

Name: Aditi Dadariya

Batch code: LISUM21

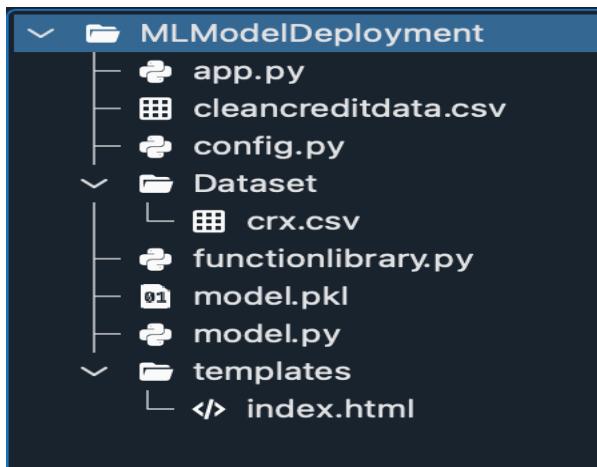
Submission date: May 28, 2023

Submitted to: Data Glacier (Week 4: Deployment on Flask)

Data Description: The dataset downloaded from Kaggle has application for credit card approval [1]. All the attribute names have been changed to meaningless name to maintain the confidentiality of data. It is mentioned that the dataset contains missing values, however as examined during pre-processing, there are no missing values, instead there are special characters found in the dataset. The datatype is different for attributes in the dataset. These are continuous, nominal with small number of values and nominal with large number of values [1].

Steps followed for model development and deployment on flask:

1. Created a new project i.e. ML Model Deployment in Spyder.



2. Created few python files i.e. model.py, app.py, config.py, functionlibrary.py, index.html
3. Created a folder 'Dataset' where the dataset file 'crx.csv' has been placed.
4. Config.py file: This file has all the global variables that has been utilised in functionlibrary.py file and model.py file.

```
#!/usr/bin/env python3
# -*- coding: utf-8 -*-
"""
Created on Tue May 23 00:06:10 2023
@author: aditidadariya
"""

##### VARIABLE DECLARATION #####
import os

# location_of_file variable is defined to store the dataset file location
absolute_path = os.path.dirname(__file__)
relative_path = "Dataset/crx.csv"
location_of_file = os.path.join(absolute_path, relative_path)

# special_char variable is defined to store any special characters
special_char = "?"

# Define models list to store the model object
models = []

# Define names list to store the name of models
names = []

# Define results list to store the accuracy score
results = []

# Define basicscore list to store the name and accuracy score
basic_score = []

# Define basicscore list to store the name and accuracy score
score = []

# Define finalresultslist to store cross validation score
final_results = []

# Defined a randstate variable to store the input for random_state in train_test_split function later
rand_state = [1,3,5,7]
```

5. functionlibrary.py file: This file contains all the functions that has been utilised in model.py file to run the data analysis, visualization and modelling.

```

  × config.py × functionlibrary.py × model.py × index.html × app.py
1  #!/usr/bin/env python3
2  #-*- coding: utf-8 -*-
3  """
4  Created on Tue May 23 00:07:18 2023
5  @author: aditidadariya
6  """
7
8  ##### Importing necessary libraries #####
9
10 import pandas as pd
11 from matplotlib import pyplot
12 import seaborn as sns
13 # Import matplotlib.pyplot to draw and save plots
14 import matplotlib.pyplot as plt
15 # Import LabelEncoder to change the data types of attributes
16 from sklearn.preprocessing import LabelEncoder
17 # Import parameters.py file to get all the variables here
18 #import config.py
19 from config import *
20 # Import Decision Tree Classifier
21 from sklearn.tree import DecisionTreeClassifier
22 # Import LinearDiscriminantAnalysis
23 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
24 # Import train_test_split function
25 from sklearn.model_selection import train_test_split
26 # Import cross_val_score function
27 from sklearn.model_selection import cross_val_score
28 from sklearn.model_selection import cross_val_predict
29 # Import StratifiedKFold function
30 from sklearn.model_selection import StratifiedKFold
31 # Import GridSearchCV function
32 from sklearn.model_selection import GridSearchCV
33 # Import scikit-learn metrics module for accuracy calculation
34 from sklearn import metrics
35 # Import confusion_matrix function
36 from sklearn.metrics import confusion_matrix
37 # Import classification_report function
38 from sklearn.metrics import classification_report
39 # Import ConfusionMatrixDisplay function to plot confusion matrix
40 from sklearn.metrics import ConfusionMatrixDisplay
41 # Import render from graphviz to convert dot file to png
42 from graphviz import render
43 # Import RandomUnderSampler, SMOTE, Pipeline to balance the dataset
44 from imblearn.under_sampling import RandomUnderSampler
45 from imblearn.over_sampling import SMOTE
46 from imblearn.pipeline import Pipeline
47
48 ###### Function declarations #####
49
50
51 # Newline gives a new line
52 def Newline():
53     print("\r\n")
54
55 # ClearLists removes all values from lists
56 def ClearLists():
57     models.clear()
58     names.clear()
59     results.clear()
60     basic_score.clear()
61     score.clear()
62     final_results.clear()
63
64
65
66 # ReadDataset function reads the dataset csv file and store it in a dataframe
67 # location_of_file is variable that stores the file location and is defined in config.py file
68 def ReadDataset(location_of_file):
69     # Read the data file from specified location
70     df = pd.read_csv(location_of_file,header=None)
71     # Specifying the column names into the dataset using dataframe, as there are no columns already specified in the dataset
72     df.columns = ['A1', 'A2', 'A3', 'A4', 'A5', 'A6', 'A7', 'A8', 'A9', 'A10', 'A11', 'A12', 'A13', 'A14', 'A15', 'A16']
73     return df
74
75
76 # for loop will iterate over all the columns in dataframe to find and drop the rows with special characters [4]
77 def RemoveSpecialChar(special_char,df):
78     for eachcol in df:
79         # Drop rows where "?" is displayed
80         df.drop(df.loc[df[eachcol] == special_char].index, inplace=True)
81     return df
82

