

## **Predicting Length of Stay at Hospitals using MIMIC data**

### **Group 11:**

**Jaya Sindhu Yannam**

**Shikhar Bharat Jain**

**Priyank Nilesh Gandhi**

**Siddharth Nandargi**

### **Member contribution statement:**

*Shikhar, Jaya, Priyank and Siddharth conceived of the presented idea. Jaya conducted the initial Literature Survey of existing work related to the chosen topic. Priyank and Siddharth performed an analysis of the existing data mining methods and implemented in Python. Shikhar performed the Regression algorithm and concluded that the model performance is not so great. All team members mutually agreed to proceed with implementing classification algorithms to improve the prediction accuracy. Shikhar implemented SMOTE to balance the imbalance dataset and proceeded to implement classification algorithms. Jaya and Priyank performed the EDA analysis and created visualizations which helped to understand the important features/attributes. Shikhar and Siddharth worked on implementing the Recursive Feature Selection methodology. Jaya, Priyank and Siddharth implemented 3 classification algorithm on the sampled dataset in Weka software. Shikhar implemented 2 classification algorithms using Python. All members discussed the results and contributed to the final report development.*

## I. Introduction

The length of stay (LOS) is the period a patient is admitted to a hospital or other similar medical facility. Moreover, LOS is correlated with disease severity and mortality. Our goal is to provide an initial length of stay prediction of new coming patients using a numerical and categorical database as a starting point. By using data mining algorithms, this prediction is made achievable. The data mining techniques are employed to make sure the procedure is reliable and lessen the effect of certainty. Moreover, studies have found patients that, at an emergency department (ED) are associated with a longer LOS and patients who develop further complications in intensive care units (ICU) have a longer LOS beforehand at the ED. To optimize healthcare planning and resources while minimizing the impact of uncertainty, the LOS must be predicted.

One of the most crucial medical resources is hospital beds, which are frequently regarded as an important indicator of hospital performance and an accurate predictor of the state of the local healthcare system's development. Patients' length of stays and hospital costs are strongly related because most hospitals have a limited supply of beds. Therefore, cutting the length of stay can save medical costs and the social burden of medical care while simultaneously increasing the turnover of inpatients. Finding the pertinent risk variables associated with patient recovery and length of stay is crucial for the medical system. As a result, one of the biggest issues facing hospital management right now is how to better allocate hospital beds and address the bed shortage.

The issue of bed allocation has not been fully solved, and there are still a lot of hospitals that need beds due to a lack of medical resources in hospitals and a number of unpredictable events during treatment. Therefore, estimating the length of stay precisely will aid in logical bed allocation and boost bed usage. An important measure of hospital management is how long patients stay there. Its prediction calls for the use of statistical techniques to summarize, analyze, and research its distribution law and change rule as well as the use of machine learning algorithms to create models that forecast the length of hospital stay. Studies contributing to LOS have regularly appeared in the literature. One study conducted determine the factors affecting LOS in the public hospital in Iran. Demonstrated that, an increase in the age would lead to increase in average LOS, Secondly, the average LOS of men is longer than women. Data mining algorithms such as Artificial Neural Network (ANN), Support Vector Machine (SVM), and Decision trees were used for this study.

## **II. Methods**

### **a. Data source**

A crucial component of this project is to be able to analyze the large amount of data contained in Electronic Medical Records (EMRs). Our dataset of choice was the MIMIC III database which houses over 50,000 actual records of people that have been admitted to ICU units between 2001 and 2012. As one would expect the MIMIC database comprises of anonymized EMRs with the data stored in a relational database format. The Data tables that we have used for our analysis and predictions are as follows:

- A) **PATIENTS**: This table has 42,520 number of instances, and links to ADMISSIONS table as well as ICUSTAYS on “Subject\_ID”
- B) **ICUSTAYS**: ICUSTAYS tables contains 61,532 rows in it and has links to PATIENTS table on “Subject\_ID”, and ADMISSIONS table on “HADM\_ID”.
- C) **ADMISSIONS**: The number of rows that the ADMISSIONS table contains is 58,976. And this table has links to PATIENTS table on “Subject\_ID” .
- D) **DIAGNOSES\_ICD**: This table has 651,047 rows, and it has links to PATIENTS on “Subject\_ID”, ADMISSIONS on “HADM\_ID”.

### **b. Data preprocessing:**

Data preprocessing involves various steps that need to be followed to prepare the data appropriately for the model building and training. These steps depend on the type of dataset we are working on, and this varies from data to data.

Firstly, we checked the data for any null values present. Upon checking that, we found that 3% of the data is missing so we removed those instances from the dataset. The length of stay variable (LOS) represented the number of days a patient stayed in the hospital. Upon examining this variable, it was found that few of the data values were negative, and the number of days cannot be negative, so we removed such values from the data tables, as a part of data cleaning process.

The admission type variable had 4 types, but it seems that “Emergency” and “Urgent” are one and the same, so the “Urgent” values were replaced with “Emergency” and the value counts were normalized.

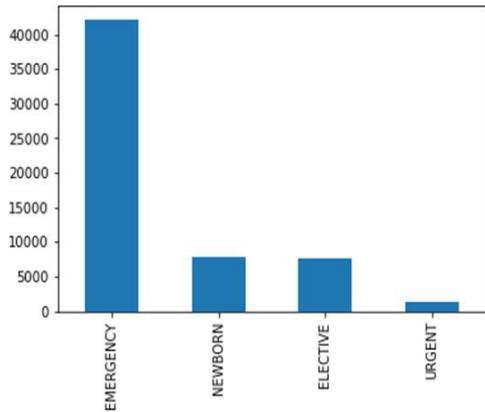


Fig. 1 Frequency of Admission type

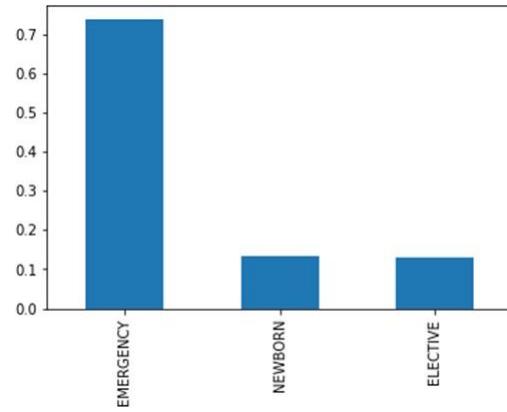


Fig. 2 After converting  
urgent category to emergency

As LOS is the target variable, it was important to check the distribution of this variable across the dataset. A histogram was generated, along with a frequency curve.

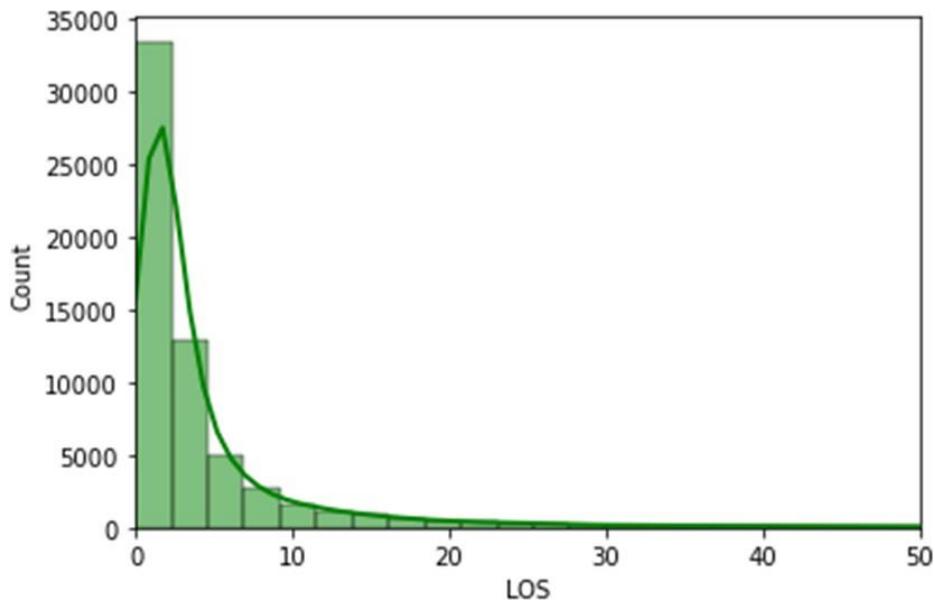


Fig. 3 Distribution of LOS variable

The above histogram suggests that the data values in LOS variable are not normally distributed, and the distribution is heavily skewed towards the right. Hence, we can conclude from this that most of the patients admitted had their length of stay in the hospital to be less than 10 days.

There was no column containing the ages of patients. So, we created a new column containing Ages of patient by using data values from 2 other columns like DOB from PATIENTS Table and ADMITTIME from ADMISSIONS table

A boxplot for Age variable was also generated, and a table containing descriptive analysis of Age was prepared.

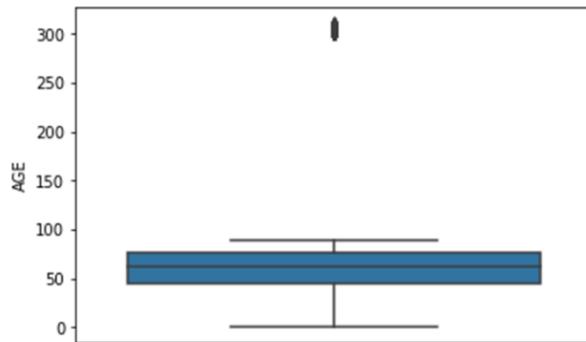


Fig. 4 Boxplot for Age variable

count	61532.000000
mean	64.925951
std	56.951788
min	0.000000
25%	44.374658
50%	62.112329
75%	76.101370
max	311.767123

Fig. 5 Descriptive analysis of Age

This boxplot helped in detecting the outliers in the age variable. And on examining these outliers, it was found that these outliers were irrelevant for the analysis and prediction but were present in a very small negligible amount across the data.

For regression analysis, min-max scaler was used to scale the data and for classification technique, we binarized the response variable where 0 represents LOS=< 10 and 1 represents LOS>10 days (about 1 and a half weeks).

Distribution of the death variable based on the genders was generated to check what gender type had most deaths and vice versa.

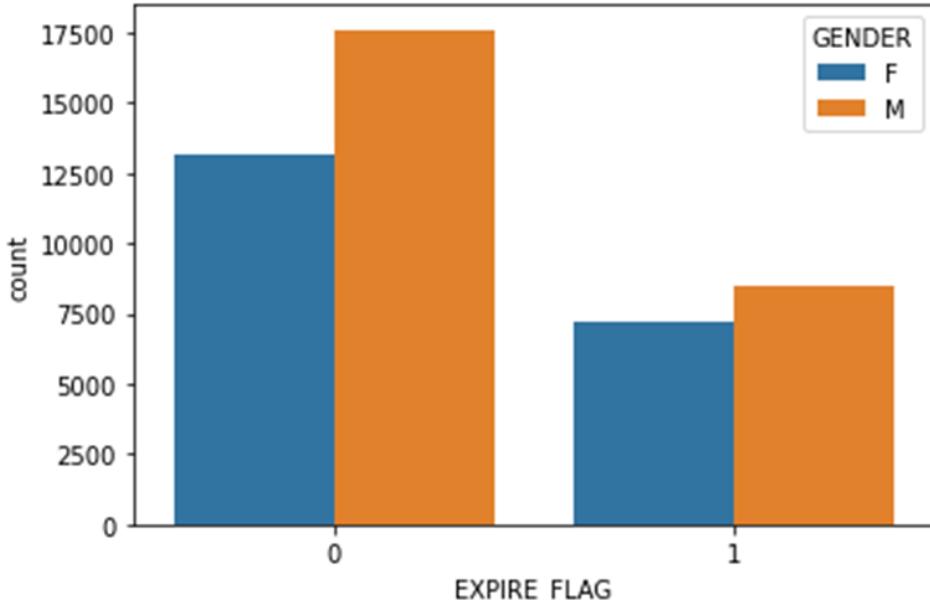


Fig. 6 Distribution of death variable based on gender

Here it can be seen that Males were the most who died and the same can be seen for alive. This indicates that the number of females admitted were less than that of male patients.

