=1,inplace=True)
input_data_test = np.array(input_data_test)
y_pred = rf.predict(input_data_test)
yhat_saved = pd.DataFrame(y_pred)

peek['82'] = yhat_saved
peek.to_csv('C:/Users/Manish/Desktop/RF_Test_result_image01_nonooverlap.csv', index=False)

```

Fig. 75: Random forest implementation for overlapping and non overlapping dataset

sensitivity are there for confusion matrix. I am considering accuracy and precision for this task. Figure 74 shows the qualitative measures result.

From the figure it can be seen that image01 overlapping data have accuracy of 72.5% and precision of 69.38%. While for non overlapping data, accuracy is 67.5% and precision is 67.6%. So, in our case overlapping data performed better than non overlapping data because more data points are available for it to train.

XV. TASK 14:IMPLEMENTING RANDOM FOREST OR DEEP LEARNING

14.1 Implement the random forest or the sequential (simple deep learning) technique as two-class and three-class classifiers using the training sets of the overlapping and non-overlapping feature vectors that you generated.

I have implemented random forest classifier for this task. It is built as a 2-class and a 3-class classifier. For this, I have first split the training dataset into features and labels and gave it as an input to the classifier to train. After training, I split the test set features and gave it to the model to predict the labels. The code for Random forest can be seen in Figure 75.

14.2 Apply the trained model to the test sets of all the categories of datasets and add the predicted labels next to their actual labels in the corresponding spreadsheets.

After predicting the labels using the step discussed above, the labels are stored into a new csv file by converting into a data frame and saving the data frame as csv file. The output

63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81
23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41
177	19	149	179	157	95	256	145	257	189	178	165	144	247	158	278	145	5	
113	125	137	149	161	173	185	197	209	221	233	245	257	269	281	293	305	317	
150	158	150	139	171	142	68	46	60	56	154	154	135	143	143	151	52	43	9
9	12	15	16	15	13	33	44	44	10	16	11	13	18	28	29	46	41	
206	199	175	166	166	160	182	184	176	198	202	214	187	278	171	182	164	185	9
217	200	200	207	218	218	220	222	218	251	253	228	225	221	218	258	219	251	1
213	221	231	234	226	229	230	236	236	209	209	216	225	229	229	225	237	237	1
129	139	139	139	139	139	139	139	139	139	139	139	139	139	139	139	139	139	
159	130	189	186	217	127	113	159	164	113	215	205	182	212	193	227	184	137	9
106	129	143	135	97	100	126	116	96	121	124	116	117	112	83	115	148	106	
132	121	121	121	121	121	121	121	121	121	121	121	121	121	121	121	121	121	
54	59	63	137	111	70	59	98	81	55	68	71	137	134	116	112	134	95	
72	196	199	200	200	202	202	202	202	75	98	96	201	205	207	205	85	64	1
156	158	160	94	71	65	54	57	51	132	156	125	125	125	125	69	81	57	54
197	197	200	200	201	204	157	159	159	201	201	201	201	201	201	201	201	201	
82	115	121	153	171	166	150	145	150	84	135	152	149	179	123	140	148	133	
54	65	92	88	66	57	53	49	50	56	64	94	66	62	74	54	54	55	52
179	196	196	196	196	196	196	196	196	196	196	196	196	196	196	196	196	196	
159	136	129	129	129	129	129	129	129	111	86	166	148	116	122	113	119	181	91

Fig. 76: Predicted labels added in 83rd column for overlapping dataset

```

#creating confusion matrix
CC_test = confusion_matrix(test_y, y_pred)
TN = CC_test[1,1]
FP = CC_test[1,0]
FN = CC_test[0,1]
TP = CC_test[0,0]

FPFN = FP+FN
TPTN = TP+TN

print(CC_test)
plot_confusion_matrix(rF, input_data_test_np, test_y)
plt.show()

