



ANALYSIS OF SURVIVAL OF LUNG CANCER PATIENTS

ST 4052 – Data Analysis Final Project Group 4

M. D. D. De Costa – s14533

S. P. P. M. Sudasinghe – s14510

K. G. S. K. Wickramasinghe - s14594

Abstract

This report is based on data recorded of patients who faced to the Cardio-thoracic surgery. The objective is to predict survival status of cancer patients who faced to the Cardio-thoracic surgery using the machine learning. The software package that we used to do this analysis was Python. We deployed many machine learning models such as Logistic Regression, Decision Tree, K Nearest Neighbors, Support Vector Machine, Random Forest, and Voting Classifier. Considering accuracy, interpretability and computational cost, we decided to go with the Logistic Regression model. The next step of this analysis is to integrate the predictive model to a data product.

Table of Contents

Abstract.....	1
Introduction	3
Objective	3
About the Dataset.....	3
Descriptive Analysis	4
Survival of the patient after 1 year from Surgery	4
Diagnosis type of the cancer.....	4
Distributons of 'in_vol' and 'out_vol'	5
Diabetes Vs Survival	5
Other disease Vs Survival.....	5
Smoking before surgery	6
Age Vs Survival	6
Size of the Tumor	6
Correlation - Association plot	7
Important Results of the Advanced Analysis	7
Issues encountered and proposed solutions.....	9
Discussion and conclusion	10
Appendix	11

List of Figures

Figure 1	4
Figure 2	4
Figure 3	4
Figure 4	5
Figure 5	5
Figure 6	5
Figure 7	6
Figure 8	6
Figure 9	6
Figure 10	6
Figure 11	6
Figure 12	7
Figure 13	9

Introduction

Worldwide, lung cancer is the most common malignancy and the most common cause of cancer deaths in the past few decades thus its five-year survival rate (17.8%) was much lower than that of other leading cancers. Owing to the high fatality rate, it remains to be an important public health issue. So even when considering Poland, lung cancer has been found as the leading cause of cancer death in both men and women. Cardio-thoracic surgery is a surgical instrument for the organs inside the thorax (the chest) which in general involves the treatments for conditions of heart and lungs. Since these organs are severely important, surgeries involving these organs have a high risk. Cardio-thoracic surgery is frequently used to assess or repair lungs affected by cancer, trauma or pulmonary disease.

Objective

To predict survival status of cancer patients who faced to the Cardio-thoracic surgery using the machine learning. The proposed solution should be able to identify and interpret risk factors in classifying survival status of lung cancer patients which can help assist the medical practitioner in their diagnosis.

About the Dataset

The data is dedicated to classification problem related to the cancer patients after a Cardio-thoracic surgery in which there are two classes class 1 and 2, where 1 indicates the death of the patient within one year after the surgery and 2 indicates the patients who survive. The data was collected from a Surgery Centre for patients who underwent major lung resections for primary lung cancer in the years 2007 to 2011.

Variable	Description
id	Id of the patient
diagnose	Diagnosis type of the cancer
in_vol	Maximal volume of gas that can be exhaled from lung (get into lung)
out_vol	Expiratory Volume per second (get out from lung)
pain	Pain before cancer (T- yes,F- no)
Haemoptysis	Haemoptysis before surgery (T- yes,F- no) (A symptom with lung cancer patients)
short_breath	Shortness of breath of patient before
cough	Cough before surgery (T- yes,F- no)

weakness	Weakness before surgery (T- yes,F- no)
tumour_size	Size of the Tumour
diabetes	Diabetes before surgery (T- yes,F- no)
breathing_difficulty	Breathing difficulty before surgery (T- yes,F- no)
other_disease	Other diseases before surgery (T- yes,F- no)
smoking	Smoking before surgery (T- yes,F- no)
asthma	Asthma before surgery (T- yes,F- no)
age	Age of the patient
survive	Survival of the patient after 1 year from Surgery (1- Died, 2- Survived)

Descriptive Analysis

Survival of the patient after 1 year from Surgery

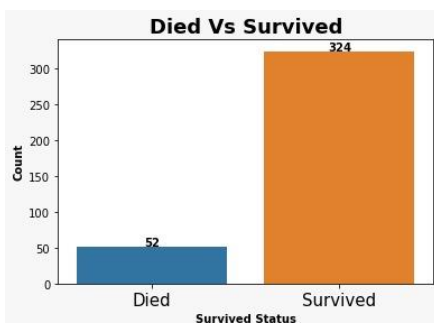


Figure 1

Out of the 376 patients, 52 patients died in the 1 year time and 324 survived, which is a 13.8% death rate. This is the response variable of the dataset. As you can see, this variable is highly unbalanced. Therefore, we should be careful when fitting models to this data.

Diagnosis type of the cancer

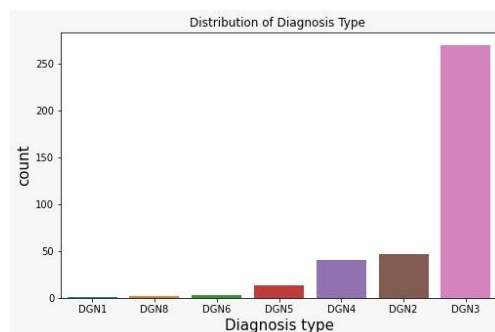


Figure 2

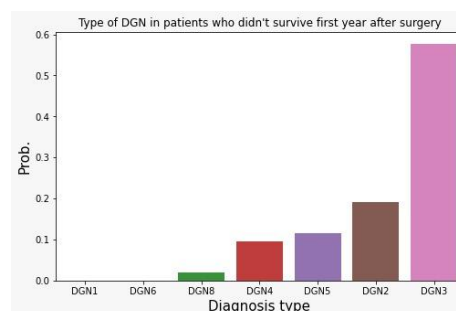


Figure 3

Referring to the first graph, for diagnosis, the large majority of patients are in category DGN3. The other categories are relatively small while category 2,4 and 5 should be considered for their counts in that order. Then the graph in the right side shows the type of diagnosis in patients who didn't survive first year after surgery. As you can see most of the died patients belongs to the category DGN3.

Distributions of 'in_vol' and 'out_vol'

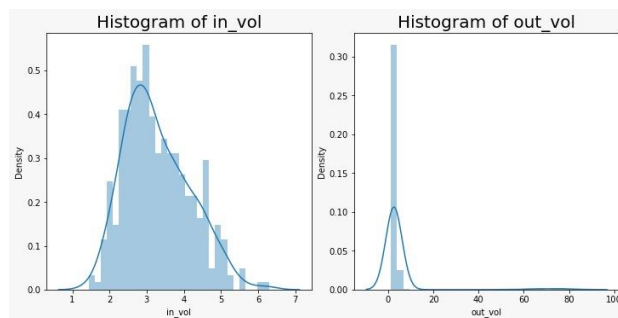


Figure 4

According to the first histogram the average maximal volume of gas that can be exhaled by a patient is around 3. Since this is a right skewed graph, we can say that there were few patients who can exhale more volume of gas. Then from the second graph we can see that majority of the patients lie between 0 and 20 expiratory volume per second.

Diabetes Vs Survival

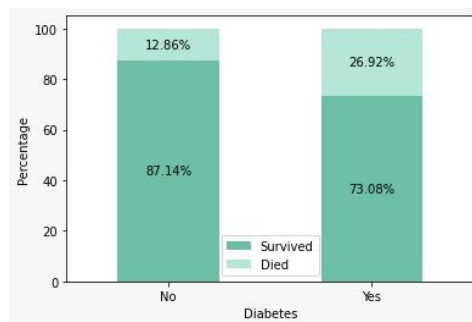


Figure 5

It seems that patients who have diabetes before surgery die more than who are not having diabetes.

Other disease Vs Survival

As you can see, the patients who have other disease have a less chance to be survived than the other patients.