```

```

83     # Encoder function transforms the character and object value in dataframe [5]
84     def Encoder(df):
85         # columnsToEncode is to get a list of columns having datatype as category and object
86         columnsToEncode = list(df.select_dtypes(include=['category','object']))
87         # le is an object of LabelEncoder with no parameter
88         le = LabelEncoder()
89         # for loop will iterate over the columns.
90         # it will first try to use LabelEncoder to fit_transform
91         # and if there are any error than the except will be executed
92         for feature in columnsToEncode:
93             try:
94                 df[feature] = le.fit_transform(df[feature])
95             except:
96                 print('Error encoding '+feature)
97         return df
98
99
100
101    # Define BasicModel function to evaluates the models
102    # displays the confusion matrix and plot it [13] [14]
103    # displays the classification_report
104    def BasicModel(models,X_train,Y_train,X_test,Y_test):
105        # Evaluate each model in turns
106        for name, model in models:
107            # Train the model
108            modelfit = model.fit(X_train,Y_train)
109            # Predict the response for test dataset
110            Y_predict = modelfit.predict(X_test)
111            # Store the accuracy in results
112            results.append(metrics.accuracy_score(Y_test, Y_predict))
113            # Store the model name in names
114            names.append(name)
115            # Print the prediction of test set
116            print('On %s Accuracy is: %f' % (name, metrics.accuracy_score(Y_test, Y_predict)*100))
117            # Store the name and accuracy in basic_score list
118            basic_score.append({'Model Name": name, "Accuracy": metrics.accuracy_score(Y_test, Y_predict)*100})
119            # Print Confusion Matrix and Classification Report
120            print(confusion_matrix(Y_test, Y_predict))
121            print(classification_report(Y_test, Y_predict))
122            # Plot Confusion Matrix [13] [14]
123            cm = confusion_matrix(Y_test, Y_predict, labels=modelfit.classes_)
124            cmdisp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=modelfit.classes_)
125            cmdisp.plot()
126            plt.show()
127        return basic_score
128
129
130    # Define BuildModelRS function to evaluate the models based on the random states [8-12]
131    # rand_state variable has been defined in the config.py file with values 1,3,5,7
132    # stores the Model Name, Random State and Accuracy of each model in score list
133    def BuildModelRS(models,rand_state,dfcreditdata,dfcredittarget):
134        # Evaluate each model in turn
135        for name, model in models:
136            # for loop will train and predict the decision tree model on different random states
137            for n in rand_state:
138                # The training set and test set has been spilted using the feature and target dataframes with different random_state
139                X_train, X_test, Y_train, Y_test = train_test_split(dfcreditdata, dfcredittarget, test_size=0.3, random_state=n)
140                # Train Decision Tree Classifier
141                modelfit = model.fit(X_train,Y_train)
142                # Predict the response for test dataset
143                Y_predict = modelfit.predict(X_test)
144                # Store the accuracy in results
145                results.append(metrics.accuracy_score(Y_test, Y_predict))
146                # Store the model name in names
147                names.append(name)
148                # Store the Model Name, Random State and Accuracy into score list
149                score.append({'Model Name": name, "Random State": int(n), "Accuracy": metrics.accuracy_score(Y_test, Y_predict)*100})
150
151
152    # Define function BuildModelBalCV to evaluate models on the balanced data and utilize cross validation
153    def BuildModelBalCV(models,X,Y):
154        # Evaluate each model in turn
155        for name, model in models:
156            # Define StratifiedKFold [17] [18]
157            skfold = StratifiedKFold(n_splits=10, random_state= None, shuffle=True)
158            # Get the X and Y using StratifiedKFold
159            skfold.get_n_splits(X,Y)
160            # Evaluate each model with cross validation
161            cv_results = cross_val_score(model, X, Y, cv=skfold, scoring='accuracy')
162            model = model.fit(X,Y)
163            # Store the accuracy in results
164            results.append(cv_results)
165            # Store the model name in names
166            names.append(name)
167            # Print the results
168            print('On %s: Mean is %f and STD is %f' % (name, cv_results.mean()*100, cv_results.std()))
169            # Store the Model Name, Mean and STD into score list
170            score.append({'Model Name": name, "Mean": cv_results.mean()*100, "STD": cv_results.std()})
171
172        return score
173        return results
174        return names

```

```

174 # Define BuildFinalModel function to evaluate models with their hyper-parameters
175 def BuildFinalModel(models,X,Y):
176     # Evaluate each model in turn
177     for name, model in models:
178         # Define StratifiedKFold
179         skfold = StratifiedKFold(n_splits=10, random_state= None, shuffle=True)
180         # Get the X and Y using StratifiedKFold
181         skfold.get_n_splits(X,Y)
182         # Evaluate each model with cross validation [21]
183         cv_results = cross_val_score(model, X, Y, cv=skfold, scoring='accuracy')
184         model = model.fit(X.values,Y)
185         # Perform cross-validation and obtain predictions
186         #predictions = cross_val_predict(model, X, Y, cv=5)           #predictions = model.fit(X, Y).predict(X)
187         # Store the cross validationscore into results
188         final_results.append(cv_results)
189         # Store model name into names
190         names.append(name)
191         # Print the results
192         #print('On %s: Mean is %f and STD is %f' % (name, cv_results.mean(), cv_results.std()))
193         # Store the Model Name, Mean and STD into score list
194         score.append({"Model Name": name, "Mean": cv_results.mean(), "STD": cv_results.std()})
195         #print('predictions are ')
196         #print(predictions)
197     return model, score, final_results
198 #return model, score, final_results, predictions
199
200
201

```

6. model.py file: This file has the data analysis, visualization and modelling process.
- Imported all the libraries.

```

8 ###### STEP 1: IMPORTING THE NECESSARY LIBRARIES #####
9
10
11 # Load all the libraries that will be utilized through the code below
12 import pandas as pd
13 #from pandas import read_csv
14 from matplotlib import pyplot
15 import seaborn as sns
16 # Import matplotlib.pyplot to draw and save plots
17 import matplotlib.pyplot as plt
18 # Import Decision Tree Classifier
19 from sklearn.tree import DecisionTreeClassifier
20 # Import LinearDiscriminantAnalysis
21 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
22 # Import train_test_split function
23 from sklearn.model_selection import train_test_split
24 # Import cross_val_score function
25 from sklearn.model_selection import cross_val_score
26 from sklearn.model_selection import cross_val_predict
27 # Import StratifiedKFold function
28 from sklearn.model_selection import StratifiedKFold
29 # Import GridSearchCV function
30 from sklearn.model_selection import GridSearchCV
31 # Import scikit-learn metrics module for accuracy calculation
32 from sklearn import metrics
33 # Import confusion_matrix function
34 from sklearn.metrics import confusion_matrix
35 # Import classification_report function
36 from sklearn.metrics import classification_report
37 # Import ConfusionMatrixDisplay function to plot confusion matrix
38 from sklearn.metrics import ConfusionMatrixDisplay
39 # Import render from graphviz to convert dot file to png
40 from graphviz import render
41 # Import RandomUnderSampler, SMOTE, Pipeline to balance the dataset
42 from imblearn.under_sampling import RandomUnderSampler
43 from imblearn.over_sampling import SMOTE
44 from imblearn.pipeline import Pipeline
45 #import config.py
46 from config import *
47 #import functionlibrary.py
48 from functionlibrary import *
49 import time
50