Later, to increase the interpretability for each ICD-9 code we mentioned the corresponding diseases and plotted it using a bar graph.

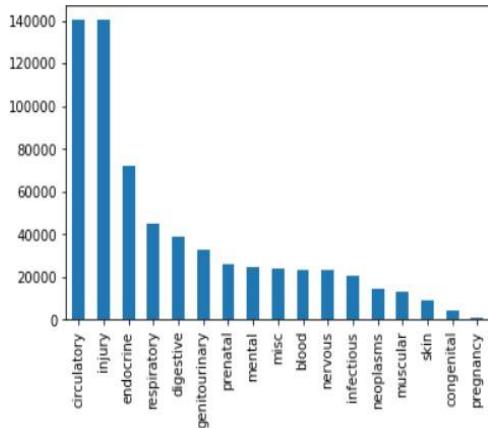


Fig. 7 ICD-9 Codes and diseases

	ROW_ID	SUBJECT_ID	HADM_ID	SEQ_NUM	ICD9_CODE	codetype
0	1297	109	172335	1.0	40301	circulatory
1	1298	109	172335	2.0	486	respiratory
2	1299	109	172335	3.0	58281	genitourinary
3	1300	109	172335	4.0	5855	genitourinary
4	1301	109	172335	5.0	4254	circulatory

Fig.8 ICD-9 Code Diseases

From the bar graph, it is evident that circulatory diseases are more prominent in ICU patients.

### **c. Feature selection process**

There are various techniques that we can use for selecting features from the dataset to get accurate predictions after training the model. There needs to maintain a balance between model complexity and model simplicity.

Too many features make a model complex, while less features make the model very simple, which would hinder the results. So, selecting appropriate number of features is very crucial.

For this we have used “Recursive Feature Elimination” technique.

In this, the least important features are eliminated at each feature elimination step, for this feature importance metrics are used, to select the least important features for model building.

In general, feature selection techniques are used to get rid of unnecessary and redundant attributes that don't contribute much to the predictive model's accuracy or, in certain cases, even make it less accurate. The accuracy of the prediction model may increase because of feature selection, but it will also be faster and more cost-effective because the model will be less complex.

We chose all the variables such as Gender\_M, nervous, neoplasams, misc, pregnancy, prenatal, respiratory, skin, age, muscular, mental, Gender\_F, circulatory, blood, injury, congenital, digestive, endocrine, genitourinary, infectious, ICU, NICU etc. Removing all death related columns because we are finding length of stay and these patients were never discharged.

We also used Weka to see what attributes Weka suggests based on the Attribute Rankings. Classifier Feature Evaluator was used by using Attribute Ranking.

Attribute selection output	
<b>Ranked attributes:</b>	
0	42 GENDER_M
0	14 nervous
0	13 neoplasms
0	11 misc
0	15 pregnancy
0	16 prenatal
0	17 respiratory
0	18 skin
0	19 age
0	12 muscular
0	10 mental
0	41 GENDER_F
0	3 circulatory
0	2 blood
0	9 injury
0	4 congenital
0	5 digestive
0	6 endocrine
0	7 genitourinary
0	8 infectious
0	20 ICU
0	21 NICU

Fig.9 Attribute Selection in Weka

```

0 22 ADMISSION_TYPE_ELECTIVE
0 35 ADMISSION_LOCATION_referral from other facilities
0 34 ADMISSION_LOCATION_EMERGENCY ROOM ADMIT
0 23 ADMISSION_TYPE_EMERGENCY
0 36 ADMISSION_LOCATION_transfer from other facilities
0 37 MARITAL_STATUS_DIVORCED
0 38 MARITAL_STATUS_MARRIED
0 39 MARITAL_STATUS_SINGLE
0 40 MARITAL_STATUS_WIDOWED
0 33 ADMISSION_LOCATION_** INFO NOT AVAILABLE **
0 32 ETHNICITY_WHITE
0 31 ETHNICITY_Unknown/other
0 26 INSURANCE_Private
0 24 ADMISSION_TYPE_NEWBORN
0 25 INSURANCE_Government
0 27 INSURANCE_Self Pay
0 30 ETHNICITY_HISPANIC/LATINO
0 28 ETHNICITY_Asian
0 29 ETHNICITY_BLACK/AFRICAN AMERICAN
0 1 ID

```

Fig.10 Attribute Selection in Weka

#### d. Description of dependent and independent variables:

A dependent variable here is Length of Stay.i.e., **LOS**. We are predicting the LOS variable of a patient in ICU. The histogram that we generated (Fig.) for LOS depicted that the LOS variable is not normally distributed and has most of its data values under 10 days (about 1 and a half weeks).

For classification problem, we have converted this variable into a binary variable, where 1 means LOS>10 days, and 0 means LOS<=10 days.

There are numerous independent variables present in the dataset which help in predicting the data for LOS variable. These variables differ based on the type of problem. For classification problems there are a different set of variables than in the regression problem.

Various dummy variables were also generated in the process, to help for better model building and prediction.

	blood	circulatory	congenital	digestive	endocrine	genitourinary	infectious	injury	mental	misc	muscular	neoplasms	nervous	pregnancy	prenatal	resp
0	0	2	0	2	5	2	0	2	0	0	0	0	0	2	0	0
1	1	2	0	4	0	0	1	0	0	1	0	0	0	0	0	0
2	0	0	0	0	1	0	0	2	1	1	0	1	0	0	0	0
3	0	1	0	2	0	0	0	1	0	0	0	0	0	0	0	0
4	1	7	0	0	3	0	0	7	0	0	0	0	0	0	0	0

Fig.11 Creation of Dummy Variables

```
diag_dummy = diag_dummy.join(diag['HADM_ID'], how="outer")
diag_dummy.head()
```

	blood	circulatory	congenital	digestive	endocrine	genitourinary	infectious	injury	mental	misc	muscular	neoplasms	nervous	pregnancy	prenatal	res
0	0	2	0	2	5	2	0	2	0	0	0	0	0	2	0	0
1	1	2	0	4	0	0	1	0	0	1	0	0	0	0	0	0
2	0	0	0	0	1	0	0	2	1	1	0	1	0	0	0	0
3	0	1	0	2	0	0	0	1	0	0	0	0	0	0	0	0
4	1	7	0	0	3	0	0	7	0	0	0	0	0	0	0	0

Fig.12 Creation of Dummy Variables

Descriptive statistics of all the variables is as follows.

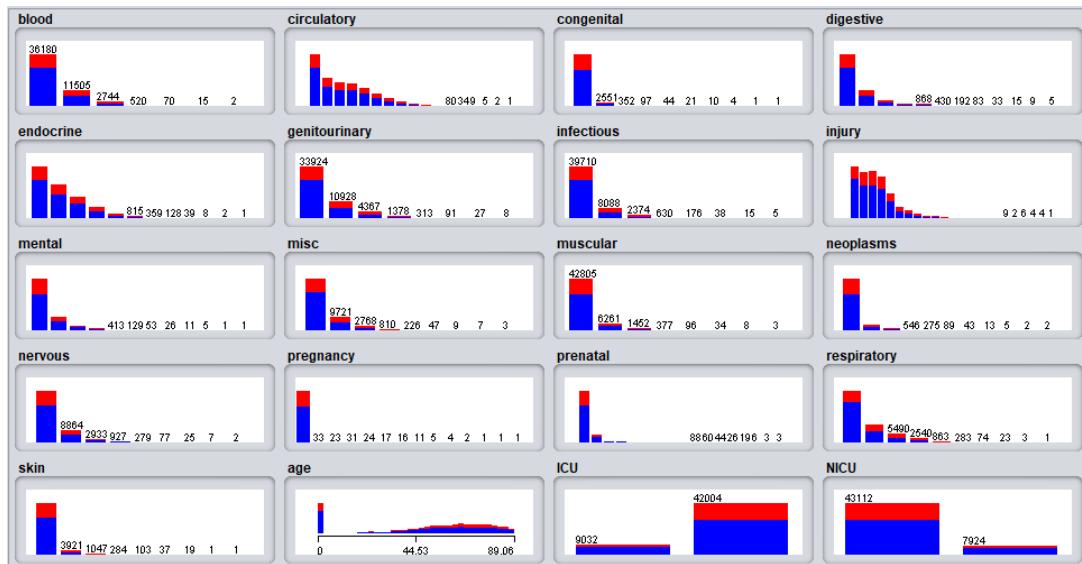


Fig.13 Statistical analysis for Independent Variables

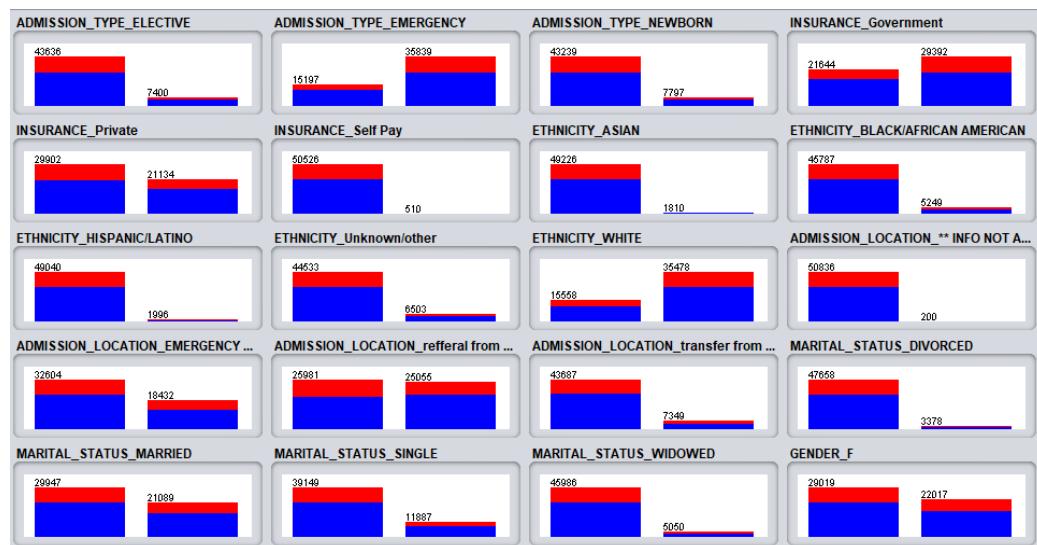


Fig.14 Statistical analysis for Independent Variables

### **III. Results and Discussion**

Random Forest is a versatile machine learning method capable of performing both regression and classification model. The random forest starts with a standard machine learning technique called “decision trees”. In a decision tree, an input is entered at the top and as it traverses down the tree the data gets put into smaller and smaller sets or classes. The random forest goes a step beyond by combining multiple trees to produce the mode class, in the case of classification, or, in the case of regression, a mean prediction by averaging the results of the individual trees. In addition, random forests correct for decision trees’ habit of over fitting to their training set. When a new input is entered into the system, it is run down all the trees. In the case of a numerical variable, the result may either be an average or weighted average of all the terminal nodes that are reached. However, if the input variable is categorical, then the result will be a voting majority.