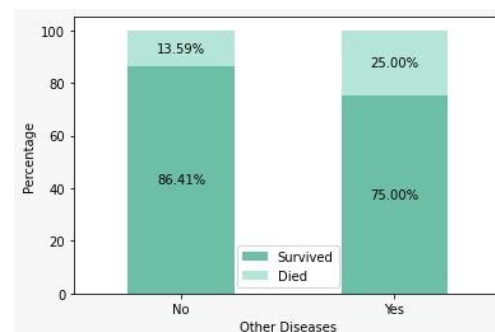


Figure 6

Smoking before surgery

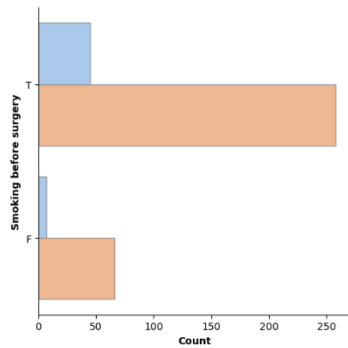


Figure 7

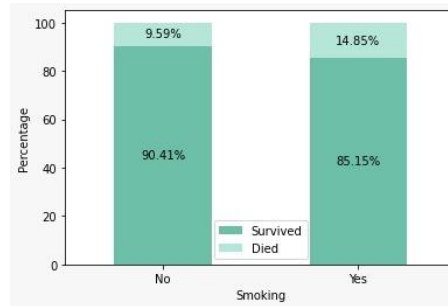


Figure 8

It seems that most of the cancer patients having smoking habits before the surgery. Also, we can see that people who smoked are died after the surgery than who are not having smoking habit. However the difference between died smokers and died non-smokers is considerably small which is around 5%.

Age Vs Survival

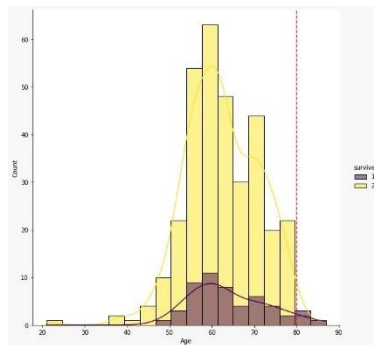


Figure 9

Here we drew two distributions of age, yellow one which represents people who survived and the brown one, people who do not survived. As you can see, patients who are more than 80 years of age has a very high chance of not surviving.

Size of the Tumor

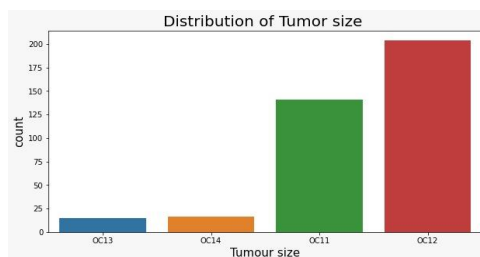


Figure 10

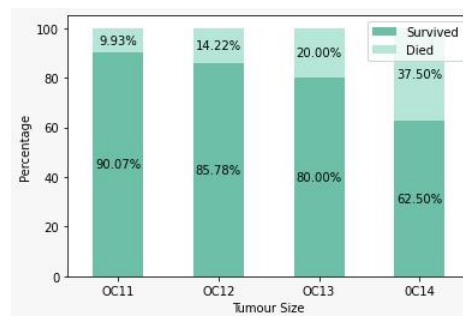


Figure 11

There are 4 sizes of tumor, OC11, OC12, OC13, OC14. Among them most of the patients have OC12 size of tumor and then OC11 size of tumor. The proportion of dead to live generally increases with the tumor size ranging from OC11 to OC14, indicating the higher tumor size correlates to higher chance of death even with surgery.

Correlation - Association plot

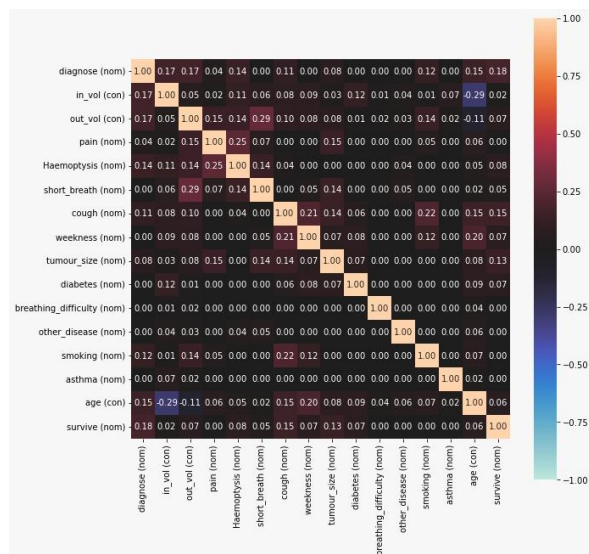


Figure 12

As you can see in the graph below, we can't see any heavy correlations between these variables. However, there is a small negative correlation between age and Maximal volume of gas that can be exhaled from lungs which is true in real life. And there is a small positive correlation between short breath and expiratory Volume per second (vol_out). Therefore, multicollinearity exists in between these variables.

Important Results of the Advanced Analysis

So now that we have thoroughly explored the patterns and trends in the data set, the next step is to utilize machine learning models to see if how well we can predict the target variable, Survive, with the feature variables.

As we mentioned in the descriptive analysis, the response variable "survive" is highly unbalanced. Therefore, we used the SMOTE (Synthetic Minority Over-sampling Technique) technique in order to balance the dataset.

Also, as we saw in our descriptive analysis, the dataset contains outlier observations. Therefore, we removed those outliers and fitted the models.

For the advanced analysis, we used some mainstream models which are Logistic Regression, Decision Tree, K Nearest Neighbors, Support Vector Machine, Random Forest, and Voting Classifier. Logistic Regression, KNN, SVC, and Voting Classifier did not have an overfitting problem. Random Forest model and decision tree models had an overfitting problem, so we used hyper parameter tuning to get rid of that.

Given below are the results of the models fitted, before balancing the data.

		Logistic Regression	Decision Tree	KNN	SVM	Random Forest	Voting Classifier
ROC AUC	Train	52.7 %	73.8 %	52.5 %	87.3 %	85.1 %	69.6 %
	Test	49.4 %	48.1 %	49.4 %	50.0 %	44.8 %	50.0 %
Accuracy	Train	84.2 %	91.4 %	85.0 %	87.2 %	90.8 %	90.6 %
	Test	84.4 %	82.2 %	84.4 %	85.6 %	76.7 %	85.6 %
Precision	Train	85.2 %	91.2 %	85.1 %	97.4 %	95.6 %	89.9 %
	Test	85.4 %	85.1 %	85.4 %	85.6 %	84.1 %	85.6 %
Recall	Train	98.4 %	99.3 %	99.7 %	87.2 %	93.4 %	100.0 %
	Test	98.7 %	96.1 %	98.7 %	100.0 %	89.6 %	100.0 %
F1 Score	Train	91.3 %	95.1 %	91.8 %	92.0 %	94.5 %	94.7 %
	Test	91.6 %	90.2 %	91.6 %	92.2 %	86.8 %	92.2 %

Table 1

Even though the results were good, and it has a very good predictive power, we used SMOTE balancing technique as mentioned above. Then we got the following results.

		Logistic Regression	Decision Tree	KNN	SVM	Random Forest	Voting Classifier
ROC AUC	Train	54.6 %	89.0 %	75.1 %	87.3 %	85.1 %	87.3 %
	Test	43.5 %	48.6 %	48.5 %	53.8 %	44.8 %	42.2 %
Accuracy	Train	72.5 %	92.5 %	67.8 %	87.2 %	90.8 %	87.2 %
	Test	74.4 %	66.7 %	55.6 %	86.7 %	76.7 %	66.7 %
Precision	Train	86.0 %	96.9 %	96.1 %	97.4 %	95.6 %	97.4 %
	Test	83.8 %	85.1 %	84.9 %	86.5 %	84.1 %	83.1 %
Recall	Train	80.6 %	94.1 %	64.5 %	87.2 %	93.4 %	87.2 %
	Test	87.0 %	85.1 %	58.4 %	100.0 %	89.6 %	76.6 %
F1 Score	Train	83.2 %	95.5 %	77.2 %	92.0 %	94.5 %	92.0 %
	Test	85.4 %	79.2 %	69.2 %	92.8 %	86.8 %	79.7 %

Table 2

By comparison of F1 scores of all the models, we could notice that applying SMOTE methodology results in poor predictability of data. So, we made the decision to go along with models fitted to original data.

Looking at those models, we realized that Logistic Regression, KNN, SVM and Voting Classifier performs exceptionally well (Decision Tree and Random Forest have a big difference in F1 scores, which indicates that the models are overfitted). Out of them, considering interpretability and computational cost, we decided to go with the Logistic Regression model.