#Qualitative measures of performance calculation using Confusion matrix
Accuracy = 1/(1+(FPFN/TPTN))
print("Test_Accuracy_Score:",Accuracy)
Precision = 1/(1+(FP/TP))
print("Test_Precision_Score:",Precision)
Sensitivity = 1/(1+(FN/TP))
print("Test_Sensitivity_Score:",Sensitivity)
Specificity = 1/(1+(FP/TN))
print("Test_Specificity_Score:",Specificity)

```

Fig. 77: Confusion matrix for Random Forest

of the csv file is as shown in the Figure 76. The output of the random forest classifier is saved in a pdf file.

14.3 Construct confusion matrices using the actual and predicted label in the 82nd and 83rd columns in the testing datasets for all the four categories of datasets.

Creation of confusion matrix is as shown in Figure 77. The output of the confusion matrix for 2class and 3class of overlapping and non overlapping data is shown from Figure 78 to Figure 81. From the output of the confusion matrix it is clear that the random forest is predicting good results because the diagonal elements have higher values compared to other elements of the matrix.

Compared to non overlapping, overlapping dataset gives better results. It is because the model has more points to train before testing.

14.4 Adapt the same two measures used in Task 2 to quantify the performance of the models.

After predicting the labels and giving it as input to confusion matrix, I have selected the accuracy and precision as my qualitative measures to quantify the performance of model.

From Figure 82 and Figure 83 it is clear that 3-class have better accuracy and precision as compared to 2-class. For 2-class non overlapping, accuracy is 89% and precision is 89.26%. For 3-class non overlapping, accuracy is 90.71% and precision is 91.55%.

For 2-class overlapping, accuracy is 98.75% and precision is 98.32%. For 3-class overlapping, accuracy is 98.40% and

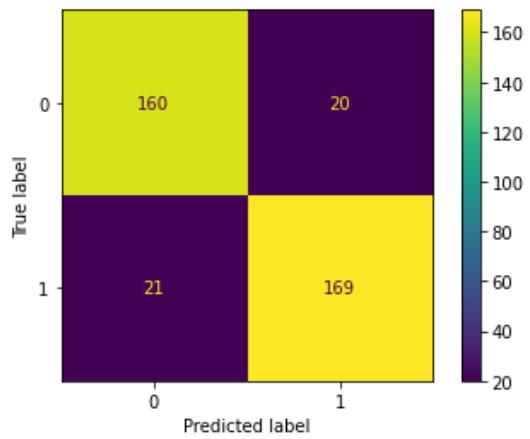


Fig. 78: Confusion matrix for 2class non overlap

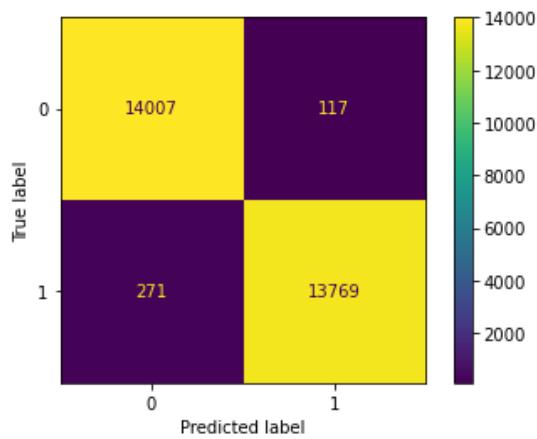


Fig. 79: Confusion matrix for 2class overlap

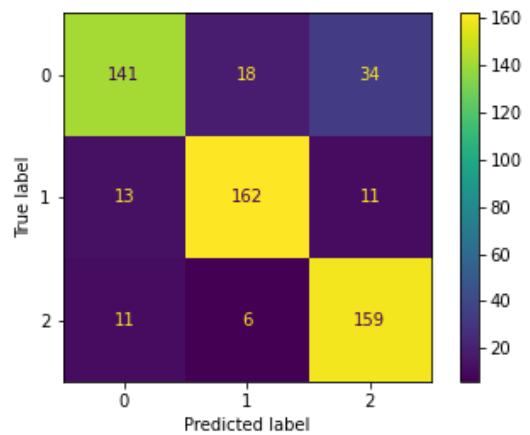


Fig. 80: Confusion matrix for 3class non overlap

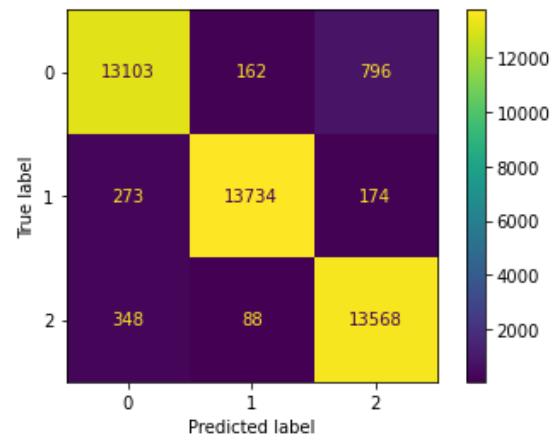


Fig. 81: Confusion matrix for 3class overlap