Initially we implemented the Regression technique on our dataset after splitting it into 80% test and 20% train data, and the following results were obtained. Here, we have used the Decision Tree and Random Forest models for the regression tasks. Decision Trees and Random Forest models can be used for both – Classification and Regression. The metrics used here are R2 (R-square) and MSE (Mean Square Error) rate. R square explains how well the regression model explains the observed data which is also known as goodness of fit measure for linear regression models. R-square value measures the strength of relationship between your model and the dependent variable on a convenient scale of 0-100% or 0 to 1. We note the our R-square values obtained for the two models implemented yield a value close to 0.5 which is considered to be a good fit. Higher the R-square value does not necessarily mean that the model fits well, but it means the model is overfitting. It is also noted that the differences between the observed and the predicted variables are small and unbiased.

Test Set

Model	R2	MSE
Decision Tree	0.3982	0.0011
Random Forest	0.4070	0.0010

Table 2

Mean Square Error (MSE) signifies the average square residual. MSE represents the error of the predicted model created based on the given set of observations in the sample. To put it in other words, MSE shows the cost incurred while making each prediction. There is no correct value for MSE. Simply put, the lower the value the better and 0 means the model is perfect. Since there is no correct answer, the MSE's basic value is in selecting one prediction model over another. MSE is a metric to know how close you got to the perfect fit. The model is often fitted by minimizing MSE, so you'd like to know how close you got to ZERO in the end. You can use this metric later to compare models, those with smaller MSE fit better to the data.

Next, we proceeded onto implementing Classification algorithms using *sklearn* library in Python where we considered the same split of dataset as we did while implementing Regression techniques. We implemented the Decision Tree and Logistic Regression as a classifier and obtained the following results.

Model	Accuracy	Sensitivity	Specificity
Decision Tree	0.7442	0.7157	0.7728
Logistic Regression	0.6402	0.6768	0.6036

Model	Accuracy	Sensitivity	Specificity
Decision Tree	0.7200	0.9405	0.7760
Logistic Regression	0.5643	0.4777	0.6033

Fig 15. Accuracy Metrics

We achieved a F1score of 0.793 for Decision tree classifier algorithm while for Logistic Regression we achieved 0.394.

We considered Accuracy, Sensitivity, F1 Score and Specificity as metrics to evaluate the classification models. Our sensitivity describes how well our test catches all our positive cases. Sensitivity is calculated by dividing the number of true-positive results by the total number of positives (which include false positives).

Our specificity describes how well our test classifies negative cases as negatives. Specificity is calculated by dividing the number of true-negative results by the total number of negatives (which include false negatives).

Our study involves implementing 3 other classification models for the sake of comparison through which we could infer the best model that suits our dataset. For this, we used the Weka software and the results for the classification algorithms implemented are as follows. It is to be noted that the Weka software could not process algorithms on the complete dataset involved in our study, hence we have performed a random sampling in Weka, and considered 15% of data for implementing these algorithms. After sampling, the total instance count was 7655.

The three classification algorithms implemented in Weka were:

1. Naïve Bayes - every pair of features being classified is independent of each other
2. J48 - J48 is a machine learning decision tree classification algorithm based on Iterative Dichotomiser 3. It is very helpful in examine the data categorically and continuously
3. A random forest is a meta estimator that fits several decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting

Classification Algorithm on Test Set				
Model	Accuracy	AUC	F measure	Software
Naïve Bayes	72%	0.757	0.721	Weka
J48	71%	0.660	0.710	Weka
Random Forest	74%	0.779	0.734	Weka

Apart from Accuracy as a metric, we have considered AUC which stands for Area under the ROC curve. AUC is used for model comparison. AUC value of 1 infers a perfect model which might be a true case of overfitting and a value under 0.5 is not so favored. Our results have AUC values greater than 0.5 and less than 1 which concludes a good fitting model. Amongst the three models,

Random Forest performs the best basis on accuracy and AUC values. Next in order is Naïve Bayes and the least performing model is J48.

### **Strengths and limitations of your study**

1. Mimic dataset does not provide any data for patients aged between 7 to 13 years.
2. Weka software could not handle large datasets and we sampled the data to perform Classification for few of the approaches.
3. MIMIC contains a vast data set and the tables in it. A careful study of all the tables and variables needs to be conducted to join the tables to form the required dataset to build the model. Extracting and joining tables is a tedious process.
4. Classification models perform better than the Regression models which has been proved with the metrics chosen while implementing models.
5. Our data does not include any bias as we did not focus on a particular race/ethnicity, nor patients with certain diseases, nor based on the resources of the hospitals, etc.
6. We followed the best Feature selection methodology as we note the models have given good accuracy and specificity values, which confirms our model fits the data best and does not lead to overfitting.
7. Initially our dataset was imbalanced which would yield to biased and inaccurate models. To improve this, we performed SMOTE (Synthetic Minority Over-Sampling technique) and balanced our dataset to ensure the models are accurate which was done by upsampling the minority class.

## **IV. Conclusion**

The scope of our project showcases how Classification algorithms improves the drawbacks of Regression on the MIMIC dataset. We can infer from the scores of Specificity and Sensitivity obtained from the Classification and Regression models in Python that the models we built are good for predicting the LOS. Predicting LOS is an example of how EHR data can be potentially used for good by applying data mining techniques. The MIMIC database offered surprisingly good depth and detail related to medical admissions which enabled me to create a hospital length-of-stay prediction model that considered a lot of interesting input features. The most surprising aspect of this work was how the patient ICD-9 diagnoses played a more important role than age when predicting the length-of-stay.

### **Future work:**

This way of predicting Length of Stay outcome in patients admitted in hospitals can be used effectively. We used random forest models to predict LOS at the ED since this model is well suited for treating data with complex interactions between variables and other non-linear effects. Models based on Deep Neural Networks are another option that could be explored in further studies.

## References:

- 1) Nouaouri, A. Samet and H. Allaoui, "Evidential data mining for length of stay (LOS) prediction problem," 2015 IEEE International Conference on Automation Science and Engineering (CASE), 2015, pp. 1415-1420, doi: 10.1109/CoASE.2015.7294296.
- 2) Hachesu PR, Ahmadi M, Alizadeh S, Sadoughi F. Use of data mining techniques to determine and predict length of stay of cardiac patients. *Healthc Inform Res*. 2013 Jun;19(2):121-9.doi: 10.4258/hir.2013.19.2.121.Epub 2013 Jun 30. PMID: 23882417; PMCID: PMC3717435
- 3) Rathor, R., & Agarkar, P. (2015). LOSH Prediction using Data Mining. *International Journal of Computer Applications*, 119.
- 4) Johnson AE;Pollard TJ;Shen L;Lehman LW;Feng M;Ghassemi M;Moody B;Szolovits P;Celi LA;Mark RG; (n.d.). Mimic-III, a freely accessible Critical Care Database. Scientific data. Retrieved November 29, 2022, from <https://pubmed.ncbi.nlm.nih.gov/27219127/>
- 5) Yong Chen, "Prediction and Analysis of Length of Stay Based on Nonlinear Weighted XGBoost Algorithm in Hospital", *Journal of Healthcare Engineering*, vol. 2021, Article ID 4714898, 9 pages, 2021. <https://doi.org/10.1155/2021/4714898>
- 6) Chrusciel, J., Girardon, F., Roquette, L. et al. The prediction of hospital length of stay using unstructured data. *BMC Med Inform Decis Mak* 21, 351 (2021). <https://doi.org/10.1186/s12911-021-01722-4>

## Appendix: SQL/STATA/Python/R codes, WEKA screen shots, etc.

**Classifier**

Choose **RandomForest** -P 100 -I 100 -num-slots 1 -K 0 -M 1.0 -V 0.001 -S 1

Test options

- Use training set
- Supplied test set
- Cross-validation Folds 10
- Percentage split % 80

[More options...](#)

(Nom) class\_los\_label

Start Stop

Result list (right-click for options)

- 16:53:20 - bayes.NaiveBayes
- 17:24:37 - trees.RandomForest
- 17:28:41 - trees.J48

Classifier output

==== Evaluation on test split ===

Time taken to test model on test split: 0.68 seconds

==== Summary ===

	Correctly Classified Instances	1139	74.3958 %
Incorrectly Classified Instances	392	25.6042 %	
Kappa statistic	0.3404		
Mean absolute error	0.3415		
Root mean squared error	0.4124		
Relative absolute error	78.9737 %		
Root relative squared error	88.5484 %		
Total Number of Instances	1531		

==== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
0.904	0.600	0.764	0.904	0.828	0.360	0.779	0.880	0	
0.400	0.096	0.661	0.400	0.499	0.360	0.779	0.629	1	
Weighted Avg.	0.744	0.439	0.731	0.744	0.723	0.360	0.779	0.800	

==== Confusion Matrix ===

a	b	<-- classified as
944	100	a = 0
292	195	b = 1

**Weka Explorer**

Preprocess **Classify** Cluster Associate Select attributes Visualize

Classifier

Choose **J48** -C 0.25 -M 2

Test options

- Use training set
- Supplied test set
- Cross-validation Folds 10
- Percentage split % 80

[More options...](#)

(Nom) class\_los\_label

Start Stop

Result list (right-click for options)

- 16:53:20 - bayes.NaiveBayes
- 17:24:37 - trees.RandomForest
- 17:28:41 - trees.J48

Classifier output

==== Evaluation on test split ===

Time taken to test model on test split: 0.5 seconds

==== Summary ===

	Correctly Classified Instances	1100	71.8485 %
Incorrectly Classified Instances	431	28.1515 %	
Kappa statistic	0.3155		
Mean absolute error	0.3244		
Root mean squared error	0.4897		
Relative absolute error	75.0083 %		
Root relative squared error	105.1492 %		
Total Number of Instances	1531		

==== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area
0.839	0.540	0.769	0.839	0.803	0.319	0.660	0.760	
0.460	0.161	0.571	0.460	0.510	0.319	0.660	0.466	
Weighted Avg.	0.718	0.419	0.706	0.718	0.709	0.319	0.660	0.667

