

Configuration Manual

MSc Research Project Data Analytics

Forename Surname Student ID: x19213590

School of Computing National College of Ireland

Supervisor:

Dr. Majid Latifi

National College of Ireland



MSc Project Submission Sheet

School of Computing

Student Name:	Silky Jain
Student ID:	x19213590
Programme:	MSc in Data Analytics Year:2021
Module:	MSc Research Project
Lecturer:	Dr. Majid Latifi
Date:	
Project Title:	Prediction of Length of Stay and Hospital Readmission for Diabetic Patients

Word Count:1405...... **Page Count:**14....

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

<u>ALL</u> internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:	Silky Jain

Date:23/09/2021.....

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST

Attach a completed copy of this sheet to each project (including multiple	
copies)	
Attach a Moodle submission receipt of the online project	
submission, to each project (including multiple copies).	
You must ensure that you retain a HARD COPY of the project, both	
for your own reference and in case a project is lost or mislaid. It is not	
sufficient to keep a copy on computer.	

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Silky Jain Student ID: x19213590

1 Introduction

The purpose of building this document is to present the implementation of the project in a concise and structured manner such that it can be replicated if required. This project aimed to build a model for the prediction of length of stay and hospital readmission using a comprehensive methodological approach involving data cleaning and data pre-processing. Another objective of this project was to identify the most influencing factors for predicting length of stay and hospital readmission. The stacked ensemble learning model is proposed which is then compared with the other base classifiers and existing models on the same field. The tools and techniques used in the project are specified in this manual.

2 System Specifications

Below are the hardware requirements which is required for running the experiment and executing code smoothly.

Operating System	Windows 10 Home
RAM	8.0 GB
Har Disk Space	100 GB Minimum
Processor	Intel(R) Core(TM) i7-9750H CPU @
	2.60GHz 2.59 GHz

3 Tools/Technology

For implementing this project, Python programming language was used with Integrated Development Environment (IDE) as Jupyter Notebook working on the Anaconda platform. The specific versions of the respective platform/language are mentioned below:

Programming Language	Python 3.8.3
IDE	Jupyter Notebook v. 6.0.3
Platform	Anaconda v. 4.9.2
Tools	Microsoft Excel, Overleaf, TeXstudio
Web Browser	Google Chrome

4 **Pre-requisites software setup**

The first step for the execution of this project is the installation of the required platform and languages.

- Python is installed using the link¹.
- Anaconda has been installed using this link².
- After execution of the project, the results are visualized in the Jupyter Notebook using the libraries such as MatPlotLib, seaborn, and Plotly.

5 Data Collection

Data for this project is extracted from the UCI Machine learning repository³ and data description is in (Strack *et al.*, 2014).

UCI Machine Learn Center for Machine Learn	ning Reposito	а ry			
Check out the <u>beta versio</u> Diabetes 130-I Download, <u>Data Folde</u>	n of the new UCI Machin US hospitals r. Data Set Descriptio	e Learning Repository v for years 19	ve are cur 99-20	rently testing! <u>Contact</u>	<u>us</u> if you hav
Abstract: This data has been	prepared to analyze factors	related to readmission as	well as ot	her outcomes pertaining t	to patients with
Data Set Characteristics:	Multivariate	Number of instances:	100000	Area:	Life
Attribute Characteristics:	Integer	Number of Attributes:	55	Date Donated	2014-05-03
Associated Tasks:	Classification, Clustering	Missing Values?	Yes	Number of Web Hits:	362225
Source: The data are submitted on beh Cios (kcios "@' vcu.edu), Jon I	nalf of the Center for Clinical DeShazo (j <u>pdeshazo "@' vcr</u>	and Translational Resear <u>J.edu</u>), and Beata Strack (;	ch, Virginia strackb *@	a Commonwealth Univers <u>y vcu.edu</u>). This data is a	ity, a recipient de-identified a
Data Set Information:					
The dataset represents 10 yea satisfied the following criteria. (1) It is an inpatient encounter (2) It is a diabetic encounter, It (3) The length of stay was at le (4) Laboratory tests were perf (5) Medications were administr	ars (1999-2008) of clinical ca (a hospital admission). nat is, one during which any sast 1 day and at most 14 da mmed during the encounter. ered during the encounter.	re at 130 US hospitals and kind of diabetes was enter ays.	d integrate	d delivery networks. It ind	cludes over 50

6 Implementation

There are various processes involved in the implementation of this project which is given below:

6.1 Data Preparation and Storage

- Extracting the CSV file from the UCI machine learning repository to the local system.
- Importing the libraries in Jupyter Notebook is shown in Figure below.