```
Test_Accuracy_Score: 0.8891891891891891
Test_Precision_Score: 0.8926553672316384
Test_Sensitivity_Score: 0.8777777777777777
Test_Specificity_Score: 0.8999999999999999
```

Fig. 82: Random Forest Measure of 2 class non overlapping dataset

precision is 98%. So it seems clear that overlapping data performs better with testing as compared to non overlapping data.

XVI. TASK 15:EVALUATION OF THE LEARNING MODELS

15.1 Find some built-in measures that are available in software libraries like metrics.accuracy_score in Python's sklearn environment. Use such two measures to compare the performance of the models.

For this task, I have used sklearn.metrics for Accuracy, Precision, Sensitivity as accuracy_score, precision_score. These are built in measures for calculation. Figure 86 shows inbuilt measure performance for lasso regression for both overlapping and non overlapping dataset and Figure 87 - Figure 87 shows Random Forest qualitative measures.

```
.... print( test_specificity_score, s)
Test_Accuracy_Score: 0.907185628742515
Test_Precision_Score: 0.9155844155844156
Test_Sensitivity_Score: 0.8867924528301887
Test_Specificity_Score: 0.9257142857142858
```

Fig. 83: Random Forest Measure of 3 class non overlapping dataset

```
Test_Accuracy_Score: 0.9875372816361312
Test_Precision_Score: 0.9832292470703811
Test_Sensitivity_Score: 0.9920702350608893
Test_Specificity_Score: 0.9829772079772079
```

Fig. 84: Random Forest Measure of 2 class overlapping dataset

```
Test_Accuracy_Score: 0.984049574655324
Test_Precision_Score: 0.9795903110047848
Test_Sensitivity_Score: 0.9877874104787033
Test_Specificity_Score: 0.9805097451274363
```

Fig. 85: Random Forest Measure of 3 class overlapping dataset

```
Overlapping measures Lasso Regression
BuiltIn_Accuracy: 0.7341996875443829
BuiltIn_Precision: 0.7224107506447672
BuiltIn_Sensitivity (recall): 0.7581196581196581

In [128]: print("Nonoverlapping measures Lasso Regression")
...: from sklearn import metrics
...: print("BuiltIn_Accuracy:",metrics.accuracy_score)
...: print("BuiltIn_Precision:",metrics.precision_score)
...: print("BuiltIn_Sensitivity (recall):",metrics.recall_score)

Nonoverlapping measures Lasso Regression
BuiltIn_Accuracy: 0.6756756756756757
BuiltIn_Precision: 0.6767676767676768
BuiltIn_Sensitivity (recall): 0.7052631578947368
```

Fig. 86: Lasso Regression builtin measures for overlapping and non overlapping dataset

- 15.2 Describe the quantitative differences between the built-in measures and the confusion-matrix-based measures that you used in the analysis to compare the models with the four categories of datasets.

For this task, I have created a table as shown in Figure 89 with values for built-in measures and measures derived from confusion matrix. From the table it is clear that there is not much difference between the values. While confusion matrix measures are slightly higher than inbuilt measures in case of 3-class classification than the built-in measures.

- 15.3 Compare all the results (qualitative measures) and de-