==== Confusion Matrix ===

a	b	<-- classified as
876	168	a = 0
263	224	b = 1

Status OK

48°F Mostly clear 

```

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

import warnings
warnings.filterwarnings("ignore")

#diag_data=pd.read_csv("C:\\Users\\shikh\\Downloads\\DIAGNOSES_ICD.csv")
#adm_data=pd.read_csv("C:\\Users\\shikh\\Downloads\\ADMISSIONS.csv")
#pat_data=pd.read_csv("C:\\Users\\shikh\\Downloads\\PATIENTS.csv")
#icu_data=pd.read_csv("C:\\Users\\shikh\\Downloads\\ICUSTAYS.csv")

adm_data = pd.read_csv('C:\\\\Users\\\\Siddi\\\\OneDrive\\\\Documents\\\\MIMIC Datasets\\\\mimic-iii-clinical-database-1.4\\\\ADMISSIONS.csv')
diag_data = pd.read_csv('C:\\\\Users\\\\Siddi\\\\OneDrive\\\\Documents\\\\MIMIC Datasets\\\\mimic-iii-clinical-database-1.4\\\\DIAGNOSES_ICD.csv')
icu_data=pd.read_csv('C:\\\\Users\\\\Siddi\\\\OneDrive\\\\Documents\\\\MIMIC Datasets\\\\mimic-iii-clinical-database-1.4\\\\ICUSTAYS.csv')
pat_data=pd.read_csv('C:\\\\Users\\\\Siddi\\\\OneDrive\\\\Documents\\\\MIMIC Datasets\\\\mimic-iii-clinical-database-1.4\\\\PATIENTS.csv')

diag_data.head()

ROW_ID SUBJECT_ID HADM_ID SEQ_NUM ICD9_CODE
----- ----- ----- ----- -----
1 print(adm_data.info())
2 print("-----")
3 print(icu_data.info())
4 print("-----")
5 print(pat_data.info())
6 print("-----")
7 print(diag_data.info())

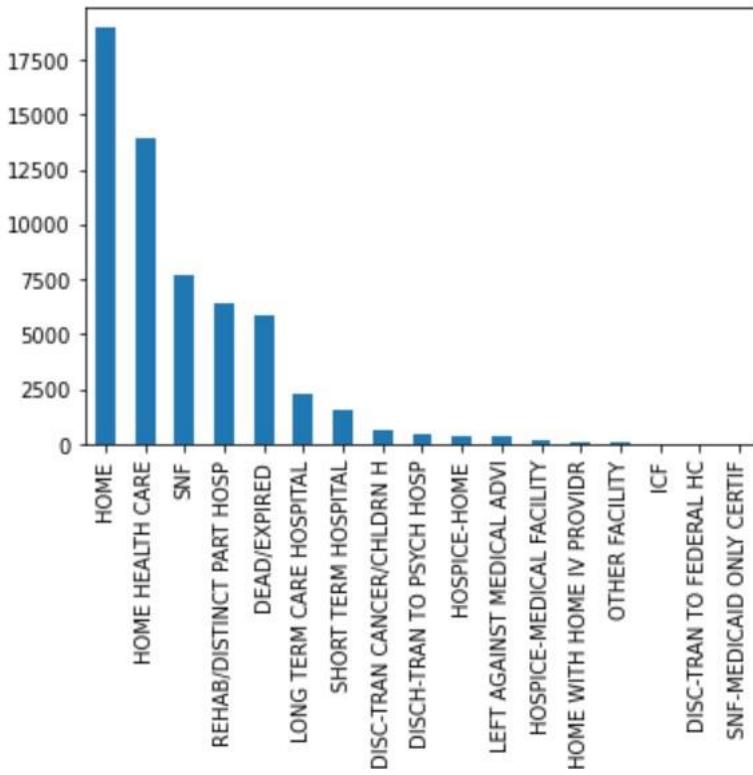
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 58976 entries, 0 to 58975
Data columns (total 19 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   ROW_ID          58976 non-null    int64  
 1   SUBJECT_ID      58976 non-null    int64  
 2   HADM_ID         58976 non-null    int64  
 3   ADMITTIME       58976 non-null    object 
 4   DISCHTIME       58976 non-null    object 
 5   DEATHTIME       5854 non-null     object 
 6   ADMISSION_TYPE  58976 non-null    object 
 7   ADMISSION_LOCATION 58976 non-null    object 
 8   DISCHARGE_LOCATION 58976 non-null    object 

```

```

1 adm_data.DISCHARGE_LOCATION.value_counts().plot(kind="bar")
2 plt.show()

```



```

In [30]: 1 # HOME WITH HOME IV PROVIDR and HOSPICE home is also a HOME with health care so replacing its value
In [31]: 1 adm_data.DISCHARGE_LOCATION=adm_data.DISCHARGE_LOCATION.apply(lambda x: str(x).replace("HOME WITH HOME IV PROVIDR","HOME HEA
<          >
In [32]: 1 adm_data.DISCHARGE_LOCATION=adm_data.DISCHARGE_LOCATION.apply(lambda x: str(x).replace("HOSPICE-HOME","HOME HEALTH CARE"))
In [33]: 1 # moreover REHAB/DISTINCT PART HOSP, LONG TERM CARE HOSPITAL, SHORT TERM HOSPITAL, DISC-TRAN CANCER/CHLDRN H ,
2 #DISCH-TRAN TO PSYCH HOSP etc are sent to some other facilities
In [34]: 1 a=["SNF","REHAB/DISTINCT PART HOSP","LONG TERM CARE HOSPITAL","SHORT TERM HOSPITAL","DISC-TRAN CANCER/CHLDRN H ","DISCH-TRAN
<          >
In [35]: 1 for i in a:
2     adm_data.DISCHARGE_LOCATION=adm_data.DISCHARGE_LOCATION.apply(lambda x: str(x).replace(i,"sent to other facilities"))
3
In [36]: 1 adm_data.DISCHARGE_LOCATION.value_counts()
Out[36]: sent to other facilities
HOME
HOME HEALTH CARE
DEAD/EXPIRED
LEFT AGAINST MEDICAL ADVI

```

```

: 1 adm_data.INSURANCE.value_counts()

: Medicare      28214
: Private       22582
: Medicaid       5785
: Government     1783
: Self Pay        611
: Name: INSURANCE, dtype: int64

: 1 # medicare and medicaid are services offered by government hence converting them

: 1 x=["Medicare","Medicaid"]

: 1 for i in x:
: 2     adm_data.INSURANCE=adm_data.INSURANCE.apply(lambda x: str(x).replace(i,"Government"))

: 1 adm_data.INSURANCE.value_counts().plot(kind="bar")
: 2 plt.show()

```

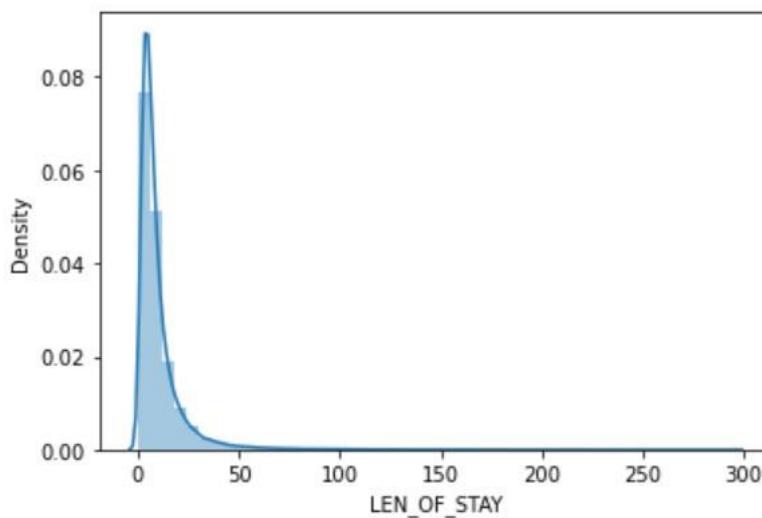


```

: 1 adm__data=adm_data[adm_data.HOSPITAL_EXPIRE_FLAG==0]

: 1 sns.distplot(adm__data.LEN_OF_STAY)
: 2 plt.show()

```



```

: 1 # we can Length of stay is right skewed

: 1 # List of columns no Longer required
: 2
: 3 m=["ROW_ID", "DTSCHTTMF", "DEATHTTMF", "DRGTTMF", "EDOUITTMF", ""

```

```

]: 1 diag_data["codetype"] = diag_data.ICD9_CODE
  2 diag_data["codetype"] = diag_data.ICD9_CODE[~diag_data.ICD9_CODE.str.contains("[a-zA-Z]").fillna(False)]
```

```

]: 1 diag_data.head()
```

```

]: 1
  ROW_ID SUBJECT_ID HADM_ID SEQ_NUM ICD9_CODE codetype
  0    1297      109   172335     1.0    40301    40301
  1    1298      109   172335     2.0     486     486
  2    1299      109   172335     3.0    58281    58281
  3    1300      109   172335     4.0    5855     5855
  4    1301      109   172335     5.0    4254     4254
```

```

]: 1 diag_data["codetype"].fillna(value='999', inplace=True)
```

```

]: 1 diag_data['codetype'] = diag_data['codetype'].str.slice(start=0, stop=3, step=1)
  2 diag_data['codetype'] = diag_data['codetype'].astype(int)
```

```

]: 1 code_ranges = [(1, 140), (140, 240), (240, 280), (280, 290), (290, 320), (320, 390),(390, 460), (460, 520), (520, 580), (580
  2           (710, 740), (740, 760), (760, 780), (780, 800), (800, 1000), (1000, 2000)]
```

```

]: 1 code_ranges = [(1, 140), (140, 240), (240, 280), (280, 290), (290, 320), (320, 390),(390, 460), (460, 520), (520, 580), (580
  2           (710, 740), (740, 760), (760, 780), (780, 800), (800, 1000), (1000, 2000)]
```

```

]: 1 diag = {0: 'infectious', 1: 'neoplasms', 2: 'endocrine', 3: 'blood', 4: 'mental', 5: 'nervous', 6: 'circulatory', 7: 'respira
  2           8: 'digestive', 9: 'genitourinary', 10: 'pregnancy', 11: 'skin', 12: 'muscular', 13: 'congenital', 14: 'prenatal', 1
  3           16: 'injury', 17: 'misc'}
```

```

]: 1 for i, j in enumerate(code_ranges):
  2     diag_data['codetype'] = np.where(diag_data['codetype'].between(j[0],j[1]),
  3                                         i,diag_data['codetype'])
```

```

]: 1 diag_data.head()
```

```

]: 1
  ROW_ID SUBJECT_ID HADM_ID SEQ_NUM ICD9_CODE codetype
  0    1297      109   172335     1.0    40301    6
  1    1298      109   172335     2.0     486     7
  2    1299      109   172335     3.0    58281    9
  3    1300      109   172335     4.0    5855     9
  4    1301      109   172335     5.0    4254     8
```

```
0]: 1 diag_data['codetype'] = diag_data['codetype'].replace(diag)
```

```
1]: 1 diag_data.head()
```

1]:

ROW_ID	SUBJECT_ID	HADM_ID	SEQ_NUM	ICD9_CODE	codetype	
0	1297	109	172335	1.0	40301	circulatory
1	1298	109	172335	2.0	486	respiratory
2	1299	109	172335	3.0	58281	genitourinary
3	1300	109	172335	4.0	5855	genitourinary
4	1301	109	172335	5.0	4254	circulatory

```
2]: 1 diag_data.codetype.value_counts().plot(kind="bar")
2 plt.show()
```



```
[1]: diag = diag_data.groupby('HADM_ID')['codetype'].apply(list).reset_index()
[2]: diag.head()
```

] :