So, the final model that we will be using for our data product is the Logistic Regression model without SMOTE.

Then we drew a feature importance plot to the model selected.

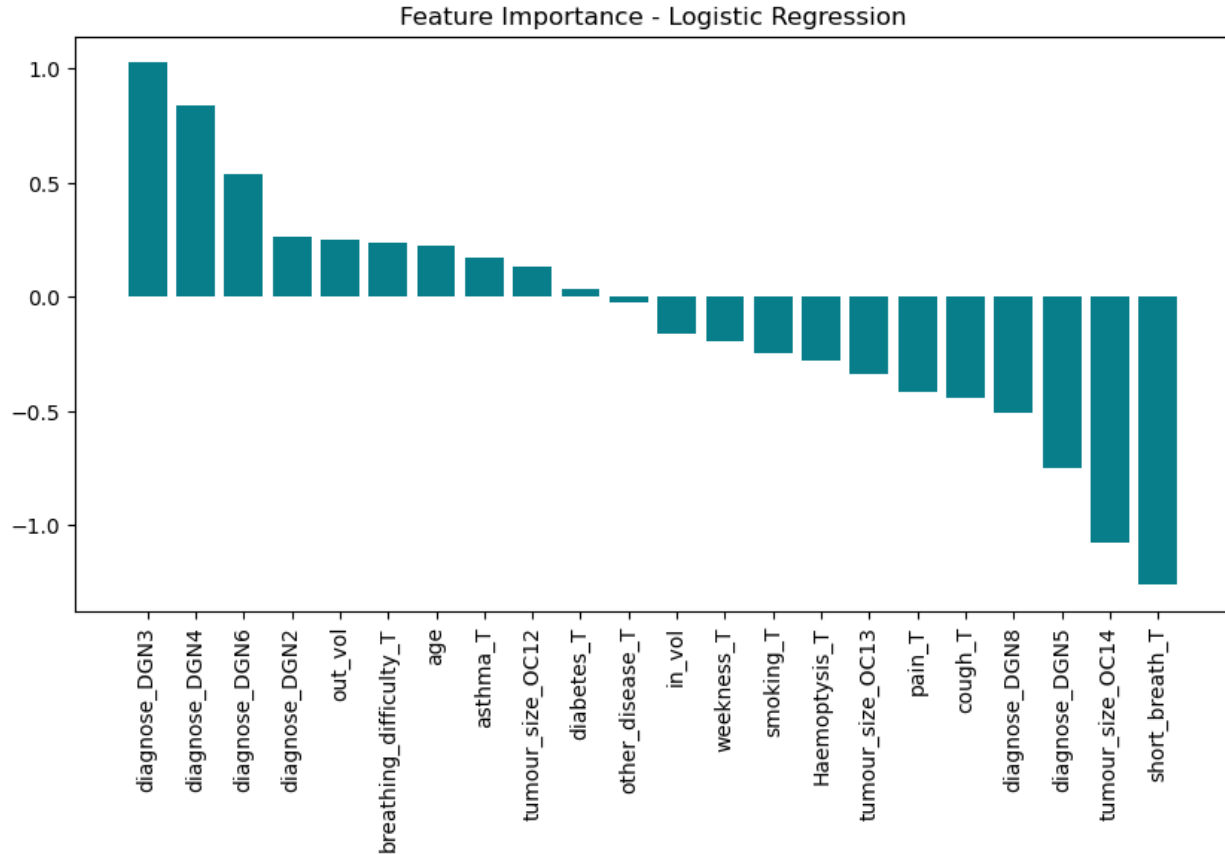


Figure 13

As you can see, all the variables are more or less important to the model. Therefore we decided to not to remove any variable.

Issues encountered and proposed solutions

Issues	Solution
The dataset is highly unbalanced	Using oversampling technique SMOTE (Synthetic Minority Over-sampling Technique) from python library imblearn.
our dataset contained less number of records and due to that it was not an easy task to interpret graphs using counts.	To overcome this, we took percentages to draw graphs and interpret them.

Discussion and conclusion

This study was conducted to predict survival status of cancer patients who faced to the Cardio-thoracic surgery using the machine learning. According to the initial descriptive analysis our target variable 'survive' was highly unbalanced. Therefore, we used SMOTE in order to balance the dataset. Then the Logistic Regression, Decision Tree, K Nearest Neighbors, Support Vector Machine were fitted to the dataset, and further to increase F1 score random forest and Voting Classifier models were fitted with hyperparameter tuning techniques. Looking at those models, we realized that Logistic Regression, KNN, SVM and Voting Classifier performs exceptionally well. Out of them, considering interpretability and computational cost, we decided to go with the Logistic Regression model. So, the final model that we will be using for our data product is the Logistic Regression model without SMOTE.