```
Random Forest 2 class overlap output
sklearn.metrics Accuracy 0.9875372816361312
sklearn.metrics precision 0.9919499748436714
sklearn.metrics sensitivity 0.982977207977208
Random Forest 2 class non overlap output
sklearn.metrics Accuracy 0.8891891891892
sklearn.metrics precision 0.8860103626943006
sklearn.metrics sensitivity 0.9
```

Fig. 87: Random Forest builtin measures for 2-class overlapping and non overlapping data

```
Random Forest 3 class non overlap output
sklearn.metrics Accuracy 0.841726618705036
sklearn.metrics precision 0.882520205490658
sklearn.metrics sensitivity 0.8832902187511079
Random Forest 3 class overlap output
sklearn.metrics Accuracy 0.9587426326129665
sklearn.metrics precision 0.9692900112056182
sklearn.metrics sensitivity 0.9690416080012648
```

Fig. 88: Random Forest builtin measures for 3-class overlapping and non overlapping data

	RANDOM FOREST		
	Buit-in measures	Confusion matrix	
	Accuracy	Precision	Accuracy
012 overlap	95.68	96.78	98.48
012 nonoverlap	84.53	88.52	90.53
01 overlap	98.62	98.62	98.62
01 non overlap	88.91	88.93	88.91

	Logistic Regression			
	Buit-in measures		Confusion matrix	
	Accuracy	Precision	Accuracy	Precision
012 overlap	45.88	23.13	68.63	63.29
012 nonoverlap	48.2	24.46	70.71	65.43
01 overlap	73.41	72.24	73.41	72.24
01 non overlap	67.56	67.67	67.56	67.67

Fig. 89: Output of builtin measures and confusion matrix based measures for RF and LR

termine which pair (a model and a category of feature vectors) is superior among the ones that you considered. Use your data and findings to support this result.

In my opinion Random Forest performs better than Logistic Regression for all category of feature vectors. Random Forest performed better in overlapping feature space then non overlapping feature space. For 2-class overlapping dataset, it gives 98.62% accuracy while for 3-class overlapping dataset, it gives 98.48 % of accuracy. So random forest for overlapping dataset is the model-category pair which is superior among all. The table in Figure 89 is evident of that.

XVII. ASSIGNMENT 3: CLASSIFIERS ON A CONCEPTUALIZED BIG DATA SYSTEM

XVIII. TASK 16: MAKING SURE ALL THE TASKS OF ASSIGNMENT 1 AND ASSIGNMENT 2 ARE COMPLETED

- 16.1 Complete all the tasks stated in assignment 1 towards developing a set of suitable feature spaces

I have completed all the tasks assigned in assignment 1. I have generated feature spaces for overlapping and non overlapping feature vector for 2-class and 3-class classification.

- 16.2 Complete all the tasks stated in assignment 2 that perform the application of Random Forest and Deep Learning algorithm using the feature spaces

For this task, I have make use of Random Forest and Deep learning method on the feature spaces generated in assignment 1. I have used these algorithms on both overlapping and non overlapping feature spaces and calculate the time and accuracy, suggest which works best on 3-class and 2-class classification for both overlapping and non overlapping blocks.

- 16.3 Running dataset on conceptualized big data system and a local system, find the differences

For this part, I am planning to generate the feature spaces that are large enough to execute on local machine but can be managed on big data system. I am planning to generate a feature space with around 1 million data points. After executing, I will be calculating the time and accuracy of the algorithm implemented in both the systems.

XIX. TASK 17: UNDERSTANDING AND THE 10-FOLD VALIDATION TECHNIQUE PRESENTED IN CHAPTER 8 OF THE TEXTBOOK

- 17.1 Read chapter 8 and explain the concept of n-fold validation in your own words. You must clearly explain the effect of different values of n on seen-data and unseen-data problems.

n-fold validation technique is a block based technique, where the dataset is divided into n equal parts. Each part is disjoint from each other. Let's try with n = 10; which is 10 fold cross validation technique.

In 10 fold cross validation technique, initially the model is fed with all the blocks except the 10th block. The 10th block is used for testing. In the next step, all the blocks are used for training except the 9th block. It is used for testing the model.

As we can see, the 9th block was seen in the first step, but the trained model in the second step didn't see 9th block. As a result the testing is done on both seen and unseen data. This is known as cross validation. This process is repeated for all other blocks and accuracy is calculated, take the average of it and we get the final accuracy. Based on accuracy, one final model is selected or retraining is done.