HADM_ID	codetype
0	[endocrine, nervous, genitourinary, digestive,...
1	[digestive, blood, infectious, digestive, circ...
2	[respiratory, respiratory, respiratory, neopla...
3	[digestive, digestive, injury, respiratory, ci...
4	[circulatory, injury, circulatory, endocrine, ...

```
[1]: 1 diag_dummy = pd.get_dummies(diag['codetype'].apply(pd.Series).stack()).sum(level=0)
2 diag_dummy.head()
```

] :

```
1 df_age_min = model_data[['SUBJECT_ID', 'ADMITTIME']].groupby('SUBJECT_ID').min().reset_index()
2 df_age_min.columns = ['SUBJECT_ID', 'ADMIT_MIN']
3 df_age_min.head()
```

	SUBJECT_ID	ADMIT_MIN
0	2	2138-07-17 19:04:00
1	3	2101-10-20 19:08:00
2	4	2191-03-16 00:28:00
3	5	2103-02-02 04:31:00
4	6	2175-05-30 07:15:00

```
1 model_data= model_data.merge(df_age_min, how='outer', on='SUBJECT_ID')

1 model_data['DOB'] = pd.to_datetime(model_data['DOB']).dt.date

1 model_data['ADMIT_MIN'] = pd.to_datetime(model_data['ADMIT_MIN']).dt.date

1 model_data['age'] = model_data.apply(lambda x: (x['ADMIT_MIN'] - x['DOB']).days/365, axis=1)

1 model_data=model_data[model_data.age<100]
```

```
1 model_data=model_data[model_data.age<100]
```

```
1 model_data.head()
```

	SUBJECT_ID	HADM_ID	ADMITTIME	ADMISSION_TYPE	ADMISSION_LOCATION	DISCHARGE_LOCATION	INSURANCE	MARITAL_STATUS	ETHNICITY	
0	22	165315	2196-04-09 12:26:00	EMERGENCY	EMERGENCY ROOM ADMIT	sent to other facilities	Private	MARRIED	WHITE	BEI
1	23	152223	2153-09-03 07:15:00	ELECTIVE	referral from other facilities	HOME HEALTH CARE	Government	MARRIED	WHITE	CORC DISEA ART
2	23	124321	2157-10-18 19:34:00	EMERGENCY	transfer from other facilities	HOME HEALTH CARE	Government	MARRIED	WHITE	
3	24	161859	2139-06-06 16:14:00	EMERGENCY	transfer from other facilities	HOME	Private	SINGLE	WHITE	
4	25	129635	2160-11-02 02:06:00	EMERGENCY	EMERGENCY ROOM ADMIT	HOME	Private	MARRIED	WHITE	ACU

5 rows × 32 columns

```
1 icu_data.head()
```

```
1 model_data.ICU=model_data.ICU.astype("object")
1 model_data.NICU=model_data.NICU.astype("object")
1 drop=["SUBJECT_ID","HADM_ID","ADMITTIME","DOB","ADMIT_MIN"]
1 model_data.drop(drop,axis=1,inplace=True)
1 model_data.drop("DIAGNOSIS",axis=1, inplace=True)
1 model_data.drop("DISCHARGE_LOCATION",axis=1, inplace=True)
1 model_data=pd.get_dummies(model_data, columns=['ADMISSION_TYPE', 'INSURANCE', 'ETHNICITY', 'ADMISSION_LOCATION', 'MARITAL_STAT
<   >
1 model_data.isnull().sum()

LEN_OF_STAY          7858
blood                7858
circulatory          7858
congenital            7858
digestive             7858
endocrine             7858
genitourinary         7858
```

```

1 from sklearn.model_selection import train_test_split
2
3 from sklearn.tree import DecisionTreeRegressor
4
5 dt = DecisionTreeRegressor(random_state=42, max_depth=4, min_samples_leaf=10)
6
7 np.random.seed(0)
8 df_train, df_test = train_test_split(model_data, train_size=0.8, test_size=0.2, random_state=100, shuffle=True)
9
10 #(df_train)
11 #df_test
12 #df_train.to_csv('train2.csv')
13 df_test.to_csv('test2.csv')
14
15 from sklearn.preprocessing import MinMaxScaler
16 scaler = MinMaxScaler()
17
18 df_train['LEN_OF_STAY'] = scaler.fit_transform(df_train[['LEN_OF_STAY']])
19 df_test['LEN_OF_STAY'] = scaler.transform(df_test[['LEN_OF_STAY']])
20
21 y_train = df_train.pop("LEN_OF_STAY")
22 X_train = df_train
23
24 X_test.shape, X_train.shape
25
26 ((10208, 41), (40828, 41))
27
28
29 1 print((df_train))
30 2 #print(len(df_test))
31
32
33      blood  circulatory  congenital  digestive  endocrine  genitourinary \
34 13105    0.0        0.0        0.0        0.0        0.0        0.0
35 46351    0.0        1.0        0.0        1.0        1.0        1.0
36 6283     0.0        0.0        0.0        0.0        3.0        0.0
37 4808     1.0        3.0        0.0        0.0        2.0        0.0
38 344      1.0        6.0        0.0        0.0        2.0        0.0
39 ...
40 ...
41 16304    0.0        2.0        0.0        0.0        1.0        0.0
42 79       0.0        3.0        0.0        1.0        0.0        1.0
43 12119    0.0        0.0        0.0        0.0        0.0        0.0
44 14147    0.0        0.0        0.0        3.0        2.0        0.0
45 38408    0.0        1.0        0.0        0.0        1.0        1.0
46
47
48      infectious  injury  mental  misc  ...
49 13105      0.0      3.0      0.0      0.0  ...
50 46351      0.0      3.0      0.0      0.0  ...
51 6283       1.0      2.0      1.0      0.0  ...
52 4808       0.0      1.0      0.0      0.0  ...
53 344        0.0      1.0      0.0      0.0  ...
54 ...
55 ...

```

```
In [149]: 1 dt.fit(X_train, y_train)

Out[149]: 
DecisionTreeRegressor
DecisionTreeRegressor(max_depth=4, min_samples_leaf=10, random_state=42)
```

```
In [150]: 1 y_train_pred = dt.predict(X_train)
```

```
In [151]: 1 from sklearn.metrics import r2_score
```

```
In [152]: 1 r2_score(y_train, y_train_pred)
```

```
Out[152]: 0.31064417653789667
```

```
In [153]: 1 from sklearn.metrics import mean_squared_error
2 mean_squared_error(y_train, y_train_pred)
```

```
Out[153]: 0.0012233773256710109
```

```
In [154]: 1 y_test_pred = dt.predict(X_test)
```

```
In [155]: 1 r2_score(y_test, y_test_pred)
```

```
Out[155]: 0.3601967080069599
```

```
]: 1 y_test_pred = dt.predict(X_test)
```

```
]: 1 r2_score(y_test, y_test_pred)
```

```
]: 0.3601967080069599
```

```
]: 1 from sklearn.metrics import mean_squared_error
2 mean_squared_error(y_test, y_test_pred)
```

```
]: 0.0011801236024813495
```

```
]: 1 from sklearn.model_selection import GridSearchCV
2
3 model = DecisionTreeRegressor()
4
5 gs = GridSearchCV(model,
6                     param_grid = {'max_depth': range(1, 11),
7                                   'min_samples_split': range(10, 60, 10)},
8                     cv=5,
9                     n_jobs=1,
10                    scoring='neg_mean_squared_error')
11
12 gs.fit(X_train, y_train)
13
14 print(gs.best_params_)
15 print(gs.best_score_)
```

```
8             cv=5,
9             n_jobs=1,
10            scoring='neg_mean_squared_error')
11
12 gs.fit(X_train, y_train)
13
14 print(gs.best_params_)
15 print(-gs.best_score_)
```

```
{'max_depth': 8, 'min_samples_split': 40}
0.001185003465366101
```

```
58]: 1 new_model = DecisionTreeRegressor(max_depth=8,
2                                         min_samples_split=40)
3
4 new_model.fit(X_train, y_train)
```

```
58]: ▾ DecisionTreeRegressor
DecisionTreeRegressor(max_depth=8, min_samples_split=40)
```

```
59]: 1 y_train_pred = new_model.predict(X_train)
```

```
60]: 1 r2_score(y_train, y_train_pred)
```

```
60]: 0.408334998845035
```

```
0.408334998845035
```

```
1 mean_squared_error(y_train, y_train_pred)
```

```
0.0010500086053830057
```

```
1 y_test_pred = new_model.predict(X_test)
```

```
1 r2_score(y_test, y_test_pred)
```

```
0.39821902012769683
```

```
1 mean_squared_error(y_test, y_test_pred)
```

```
0.0011099910656280029
```

```
1 from sklearn.model_selection import GridSearchCV
2 from sklearn.ensemble import RandomForestRegressor
3 # Create the parameter grid based on the results of random search
4 param_grid = {
5     'bootstrap': [True],
6     'max_depth': [80, 90, 100, 110],
7     'max_features': [2, 3],
8     'min_samples_leaf': [3, 4, 5],
9     'min_samples_split': [8, 10, 12]
```

```

o     'min_samples_leaf': [2, 4, 7],
9     'min_samples_split': [8, 10, 12],
10    'n_estimators': [100, 200, 300, 1000]
11 }
12 # Create a based model
13 rf = RandomForestRegressor()
14 # Instantiate the grid search model
15 grid_search = GridSearchCV(estimator = rf, param_grid = param_grid,
16                           cv = 3, n_jobs = -1, verbose = 2)

1 #grid_search.fit(X_train, y_train)
2 #grid_search.best_params_

1 rf = RandomForestRegressor(bootstrap= True,max_depth=100,max_features= 3,
2                               min_samples_leaf= 3,min_samples_split= 10,n_estimators= 300)

1 rf.fit(X_train, y_train)

▼
      RandomForestRegressor
RandomForestRegressor(max_depth=100, max_features=3, min_samples_leaf=3,
                      min_samples_split=10, n_estimators=300)

1 y_train_pred = rf.predict(X_train)

:
:   1 y_train_pred = rf.predict(X_train)
:
:   1 r2_score(y_train, y_train_pred)
:
: 0.4991058021742466

:
:   1 mean_squared_error(y_train, y_train_pred)
:
: 0.0008889206173709558

:
:   1 y_test_pred = rf.predict(X_test)
:
:   1 r2_score(y_test, y_test_pred)
:
: 0.4070452076206904

:
:   1 mean_squared_error(y_test, y_test_pred)
:
: 0.0010937110740887898

:
:   1 model_data.head()
:
:
LEN_OF_STAY blood circulatory congenital digestive endocrine genitourinary infectious injury mental ... ADMISSION_LOCATION** ADMISSION_LOCATION**  

INFO NOT AVAILABLE**

```

```
: 1 # we can see data data is imbalanced
:
: 1 X=class_data.drop(["class_los_label"],axis=1)
:
: 1 y=class_data["class_los_label"]
:
: 1 X_train,X_test,y_train,y_test=train_test_split(X, y, train_size=0.8, test_size=0.2,random_state=100, shuffle=True)
:
: 1 y_train.value_counts(normalize=True)
:
: 0    0.693372
: 1    0.306628
: Name: class_los_label, dtype: float64
:
: 1 y_test.value_counts(normalize=True)
:
: 0    0.689753
: 1    0.310247
: Name: class_los_label, dtype: float64
:
: 1 # handling imbalance problem
: 2
: 3 # SMOTE
: 4 from imblearn.over_sampling import SMOTE
```

```
:
3 # SMOTE
4 from imblearn.over_sampling import SMOTE
5 smt = SMOTE(random_state=27)
6 X_resampled_smt, y_resampled_smt = smt.fit_resample(X_train, y_train)
7 len(X_resampled_smt)
```