```

- Read the Credit Card Approval dataset “crx.csv” located in Dataset folder.

```

51 ##### STEP 2: READING THE DATA #####
52
53
54
55 # Calling ReadDataset function from functionlibrary.py file.
56 # It reads the dataset file [1] and store it in a dataframe along with its attributes
57 # location_of_file is defined in the config.py file
58 dfcredit = ReadDataset(location_of_file)
59 print(dfcredit.head(5))
60 # Print the number of rows and columns present in dataset
61 print("The dataset has Rows {} and Columns {}".format(dfcredit.shape[0], dfcredit.shape[1]))
62
63 # Give a new line to clearly format the output
64 Newline()
65
66 # Output:
67 #   A1    A2    A3 A4 A5 A6 A7    A8 A9 A10   A11 A12 A13   A14 A15 A16
68 # 0 b 30.83 0.000 u g w v 1.25 t t 1 f g 00202 0 +
69 # 1 a 58.67 4.460 u g q h 3.04 t t 6 f g 00043 560 +
70 # 2 a 24.50 0.500 u g q h 1.50 t f 0 f g 00280 824 +
71 # 3 b 27.83 1.540 u g w v 3.75 t t 5 t g 00100 3 +
72 # 4 b 20.17 5.625 u g w v 1.71 t f 0 f s 00120 0 +
73 # The dataset has Rows 690 and Columns 16
74

```

- c. Pre-processed the dataset by verifying the missing values, removing the special character “?”, converting the datatypes of variables and encoding the character values to int values.

```

75 ##### STEP 3: PRE-PROCESSING THE DATA #####
76
77 ===== 3.1. Understand the data =====
78
79 # Summarizing the data statistically by getting the mean, std and other values for all the columns
80 print(dfcredit.describe())
81
82 # From the output it is clear that the mean, std and other analysis are calculated
83 # only on columns A3, A8, A11, A15, which specifies that other columns are having
84 # some unnecessary data.
85 # To get rid of missing values, Step 3.2 is performed
86
87 # Give a new line to clearly format the output
88 Newline()
89
90 #           A3          A8          A11          A15
91 #count  690.000000  690.000000  690.000000  690.000000
92 #mean   4.758125  2.223406  2.400000  1017.385507
93 #std    4.781835  3.346513  4.86294  5210.12598
94 #min    0.000000  0.000000  0.000000  0.000000
95 #25%   1.000000  0.165000  0.000000  0.000000
96 #50%   2.750000  1.000000  0.000000  5.000000
97 #75%   7.287500  2.625000  3.000000  395.500000
98 #max   28.000000  28.500000  67.00000  100000.00000
99
100 ===== 3.2. Find the missing values =====
101
102 # Printing the missing values (NaN)
103 print(dfcredit.isnull().sum())
104
105 # Give a new line to clearly format the output
106 Newline()
107
108 # Output:
109 #A1    0
110 #A2    0
111 #A3    0
112 #A4    0
113 #A5    0
114 #A6    0
115 #A7    0
116 #A8    0
117 #A9    0
118 #A10   0
119 #A11   0
120 #A12   0
121 #A13   0
122 #A14   0
123 #A15   0
124 #A16   0
125 #dtype: int64
126
127 # The output shows that there are no missing values (NaN) available in the data
128
129 # To get a clear view of this kind of data, dataset is verified manually and found that
130 # there are "?" special character specified in many columns
131 # Therefore Step 3.3 will find the rows having "?" and drop the respective row from the dataset
132
133 ===== 3.3. Find and drop unnecessary values =====
134
135 # Finding the special character "?" in the dataset and dropping the rows, to get a clean dataset [4]
136
137 # Calling RemoveSpecialChar function from functionlibrary.py file to remove all the special characters from dataset
138 dfcredit = RemoveSpecialChar(special_char,dfcredit)
139
140 # Summarizing the dataset by printing the number of rows and columns present in dataset
141 print("The updated dataset has Rows {} and Columns {}".format(dfcredit.shape[0], dfcredit.shape[1]))
142
143 # Give a new line to clearly format the output
144 Newline()
145
146 # Summarizing the data statistically again by getting the mean, std and other values for all the columns
147 print(dfcredit.describe())
148
149 # The output again shows that mean, std etc is done on attributes A3, A8, A11 and A15 only, instead of all the attributes
150
151 # Give a new line to clearly format the output
152 Newline()
153
154 # Verify the data types of all attributes in the dataset, to verify the difference between the datatypes
155 print(dfcredit.info())
156
157 # The information shows that the datatype of attributes are of int, float and object types
158 # The data needs to be converted to have consistency among all the attributes,
159 # because the dataset does not have a clear instruction of what all the attributes signifies to.
160 # Therefore, it becomes necessary to have the data of same date types
161
162 # Give a new line to clearly format the output
163 Newline()
164