¹ https://www.python.org/downloads/release/python-383/

² https://anaconda.org/conda-forge/conda/files?version=4.9.2

³ https://archive.ics.uci.edu/ml/datasets/diabetes+130-us+hospitals+for+years+1999-2008



• Reading data from CSV files and storing it in the form of a dataframe.

Reading data from csv files

In [3]:	1 2 3 4 5 6 7 8	data = pd def displ with display_a IDs_mappi display_a ∢	<pre>.read_csv(' ay_all(data pd.option_display(data ll(data.hea ng = pd.rea ll(IDs_mapp)</pre>	"C:/Users/sil a): context("disp a) ad()) ad_csv("C:/Us bing.head(67)	ky/Desk lay.max ers/sil)	top/[_row' ky/De	DAPA - p ", 100, esktop/D	redictive analy "display.max_co APA - predictiv	ysis\dataset_diabetes plumns", 100): ve analysis\dataset_c	\dataset_diabetes	\diabetic_data iabetes\IDs_ma	.csv") pping.
		encounter_id	patient_nbr	race	gender	age	weight	admission_type_id	discharge_disposition_id	admission_source_id	time_in_hospital	payer
	0	2278392	8222157	Caucasian	Female	[0- 10)	?	6	25	1	1	
	1	149190	55629189	Caucasian	Female	[10- 20)	?	1	1	7	3	
	2	64410	86047875	AfricanAmerican	Female	[20- 30)	?	1	1	7	2	
	3	500364	82442376	Caucasian	Male	[30- 40)	?	1	1	7	2	
	4	16680	42519267	Caucasian	Male	[40- 50)	?	1	1	7	1	
		admissio	on_type_id				descripti	ion				
	n		1				Emerner	nev				

6.2 Exploratory Data Analysis

Before any cleaning and transformation of data, it is important to understand the data and to be aware of the features to focus on to undertake data cleaning and pre-processing.





The above figure shows the plot showing the number of patients who were readmitted or not from the dataset.



The figure shows the plot of the count of patients based on the gender and age group which shows that there is an almost equal number of patients in terms of gender and there are a greater number of patients observed from the age 65 to 85 years.



These figures show the readmission of patients concerning admission type, so there are high cases of emergencies. It also shows the number of days patients stay at the hospital where most of the patients stayed for 1 to 3 days while admitted.

6.3 Data Pre-processing

Data cleaning is an important step in model building to get optimal performance. The steps involved in cleaning are converting the datatypes of the variables, checking and removal of duplicate rows, dropping columns with high multicollinearity, checking the columns having missing values and replacing them mean or other values, or dropping the columns if the percentage of the missing value is high as shown in Figures below.