56618

```
1 # Importing decision tree classifier
2 from sklearn.tree import DecisionTreeClassifier
```

```
1 import sklearn
2 from sklearn.feature_selection import RFE
```

```
1 dtree = DecisionTreeClassifier()
```

```
1
2 # Importing RFE
3
4 # Instantiate RFE with 15 columns
5 rfe = RFE(20)
6
7 # Fit the rfe model with train set
8 rfe = rfe.fit(X_resampled_smt, y_resampled_smt)
```

```
1 dtree = DecisionTreeClassifier()
2 # Importing RFE
3 from sklearn.feature_selection import RFE
4
5 # Instantiate RFE with 15 columns
6 rfe = RFE(dtree,n_features_to_select=20)
7
8 # Fit the rfe model with train set
9 rfe = rfe.fit(X_resampled_smt, y_resampled_smt)
```

```
1 # RFE selected columns
2 rfe_cols = X_resampled_smt.columns[rfe.support_]
3 print(rfe_cols)
```

```
Index(['blood', 'circulatory', 'digestive', 'endocrine', 'genitourinary',
       'infectious', 'injury', 'mental', 'misc', 'muscular', 'neoplasms',
       'nervous', 'prenatal', 'respiratory', 'skin', 'age', 'ETHNICITY_WHITE',
       'ADMISSION_LOCATION_refferal from other facilities',
       'MARITAL_STATUS_MARRIED', 'GENDER_F'],
      dtype='object')
```

```
1 # Importing Libraries for cross validation
```

```
1 # Importing libraries for cross validation
2 from sklearn.model_selection import KFold
3 from sklearn.model_selection import cross_val_score
4 from sklearn.model_selection import GridSearchCV
5 # Create the parameter grid
6 param_grid = {
7     'max_depth': range(5, 15, 5),
8     'min_samples_leaf': range(50, 150, 50),
9     'min_samples_split': range(50, 150, 50),
10 }
11
12
13 # Instantiate the grid search model
14 dtree = DecisionTreeClassifier()
15
16 grid_search = GridSearchCV(estimator = dtree,
17                             param_grid = param_grid,
18                             scoring= 'recall',
19                             cv = 5,
20                             verbose = 1)
21
22 # Fit the grid search to the data
23 grid_search.fit(X_resampled_smt[rfe_cols],y_resampled_smt)
```

Fitting 5 folds for each of 8 candidates, totalling 40 fits

```
▶ GridSearchCV
```

```
Fitting 5 folds for each of 8 candidates, totalling 40 fits
:
:   GridSearchCV
:     estimator: DecisionTreeClassifier
:       DecisionTreeClassifier
```

```
:
:   # cv results
:   cv_results = pd.DataFrame(grid_search.cv_results_)
:   cv_results
```

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_max_depth	param_min_samples_leaf	param_min_samples_split	params	split
0	0.155799	0.012851	0.036331	0.006940	5	50	50	{'max_depth': 5, 'min_samples_leaf': 50, 'min...}	
1	0.225724	0.047187	0.052582	0.011328	5	50	100	{'max_depth': 5, 'min_samples_leaf': 50, 'min...}	
2	0.165192	0.038695	0.036868	0.015044	5	100	50	{'max_depth': 5, 'min_samples_leaf': 100, 'min...}	
3	0.134946	0.007306	0.031133	0.002559	5	100	100	{'max_depth': 5, 'min_samples_leaf': 100, 'min...}	

```
:
:   print(grid_search.best_estimator_)
```

```
DecisionTreeClassifier(max_depth=5, min_samples_leaf=50, min_samples_split=50)
```

```
:
:   dt = DecisionTreeClassifier(max_depth=5, min_samples_leaf=50, min_samples_split=50)
:   dt.fit(X_resampled_smt[rfe_cols], y_resampled_smt)
```

```
:
:   DecisionTreeClassifier
```

```
DecisionTreeClassifier(max_depth=5, min_samples_leaf=50, min_samples_split=50)
```

```
:
:   # Predictions on the train set
:   y_train_pred = dt.predict(X_resampled_smt[rfe_cols])
```

```
:
:   # Confusion matrix
:   from sklearn import metrics
:   from sklearn.metrics import confusion_matrix
:   confusion = metrics.confusion_matrix(y_resampled_smt, y_train_pred)
:   print(confusion)
```

```
[[21855  6454]
 [ 8015 20294]]
```

```
[[21855 6454]
 [ 8015 20294]]
```

```
1 TP = confusion[1,1] # true positive
2 TN = confusion[0,0] # true negatives
3 FP = confusion[0,1] # false positives
4 FN = confusion[1,0] # false negatives

1 # Accuracy
2 print("Accuracy:-",metrics.accuracy_score(y_resampled_smt, y_train_pred))
3
4 # Sensitivity
5 print("Sensitivity:-",TP / float(TP+FN))
6
7 # Specificity
8 print("Specificity:-", TN / float(TN+FP))
```

```
Accuracy:- 0.744445229432336
```

```
Sensitivity:- 0.7168744922109577
```

```
Specificity:- 0.7720159666537144
```

```
1 # Prediction on the test set
2 y_test_pred = dt.predict(X_test[rfe_cols])
```

```
1 " " " " "
```

```
: 1 # Prediction on the test set
2 y_test_pred = dt.predict(X_test[rfe_cols])
```

```
: 1 # Confusion matrix
2 confusion = metrics.confusion_matrix(y_test, y_test_pred)
3 print(confusion)
```

```
[[5460 1581]
```

```
[1268 1899]]
```

```
: 1 P = confusion[1,1] # true positive
2 TN = confusion[0,0] # true negatives
3 FP = confusion[0,1] # false positives
4 FN = confusion[1,0] # false negatives
```

```
: 1 # Accuracy
2 print("Accuracy:-",metrics.accuracy_score(y_test, y_test_pred))
3
```

```
: 4 # Sensitivity
5 print("Sensitivity:-",TP / float(TP+FN))
6
```

```
: 7 # Specificity
8 print("Specificity:-", TN / float(TN+FP))
9
```

```
: 10 #F1 Score
11 from sklearn.metrics import f1_score
```

```
Accuracy:- 0.744445229432336
Sensitivity:- 0.7168744922109577
Specificity:- 0.7720159666537144
```

```
1 # Prediction on the test set
2 y_test_pred = dt.predict(X_test[rfe_cols])
```

```
1 # Confusion matrix
2 confusion = metrics.confusion_matrix(y_test, y_test_pred)
3 print(confusion)
```

```
[[5460 1581]
 [1268 1899]]
```

```
1 P = confusion[1,1] # true positive
2 TN = confusion[0,0] # true negatives
3 FP = confusion[0,1] # false positives
4 FN = confusion[1,0] # false negatives
```

```
1 # Accuracy
2 print("Accuracy:-",metrics.accuracy_score(y_test, y_test_pred))
3
4 # Sensitivity
```

```
1 # Accuracy
2 print("Accuracy:-",metrics.accuracy_score(y_test, y_test_pred))
3
4 # Sensitivity
5 print("Sensitivity:-",TP / float(TP+FN))
6
7 # Specificity
8 print("Specificity:-", TN / float(TN+FP))
9
10 #F1 Score
11 from sklearn.metrics import f1_score
12 print("F1_Score",f1_score(y_test, y_test_pred, pos_label="0"))
13
```

```
Accuracy:- 0.7209051724137931
Sensitivity:- 0.9411928392542436
Specificity:- 0.7754580315296122
F1_Score 0.7930859176410777
```

```
1 ## Logistic regression

1 ##### Importing stats model
2 import statsmodels.api as sm

1 X_resampled_smt.ICU=X_resampled_smt.ICU.astype("int")

1 X_resampled_smt.NICU=X_resampled_smt.NICU.astype("int")

1 X_resampled_smt.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 56618 entries, 0 to 56617
Data columns (total 41 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   blood            56618 non-null   float64
 1   circulatory     56618 non-null   float64
 2   congenital       56618 non-null   float64
 3   digestive        56618 non-null   float64
 4   endocrine        56618 non-null   float64
 5   genitourinary    56618 non-null   float64
 6   infectious       56618 non-null   float64
 7   injury            56618 non-null   float64
 8   mental            56618 non-null   float64
```

```

1 y_resampled_smt=y_resampled_smt.astype("int")

1 lreg=sm.GLM(y_resampled_smt,(sm.add_constant(X_resampled_smt)),family=sm.families.Binomial())

1 # Fit the model
2 lreg=lreg.fit().summary()

1 lreg

```

Generalized Linear Model Regression Results

<b>Dep. Variable:</b>	class_los_label	<b>No. Observations:</b>	56618
<b>Model:</b>	GLM	<b>Df Residuals:</b>	56576
<b>Model Family:</b>	Binomial	<b>Df Model:</b>	41
<b>Link Function:</b>	Logit	<b>Scale:</b>	1.0000
<b>Method:</b>	IRLS	<b>Log-Likelihood:</b>	-23895.
<b>Date:</b>	Sat, 10 Dec 2022	<b>Deviance:</b>	47791.
<b>Time:</b>	21:25:07	<b>Pearson chi2:</b>	4.52e+04
<b>No. Iterations:</b>	29	<b>Pseudo R-squ. (CS):</b>	0.4185
<b>Covariance Type:</b>	nonrobust		

```
1 lreg
```

Generalized Linear Model Regression Results

<b>Dep. Variable:</b>	class_los_label	<b>No. Observations:</b>	56618
<b>Model:</b>	GLM	<b>Df Residuals:</b>	56576
<b>Model Family:</b>	Binomial	<b>Df Model:</b>	41
<b>Link Function:</b>	Logit	<b>Scale:</b>	1.0000
<b>Method:</b>	IRLS	<b>Log-Likelihood:</b>	-23895.
<b>Date:</b>	Sat, 10 Dec 2022	<b>Deviance:</b>	47791.
<b>Time:</b>	21:25:07	<b>Pearson chi2:</b>	4.52e+04
<b>No. Iterations:</b>	29	<b>Pseudo R-squ. (CS):</b>	0.4185
<b>Covariance Type:</b>	nonrobust		

	coef	std err	z	P> z	[0.025	0.975]
const	144.6232	7.86e+04	0.002	0.999	-1.54e+05	1.54e+05
blood	0.0625	0.019	3.301	0.001	0.025	0.100
circulatory	0.1111	0.006	17.518	0.000	0.099	0.123
congenital	0.2610	0.037	7.109	0.000	0.189	0.333
digestive	0.2263	0.011	21.044	0.000	0.205	0.247