Appendix

<pre> 1. import pandas as pd 2. import numpy as np 3. import matplotlib.pyplot as plt 4. import seaborn as sns 5. import plotly.express as px 6. from sklearn.model_selection import train_test_split 7. from imblearn.over_sampling import SMOTE, ADASYN 8. from sklearn.metrics import accuracy_score, recall_score, precision_score, f1_score 9. from sklearn.preprocessing import Robust- Scaler 10. from sklearn.ensemble import Random- ForestClassifier 11. from sklearn.model_selection import cross_val_score 12. from sklearn.model_selection import GridSearchCV 13. from sklearn.ensemble import Random- ForestClassifier 14. from sklearn.linear_model import LogisticRegression 15. from sklearn.metrics import classifica- tion_report </pre>	<pre> 16. from sklearn.metrics import roc_auc_score 17. from sklearn.metrics import precision_score 18. from sklearn.metrics import recall_score 19. from sklearn.metrics import accuracy_score 20. from sklearn.metrics import confusion_ma- trix 21. from sklearn.metrics import f1_score 22. from sklearn.tree import DecisionTreeClas- sifier 23. from sklearn.model_selection import GridSearchCV 24. from sklearn.neighbors import KNeigh- borsClassifier 25. from sklearn.ensemble import VotingClassi- fier 26. from sklearn import svm 27. 28. data = pd.read_csv('ThoracicSurgery.csv') 29. data.head() 30. data.isna().sum() 31. data.duplicated().sum() 32. data.shape 33. data.columns 34. cats = ['diagnose','pain', 'Haemoptysis', </pre>	<pre> 35. 'short_breath', 'cough', 'weakness', 'tumour_size', 'diabetes', 36. 'breathing_difficulty', 'other_dis- ease', 'smoking', 'asthma','survive'] 37. nums = ['in_vol', 'out_vol','age'] 38. ## Descriptive Analysis 39. plt.figure(figsize=(15, 7)) 40. for i in range(0, len(nums)): 41. plt.subplot(1, 3, i+1) 42. sns.box- plot(y=data[nums[i]],color='green',ori- ent='v') 43. plt.tight_layout() 44. 45. Q1 = data['age'].quantile(0.25) 46. Q3 = data['age'].quantile(0.75) 47. IQR = Q3 - Q1 48. filter = (data['age'] >= Q1 - 1.5 * IQR) & (data['age'] <= Q3 + 1.5 * IQR) 49. data = data.loc[filter] 50. 51. Q1 = data['in_vol'].quantile(0.25) 52. Q3 = data['in_vol'].quantile(0.75) 53. IQR = Q3 - Q1 </pre>
<pre> 54. filter = (data['in_vol'] >= Q1 - 1.5 * IQR) & (data['in_vol'] <= Q3 + 1.5 * IQR) 55. data = data.loc[filter] 56. 57. Q1 = data['out_vol'].quantile(0.25) 58. Q3 = data['out_vol'].quantile(0.75) 59. IQR = Q3 - Q1 60. filter = (data['out_vol'] >= Q1 - 1.5 * IQR) & (data['out_vol'] <= Q3 + 1.5 * IQR) 61. data = data.loc[filter] 62. 63. px.histo- gram(data,x='age',nbins=40,color='survive') 64. 65. fig, ax1 = plt.subplots(1,1) 66. g = sns.countplot(x=data['survive'],ax=ax1) 67. plt.title('Deaths vs Survived',font- size=22,weight='bold') 68. plt.ylabel('Count', fontsize = 20, weight='bold') 69. g.set_xticklabels(['Died', 'Survived'], fontsize = 18) 70. plt.show() 71. </pre>	<pre> 72. Diagnose = ['Diagno- sis_type':['DGN1','DGN2','DGN3','DGN4','DGN5', 'DGN6','DGN8'], 73. 'Count':[1,47,270,40,13,3,2]} 74. df_diagnosis = pd.DataFrame(Diagnosis) 75. df_diagnosis = df_diagnosis.sort_val- ues(["Count"], ascending=True) 76. df_diagnosis 77. 78. plt.figure(figsize=(8,5)) 79. sns.barplot(x="Diagnosis_type", y="Count",data = df_diagnosis, ci = None) 80. plt.title('Distribution of Diagnosis Type ') 81. plt.xlabel('Diagnosis type', fontsize = 15) 82. plt.ylabel('count',fontsize = 15) 83. plt.show() 84. 85. pd_df=train[train['survive'] == 1] 86. pd_df['diagnosis'].value_counts() 87. b=pd_df['diagnose'].value_counts().sum() 88. Diagnosis1 = {'Diagno- sis_type':['DGN1','DGN2','DGN3','DGN4','DGN5', 'DGN6','DGN8'], 89. 'Prob':[0,10/b,30/b,5/b,6/b,0,1/b]} </pre>	<pre> 90. df1_diagnosis = pd.DataFrame(Diagnosis1) 91. df1_diagnosis = df1_diagnosis.sort_val- ues(["Prob"], ascending=True) 92. df1_diagnosis 93. 94. # Type of DGN in patients who didn't survive first year after surgery 95. plt.figure(figsize=(8,5)) 96. sns.barplot(x="Diagnosis_type", y="Prob",data = df1_diagnosis, ci = None) 97. plt.title("Type of DGN in patients who didn't survive first year after surgery") 98. plt.xlabel('Diagnosis type', fontsize = 15) 99. plt.ylabel('Prob.',fontsize = 15) 100. plt.show() 101. ## Distributions of 'FVC' and 'FEV1' 102. def plot_hist(col, bins=30, title="",xla- bel="",ax=None): 103. sns.distplot(col, bins=bins,ax=ax) 104. ax.set_title(f'Histogram of {ti- tle}',fontsize=20) 105. ax.set_xlabel(xlabel) 106. 107. fig, axes = plt.subplots(1,2,fig- size=(10,5),constrained layout=True) </pre>
<pre> 108. plot_hist(train.in_vol, 109. title='in_vol', 110. xlabel='in_vol', 111. ax=axes[0]) 112. plot_hist(train.out_vol, 113. bins=30, 114. title='out_vol', 115. xlabel='out_vol', 116. ax=axes[1]) 117. ## Dependency of diabetes in determining survival of a patient 118. pd.crosstab(train['diabetes'], train['sur- vive']).apply(lambda r: (r/r.sum())*100, axis=1) 119. rating = pd.DataFrame({ 120. 'Survived': [87.14, 73.08], 121. 'Died': [12.86, 26.92] 122. Class = ["No","Yes"] 123. ax = rating.plot(stacked=True, kind='bar', color = ['#6dbfa7', '#b5e6d7']) 124. for bar in ax.patches: 125. height = bar.get_height() 126. width = bar.get_width() 127. x = bar.get_x() 128. y = bar.get_y() </pre>	<pre> 129. label_text = str('{:.2f}').for- mat(height) + '%' 130. label_x = x + width / 2 131. label_y = y + height / 2 132. ax.text(label_x, label_y, label_text, ha='center', va='center') 133. 134. ax.set_xticklabels(Class,rotation='hori- zontal') 135. plt.ylabel('Percentage') 136. plt.xlabel('Diabetes') 137. plt.show() 138. ## Dependency of having PAD in determin- ing survival of a patient 139. pd.crosstab(train['other_disease'], train['survive']).apply(lambda r: (r/r.sum())*100, axis=1) 140. rating = pd.DataFrame({ 141. 'Survived': [86.41, 75], 142. 'Died': [13.59, 25] 143. }) 144. Class = ["No","Yes"] 145. ax = rating.plot(stacked=True, kind='bar', color = ['#6dbfa7', '#b5e6d7']) 146. </pre>	<pre> 147. for bar in ax.patches: 148. height = bar.get_height() 149. width = bar.get_width() 150. x = bar.get_x() 151. y = bar.get_y() 152. label_text = str('{:.2f}').for- mat(height) + '%' 153. label_x = x + width / 2 154. label_y = y + height / 2 155. ax.text(label_x, label_y, label_text, ha='center', va='center') 156. 157. ax.set_xticklabels(Class,rotation='hori- zontal') 158. plt.ylabel('Percentage') 159. plt.xlabel('Other Diseases') 160. plt.show() 161. ## Dependency of smoking in determining the survival of a patient. 162. pd.crosstab(train['smoking'], train['sur- vive']).apply(lambda r: (r/r.sum())*100, axis=1) 163. rating = pd.DataFrame({ 164. 'Survived': [90.41, 85.15], 165. 'Died': [9.59, 14.85] </pre>

```

166. })
167. Class = ["No", "Yes"]
168. ax = rating.plot(stacked=True, kind='bar',
169. color = ['#6dbfa7', '#b5e6d7'])
169. for bar in ax.patches:
170.     height = bar.get_height()
171.     width = bar.get_width()
172.     x = bar.get_x()
173.     y = bar.get_y()
174.     label_text = str(':.2f').for-
175.     mat(height)) + '%'
176.     label_x = x + width / 2
177.     label_y = y + height / 2
178.     ax.text(label_x, label_y, label_text,
179.     ha='center', va='center')
180. plt.ylabel('Percentage')
181. plt.xlabel('Smoking')
182. plt.show()
183. # ## Age & survival
184. g=sns.displot(data=train, x="age",
185. hue='survive', kde=True,
186. height=8, palette='viridis')

186. plt.axvline(x=80, ls='--', color='red')
187. plt.xlabel('Age')
188. plt.ylabel('Count')
189. plt.show()
190. # ## Size of the Tumour
191. train['tumour_size'].value_counts()
192. tumour = {'Tu-
193.     mour_size': ['OC13', 'OC14', 'OC11', 'OC12'],
194.     'Count': [15, 16, 141, 204]}
195. df_tumour = pd.DataFrame(tumour)
196. df_tumour = df_tumour.sort_val-
197.     ues(["Count"], ascending=True)
198. plt.figure(figsize=(10,5))
199. sns.barplot(x="Tumour size",
200. y="Count", data = df_tumour, ci = None)
201. plt.title('Distribution of Tumor size',
202.     fontsize = 20)
203. plt.xlabel('Tumour size', fontsize = 15)
204. plt.ylabel('count', fontsize = 15)
205. plt.show()
206. pd.crosstab(train['tumour_size'],
207. train['survive']).apply(lambda r:
208.     (r/r.sum())*100, axis=1)

209. rating = pd.DataFrame({
210.     'Survived': [ 90.07, 85.78, 80.00,
211.     62.50],
212.     'Died': [ 9.93, 14.22, 20.00, 37.50]
213. })
214. Class = ["OC11", "OC12", "OC13", "OC14"]
215. ax = rating.plot(stacked=True, kind='bar',
216. color = ['#6dbfa7', '#b5e6d7'])
217. for bar in ax.patches:
218.     height = bar.get_height()
219.     width = bar.get_width()
220.     x = bar.get_x()
221.     y = bar.get_y()
222.     label_text = str(':.2f').for-
223.     mat(height)) + '%'
224.     label_x = x + width / 2
225.     label_y = y + height / 2
226.     ax.text(label_x, label_y, label_text,
227.     ha='center', va='center')
228. ax.set_xticklabels(Class, rotation='hori-
229.     zontal')
230. plt.ylabel('Percentage')
231. plt.xlabel('Tumour Size')
232. plt.show()

225. cats = ['diagnose', 'pain', 'Haemoptysis',
226.     'short_breath', 'cough', 'week-
227.     ness', 'tumour_size', 'diabetes',
228.     'breathing_difficulty', 'other_dis-