With different values of n, the accuracy and error in each step will change. For large value of n, we will get better accuracy because our data will get more points to train. But optimal data size for training is also to be determined to avoid over fitting problem. Along with the value of number of blocks, the size of blocks also matters. If we divide a small data into many blocks, each block contains very few data points. If we train and validate it on such blocks, it won't be accurate. So size of block also matters. We need to define an appropriate values of the size of blocks.

In our case, we will be dividing the data into unequal block sizes and performing training and testing on it.

- 17.2 Now connect your data set to the 10-fold validation concept that is explained in Figure 8.6 by identifying the number of observations in each fold and the statistical differences between each fold.

I have generated my dataset using 2 types of fruit images: apples and bananas. I have selected 5 images of each fruit type, so in total 10 images. I then divided my data into 10 blocks and apply the 10 fold validation technique that we discussed above on local machine and on big data system.

For dividing the data, I have plotted the histogram of a feature to randomly divide the data into blocks. The detailed information is explained in further steps.

- 17.3 Make sure to draw your conceptualized big data system and write down the data characteristics of a subdomain that is delivered to a node in the conceptualized cluster.

Our big data system contains 9 nodes. Each node is feed with training and test data. Training data contains 8 blocks while test data contains 1 block. The size of each block is unique.

Each block contains data points within a certain range. This is the data characteristic of each block that is feed. The

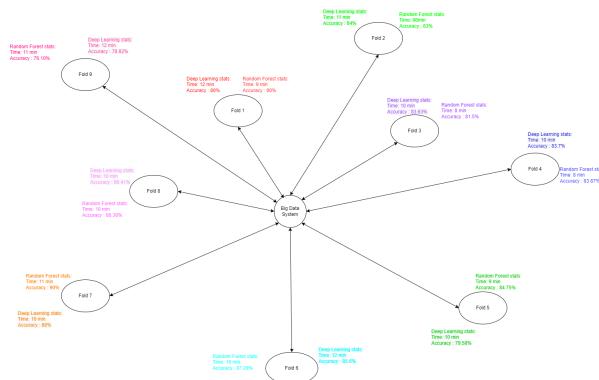


Fig. 90: Conceptualized Big Data System

first block contains data points within the range of (0,25.5), second block contains in range of (25.5,51). The third block contains in the range of (51,76.5), fourth block in the range of (76.5,102), fifth block in the range of (102,127.5), sixth block in the range of (127.5,153). The seventh block in contains in the range of (153,178.5), eighth block in the range of (178.5,204), ninth block within the range of (204,255).

The diagram of our conceptualized big data system can be seen in Figure 90. Our big data system has one main system and 9 sub nodes. These sub nodes are basically the machines on which we will run each folds simultaneously. Figure 90 shows the big data system when run for Random Forest and Deep Learning model. The accuracy and time taken for each model is mentioned on the nodes. The final time taken and accuracy using each technique is calculated by taking mean of time and mean of accuracies. It is an example of Map-Reduce where data is divided among the nine clusters so processing of big data becomes efficient. The key value pair for map reduce is mentioned as: key is the mean of the range of each bin and the value is the data points contained in that bin of histogram.

The key value pair for a block is defined as: the key contains the mean of range of bin. For example, out first key will be $(0+25.5)/2 = 12.75$. For second block it will be $(25.5 + 51)/2 = 38.25$. Similarly is the method for all other blocks/splits. For the value contained within the blocks, it contains all the rows that fall within that range. For example, block_1 contains the points whose feature23 values lie within the range of 0 to 25.5. Such is the key-value pair for all other blocks/splits.

- 17.3 Then create a data frame (or a feature space) for each node and save it as a spreadsheet (.csv file). Do the same for every node and save the spreadsheets for performing Task 3.

For this task, I have merged the data frames to create a big csv files with block-1 to block-8 in it and block-9 is in the testing. I did the same for all other blocks, keeping one block in the testing and combining all the remaining blocks. Then I fed the nodes with that train and test data.

First node is fed with data that contains block-9 as testing and remaining blocks in the training. Second node is fed with