Missing Values

<pre>2 #In the dataset missing values are represented as '?' sign 3 for col in df.columns: 4 if df[col].dtype == object: 5 print(col,df[col][df[col] == '?'].count()) race 2273 gender 0 age 0 weight 98569 payer_code 40256 medical_specialty 49949 diag_1 21 diag_2 358 diag_3 1423 max_glu_serum 0 A1Cresult 0 metformin 0 repaglinide 0 nateglinide 0 chlorpropamide 0 glimepiride 0 acetobexamide 0 glipizide 0 glyburide 0 tolbutamide 0</pre>	In [99]]: 1 #Checking for missing values in dataset										
<pre>3 for col in df.columns: 4 if df[col].dtype == object: 5 print(col,df[col][df[col] == '?'].count()) race 2273 gender 0 age 0 weight 98569 payer_code 40256 medical_specialty 49949 diag_1 21 diag_2 358 diag_3 1423 max_glu_serum 0 A1Cresult 0 metformin 0 repaglinide 0 nateglinide 0 chlorpropamide 0 glimepiride 0 acetobexamide 0 glipizide 0 glyburide 0 tolbutamide 0</pre>		2 #In the dataset missing values are represented as '?' sign										
<pre>4 5 if df[col].dtype == object: 5 print(col,df[col][df[col] == '?'].count()) race 2273 gender 0 age 0 weight 98569 payer_code 40256 medical_specialty 49949 diag_1 21 diag_2 358 diag_3 1423 max_glu_serum 0 A1Cresult 0 metformin 0 repaglinide 0 nateglinide 0 chlorpropamide 0 glimepiride 0 acetobexamide 0 glipizide 0 glyburide 0 tolbutamide 0</pre>		3 for col in df.columns:										
<pre>5 print(col,df[col][df[col] == '?'].count()) race 2273 gender 0 age 0 weight 98569 payer_code 40256 medical_specialty 49949 diag_1 21 diag_2 358 diag_3 1423 max_glu_serum 0 A1Cresult 0 metformin 0 repaglinide 0 chlorpropamide 0 glimepiride 0 acetohexamide 0 glipizide 0 glyburide 0 tolbutamide 0</pre>		<pre>4 if df[col].dtype == object:</pre>										
race 2273 gender 0 age 0 weight 98569 payer_code 40256 medical_specialty 49949 diag_1 21 diag_2 358 diag_3 1423 max_glu_serum 0 A1Cresult 0 metformin 0 repaglinide 0 nateglinide 0 chlorpropamide 0 glimepiride 0 acetohexamide 0 glipizide 0 glyburide 0 tolbutamide 0		<pre>5 print(col,df[col][df[col] == '?'].count())</pre>										
gender 0 age 0 weight 98569 payer_code 40256 medical_specialty 49949 diag_1 21 diag_2 358 diag_3 1423 max_glu_serum 0 A1Cresult 0 metformin 0 repaglinide 0 nateglinide 0 chlorpropamide 0 glimepiride 0 acetohexamide 0 glipizide 0 glyburide 0 tolbutamide 0		race 2273										
age 0 weight 98569 payer_code 40256 medical_specialty 49949 diag_1 21 diag_2 358 diag_3 1423 max_glu_serum 0 A1Cresult 0 metformin 0 repaglinide 0 nateglinide 0 chlorpropamide 0 glimepiride 0 acetohexamide 0 glipizide 0 glyburide 0 tolbutamide 0		gender Ø										
<pre>weight 98569 payer_code 40256 medical_specialty 49949 diag_1 21 diag_2 358 diag_3 1423 max_glu_serum 0 A1Cresult 0 metformin 0 repaglinide 0 nateglinide 0 chlorpropamide 0 glimepiride 0 acetohexamide 0 glipizide 0 tolbutamide 0</pre>		age 0										
payer_code 40256 medical_specialty 49949 diag_1 21 diag_2 358 diag_3 1423 max_glu_serum 0 A1Cresult 0 metformin 0 repaglinide 0 nateglinide 0 chlorpropamide 0 glimepiride 0 acetohexamide 0 glipizide 0 glyburide 0 tolbutamide 0		weight 98569										
<pre>medical_specialty 49949 diag_1 21 diag_2 358 diag_3 1423 max_glu_serum 0 A1Cresult 0 metformin 0 repaglinide 0 nateglinide 0 glimepiride 0 acetohexamide 0 glipizide 0 glyburide 0 tolbutamide 0</pre>		payer_code 40256										
diag_1 21 diag_2 358 diag_3 1423 max_glu_serum 0 A1Cresult 0 metformin 0 repaglinide 0 nateglinide 0 chlorpropamide 0 glimepiride 0 acetohexamide 0 glipizide 0 glyburide 0 tolbutamide 0		medical_specialty 49949										
diag_2 358 diag_3 1423 max_glu_serum 0 A1Cresult 0 metformin 0 repaglinide 0 nateglinide 0 glimepiride 0 acetohexamide 0 glipizide 0 glyburide 0 tolbutamide 0		diag_1 21										
alag_3 1423 max_glu_serum 0 A1Cresult 0 metformin 0 repaglinide 0 nateglinide 0 glimepiride 0 acetohexamide 0 glipizide 0 glyburide 0 tolbutamide 0		diag_2 358										
max_glu_serum 0 AlCresult 0 metformin 0 repaglinide 0 nateglinide 0 glimepiride 0 acetohexamide 0 glipizide 0 glyburide 0 tolbutamide 0		diag_3 1423										
Michaelini o metformin o repaglinide o nateglinide o glimepiride o acetohexamide o glipizide o glyburide o tolbutamide o		max_giu_serum 0										
repaglinide 0 nateglinide 0 chlorpropamide 0 glimepiride 0 acetohexamide 0 glipizide 0 glyburide 0 tolbutamide 0		Alcresult 0										
nateglinide 0 chlorpropamide 0 glimepiride 0 acetohexamide 0 glipizide 0 glyburide 0 tolbutamide 0		nenaglinide 0										
chlorpropamide 0 glimepiride 0 acetohexamide 0 glipizide 0 glyburide 0 tolbutamide 0		repagninide 0										
glimepiride 0 acetohexamide 0 glipizide 0 glyburide 0 tolbutamide 0		chlorpronamide 0										
acetohexamide 0 glipizide 0 glyburide 0 tolbutamide 0		glimepiride 0										
glipizide 0 glyburide 0 tolbutamide 0		acetohexamide 0										
glyburide 0 tolbutamide 0		glipizide 0										
tolbutamide 0		glyburide 0										
		tolbutamide 0										
	in [100]:	<pre>1 # gender was coded differently so we use a custom count for this one 2 print('gender', df['gender'][df['gender'] == 'Unknown/Invalid'].count())</pre>										
<pre>[100]: 1 # gender was coded differently so we use a custom count for this one 2 print('gender', df['gender'][df['gender'] == 'Unknown/Invalid'].count())</pre>		gender 3										
<pre>[100]: 1 # gender was coded differently so we use a custom count for this one 2 print('gender', df['gender'][df['gender'] == 'Unknown/Invalid'].count()) gender 3</pre>	[101]	1 #deepping columns with lange number of missing values										