229.     ease', 'smoking', 'asthma', 'survive']
230. nums = ['in_vol', 'out_vol', 'age']
231. plt.figure(figsize=(10,8))
232. plt.title('Correlation matrix - Quantita-
233.     tive variables')
234. sns.heatmap(train[nums].corr(), cbar=True,
235.     fmt='.1f', annot=True, annot_kws={'size':8},
236.     cmap="crest")
237. plt.show()
238. # from dython - all the variables
239. from dython.nominal import associations
240. associations(train.drop('id', axis=1), nom-
241.     inal_columns=cats, numerical_columns=nums,
242.     mark_columns=True,
243.     fmt='.2f', cmap="crest", figsize=(10, 10))
244. data = pd.get_dummies(data, col-
245.     umns=cats, drop_first=True)
246. rs = RobustScaler()

242. data[nums] = rs.fit_transform(data[nums])
243. # Train test split
244. X = data.drop(['survive_2', 'id'], axis=1)
245. y = data['survive_2']
246. X_train, X_test, y_train, y_test =
247.     train_test_split(X, y, test_size = 0.2, ran-
248.     dom_state = 0)
249. # Oversampling
250. sm = SMOTE()
251. X_train_sm, y_train_sm =
252.     sm.fit_resample(X_train, y_train)
253. # Advanced analysis
254. rfc = RandomForestClassifier()
255. rfc.fit(X_train, y_train)
256. y_pred = rfc.predict(X_train)
257. print('LogReg: ROC AUC = ', str(round(roc_auc_score(y_train,
258.     y_pred)*100, 1)), '%')
259. print('LogReg: Precision = ', str(round(precision_score(y_train,
260.     y_pred)*100, 1)), '%')
261. print('LogReg: Recall = ', str(round(recall_score(y_train,
262.     y_pred)*100, 1)), '%')
263. print('LogReg: Accuracy = ', str(round(accuracy_score(y_train,
264.     y_pred)*100, 1)), '%')
265. print('LogReg: F1-Score = ', str(round(f1_score(y_train,
266.     y_pred)*100, 1)), '%')
267. confusion_matrix(y_train, y_pred)
268. rfc = rfc.predict(X_test)
269. print('LogReg: ROC AUC = ', str(round(roc_auc_score(y_test,
270.     rfc.predict(X_test)*100, 1)), '%')
271. print('LogReg: Precision = ', str(round(precision_score(y_test,
272.     rfc.predict(X_test)*100, 1)), '%')
273. print('LogReg: Recall = ', str(round(recall_score(y_test,
274.     rfc.predict(X_test)*100, 1)), '%')
275. print('LogReg: Accuracy = ', str(round(accuracy_score(y_test,
276.     rfc.predict(X_test)*100, 1)), '%')
277. print('LogReg: F1-Score = ', str(round(f1_score(y_test,
278.     rfc.predict(X_test)*100, 1)), '%')

271. print('LogReg: ROC AUC = ', str(round(roc_auc_score(y_test,
272.     pred)*100, 1)), '%')
273. print('LogReg: Precision = ', str(round(precision_score(y_test,
274.     pred)*100, 1)), '%')
275. print('LogReg: Recall = ', str(round(recall_score(y_test,
276.     pred)*100, 1)), '%')
277. print('LogReg: Accuracy = ', str(round(accuracy_score(y_test,
278.     pred)*100, 1)), '%')
279. print('LogReg: F1-Score = ', str(round(f1_score(y_test,
280.     pred)*100, 1)), '%')
281. confusion_matrix(y_test, pred)
282. # ## Hyperparameter Tuning
283. rfc=RandomForestClassifier(ran-
284.     dom_state=42)
285. param_grid = {
286.     'n_estimators': [200, 500],
287.     'max_features': ['auto', 'sqrt',
288.     'log2'],
289.     'max_depth': [4, 5, 6, 7, 8],
290.     'criterion': ['gini', 'entropy']}
291. cv_rfc = GridSearchCV(estimator=rfc,
292.     param_grid=param_grid, cv=5)
293. cv_rfc.fit(X_train, y_train)
294. cv_rfc.best_params_

287. rfcl=RandomForestClassifier(ran-
288.     dom_state=42, n_estimators= 200, max_depth=4,
289.     criterion='gini')
290. rfcl.fit(X_train, y_train)
291. # Predict the train data
292. pred = rfcl.predict(X_train)
293. print('LogReg: ROC AUC = ', str(round(roc_auc_score(y_train,
294.     pred)*100, 1)), '%')
295. print('LogReg: Precision = ', str(round(precision_score(y_train,
296.     pred)*100, 1)), '%')
297. print('LogReg: Recall = ', str(round(recall_score(y_train,
298.     pred)*100, 1)), '%')
299. print('LogReg: Accuracy = ', str(round(accuracy_score(y_train,
300.     pred)*100, 1)), '%')
301. print('LogReg: F1-Score = ', str(round(f1_score(y_train,
302.     pred)*100, 1)), '%')
303. confusion_matrix(y_train, pred)
304. # Predict the test data
305. pred = rfcl.predict(X_test)
306. print('LogReg: ROC AUC = ', str(round(roc_auc_score(y_test,
307.     pred)*100, 1)), '%')
308. print('LogReg: Precision = ', str(round(precision_score(y_test,
309.     pred)*100, 1)), '%')
310. print('LogReg: Recall = ', str(round(recall_score(y_test,
311.     pred)*100, 1)), '%')
312. print('LogReg: Accuracy = ', str(round(accuracy_score(y_test,
313.     pred)*100, 1)), '%')
314. print('LogReg: F1-Score = ', str(round(f1_score(y_test,
315.     pred)*100, 1)), '%')
316. confusion_matrix(y_test, pred)

301. print('LogReg: ROC AUC = ', str(round(roc_auc_score(y_test,
302.     pred)*100, 1)), '%')
303. print('LogReg: Precision = ', str(round(precision_score(y_test,
304.     pred)*100, 1)), '%')
305. print('LogReg: Recall = ', str(round(recall_score(y_test,
306.     pred)*100, 1)), '%')
307. print('LogReg: Accuracy = ', str(round(accuracy_score(y_test,
308.     pred)*100, 1)), '%')
309. print('LogReg: F1-Score = ', str(round(f1_score(y_test,
310.     pred)*100, 1)), '%')
311. confusion_matrix(y_test, pred)
312. # ## Smote
313. rfc=RandomForestClassifier(ran-
314.     dom_state=42)
315. param_grid = {
316.     'n_estimators': [200, 500],
317.     'max_features': ['auto', 'sqrt',
318.     'log2'],
319.     'max_depth': [4, 5, 6, 7, 8],
320.     'criterion': ['gini', 'entropy']}
321. cv_rfc = GridSearchCV(estimator=rfc,
322.     param_grid=param_grid, cv=5)
323. cv_rfc.fit(X_train, y_train)
324. cv_rfc.best_params_

315. 'criterion': ['gini', 'entropy']
316. CV_rfc = GridSearchCV(estimator=rfc,
317.     param_grid=param_grid, cv=5)
318. CV_rfc.fit(X_train_sm, y_train_sm)
319. CV_rfc.best_params_
320. rfcl=RandomForestClassifier(ran-
321.     dom_state=42, n_estimators= 500, max_depth=8,
322.     criterion='entropy')
323. rfcl.fit(X_train_sm, y_train_sm)
324. # Predict the train data
325. pred = rfcl.predict(X_train)
326. print('LogReg: ROC AUC = ', str(round(roc_auc_score(y_train,
327.     pred)*100, 1)), '%')
328. print('LogReg: Precision = ', str(round(precision_score(y_train,
329.     pred)*100, 1)), '%')
330. print('LogReg: Recall = ', str(round(recall_score(y_train,
331.     pred)*100, 1)), '%')
332. print('LogReg: Accuracy = ', str(round(accuracy_score(y_train,
333.     pred)*100, 1)), '%')
334. print('LogReg: F1-Score = ', str(round(f1_score(y_train,
335.     pred)*100, 1)), '%')
336. confusion_matrix(y_train, pred)
337. # Predict the test data
338. pred = rfcl.predict(X_test)
339. print('LogReg: ROC AUC = ', str(round(roc_auc_score(y_test,
340.     pred)*100, 1)), '%')
341. print('LogReg: Precision = ', str(round(precision_score(y_test,
342.     pred)*100, 1)), '%')
343. print('LogReg: Recall = ', str(round(recall_score(y_test,
344.     pred)*100, 1)), '%')
345. print('LogReg: Accuracy = ', str(round(accuracy_score(y_test,
346.     pred)*100, 1)), '%')
347. print('LogReg: F1-Score = ', str(round(f1_score(y_test,
348.     pred)*100, 1)), '%')
349. confusion_matrix(y_test, pred)

330. print('LogReg: F1-Score = ', str(round(f1_score(y_train,
331.     pred)*100, 1)), '%')
332. print('LogReg: ROC AUC = ', str(round(roc_auc_score(y_train,
333.     pred)*100, 1)), '%')
334. print('LogReg: Precision = ', str(round(precision_score(y_train,
335.     pred)*100, 1)), '%')
336. print('LogReg: Recall = ', str(round(recall_score(y_train,
337.     pred)*100, 1)), '%')
338. print('LogReg: Accuracy = ', str(round(accuracy_score(y_train,
339.     pred)*100, 1)), '%')
340. print('LogReg: F1-Score = ', str(round(f1_score(y_train,
341.     pred)*100, 1)), '%')
342. confusion_matrix(y_train, pred)
343. # Logistic Regression
344. logreg = LogisticRegression(ran-
345.     dom_state=42, solver = 'liblinear')
346. logreg.fit(X_train, y_train)
347. # Predict the train data
348. ytrain_predicted = logreg.predict(X_train)
349. print('LogReg: ROC AUC = ', str(round(roc_auc_score(y_train,
350.     ytrain_predicted)*100, 1)), '%')
351. print('LogReg: Precision = ', str(round(precision_score(y_train,
352.     ytrain_predicted)*100, 1)), '%')
353. print('LogReg: Recall = ', str(round(recall_score(y_train,
354.     ytrain_predicted)*100, 1)), '%')
355. print('LogReg: Accuracy = ', str(round(accuracy_score(y_train,
356.     ytrain_predicted)*100, 1)), '%')
357. print('LogReg: F1-Score = ', str(round(f1_score(y_train,
358.     ytrain_predicted)*100, 1)), '%')
359. confusion_matrix(y_train, ytrain_predicted)
360. # Predict the test data
361. ytest_predicted = logreg.predict(X_test)
362. print('LogReg: ROC AUC = ', str(round(roc_auc_score(y_test,
363.     ytest_predicted)*100, 1)), '%')
364. print('LogReg: Precision = ', str(round(precision_score(y_test,
365.     ytest_predicted)*100, 1)), '%')
366. print('LogReg: Recall = ', str(round(recall_score(y_test,
367.     ytest_predicted)*100, 1)), '%')
368. print('LogReg: Accuracy = ', str(round(accuracy_score(y_test,
369.     ytest_predicted)*100, 1)), '%')
370. print('LogReg: F1-Score = ', str(round(f1_score(y_test,
371.     ytest_predicted)*100, 1)), '%')
372. confusion_matrix(y_test, ytest_predicted)