	<b>ETHNICITY_BLACK/AFRICAN AMERICAN</b>	-29.4138	2.89e+04	-0.001	0.999	-5.66e+04	5.66e+04
	<b>ETHNICITY_HISPANIC/LATINO</b>	-29.3425	2.89e+04	-0.001	0.999	-5.66e+04	5.66e+04
	<b>ETHNICITY_Unknown/other</b>	-29.2404	2.89e+04	-0.001	0.999	-5.66e+04	5.66e+04
	<b>ETHNICITY_WHITE</b>	-29.1738	2.89e+04	-0.001	0.999	-5.66e+04	5.66e+04
	<b>ADMISSION_LOCATION_** INFO NOT AVAILABLE **</b>	-30.5316	2.73e+04	-0.001	0.999	-5.35e+04	5.34e+04
	<b>ADMISSION_LOCATION_EMERGENCY ROOM ADMIT</b>	-29.0338	2.73e+04	-0.001	0.999	-5.35e+04	5.34e+04
	<b>ADMISSION_LOCATION_refferral from other facilities</b>	-29.5240	2.73e+04	-0.001	0.999	-5.35e+04	5.34e+04
	<b>ADMISSION_LOCATION_transfer from other facilities</b>	-28.8920	2.73e+04	-0.001	0.999	-5.35e+04	5.34e+04
	<b>MARITAL_STATUS_DIVORCED</b>	-1.3159	0.061	-21.438	0.000	-1.436	-1.196
	<b>MARITAL_STATUS_MARRIED</b>	-1.1739	0.046	-25.601	0.000	-1.264	-1.084
	<b>MARITAL_STATUS_SINGLE</b>	-1.2900	0.049	-26.278	0.000	-1.386	-1.194
	<b>MARITAL_STATUS_WIDOWED</b>	-1.3066	0.057	-23.103	0.000	-1.417	-1.196
	<b>GENDER_F</b>	-29.1795	3.12e+04	-0.001	0.999	-6.12e+04	6.12e+04
	<b>GENDER_M</b>	-29.1409	3.12e+04	-0.001	0.999	-6.12e+04	6.12e+04

```
[1]: 1 # Importing Logistic regression from sklearn
2 from sklearn.linear_model import LogisticRegression
3 # Instantiate the Logistic regression
4 logreg = LogisticRegression()
```

```
1 # Importing Logistic regression from sklearn
2 from sklearn.linear_model import LogisticRegression
3 # Instantiate the Logistic regression
4 logreg = LogisticRegression()
```

```
1 # Instantiate RFE with 15 columns
2 rfe = RFE(logreg, n_features_to_select=15)
3
4 # Fit the rfe model with train set
5 rfe = rfe.fit(X_resampled_smt, y_resampled_smt)
```

```
1 rfe_cols = X_resampled_smt.columns[rfe.support_]
2 print(rfe_cols)
```

```
Index(['INSURANCE_Government', 'INSURANCE_Private', 'INSURANCE_Self Pay',
       'ETHNICITY_ASIAN', 'ETHNICITY_BLACK/AFRICAN AMERICAN',
       'ETHNICITY_HISPANIC/LATINO', 'ETHNICITY_Unknown/other',
       'ETHNICITY_WHITE', 'ADMISSION_LOCATION_** INFO NOT AVAILABLE **',
       'ADMISSION_LOCATION_EMERGENCY ROOM ADMIT',
       'ADMISSION_LOCATION_refferral from other facilities',
       'ADMISSION_LOCATION_transfer from other facilities',
       'MARITAL_STATUS_SINGLE', 'GENDER_F', 'GENDER_M'],
      dtype='object')
```

```
1 # Adding constant to X_train
```

```

1 # Adding constant to X_train
2 X_train_sm_1 = sm.add_constant(X_resampled_smt[rfe_cols])
3
4 #Instantiate the model
5 log_1 = sm.GLM(y_resampled_smt, X_train_sm_1, family=sm.families.Binomial())
6
7 # Fit the model
8 log_1 = log_1.fit()
9
10 log_1.summary()

```

: Generalized Linear Model Regression Results

<b>Dep. Variable:</b>	class_los_label	<b>No. Observations:</b>	56618				
<b>Model:</b>	GLM	<b>Df Residuals:</b>	56602				
<b>Model Family:</b>	Binomial	<b>Df Model:</b>	15				
<b>Link Function:</b>	Logit	<b>Scale:</b>	1.0000				
<b>Method:</b>	IRLS	<b>Log-Likelihood:</b>	-30318.				
<b>Date:</b>	Sat, 10 Dec 2022	<b>Deviance:</b>	60636.				
<b>Time:</b>	21:25:17	<b>Pearson chi2:</b>	4.61e+04				
<b>No. Iterations:</b>	29	<b>Pseudo R-squ. (CS):</b>	0.2705				
<b>Covariance Type:</b>	nonrobust						
		coef	std err	z	P> z	[0.025	0.975]
	const	117.6701	6.71e+04	0.002	0.999	-1.31e+05	1.32e+05
	INSURANCE_Government	-29.2458	3.84e+04	-0.001	0.999	-7.53e+04	7.52e+04
	INSURANCE_Private	-29.4039	3.84e+04	-0.001	0.999	-7.53e+04	7.52e+04
	INSURANCE_Self Pay	-30.1317	3.84e+04	-0.001	0.999	-7.53e+04	7.52e+04
	ETHNICITY_ASIAN	-30.0488	3.13e+04	-0.001	0.999	-6.13e+04	6.13e+04
	ETHNICITY_BLACK/AFRICAN AMERICAN	-29.6896	3.13e+04	-0.001	0.999	-6.13e+04	6.13e+04
	ETHNICITY_HISPANIC/LATINO	-29.7361	3.13e+04	-0.001	0.999	-6.13e+04	6.13e+04
	ETHNICITY_Unknown/other	-29.6023	3.13e+04	-0.001	0.999	-6.13e+04	6.13e+04
	ETHNICITY_WHITE	-29.5043	3.13e+04	-0.001	0.999	-6.13e+04	6.13e+04
ADMISSION_LOCATION_** INFO NOT AVAILABLE **		-31.2396	3.01e+04	-0.001	0.999	-5.9e+04	5.89e+04
ADMISSION_LOCATION_EMERGENCY ROOM ADMIT		-29.5970	3.01e+04	-0.001	0.999	-5.9e+04	5.89e+04
ADMISSION_LOCATION_referral from other facilities		-29.8224	3.01e+04	-0.001	0.999	-5.9e+04	5.89e+04
ADMISSION_LOCATION_transfer from other facilities		-29.4197	3.01e+04	-0.001	0.999	-5.9e+04	5.89e+04
MARITAL_STATUS_SINGLE		-0.3207	0.025	-13.009	0.000	-0.369	-0.272
GENDER_F		-29.5266	3.38e+04	-0.001	0.999	-6.63e+04	6.63e+04
GENDER_M		-29.4972	3.38e+04	-0.001	0.999	-6.63e+04	6.63e+04

```

1 # Check for the VIF values of the feature variables.
2 from statsmodels.stats.outliers_influence import variance_inflation_factor

1 vif = pd.DataFrame()
2 vif['Features'] = X_resampled_smt[rfe_cols].columns
3 vif['VIF'] = [variance_inflation_factor(X_resampled_smt[rfe_cols].values, i) for i in range(X_resampled_smt[rfe_cols].shape[0])]
4 vif['VIF'] = round(vif['VIF'], 2)
5 vif = vif.sort_values(by = "VIF", ascending = False)
6 vif

```

	Features	VIF
8	ADMISSION_LOCATION_** INFO NOT AVAILABLE **	1.04
2	INSURANCE_Self Pay	0.47
3	ETHNICITY_ASIAN	0.22
5	ETHNICITY_HISPANIC/LATINO	0.19
6	ETHNICITY_Unknown/other	0.10
1	INSURANCE_Private	0.08
4	ETHNICITY_BLACK/AFRICAN AMERICAN	0.08
0	INSURANCE_Government	0.07
11	GENDER_M	0.07
13	GENDER_F	0.04
7	ETHNICITY_WHITE	0.03
12	MARITAL_STATUS_SINGLE	0.01
9	ADMISSION_LOCATION_EMERGENCY ROOM ADMIT	0.00
10	ADMISSION_LOCATION_referral from other facilities	0.00
11	ADMISSION_LOCATION_transfer from other facilities	0.00

```

227]: 1 y_train_pred=log_1.predict(X_train_sm_1)
2 y_train_pred.head()

```

```

227]: 0    0.249350
1    0.378761
2    0.349263
3    0.394379
4    0.449278
dtype: float64

```

```

228]: 1 y_train_pred_final=pd.DataFrame({'class_label':y_resampled_smt.values, 'prob':y_train_pred.values})
2
3 y_train_pred_final.head(10)

```

```

228]: class_label      prob

```

```

: 1 y_train_pred_final= pd.DataFrame({'class_label':y_resampled_smt.values, 'prob':y_train_pred.values})
: 2
: 3 y_train_pred_final.head(10)

:      class_label      prob
: 0          0  0.249350
: 1          0  0.378761
: 2          0  0.349263
: 3          0  0.394379
: 4          1  0.449278
: 5          0  0.357313
: 6          0  0.364105
: 7          1  0.401435
: 8          0  0.401435
: 9          0  0.357313

: 1 ## Finding Optimal Probability Cutoff Point

: 1 # Creating columns for different probability cutoffs
: 2 prob_cutoff = [float(p/10) for p in range(10)]

: 1 ## Finding Optimal Probability Cutoff Point

: 1 # Creating columns for different probability cutoffs
: 2 prob_cutoff = [float(p/10) for p in range(10)]
: 3
: 4 for i in prob_cutoff:
: 5     y_train_pred_final[i] = y_train_pred_final['prob'].map(lambda x : 1 if x > i else 0)
: 6
: 7 y_train_pred_final.head()

:      class_label      prob  0.0  0.1  0.2  0.3  0.4  0.5  0.6  0.7  0.8  0.9
: 0          0  0.249350   1   1   1   0   0   0   0   0   0   0
: 1          0  0.378761   1   1   1   1   0   0   0   0   0   0
: 2          0  0.349263   1   1   1   1   0   0   0   0   0   0
: 3          0  0.394379   1   1   1   1   0   0   0   0   0   0
: 4          1  0.449278   1   1   1   1   1   0   0   0   0   0

```

```

: 1 # Creating a dataframe
: 2 from sklearn import metrics
: 3 from sklearn.metrics import confusion_matrix
: 4 cutoff_df = pd.DataFrame(columns=['probability', 'accuracy', 'sensitivity', 'specificity'])
: 5

```

```

1 # Creating a dataframe
2 from sklearn import metrics
3 from sklearn.metrics import confusion_matrix
4 cutoff_df = pd.DataFrame(columns=['probability', 'accuracy', 'sensitivity', 'specificity'])
5
6 for i in prob_cutoff:
7     cm1 = metrics.confusion_matrix(y_train_pred_final['class_label'], y_train_pred_final[i] )
8     total1=sum(sum(cm1))
9     accuracy = (cm1[0,0]+cm1[1,1])/total1
10
11    speci = cm1[0,0]/(cm1[0,0]+cm1[0,1])
12    sensi = cm1[1,1]/(cm1[1,0]+cm1[1,1])
13    cutoff_df.loc[i] = [ i ,accuracy,sensi,speci]
14 print(cutoff_df)

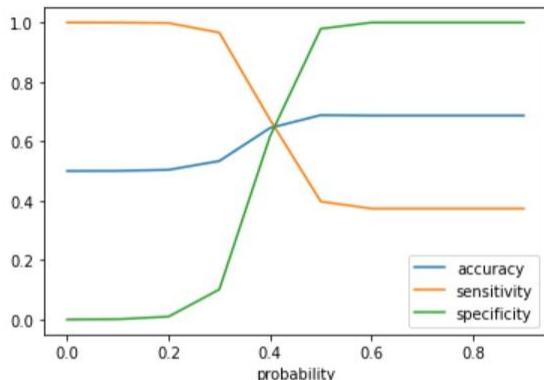
      probability   accuracy   sensitivity   specificity
0.0          0.0  0.500000  1.000000  0.000000
0.1          0.1  0.500512  0.999894  0.001130
0.2          0.2  0.504009  0.998304  0.009714
0.3          0.3  0.533664  0.966089  0.101240
0.4          0.4  0.645025  0.673143  0.616906
0.5          0.5  0.688191  0.397294  0.979088
0.6          0.6  0.686707  0.373415  1.000000
0.7          0.7  0.686707  0.373415  1.000000
0.8          0.8  0.686707  0.373415  1.000000
0.9          0.9  0.686707  0.373415  1.000000