In [101]:	1 2	<pre>#dropping columns with large number of missing values df = df.drop(['weight','payer_code','medical_specialty'], axis = 1)</pre>
In [102]:	1 2 3 4 5 6 7 8 9 10	<pre>drop_Idx = set(df[(df['diag_1'] == '?') & (df['diag_2'] == '?') & (df['diag_3'] == '?')].index) drop_Idx = drop_Idx.union(set(df['diag_2'][df['diag_2'] == '?'].index)) drop_Idx = drop_Idx.union(set(df['diag_3'] == '?'].index)) drop_Idx = drop_Idx.union(set(df['diag_3'][df['diag_3'] == '?'].index)) drop_Idx = drop_Idx.union(set(df['diag_3'][df['diag_3'] == '?'].index)) drop_Idx = drop_Idx.union(set(df['discharge_disposition_id'] == 11].index)) drop_Idx = list(set(df.index) - set(drop_Idx)) df = df.iloc[new_Idx]</pre>
In [103]:	1	<pre>df = df.drop(['citoglipton', 'examide'], axis = 1)</pre>
In [104]:	1 2 3 4 5 6	<pre>#Checking for missing values in the data for col in df.columns: if df[col].dtype == object: print(col,df[col][df[col] == '?'].count()) print('gender', df['gender'][df['gender'] == 'Unknown/Invalid'].count())</pre>
	rac gen dia dia dia dia max A1C	e 0 der 0 g_1 0 g_2 0 g_3 0 gglu_serum 0 result 0 formin 0

6.4 Data Transformation

After the data is cleaned there are steps to transform data so that it becomes easy for the machine learning algorithms to process it. The steps involved in pre-processing are feature engineering, feature encoding, feature selection, sampling, and scaling of data.

Encoding the outcome variable: The outcome we are looking at is whether the patient gets readmitted to the hospital within 30 days or not. The variable actually has < 30, > 30 and No Readmission categories. To reduce our problem to a binary classification, we combined the readmission after 30 days and no readmission into a single category

In [117]:	1 2 3	<pre>df['readmitted'] = df['readmitted'].replace('>30', 0) df['readmitted'] = df['readmitted'].replace('<30', 1) df['readmitted'] = df['readmitted'].replace('NO', 0)</pre>
In [118]:	1 2 3 4 5 6 7	<pre># Creating additional columns for diagnosis df['level1_diag1'] = df['diag_1'] df['level2_diag1'] = df['diag_1'] df['level2_diag2'] = df['diag_2'] df['level2_diag2'] = df['diag_2'] df['level1_diag3'] = df['diag_3'] df['level2_diag3'] = df['diag_3']</pre>
In [119]:	1 2 3 4 5 6 7 8 9 10 11 12	<pre>df.loc[df['diag_1'].str.contains('V'), ['level1_diag1', 'level2_diag1']] = 0 df.loc[df['diag_1'].str.contains('E'), ['level1_diag2', 'level2_diag2']] = 0 df.loc[df['diag_2'].str.contains('V'), ['level1_diag2', 'level2_diag2']] = 0 df.loc[df['diag_2'].str.contains('E'), ['level1_diag2', 'level2_diag2']] = 0 df.loc[df['diag_3'].str.contains('F'), ['level1_diag3', 'level2_diag2']] = 0 df.loc[df['diag_3'].str.contains('E'), ['level1_diag3', 'level2_diag3']] = 0 df.loc[df['diag_3'].str.contains('E'), ['level1_diag3', 'level2_diag3']] = 0 df.loc[df['diag_3'].str.contains('E'), ['level1_diag3', 'level2_diag3']] = 0 df['level1_diag1'] = df['level1_diag1'].replace('P', -1) df['level2_diag1'] = df['level2_diag2'].replace('P', -1) df['level1_diag2'] = df['level2_diag2'].replace('P', -1) df['level1_diag2'] = df['level2_diag3'].replace('P', -1) df['level2_diag3'] = df['level2_diag3'].replace('P', -1)</pre>