```



```

355. confusion_matrix(y_train, ytrain_pre-
dicted)
356. # Predict the test data
357. ytest_predicted = logreg.predict(X_test)
358. print('\n\nclassification report')
359. print(classification_report(y_test,
ytest_predicted)) # generate the precision,
recall, f-1 score, num
360. print('LogReg: ROC AUC =
',str(round(roc_auc_score(y_test, ytest_pre-
dicted)*100,1)), '%')
361. print('LogReg: Precision =
',str(round(precision_score(y_test,
ytest_predicted)*100,1)), '%')
362. print('LogReg: Recall = ',str(round(re-
call_score(y_test, ytest_predicted)*100,1)),
'%')
363. print('LogReg: Accuracy = ',str(round(ac-
curacy_score(y_test, ytest_pre-
dicted)*100,1)), '%')
364. print('LogReg: F1-Score =
',str(round(f1_score(y_test, ytest_pre-
dicted)*100,1)), '%')
365. confusion_matrix(y_test, ytest_predicted)
366. #Important variables

367. importances_lr = pd.DataFrame(data = {'At-
tribute': X_train.columns, 'Importance':
logreg.coef_[0]})
368. importances_lr = importances_lr.sort_val-
ues(by='Importance', ascending = False)
369. importances_lr
370. #plot
371. plt.figure(figsize=(10,5))
372. plt.bar(x=importances_lr['Attribute'],
height = importances_lr['Importance'], color
= '#087E8B')
373. plt.title('Feature Importance - Logistic
Regression')
374. plt.xticks(rotation='vertical')
375. plt.show()
376. # ### SMOTE
377. logreg = LogisticRegression(ran-
dom_state=42, solver = 'liblinear')
378. logreg.fit(X_train_sm, y_train_sm)
379. # Predict the train data
380. ytrain_predicted = logreg.predict(X_train)
381. print('\n\nclassification report')
382. print(classification_report(y_train,
ytrain_predicted)) # generate the precision,
recall, f-1 score, num

383. print('LogReg: ROC AUC =
',str(round(roc_auc_score(y_train,
ytrain_predicted)*100,1)), '%')
384. print('LogReg: Precision =
',str(round(precision_score(y_train,
ytrain_predicted)*100,1)), '%')
385. print('LogReg: Recall = ',str(round(re-
call_score(y_train, ytrain_pre-
dicted)*100,1)), '%')
386. print('LogReg: Accuracy = ',str(round(ac-
curacy_score(y_train, ytrain_pre-
dicted)*100,1)), '%')
387. print('LogReg: F1-Score =
',str(round(f1_score(y_train, ytrain_pre-
dicted)*100,1)), '%')
388. confusion_matrix(y_train, ytrain_pre-
dicted)
389. # Predict the test data
390. ytest_predicted = logreg.predict(X_test)
391. print('\n\nclassification report')
392. print(classification_report(y_test,
ytest_predicted)) # generate the precision,
recall, f-1 score, num

393. print('LogReg: ROC AUC =
',str(round(roc_auc_score(y_test, ytest_pre-
dicted)*100,1)), '%')
394. print('LogReg: Precision =
',str(round(precision_score(y_test,
ytest_predicted)*100,1)), '%')
395. print('LogReg: Recall = ',str(round(re-
call_score(y_test, ytest_predicted)*100,1)),
'%')
396. print('LogReg: Accuracy = ',str(round(ac-
curacy_score(y_test, ytest_pre-
dicted)*100,1)), '%')
397. print('LogReg: F1-Score =
',str(round(f1_score(y_test, ytest_pre-
dicted)*100,1)), '%')
398. confusion_matrix(y_test, ytest_predicted)
399. # ## Decision Tree
400. # ### Without hyperparameter tuning (over-
fitted)
401. dt = DecisionTreeClassifier(ran-
dom_state=101)
402. dt.fit(X_train_sm, y_train_sm)
403. #y predicted = dt.predict(xtest)
404. # Predict the train data
405. ytrain_predicted = dt.predict(X_train)

406. print('\n\nclassification report')
407. print(classification_report(y_train,
ytrain_predicted)) # generate the precision,
recall, f-1 score, num
408. print('DT: ROC AUC =
',str(round(roc_auc_score(y_train,
ytrain_predicted)*100,1)), '%')
409. print('DT: Precision = ',str(round(preci-
sion_score(y_train, ytrain_pre-
dicted)*100,1)), '%')
410. print('DT: Recall = ',str(round(re-
call_score(y_train, ytrain_pre-
dicted)*100,1)), '%')
411. print('DT: Accuracy = ',str(round(accu-
racy_score(y_train, ytrain_pre-
dicted)*100,1)), '%')
412. print('DT: F1-Score =
',str(round(f1_score(y_train, ytrain_pre-
dicted)*100,1)), '%')
413. confusion_matrix(y_train, ytrain_pre-
dicted)
414. ## Decision tree overfitted, so we used
Random Forest
415. # Predict the test data
416. ytest_predicted = dt.predict(X_test)

417. print('\n\nclassification report')
418. print(classification_report(y_test,
ytest_predicted)) # generate the precision,
recall, f-1 score, num
419. print('DT: ROC AUC =
',str(round(roc_auc_score(y_test, ytest_pre-
dicted)*100,1)), '%')
420. print('DT: Precision = ',str(round(preci-
sion_score(y_test, ytest_predicted)*100,1)),
'%')
421. print('DT: Recall = ',str(round(re-
call_score(y_test, ytest_predicted)*100,1)),
'%')
422. print('DT: Accuracy = ',str(round(accu-
racy_score(y_test, ytest_predicted)*100,1)),
'%')
423. print('DT: F1-Score =
',str(round(f1_score(y_test, ytest_pre-
dicted)*100,1)), '%')
424. confusion_matrix(y_test, ytest_predicted)
425. # ### Hyperparameter tuning
426. param_grid = {'cgp.alpha': [0.1, .01,
.001],
'max_depth' : [5, 6, 7, 8,
9],

428. 'criterion' :['gini', 'en-
tropy']
429. }
430. tree_clas = DecisionTreeClassifier(ran-
dom_state=1024)
431. grid_search = GridSearchCV(estima-
tor=tree_clas, param_grid=param_grid, cv=5,
verbose=True)
432. grid_search.fit(X_train_sm, y_train_sm)
433. grid_search.best_params_
434. # ## Original data
435. dt = DecisionTreeClassifier(ran-
dom_state=101, ccp_alpha= 0.001, criterion=
'gini', max_depth= 8)
436. dt.fit(X_train, y_train)
437. # Predict the train data
438. ytrain_predicted = dt.predict(X_train)
439. print('\n\nclassification report')
440. print(classification_report(y_train,
ytrain_predicted)) # generate the precision,
recall, f-1 score, num
441. print('DT: ROC AUC =
',str(round(roc_auc_score(y_train,
ytrain_predicted)*100,1)), '%')

442. print('DT: Precision = ',str(round(preci-
sion_score(y_train, ytrain_pre-
dicted)*100,1)), '%')
443. print('DT: Recall = ',str(round(re-
call_score(y_train, ytrain_pre-
dicted)*100,1)), '%')
444. print('DT: Accuracy = ',str(round(accu-
racy_score(y_train, ytrain_pre-
dicted)*100,1)), '%')
445. print('DT: F1-Score =
',str(round(f1_score(y_train, ytrain_pre-
dicted)*100,1)), '%')
446. confusion_matrix(y_train, ytrain_pre-
dicted)
447. # Predict the test data
448. ytest_predicted = dt.predict(X_test)
449. print('\n\nclassification report')
450. print(classification_report(y_test,
ytest_predicted)) # generate the precision,
recall, f-1 score, num
451. print('DT: ROC AUC =
',str(round(roc_auc_score(y_test, ytest_pre-
dicted)*100,1)), '%')

452. print('DT: Precision = ',str(round(preci-
sion_score(y_test, ytest_predicted)*100,1)),
'%')
453. print('DT: Recall = ',str(round(re-
call_score(y_test, ytest_predicted)*100,1)),
'%')
454. print('DT: Accuracy = ',str(round(accu-
racy_score(y_test, ytest_predicted)*100,1)),
'%')
455. print('DT: F1-Score =
',str(round(f1_score(y_test, ytest_pre-
dicted)*100,1)), '%')
456. confusion_matrix(y_test, ytest_predicted)
457. # ### SMOTE
458. dt = DecisionTreeClassifier(ran-
dom_state=101, ccp_alpha= 0.001, criterion=
'gini', max_depth= 8)
459. dt.fit(X_train_sm, y_train_sm)
460. # Predict the train data
461. ytrain_predicted = dt.predict(X_train)
462. print('\n\nclassification report')
463. print(classification_report(y_train,
ytrain_predicted)) # generate the precision,
recall, f-1 score, num

464. print('DT: ROC AUC =
',str(round(roc_auc_score(y_train,
ytrain_predicted)*100,1)), '%')
465. print('DT: Precision = ',str(round(preci-
sion_score(y_train, ytrain_pre-
dicted)*100,1)), '%')
466. print('DT: Recall = ',str(round(re-
call_score(y_train, ytrain_pre-
dicted)*100,1)), '%')
467. print('DT: Accuracy = ',str(round(accu-
racy_score(y_train, ytrain_pre-
dicted)*100,1)), '%')
468. print('DT: F1-Score =
',str(round(f1_score(y_train, ytrain_pre-
dicted)*100,1)), '%')
469. confusion_matrix(y_train, ytrain_pre-
dicted)
470. # Predict the test data
471. ytest_predicted = dt.predict(X_test)
472. print('\n\nclassification report')
473. print(classification_report(y_test,
ytest_predicted)) # generate the precision,
recall, f-1 score, num

474. print('DT: ROC AUC =
',str(round(roc_auc_score(y_test, ytest_pre-
dicted)*100,1)), '%')
475. print('DT: Precision = ',str(round(preci-
sion_score(y_test, ytest_predicted)*100,1)),
'%')
476. print('DT: Recall = ',str(round(re-
call_score(y_test, ytest_predicted)*100,1)),
'%')
477. print('DT: Accuracy = ',str(round(accu-
racy_score(y_test, ytest_predicted)*100,1)),
'%')
478. print('DT: F1-Score =
',str(round(f1_score(y_test, ytest_pre-
dicted)*100,1)), '%')
479. confusion_matrix(y_test, ytest_predicted)
480. knn = KNeighborsClassifier()
481. knn_model = knn.fit(X_train, y_train)
482. ypredicted = knn.predict(X_train)
483. print('\n\nclassification report')
484. print(classification_report(y_train,
ypredicted)) # generate the precision, re-
call, f-1 score, num

485. print('KNN: ROC AUC =
',str(round(roc_auc_score(y_train,
ypredicted)*100,1)), '%')
486. print('KNN: Precision = ',str(round(preci-
sion_score(y_train, ypredicted)*100,1)),
'%')
487. print('KNN: Recall = ',str(round(re-
call_score(y_train, ypredicted)*100,1)),
'%')
488. print('KNN: Accuracy = ',str(round(accu-
racy_score(y_train, ypredicted)*100,1)),
'%')
489. print('KNN: F1-Score =
',str(round(f1_score(y_train,
ypredicted)*100,1)), '%')
490. confusion_matrix(y_train, ypredicted)
491. ypredicted = knn.predict(X_test)
492. print('\n\nclassification report')
493. print(classification_report(y_test,
ypredicted)) # generate the precision, re-
call, f-1 score, num
494. print('KNN: ROC AUC =
',str(round(roc_auc_score(y_test,
ypredicted)*100,1)), '%')