```

```

165 #===== 3.4. Analysing and making consistent data type for all attributes =====
166
167 # Convert the column values from string/object into integer using labelencoder [5]
168
169 # In the dataset, there are mix the data types, such as values as string, float and integer
170 # To have same datatype, created Encoder function to convert the strings into integer values
171
172 # Calling the Encoder function passing the dataframe as parameter that creates a new dataframe to store the new values
173 dfcredit = Encoder(dfcredit)
174
175 # Save the new dataframe into a csv file
176 df = pd.DataFrame(dfcredit)
177 df.to_csv('cleancreditdata.csv')
178
179 # Give a new line to clearly format the output
180 Newline()
181
182 # Summarizing the data statistically again by getting the mean, std and other values on all the columns
183 print(dfcredit.describe())
184 print(dfcredit.info())
185 # Now the count, mean, std, min, max etc are shown for all attributes, this means the data is clean now
186
187 # Give a new line to clearly format the output
188 Newline()
189
190 #===== 3.5. Class distribution =====
191
192 # As mentioned in the dataset details, "A16" attribute is the class attribute [1]
193
194 # Finding the number of rows that belong to each class
195 print(dfcredit.groupby('A16').size())
196
197 # Output is showing that the class "0" is having 296 rows and class "1" is having 357 rows
198 # By this we understand that there are only 2 classes 0 and 1 available in the data
199 # which means this is a Binary Classification type.
200 # This also show the 0 is assigned to denied requests for applications asking credit card approval
201 # and 1 is assigned for approved requests for applications asking credit card approval
202 # This also states that the data is imbalanced.
203 # To understand the dataset more, Data Visualization by plotting the pair plot and heat graph is done in next step
204
205 # Give a new line to clearly format the output
206 Newline()
207

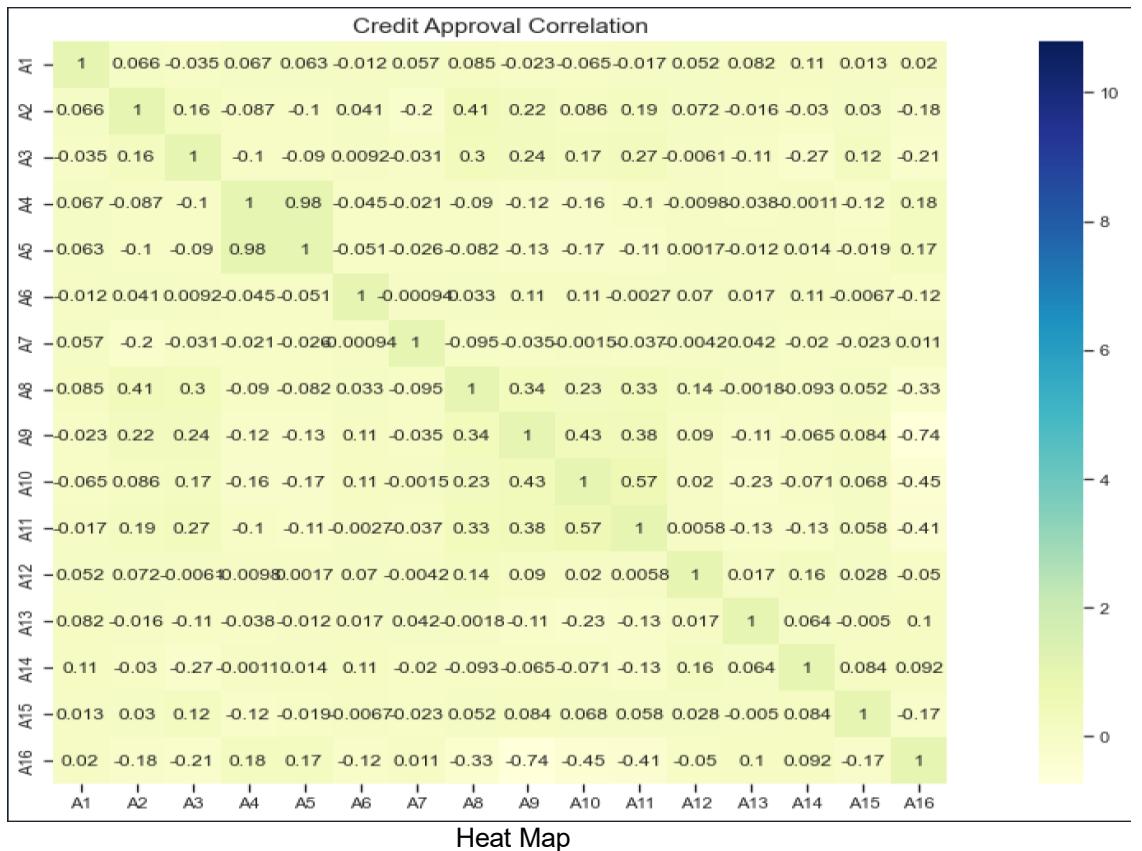
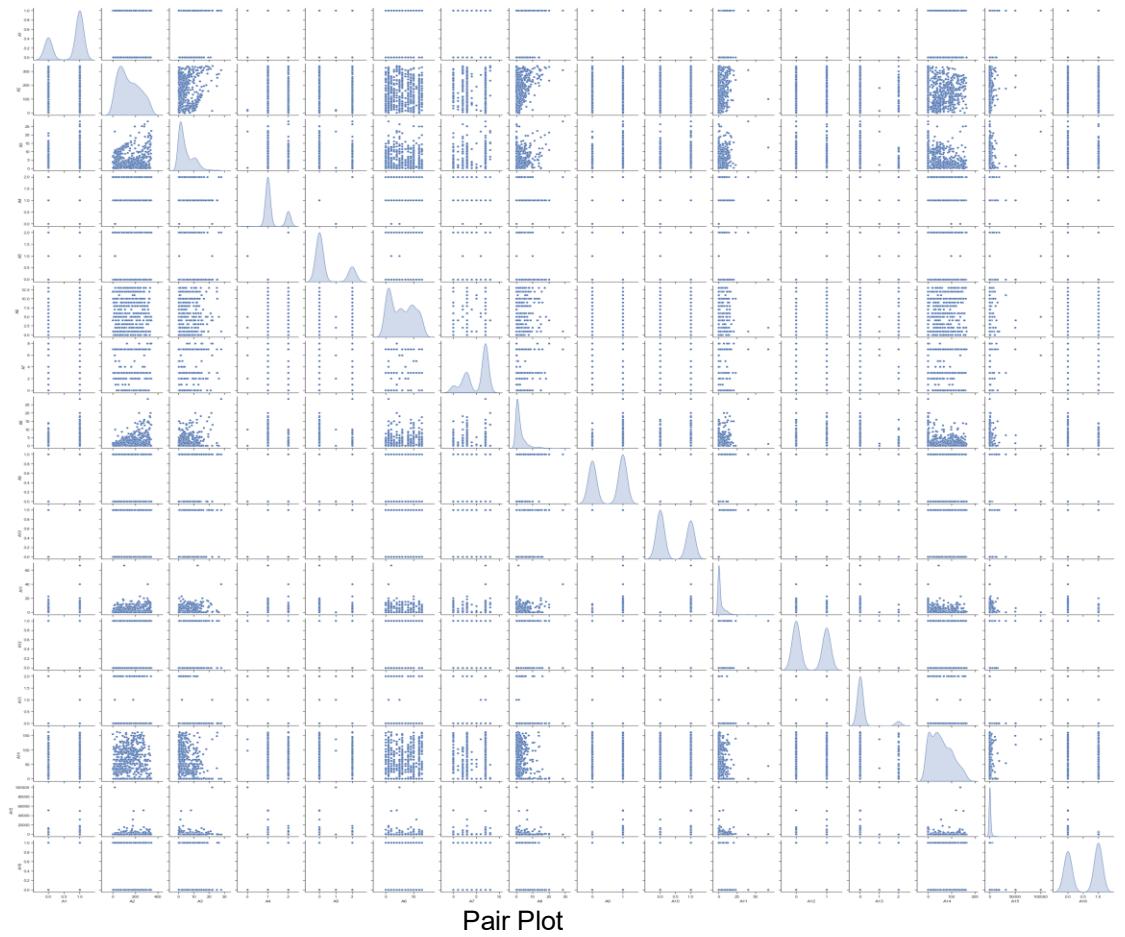
```