```

min_point = min(input_data['23'])
max_point = max(input_data['23'])
print(f'The min value of feature 23 is {min_point} and max value of feature 23 is {max_point}')

plt.hist(f, bins=10, color='white')
plt.title('Total data')
plt.xlabel('Feature 23')
plt.ylabel('frequency')
plt.show()

first_block=pd.DataFrame()
first_block = input_data[f<=25.5]
first_block.to_csv('C:/Users/Munish/Desktop/csvs/10 fold/first_block.csv', index= False)

second_block=pd.DataFrame()
second_block = input_data[f.between(25.5,51)]
second_block.to_csv('C:/Users/Munish/Desktop/csvs/10 fold/second_block.csv', index= False)

third_block=pd.DataFrame()
third_block = input_data[f.between(51,76.5)]
third_block.to_csv('C:/Users/Munish/Desktop/csvs/10 fold/third_block.csv', index= False)

fourth_block=pd.DataFrame()
fourth_block = input_data[f.between(76.5,102)]
fourth_block.to_csv('C:/Users/Munish/Desktop/csvs/10 fold/fourth_block.csv', index= False)

fifth_block=pd.DataFrame()
fifth_block = input_data[f.between(102,127.5)]
fifth_block.to_csv('C:/Users/Munish/Desktop/csvs/10 fold/fifth_block.csv', index= False)

sixth_block=pd.DataFrame()
sixth_block = input_data[f.between(127.5,153)]
sixth_block.to_csv('C:/Users/Munish/Desktop/csvs/10 fold/sixth_block.csv', index= False)

seventh_block=pd.DataFrame()
seventh_block = input_data[f.between(153,178.5)]
seventh_block.to_csv('C:/Users/Munish/Desktop/csvs/10 fold/seventh_block.csv', index= False)