```

```

1 # Plotting accuracy, sensitivity and specificity for different probabilities.
2 cutoff_df.plot('probability', ['accuracy','sensitivity','specificity'])
3 plt.show()

```



```

1 # Creating a column with name "predicted", which is the predicted value for 0.6 cutoff
2 y_train_pred_final['predicted'] = y_train_pred_final['prob'].map(lambda x: 1 if x > 0.4 else 0)
3 y_train_pred_final.head()

```

```
class_label      prob  0.0  0.1  0.2  0.3  0.4  0.5  0.6  0.7  0.8  0.9  predicted

```

```
33]: 1 # Creating a column with name "predicted", which is the predicted value for 0.6 cutoff
2 y_train_pred_final['predicted'] = y_train_pred_final['prob'].map(lambda x: 1 if x > 0.4 else 0)
3 y_train_pred_final.head()
```

```
33]:
```

	class_label	prob	0.0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	predicted
0	0	0.249350	1	1	1	0	0	0	0	0	0	0	0
1	0	0.378761	1	1	1	1	0	0	0	0	0	0	0
2	0	0.349263	1	1	1	1	0	0	0	0	0	0	0
3	0	0.394379	1	1	1	1	0	0	0	0	0	0	0
4	1	0.449278	1	1	1	1	1	0	0	0	0	0	1

```
34]: 1 # Confusion metrics
2 confusion = metrics.confusion_matrix(y_train_pred_final['class_label'], y_train_pred_final['predicted'])
3 print(confusion)

[[17464 10845]
 [ 9253 19056]]
```

```
35]: 1 TP = confusion[1,1] # true positive
2 TN = confusion[0,0] # true negatives
3 FP = confusion[0,1] # false positives
4 FN = confusion[1,0] # false negatives
```

```
: 1 # Confusion metrics
2 confusion = metrics.confusion_matrix(y_train_pred_final['class_label'], y_train_pred_final['predicted'])
3 print(confusion)

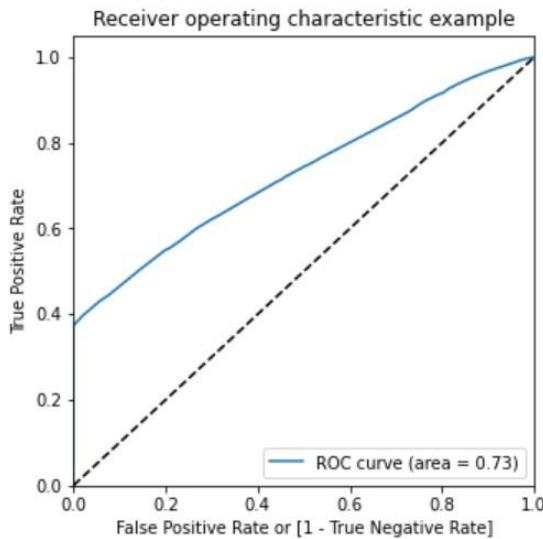
[[17464 10845]
 [ 9253 19056]]
```

```
: 1 TP = confusion[1,1] # true positive
2 TN = confusion[0,0] # true negatives
3 FP = confusion[0,1] # false positives
4 FN = confusion[1,0] # false negatives
```

```
: 1 # Accuracy
2 print("Accuracy:-",metrics.accuracy_score(y_train_pred_final['class_label'], y_train_pred_final['predicted']))
3
4 # Sensitivity
5 print("Sensitivity:-",TP / float(TP+FN))
6
7 # Specificity
8 print("Specificity:-", TN / float(TN+FP))
```

```
Accuracy:- 0.6450245504963086
Sensitivity:- 0.6731428167720513
Specificity:- 0.6169062842205659
```

```
| 18 |     return None  
|:  
| 1 | draw_roc(y_train_pred_final['class_label'], y_train_pred_final['prob'])
```



```
|:  
| 1 | X_test_log=X_test[rfe_cols]  
|:  
| 1 | X_test_sm = sm.add_constant(X_test_log)  
|:  
| 1 | y_test_pred = log_1.predict(X_test_sm)  
|:  
| 1 | y_test_pred.head()  
|:  
| 37611    0.364105  
| 41184    0.417696  
| 45658    0.366555  
| 4908     0.378761  
| 2995     0.456576  
| dtype: float64  
|:  
| 1 | y_pred_1 = pd.DataFrame(y_test_pred)  
| 2 | y_pred_1.head()  
|:  
| 0  
| 37611  0.364105  
| 41184  0.417696  
| 45658  0.366555  
| 4908   0.378761
```

```

1 y_test_pred_final.test_predicted=y_test_pred_final.test_predicted.astype("int64")
2 y_test_pred_final.class_los_label=y_test_pred_final.class_los_label.astype("int64")

1 confusion = metrics.confusion_matrix(y_test_pred_final['class_los_label'], y_test_pred_final['test_predicted'])
2 print(confusion)

[[4339 2702]
 [1728 1439]]


1 TP = confusion[1,1] # true positive
2 TN = confusion[0,0] # true negatives
3 FP = confusion[0,1] # false positives
4 FN = confusion[1,0] # false negatives

1 # Accuracy
2 print("Accuracy:-",metrics.accuracy_score(y_test_pred_final['class_los_label'], y_test_pred_final['test_predicted']))
3
4 # Sensitivity
5 print("Sensitivity:-",TP / float(TP+FN))
6
7 # Specificity
8 print("Specificity:-", TN / float(TN+FP))
9
10 #F1 Score
11 from sklearn.metrics import f1_score
12 print("F1_Score",f1_score(y_test_pred_final['class_los_label'], y_test_pred_final['test_predicted']))


1 # Accuracy
2 print("Accuracy:-",metrics.accuracy_score(y_test_pred_final['class_los_label'], y_test_pred_final['test_predicted']))
3
4 # Sensitivity
5 print("Sensitivity:-",TP / float(TP+FN))
6
7 # Specificity
8 print("Specificity:-", TN / float(TN+FP))
9
10 #F1 Score
11 from sklearn.metrics import f1_score
12 print("F1_Score",f1_score(y_test_pred_final['class_los_label'], y_test_pred_final['test_predicted']))


Accuracy:- 0.5660266457680251
Sensitivity:- 0.45437322387117146
Specificity:- 0.6162476920891918
F1_Score 0.3938149972632732

```

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier Choose **NaiveBayes**

- Test options

- Use training set
- Supplied test set Set...
- Cross-validation Folds 10
- Percentage split % 70
- [More options...](#)

(Nom) class\_los\_label

Start Stop

Result list (right-click for options)

- 21:56:07 - bayes.NaiveBayes
- 21:56:32 - bayes.NaiveBayes**

Classifier output

```
==== Evaluation on test split ====
Time taken to test model on test split: 0.03 seconds

==== Summary ====
Correctly Classified Instances      1671          72.7787 %
Incorrectly Classified Instances   625           27.2213 %
Kappa statistic                      0.349
Mean absolute error                  0.2856
Root mean squared error              0.48
Relative absolute error              65.9209 %
Root relative squared error        102.6854 %
Total Number of Instances          2296

==== Detailed Accuracy By Class ====

```

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
0.838	0.504	0.778	0.838	0.807	0.352	0.759	0.868	0	
0.496	0.162	0.593	0.496	0.540	0.352	0.759	0.597	1	
<b>Weighted Avg.</b>	0.728	0.394	0.718	0.728	0.721	0.352	0.759	0.780	

```
==== Confusion Matrix ====

```

	a	b	<-- classified as
1304	252	1	a = 0
373	367	1	b = 1

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier Choose **J48 - C 0.25 -M 2**

- Test options

- Use training set
- Supplied test set Set...
- Cross-validation Folds 10
- Percentage split % 70
- [More options...](#)

(Nom) class\_los\_label

Start Stop

Result list (right-click for options)

- 21:56:07 - bayes.NaiveBayes
- 21:56:32 - bayes.NaiveBayes
- 21:59:47 - trees.J48**

Classifier output

```
==== Evaluation on test split ====
Time taken to test model on test split: 0.03 seconds

==== Summary ====
Correctly Classified Instances      1654          72.0383 %
Incorrectly Classified Instances   642           27.9617 %
Kappa statistic                      0.3194
Mean absolute error                  0.3373
Root mean squared error              0.4885
Relative absolute error              77.862 %
Root relative squared error        104.4938 %
Total Number of Instances          2296

==== Detailed Accuracy By Class ====

```

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
0.848	0.547	0.765	0.848	0.804	0.325	0.631	0.738	0	
0.453	0.152	0.586	0.453	0.511	0.325	0.631	0.458	1	
<b>Weighted Avg.</b>	0.720	0.420	0.707	0.720	0.710	0.325	0.631	0.647	

```
==== Confusion Matrix ====

```

	a	b	<-- classified as
1319	237	1	a = 0
405	335	1	b = 1

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose **RandomForest -P 100 -I 100 -num-slots 1 -K 0 -M 1.0 -V 0.001 -S 1**

Test options

- Use training set
- Supplied test set Set...
- Cross-validation Folds
- Percentage split %

[More options...](#)

(Nom) class\_los\_label

Start Stop

Result list (right-click for options)

- 21:56:07 - bayes.NaiveBayes
- 21:56:32 - bayes.NaiveBayes
- 21:59:47 - trees.J48
- 22:00:40 - trees.RandomForest**

Classifier output

```
==== Evaluation on test split ====
Time taken to test model on test split: 0.14 seconds
==== Summary ====
Correctly Classified Instances      1730          75.3484 %
Incorrectly Classified Instances   566           24.6516 %
Kappa statistic                   0.3706
Mean absolute error               0.3416
Root mean squared error          0.4114
Relative absolute error           78.8641 %
Root relative squared error     88.0055 %
Total Number of Instances        2296

==== Detailed Accuracy By Class ====

      TP Rate  FP Rate  Precision  Recall   F-Measure  MCC    ROC Area  PRC Area  Class
      0.912    0.580    0.768     0.912    0.834     0.392   0.784    0.876     0
      0.420    0.088    0.694     0.420    0.524     0.392   0.784    0.648     1
Weighted Avg.   0.753    0.421    0.744     0.753    0.734     0.392   0.784    0.802

==== Confusion Matrix ====

      a      b  <-- classified as
1419  137 |  a = 0
429   311 |  b = 1
```