- d. Visualised the data to understand the correlation between variables by creating Pair plot and heat map as displayed below.

```

208 ##### Step 4: DATA VISUALIZATION #####
209
210 #===== 4.1. Plot Pairplot =====
211
212 # To visualize the data, pair plot has been demonstrated below [6]
213
214 # Clear the plot dimensions
215 plt.clf()
216
217 # Plot the multivariate plot using pairplot
218 sns.set(style="ticks", color_codes=True)
219 grh = sns.pairplot(dfcredit,diag_kind="kde")
220 plt.show()
221
222 # Give a new line to clearly format the output
223 Newline()
224
225 # Wait for the graph to be displayed
226 time.sleep(30)
227
228 # The plots suggests that there is a high correlation and a predictable relationship
229 # between the diagonal grouping of some pairs of attributes
230 # As stated earlier that the data is imbalanced, the graph is not clearly linear too.
231
232 #===== 4.2. Plot Heat map =====
233
234 # To visualize the data again, heat map is demonstrated below [7]
235
236 # Clear the plot dimensions
237 plt.clf()
238
239 # Plot Heat graph to visualise the correlation
240 plt.figure(figsize=(20,10))
241 matrix1=dfcredit.corr()
242 plt.title('Credit Approval Correlation', size = 15)
243 sns.heatmap(matrix1, vmax=10.8, square=True, cmap='YlGnBu', annot=True)
244
245 # Save the pairplot in png file
246 #plt.savefig("heatplot.png")
247
248 # Give a new line to clearly format the output
249 Newline()
250
251 # Wait for the graph to be displayed
252 time.sleep(30)
253
254 # As observed earlier in the pairplot, the data has high correlation and a predictable relationship
255 # between the diagonal grouping of some pairs of attributes.
256 # The same has been observed in Heat map as well.
257 # To fix the imbalance problem, I have used the balancing technique RandomUnderSampler and SMOTE
258 # in the further steps after checking the performance of Decision Tree and Linear Discriminant Analysis models
259 # using the basic HoldOut Validation (i.e. train and test split)
260

```



- e. Developed 2 machine learning algorithms (Decision Tree Classifier and Linear Discriminant Analysis) to compare between the models based on its accuracy.

```

262 ##### STEP 5: FEATURE SELECTION #####
263
264 # To implement any model on the data, we need to do feature selection first
265 # and therefore the dataset is divided into 2 parts.
266 # Firstly features (all the attributes except the target attribute) and secondly target (class attribute),
267 # The dfcredit dataframe has been split into dfcredittarget having only column A16, as this is the class attribute,
268 # and all the other attributes are taken into dfcreditdata dataframe as features
269
270 dfcreditdata = dfcredit.drop("A16", axis=1) # Features
271 dfcredittarget = dfcredit.iloc[:, -1] # Target
272
273 # None of the feature selection methods are used here because the dataset is having very less number of attributes.
274 # and therefore, the dimensionality reduction has not been performed here.
275 # Also, the attribute names have been changed to meaningless by the author [1], hence all the attributes are considered as features
276
277 ##### BUILDING DECISION TREE CLASSIFIER AND LINEAR DISCRIMINANT ANALYSIS MODELS #####
278
279 ##### STEP 6: BUILD BASE MODELS WITH HOLDOUT VALIDATION #####
280
281 # Base Model is built below using the HoldOut Validation (train test split) to evaluate the accuracy
282
283 # ===== 6.1. Split the data to train and test set =====
284
285 # Spot Check Algorithms with Decision Tree [8] [9] [10] and Linear Discriminant algorithms [11] [12]
286
287 # Data is split into training and test set using the feature and target dataframes
288 # to understand the model performance with 70% training set and 30% test set
289 X_train, X_test, Y_train, Y_test = train_test_split(dfcreditdata, dfcredittarget, test_size=0.3, random_state=None)
290
291 # ===== 6.2. Define models =====
292
293 # Clear lists
294 ▲ 294 ClearLists()
295
296 # Create model object of Decision Tree Classifier [8] [9] [10]
297 ▲ 297 models.append(('DTCLS', DecisionTreeClassifier()))
298 # Create model object of Linear Discriminant Analysis [11] [12]
299 ▲ 299 models.append(('LDA', LinearDiscriminantAnalysis()))
300
301 # ===== 6.3. Build Model and Evaluate it =====
302
303 # Calling BasicModel function
304 ▲ 304 BasicModel(models, X_train, Y_train, X_test, Y_test)
305
306 # Create a dataframe to store accuracy
307 ▲ 307 dfbasicsscore = pd.DataFrame(basic_score)
308 print(dfbasicsscore.head())
309
310 # Give a new line to clearly format the output
311 ▲ 311 Newline()
312
313 # Output:
314 # On DTCLS Accuracy is: 81.632653
315 # On LDA Accuracy is: 85.204082
316 # The output shows that the accuracy for LDA model is better than Decision Tree Classifier.
317 # To improve the performance of each model, the evaluation is done with respect to different random_state in the next step
318
319 ##### STEP 7: BUILD MODELS WITH RANDOM STATE #####
320
321 # Model is build with random state to evaluate the accuracy
322
323 # ===== 7.1. Define models =====
324
325 # Clear lists
326 ▲ 326 ClearLists()
327
328 # Create model object of Decision Tree Classifier
329 ▲ 329 models.append(('DTCLS', DecisionTreeClassifier()))
330 # Create model object of Linear Discriminant Analysis
331 ▲ 331 models.append(('LDA', LinearDiscriminantAnalysis()))
332
333 # ===== 7.2. Build Model and Evaluate it =====
334
335 # Calling BuildModelRS to evaluate the models with random states and return the accuracy
336 ▲ 336 BuildModelRS(models, rand_state, dfcreditdata, dfcredittarget)
337
338 # Create a dataframe to store accuracy
339 ▲ 339 dfrsscore = pd.DataFrame(score)
340 print(dfrsscore.head(8))
341
342 # Give a new line to clearly format the output
343 ▲ 343 Newline()
344
345 # The output shows that the accuracy changes in each random state. Therefore, data has been balanced in further steps.
346
347
348