eighth_block=pd.DataFrame()
eighth_block = input_data[f.between(178.5,204)]
eighth_block.to_csv('C:/Users/Munish/Desktop/csvs/10 fold/eighth_block.csv', index= False)

```

Fig. 91: 10 folds data split

data that contains block-8 as testing and remaining blocks in the training. Third node is fed with data that contains block-7 as testing and remaining blocks in the training. Fourth node is fed with data that contains block-6 as testing and remaining blocks in the training.

Fifth node is fed with data that contains block-5 as testing and remaining blocks in the training. Sixth node is fed with data that contains block-4 as testing and remaining blocks in the training. Seventh node is fed with data that contains block-3 as testing and remaining blocks in the training. Eighth node is fed with data that contains block-2 as testing and remaining blocks in the training. Ninth node is fed with data that contains block-1 as testing and remaining blocks in the training. All these nodes are fed with data at the same time and performing parallelization. The accuracy and time duration is calculated. The average of accuracy is done and we get the final accuracy.

There are 9 files for training and nine files for testing. It can be seen in the data folder.

XX. TASK 18: MACHINE LEARNING TASKS

18.1 Prepare a dataset with 2-class classification with non overlapping feature space. Then divide the data to perform 10 fold cross validation.

For this step, I have used 10 images in total: 5 images of apples and 5 images of banana. I have generated the data using overlapping feature vectors and have generated the combined data with more than 700,000 data points.

Further to generate 10 blocks of data, I have selected one feature: Feature 23 and calculated the min() and max() of the feature 23 to get the range of it. After I have got the range I have plotted the histogram of it with bins = 10. By doing this, I get 10 blocks of data. But the tenth block was of size much smaller than other blocks. So combined the block-9 and block-10. The code for this is shown in Figure 91. The histogram can be seen in the Figure 92. From the histogram it is clear that the size of block-10 is much smaller. So I appended block-9 with block-10 and finally get 9 blocks of data with comparative size of blocks.

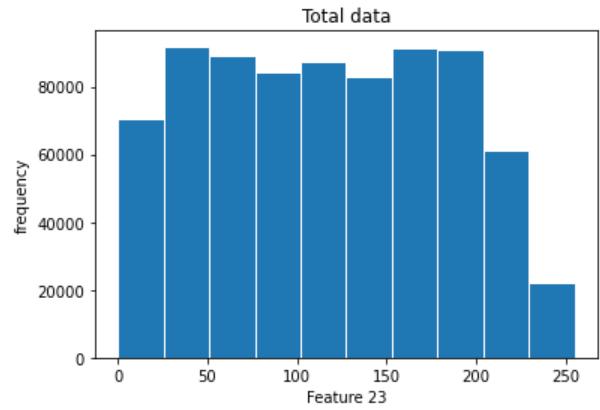


Fig. 92: Histogram of feature 23

18.2 Now apply random forest and deep learning to the entire data frame and calculate the accuracies and measure the computational time. Then do the same to each fold and calculate the time durations.

I have applied Random Forest and Deep Learning to entire dataset. Figure 93 shows the accuracies of random forest and deep learning. From the figure it can be seen that random forest took around 14 minutes to train the 80% of data that has 700,000 points with accuracy of 92.43%. While, deep learning method took 33 minutes to train 80% of the data with accuracy of 86%.

The deep learning method was trained on 10 epochs with batch size of 20. Provided more batch size and epochs, system would have performed better.

For 10 fold data, random forest took 10 minutes to train with average accuracy of 75.48%. The value of n_estimators for this task is 30. This 75% accuracy can be improved by increasing the value of n_estimators.

Deep learning method took 10 minutes with the accuracy of 87% while running on 5 epochs and batch size of 50. Provided more epoch and optimal batch size, accuracy would have increased. The details about accuracy and time can be seen in Figure 94 and Figure 95

XXI. TASK 19: COMPARING THE RESULTS ON A LOCAL SYSTEM VERSUS A BIG DATA SYSTEM

19.1 Compile all the results that you produced from the experiment that implemented the random forest and deep learning techniques on the conceptualized local system using your data set/s.

The results of implementing random forest on local system using merged data can be seen in the Figure 96 and Figure 97. For deep learning results can be seen in Figure 98. From the figure it can be seen that random forest has 94% accuracy while deep learning method has 86% accuracy. The reason for lower accuracy in deep learning could be low number of epochs and batch size and no chance to perform parallelization.

Random Forest (Complete Data)

Time	Accuracy
14 min	92.40%

Deep Learning(Complete Data)

Time	Accuracy
33 min	86%

Fig. 93: Total Data accuracies

Random Forest

Train	Test	Time	Accuracy
except 9	9 fold	9 min	80%
except 8	8 fold	8 min	83%
except 7	7 fold	8 min	81.50%
except 6	6 fold	8 min	83.87%
except 5	5 fold	8 min	84.75%
except 4	4 fold	10 min	87.28%
except 3	3 fold	11 min	90%
except 2	2 fold	10 min	88.36%
except 1	1 fold	11 min	76.10%

Fig. 94: 10 folds data accuracies of Random Forest

Deep Learning

Train	Test	Time	Accuracy
except 9	9 fold	12 min	80%
except 8	8 fold	11 min	84%
except 7	7 fold	10 min	83.63%
except 6	6 fold	10 min	83.70%
except 5	5 fold	10 min	79.58%
except 4	4 fold	12 min	85.60%
except 3	3 fold	10 min	88%
except 2	2 fold	10 min	86.41%
except 1	1 fold	12 min	78.82%

Fig. 95: 10 folds data accuracies of Deep Learning

```
In [5]: print("Our_Specificity:",metrics.specificity_score(y_test,y_hat))
Our_Accuracy_Score: 0.9240307969254877
Our_Precision_Score: 0.9202618013817295
Our_Sensitivity_Score: 0.9526603635719586
Our_Specificity_Score: 0.8836908197488496