```

```

495. print('KNN: Precision = ',str(round(precision_score(y_test, ypredicted)*100,1)), '%')
496. print('KNN: Recall = ',str(round(recall_score(y_test, ypredicted)*100,1)), '%')
497. print('KNN: Accuracy = ',str(round(accuracy_score(y_test, ypredicted)*100,1)), '%')
498. print('KNN: F1-Score = ',str(round(f1_score(y_test, ypredicted)*100,1)), '%')
499. confusion_matrix(y_test, ypredicted)
500. # ### SMOTE
501. knn = KNeighborsClassifier()
502. knn_model = knn.fit(X_train_sm, y_train_sm)
503. ypredicted = knn.predict(X_train)
504. print('\nclassification report')
505. print(classification_report(y_train, ypredicted)) # generate the precision, recall, f-1 score, num
506. print('KNN: ROC AUC = ',str(round(roc_auc_score(y_train, ypredicted)*100,1)), '%')
507. print('KNN: Precision = ',str(round(precision_score(y_train, ypredicted)*100,1)), '%')
508. print('KNN: Recall = ',str(round(recall_score(y_train, ypredicted)*100,1)), '%')
509. print('KNN: Accuracy = ',str(round(accuracy_score(y_train, ypredicted)*100,1)), '%')
510. print('KNN: F1-Score = ',str(round(f1_score(y_train, ypredicted)*100,1)), '%')
511. confusion_matrix(y_train, ypredicted)
512. ypredicted = knn.predict(X_test)
513. print('\nclassification report')
514. print(classification_report(y_test, ypredicted)) # generate the precision, recall, f-1 score, num
515. print('KNN: ROC AUC = ',str(round(roc_auc_score(y_test, ypredicted)*100,1)), '%')
516. print('KNN: Precision = ',str(round(precision_score(y_test, ypredicted)*100,1)), '%')
517. print('KNN: Recall = ',str(round(recall_score(y_test, ypredicted)*100,1)), '%')
518. print('KNN: Accuracy = ',str(round(accuracy_score(y_test, ypredicted)*100,1)), '%')
519. print('KNN: F1-Score = ',str(round(f1_score(y_test, ypredicted)*100,1)), '%')
520. confusion_matrix(y_test, ypredicted)
521. # ### Voting Classifier
522. # ### Original Data
523. voting_clf = VotingClassifier(
524. estimators=[('lr', logreg), ('rf', rfc1), ('knn', knn), ('DT', dt)],
525. voting='hard')
526. voting_clf.fit(X_train, y_train)
527. from sklearn.metrics import accuracy_score
528. for clf in (knn, voting_clf):
529.     clf.fit(X_train, y_train)
530.     y_pred = clf.predict(X_test)
531.     print(clf.__class__.__name__, accuracy_score(y_test, y_pred))
532.     print(clf.__class__.__name__, precision_score(y_test, y_pred))
533.     print(clf.__class__.__name__, recall_score(y_test, y_pred))
534.     print(clf.__class__.__name__, f1_score(y_test, y_pred))
535. ypredicted = voting_clf.predict(X_train)
536. print('\nclassification report')
537. print(classification_report(y_train, ypredicted)) # generate the precision, recall, f-1 score, num
538. print('VC: ROC AUC = ',str(round(roc_auc_score(y_train, ypredicted)*100,1)), '%')
539. print('VC: Precision = ',str(round(precision_score(y_train, ypredicted)*100,1)), '%')
540. print('VC: Recall = ',str(round(recall_score(y_train, ypredicted)*100,1)), '%')
541. print('VC: Accuracy = ',str(round(accuracy_score(y_train, ypredicted)*100,1)), '%')
542. print('VC: F1-Score = ',str(round(f1_score(y_train, ypredicted)*100,1)), '%')
543. confusion_matrix(y_train, ypredicted)
544. ypredicted = voting_clf.predict(X_test)
545. print('\nclassification report')
546. print(classification_report(y_test, ypredicted)) # generate the precision, recall, f-1 score, num
547. print('VC: ROC AUC = ',str(round(roc_auc_score(y_test, ypredicted)*100,1)), '%')
548. print('VC: Precision = ',str(round(precision_score(y_test, ypredicted)*100,1)), '%')
549. print('VC: Recall = ',str(round(recall_score(y_test, ypredicted)*100,1)), '%')
550. print('VC: Accuracy = ',str(round(accuracy_score(y_test, ypredicted)*100,1)), '%')
551. print('VC: F1-Score = ',str(round(f1_score(y_test, ypredicted)*100,1)), '%')
552. confusion_matrix(y_test, ypredicted)
553. # ### SMOTE
554. voting_clf = VotingClassifier(
555. estimators=[('lr', logreg), ('rf', rfc1), ('knn', knn), ('DT', dt)],
556. voting='hard')
557. voting_clf.fit(X_train_sm, y_train_sm)
558. for clf in (knn, voting_clf):
559.     clf.fit(X_train_sm, y_train_sm)
560.     y_pred = clf.predict(X_test)
561.     print(clf.__class__.__name__, accuracy_score(y_test, y_pred))
562.     print(clf.__class__.__name__, precision_score(y_test, y_pred))
563.     print(clf.__class__.__name__, recall_score(y_test, y_pred))
564.     print(clf.__class__.__name__, f1_score(y_test, y_pred))
565. ypredicted = voting_clf.predict(X_train)
566. print('\nclassification report')
567. print(classification_report(y_train, ypredicted)) # generate the precision, recall, f-1 score, num
568. print('VC: ROC AUC = ',str(round(roc_auc_score(y_train, ypredicted)*100,1)), '%')
569. print('VC: Precision = ',str(round(precision_score(y_train, ypredicted)*100,1)), '%')
570. print('VC: Recall = ',str(round(recall_score(y_train, ypredicted)*100,1)), '%')
571. print('VC: Accuracy = ',str(round(accuracy_score(y_train, ypredicted)*100,1)), '%')
572. print('VC: F1-Score = ',str(round(f1_score(y_train, ypredicted)*100,1)), '%')
573. confusion_matrix(y_train, ypredicted)
574. ypredicted = voting_clf.predict(X_test)
575. print('\nclassification report')
576. print(classification_report(y_test, ypredicted)) # generate the precision, recall, f-1 score, num
577. print('VC: ROC AUC = ',str(round(roc_auc_score(y_test, ypredicted)*100,1)), '%')
578. print('VC: Precision = ',str(round(precision_score(y_test, ypredicted)*100,1)), '%')
579. print('VC: Recall = ',str(round(recall_score(y_test, ypredicted)*100,1)), '%')
580. print('VC: Accuracy = ',str(round(accuracy_score(y_test, ypredicted)*100,1)), '%')
581. print('VC: F1-Score = ',str(round(f1_score(y_test, ypredicted)*100,1)), '%')
582. confusion_matrix(y_test, ypredicted)
583. # ### SVM
584. # ### Original Data
585. clf = svm.SVC(kernel='linear', C = 1.0)
586. clf.fit(X_train, y_train)
587. ypredicted = voting_clf.predict(X_train)
588. print('\nclassification report')
589. print(classification_report(y_train, ypredicted)) # generate the precision, recall, f-1 score, num
590. print('SVM: ROC AUC = ',str(round(roc_auc_score(y_train, ypredicted)*100,1)), '%')
591. print('SVM: Precision = ',str(round(precision_score(y_train, ypredicted)*100,1)), '%')
592. print('SVM: Recall = ',str(round(recall_score(y_train, ypredicted)*100,1)), '%')
593. print('SVM: Accuracy = ',str(round(accuracy_score(y_train, ypredicted)*100,1)), '%')
594. print('SVM: F1-Score = ',str(round(f1_score(y_train, ypredicted)*100,1)), '%')
595. confusion_matrix(y_train, ypredicted)
596. ypredicted = clf.predict(X_test)
597. print('\nclassification report')
598. print(classification_report(y_test, ypredicted)) # generate the precision, recall, f-1 score, num
599. print('SVM: ROC AUC = ',str(round(roc_auc_score(y_test, ypredicted)*100,1)), '%')
600. print('SVM: Precision = ',str(round(precision_score(y_test, ypredicted)*100,1)), '%')
601. print('SVM: Recall = ',str(round(recall_score(y_test, ypredicted)*100,1)), '%')
602. print('SVM: Accuracy = ',str(round(accuracy_score(y_test, ypredicted)*100,1)), '%')
603. print('SVM: F1-Score = ',str(round(f1_score(y_test, ypredicted)*100,1)), '%')
604. confusion_matrix(y_test, ypredicted)
605. # ### SMOTE
606. clf = svm.SVC(kernel='linear', C = 1.0)
607. clf.fit(X_train_sm, y_train_sm)
608. ypredicted = voting_clf.predict(X_train)
609. print('\nclassification report')
610. print(classification_report(y_train, ypredicted)) # generate the precision, recall, f-1 score, num