```

```

349 ##### STEP 8: BALANCING THE DATA #####
350 # Installed imbalance-learn library using "conda install -c conda-forge imbalanced-learn" [15]
351 # To balance the data, RandomUnderSampler and SMOTE is utilized to undersample and oversample the data along with pipeline [16]
352
353 # Define oversample with SMOTE function
354 oversample = SMOTE()
355 # Define undersample with RandomUnderSampler function
356 undersample = RandomUnderSampler()
357 # Define Steps for oversample and undersample
358 steps = [('o', oversample), ('u', undersample)]
359 # Define the pipeline with the steps
360 pipeline = Pipeline(steps = steps)
361 # Fit the features and target using pipeline and resample them to get X and Y
362 X, Y = pipeline.fit_resample(dfcreditdata, dfcredittarget)
363 # Print the shape of X and Y
364 #print("The features dataset has Rows {} and Columns {}".format(X.shape[0], X.shape[1]))
365 #print("The target dataset has Rows {} and Columns {}".format(Y.shape[0], 0))
366
367 # Give a new line to clearly format the output
368 Newline()
369
370 # By balancing the data using oversample and undersample, X and Y now have adequate amount of data available
371 # Reducing the dimensionality is not needed because the dataset is small with only 15 features
372
373 ###### STEP 9: OPTIMIZE MODEL ON BALANCED DATA USING CROSS VALIDATION #####
374
375 # StratifiedKFold Cross Validation is utilized to improve the performance as it splits the data approximately in the same percentage [17] [18]
376 # StratifiedKFold Cross Validation is used because there are 2 classes, and the split of train and test set could be properly done
377
378 # ===== 9.1. Define models =====
379
380 # Clear lists
381 ClearLists()
382
383 # Create model object of Decision Tree Classifier
384 models.append(('DTCLS', DecisionTreeClassifier()))
385 # Create model object of Linear Discriminant Analysis
386 models.append(('LDA', LinearDiscriminantAnalysis()))
387
388 # ===== 9.2. Build Model and Evaluate it =====
389
390 # Calling BuildModelBalCV function to get the accuracy, cross validation results and names of models used
391 BuildModelBalCV(models,X,Y)
392
393 # Create a dataframe to store accuracy
394 dfscore = pd.DataFrame(score)
395 print(dfscore.head())
396
397 # Give a new line to clearly format the output
398 Newline()
399
400 # Compare Algorithms and plot them in boxplot
401 pyplot.clf()
402 pyplot.boxplot(results, labels=names)
403 pyplot.title('Algorithm Comparison')
404 pyplot.show()
405
406 # Give a new line to clearly format the output
407 Newline()
408
409 # Outout:
410 # On DTCLS: Mean is 82.355243 and STD is 0.033680
411 # On LDA: Mean is 87.255477 and STD is 0.040058
412
413 # LDA is having the highest Mean value of 87.25%.
414 # Decision Tree has also performed good with Mean as 82.35%
415 # Looking at the boxplot and the Mean values, it is clear that LDA has performed better even after balancing the data
416 # along with StratifiedKFold Cross Validation by giving Mean as 87.25%
417
418 # Give a new line to clearly format the output
419 Newline()
420

```

- f. Tuned the Decision Tree Classifier and LDA models with the hyperparameters to get the best parameters.

```

420 ##### STEP 10: TUNE DECISION TREE BY GRIDSEARCHCV #####
421
422 # Tuning the Decision Tree Classifier with GridSearchCV to find the best hyperparameter [19] [20]
423
424 # Define decision tree classifier model
425 dtmodel = DecisionTreeClassifier()
426 # Define StratifiedKFold
427 skfold = StratifiedKFold(n_splits=10, random_state= None, shuffle=True)
428 # Define parameters
429 param_dict = dict()
430 param_dict = {"criterion": ["gini", "entropy"], "max_depth": [7,5,3]}
431 # Build GridSearchCV to get the accuracy
432 search = GridSearchCV(dtmodel, param_dict, scoring='accuracy', cv=skfold, n_jobs=-1)
433 # Fit the GridSearchCV
434 results = search.fit(X, Y)
435 # Summarize
436 dt_accuracy = results.best_score_
437 dt_para = results.best_params_
438 print('Decision Tree Mean Accuracy: %f' % results.best_score_)
439 print('Config: %s' % results.best_params_)
440
441 # Give a new line to clearly format the output
442 Newline()
443
444 # Output:
445 # Decision Tree Mean Accuracy: 0.865630
446 # Config: {'criterion': 'entropy', 'max_depth': 3}
447 # The best parameter for Decision Tree Classifier are 'criterion' as 'entropy', 'max_depth' as 3,
448 # using which the performance of decision tree has been increased slightly with 86.56%
449

```

```

450 ##### STEP 11: TUNE LDA BY GRIDSEARCHCV #####
451
452 # Tuning the LDA with GridSearchCV to find the best hyperparameter [19]
453
454 # Define LDA model
455 ldamodel = LinearDiscriminantAnalysis()
456 # Define StratifiedKFold
457 skfold = StratifiedKFold(n_splits=10, random_state=None, shuffle=True)
458 # Define parameters
459 grid = dict()
460 grid['solver'] = ['svd', 'lsqr', 'eigen']
461 # Build GridSearchCV to get the accuracy
462 search = GridSearchCV(ldamodel, grid, scoring='accuracy', cv=skfold, n_jobs=-1)
463 # Fit the GridSearchCV
464 results = search.fit(X, Y)
465 # Summarize
466 lda_accuracy = results.best_score_
467 lda_para = results.best_params_
468 print('LDA Mean Accuracy: %f' % results.best_score_)
469 print('Config: %s' % results.best_params_)
470 #print(lda_accuracy)
471 #print(lda_para['solver'])
472
473 # Give a new line to clearly format the output
474 Newline()
475
476 # Output:
477 # LDA Mean Accuracy: 0.871088
478 # Config: {'solver': 'svd'}
479
480 # The best parameter LDA are 'solver' as 'svd'
481 # using which the performance of LDA model has been increased slightly to 87.10%
482
483