In [5]: from sklearn import metrics
...: print("BuiltIn_Accuracy:",metrics.accuracy_score(y_test, yhat_test))
...: print("BuiltIn_Precision:",metrics.precision_score(y_test, yhat_test))
...: print("BuiltIn_Sensitivity (recall):",metrics.recall_score(y_test, yhat_test))
BuiltIn_Accuracy: 0.9240307969254877
BuiltIn_Precision: 0.9298153467377923
BuiltIn_Sensitivity (recall): 0.8836908197488496
```

Fig. 96: Random Forest on merged Data

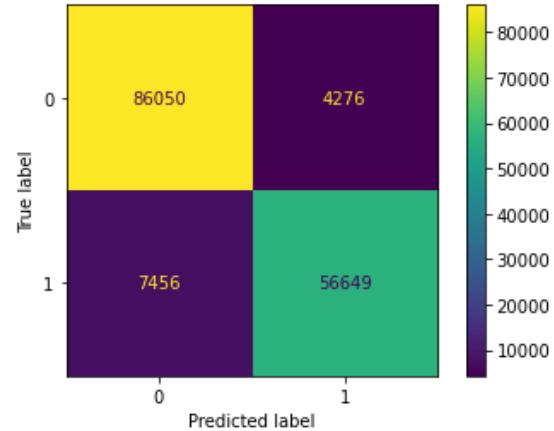


Fig. 97: Confusion matrix of Random Forest on merged data

19.2 Compile all the results that you produced from the experiment that implemented the random forest and deep learning techniques on the conceptualized big system using the 10-fold data sets.

With conceptualized big data system and 10 folds, random forest takes 10 minutes to train and validate all 10 blocks with average accuracy of 75.48%. For this task I have chosen n_estimators=50, This accuracy can be improved by increasing the number of n_estimators to 100 or 200 but it will take more time to train the data.

With deep learning method and parallelization concept, ten folds took 10 minutes with 87% accuracy on average. This accuracy can be further increased by increasing the number of epoch and taking more batch size in each epoch. For this task, I have used 5 epochs with batch size of 50. But it can be increased much more.

I will include the screenshots of each fold accuracy for random forest and deep learning in the screenshot folder. The accuracy for random forest and deep learning on 1st fold is shown in Figure 99 and Figure 100.

19.3 Compare all the results (qualitative measures) and de-

```
Our_Accuracy_Score: 0.8675913514773589
Our_Precision_Score: 0.8847060118916539
Our_Sensitivity_Score: 0.889544538671036
Our_Specificity_Score: 0.836658606972935
BuiltIn_Accuracy: 0.8675913514773588
BuiltIn_Precision: 0.8431560579145116
BuiltIn_Sensitivity (recall): 0.836658606972935
```

Fig. 98: Deep learning on merged data

```
    *** print( classification_report ) ->
Qualitative measures using Confusion matrix
Test_Accuracy_Score: 0.8005358129700377
Test_Precision_Score: 0.8189236356775624
Test_Sensitivity_Score: 0.5763128194914692
Test_Specificity_Score: 0.9277187829492394
Using inbuilt library
sklearn.metrics Accuracy 0.8005358129700377
sklearn.metrics precision 0.7942518174637716
sklearn.metrics sensitivity 0.9277187829492394
```

Fig. 99: Random forest accuracy on split_1

```
Our_Accuracy_Score: 0.7987457651553377
Our_Precision_Score: 0.7411989611888616
Our_Sensitivity_Score: 0.6821018389431056
Our_Specificity_Score: 0.8649081186925742
BuiltIn_Accuracy: 0.7987457651553377
BuiltIn_Precision: 0.8274849587491444
BuiltIn_Sensitivity (recall): 0.8649081186925742
```

Fig. 100: Deep learning accuracy on split_1

termine if the conceptualized big data system was able to run the application that was failed to be completed on the conceptualized local system.

From the tasks 19.1 and 19.2 we can see that Deep Learning works well on a big data system with Ten Folds. As you can see its accuracy increased from 86% in local system to 87% on big data system.

As for random forest, it took 8.3 minutes to run on big data system because each fold is run simultaneously to each node of the big data system. It gives the total accuracy of 75.48%. While on local machine random forest is feed with the full merged data at a time so it took 14 minutes with the accuracy of 92.40%.

So from the results it is clear that big data system helps in running the application that failed to be completed on conceptualized local system.