```

```

611. print('SVM: ROC AUC =
    ',str(round(roc_auc_score(y_train,
ypredicted)*100,1)), '%')
612. print('SVM: Precision = ',str(round(precision_score(y_train, ypredicted)*100,1)),
    '%')
613. print('SVM: Recall = ',str(round(recall_score(y_train, ypredicted)*100,1)),
    '%')
614. print('SVM: Accuracy = ',str(round(accuracy_score(y_train, ypredicted)*100,1)),
    '%')
615. print('SVM: F1-Score =
    ',str(round(f1_score(y_train,
ypredicted)*100,1)), '%')
616. confusion_matrix(y_train, ypredicted)
617. ypredicted = clf.predict(X_test)
618. print('\nclassification report')
619. print(classification_report(y_test,
ypredicted)) # generate the precision, recall, f-1 score, num
620. print('SVM: ROC AUC =
    ',str(round(roc_auc_score(y_test,
ypredicted)*100,1)), '%')

621. print('SVM: Precision = ',str(round(precision_score(y_test, ypredicted)*100,1)), '%')
622. print('SVM: Recall = ',str(round(recall_score(y_test, ypredicted)*100,1)), '%')
623. print('SVM: Accuracy = ',str(round(accuracy_score(y_test, ypredicted)*100,1)), '%')
624. print('SVM: F1-Score =
    ',str(round(f1_score(y_test,
ypredicted)*100,1)), '%')
625. confusion_matrix(y_test, ypredicted)

```