```

g. Created models with best hyperparameters along with cross validation method.

```

484 ##### STEP 12: BUILD TUNED MODEL #####
485
486 # Building the model finally on balanced data with best hyper-parameters along with StratifiedKFold cross validation [21]
487
488 # ===== 12.1. Define models =====
489
490 # Clear lists
491 ClearLists()
492
493 # Define models with the best hyperparameters found earlier
494 models.append(('DTCLS', DecisionTreeClassifier(criterion=dt_para['criterion'], max_depth=dt_para['max_depth'])))
495 models.append(('LDA', LinearDiscriminantAnalysis(solver=lda_para['solver'])))
496
497 # ===== 12.2. Build Model and Evaluate it =====
498
499 # Calling BuildFinalModel to get the accuracy after tuning
500 BuildFinalModel(models,X,Y)
501 # Create a datafram to store accuracy
502 dffinalscore = pd.DataFrame(score)
503 print(dffinalscore.head())
504
505 max_mean = dffinalscore['Mean'].max()
506 max_mean_idx = dffinalscore['Mean'].idxmax()
507 model_name = dffinalscore.loc[max_mean_idx, 'Model Name']
508
509 # Give a new line to clearly format the output
510 Newline()
511
512 print('%s has the max accuracy as %f' % (model_name, max_mean))
513
514 # Give a new line to clearly format the output
515 Newline()
516
517 # Compare Algorithms and plot them in boxplot
518 pyplot.clf()
519 pyplot.boxplot(final_results, labels=names)
520 pyplot.title('Algorithm Comparison')
521 pyplot.show()
522
523 # Give a new line to clearly format the output
524 Newline()
525
526 # Output:
527 # On DTCLS: Mean is 84.178404 and STD is 0.035138
528 # On LDA: Mean is 86.977700 and STD is 0.029924
529
530 # Looking at the Mean and Boxplot, it is clear that LDA has performed better than
531 # Decision Tree after tuning the model with hyperparameters using StratifiedKFold cross validation.
532 # However, there is not much difference in both the accuracies.
533
534

```

h. Trained the best performed model and dumped it in model.pkl file.

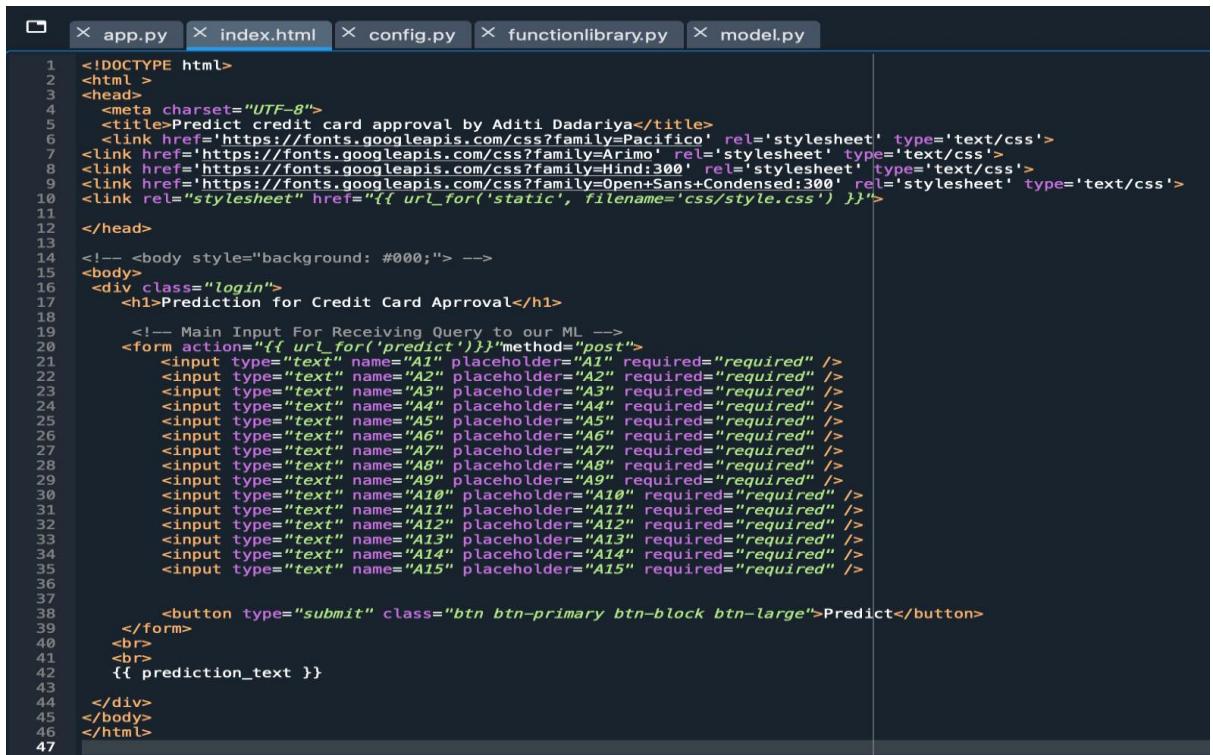
```

534 ##### STEP 12: BUILD BEST PERFORMED MODEL #####
535
536
537
538 # Clear lists
539 ClearLists()
540
541
542 # Define models with the best hyperparameters found earlier
543 if model_name == "DTCLS":
544     models.append(('DTCLS', DecisionTreeClassifier(criterion=dt_para['criterion'], max_depth=dt_para['max_depth'])))
545 elif model_name == 'LDA':
546     models.append(('LDA', LinearDiscriminantAnalysis(solver=lda_para['solver'])))
547
548 # Calling BuildFinalModel to get the prediction after tuning
549 model, score, final_results = BuildFinalModel(models,X,Y)
550
551 #$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$$
552 # Dump the model.pkl file
553 import pickle
554 pickle.dump(model, open('model.pkl', 'wb'))
555
556
557

```

i. When model.py is run, a new model.pkl file is created.

7. index.html file: This file creates a web application with 15 text boxes and a “Predict” button.

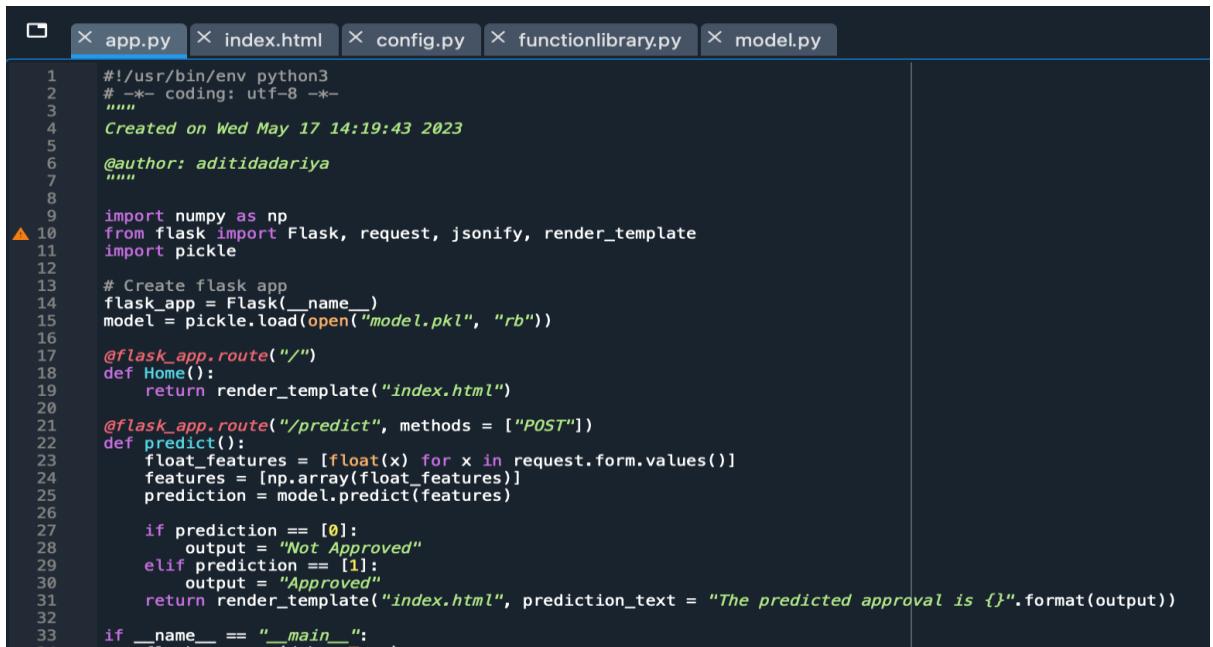


```

1  <!DOCTYPE html>
2  <html >
3  <head>
4      <meta charset="UTF-8">
5      <title>Predict credit card approval by Aditi Dadariya</title>
6      <link href='https://fonts.googleapis.com/css?family=Pacifico' rel='stylesheet' type='text/css'>
7      <link href='https://fonts.googleapis.com/css?family=Arimo' rel='stylesheet' type='text/css'>
8      <link href='https://fonts.googleapis.com/css?family=Hind:300' rel='stylesheet' type='text/css'>
9      <link href='https://fonts.googleapis.com/css?family=Open+Sans+Condensed:300' rel='stylesheet' type='text/css'>
10     <link rel="stylesheet" href="{{ url_for('static', filename='css/style.css') }}>
11
12 </head>
13
14 <!-- <body style="background: #000;"> -->
15 <body>
16     <div class="login">
17         <h1>Prediction for Credit Card Approval</h1>
18
19         <!-- Main Input For Receiving Query to our ML -->
20         <form action="{{ url_for('predict') }}" method="post">
21             <input type="text" name="A1" placeholder="A1" required="required" />
22             <input type="text" name="A2" placeholder="A2" required="required" />
23             <input type="text" name="A3" placeholder="A3" required="required" />
24             <input type="text" name="A4" placeholder="A4" required="required" />
25             <input type="text" name="A5" placeholder="A5" required="required" />
26             <input type="text" name="A6" placeholder="A6" required="required" />
27             <input type="text" name="A7" placeholder="A7" required="required" />
28             <input type="text" name="A8" placeholder="A8" required="required" />
29             <input type="text" name="A9" placeholder="A9" required="required" />
30             <input type="text" name="A10" placeholder="A10" required="required" />
31             <input type="text" name="A11" placeholder="A11" required="required" />
32             <input type="text" name="A12" placeholder="A12" required="required" />
33             <input type="text" name="A13" placeholder="A13" required="required" />
34             <input type="text" name="A14" placeholder="A14" required="required" />
35             <input type="text" name="A15" placeholder="A15" required="required" />
36
37
38             <button type="submit" class="btn btn-primary btn-block btn-large">Predict</button>
39
40         <br>
41         <br>
42         {{ prediction_text }}
43
44     </div>
45
46 </body>
47 </html>

```

8. app.py file: This file runs index.html to open a web application, uses model.pkl file to predict the output and display the output in web application. The web application is at “<http://127.0.0.1:5000/>”



```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4      Created on Wed May 17 14:19:43 2023
5
6      @author: aditidadariya
7
8
9      import numpy as np
10     from flask import Flask, request, jsonify, render_template
11     import pickle
12
13     # Create flask app
14     flask_app = Flask(__name__)
15     model = pickle.load(open("model.pkl", "rb"))
16
17     @flask_app.route("/")
18     def Home():
19         return render_template("index.html")
20
21     @flask_app.route("/predict", methods = ["POST"])
22     def predict():
23         float_features = [float(x) for x in request.form.values()]
24         features = [np.array(float_features)]
25         prediction = model.predict(features)
26
27         if prediction == [0]:
28             output = "Not Approved"
29         elif prediction == [1]:
30             output = "Approved"
31         return render_template("index.html", prediction_text = "The predicted approval is {}".format(output))
32
33     if __name__ == "__main__":
34         flask_app.run(debug=True)
35

```

```

In [8]: runfile('/Users/aditidadariya/Aditi Personal/Data Glacier - Internship/Week4/MLModelDeployment/app.py', wdir='/Users/aditidadariya/Aditi Personal/Data Glacier - Internship/Week4/MLModelDeployment')
* Running on flask app 'app'
* Debug mode: on
WARNING: This is a development server. Do not use it in a production deployment. Use a production WSGI server instead.
* Running on http://127.0.0.1:5000
* Press Ctrl+C to quit
* Restarting with watchdog (fsevents)
* Debugger is active!
* Debugger PIN: 514-666-156

```

9. Web application is below with the inputs for credit cards application provided and the output of prediction displayed.

Prediction for Credit Card Approval

A1	A2	A3	A4	A5	A6	A7	A8	A9
A10	A11	A12	A13	A14	A15	<input type="button" value="Predict"/>		

Prediction for Credit Card Approval

1	171	1.5	1	0	2	3	5.5	1
1	3	1	0	0	0	<input type="button" value="Predict"/>		

Prediction for Credit Card Approval

A1	A2	A3	A4	A5	A6	A7	A8	A9
A10	A11	A12	A13	A14	A15	<input type="button" value="Predict"/>		

The predicted approval is Not Approved

Prediction for Credit Card Approval

1	79	1	1	0	1	7	0.5	0
0	0	1	2	94	0	<input type="button" value="Predict"/>		

Prediction for Credit Card Approval

A1	A2	A3	A4	A5	A6	A7	A8	A9
A10	A11	A12	A13	A14	A15	<input type="button" value="Predict"/>		

The predicted approval is Approved

10. References:

[1] <https://archive.ics.uci.edu/ml/datasets/Credit+Approval>