Feature encoding is done to encode the values of features into fewer categories such that it becomes easier for the ML algorithms to understand data as shown in the figure.

Data Transformation

Feature Engineering

Service Utilization

Service utilization: The data contains variables for number of inpatient (admissions), emergency room visits and outpatient visits for a given patient in the previous one year. These are (crude) measures of how much hospital/clinic services a person has used in the past year. We added these three to create a new variable called service utilization (see figure below).

]: 1	<pre>1 df['service_utilization'] = df['number_outpatient'] + df['number_emergency'] + df['number_inpatient']</pre>										
:	df.head(10).T										
:		1	2	3	4	5	6	7	8	9	10
0	encounter_id	149 1 90	64410	500364	16680	35754	55842	63768	12522	15738	28236
	patient_nbr	55629189	86047875	82442376	42519267	82637451	84259809	114882984	48330783	63555939	89869032
	race	Caucasian	AfricanAmerican	Caucasian	AfricanAmerican						
	gender	Female	Female	Male	Male	Male	Male	Male	Female	Female	Female
	age	[10-20)	[20-30)	[30-40)	[40-50)	[50-60)	[60-70)	[70-80)	[80-90)	[90-100)	[40-50)
	admission_type_id	1	1	1	1	2	3	1	2	3	1
di	scharge_disposition_id	1	1	1	1	1	1	1	1	3	1
	admission_source_id	7	7	7	7	2	2	7	4	4	7
	time_in_hospital	3	2	2	1	3	4	5	13	12	9
	num_lab_procedures	59	11	44	51	31	70	73	68	33	47

The figure shows the engineering of the new feature 'service_utilization' by combining the existing feature in the dataset.

In [160]:	<pre>1 # Importing Libraries for SMOTE sampling 2 from imblearn.over_sampling import SMOTE 3 from collections import Counter 4 print('Original dataset shape {}'.format(Counter(y_train)))</pre>											
C	Original dataset shape Counter({0: 43711, 1: 4053})											
In [161]:	1 sm = SMOTE(random_state=20)											
In [162]:	1 pip installupgrade scikit-learn											
r F (F F T T T T T T T T T T T T T T T T	<pre>n) (1.19.5) Requirement already satisfied, skipping upgrade: joblib>=0.11 in c:\users\silky\anaconda3\lib\site-packages (from scikit-learn) (0.16.0) Requirement already satisfied, skipping upgrade: scipy>=0.19.1 in c:\users\silky\anaconda3\lib\site-packages (from scikit-learn) (1.5.0) Requirement already satisfied, skipping upgrade: threadpoolctl>=2.0.0 in c:\users\silky\anaconda3\lib\site-packages (from scikit t-learn) (2.1.0) Note: you may need to restart the kernel to use updated packages.</pre>											
Out[163]:					2							
-	encounter id	1/0100	64410	500364	16680	35754						
	patient nbr	55629189	86047875	82442376	42519267	82637451						
	race	Caucasian	AfricanAmerican	Caucasian	Caucasian	Caucasian						
	gender	0	0	1	1	1						
	age	2	3	4	5	6						
	admission_type_id	1	1	1	1	1						

The figure shows the SMOTE oversampling where the target variable 'readmitted' is having imbalanced data which is balanced using the sampling technique.

Features are selected based on the important feature identified by the Random Forest algorithm and then those features are pre-processed and used for building the model as shown in the figures below.



Feature Scaling

In [141]:	1 2 3 4 5	<pre># Feature Scaling datf = pd.DataFrame() datf['features'] = numerics datf['std_dev'] = datf['features'].apply(lambda x: df[x].std()) datf['mean'] = datf['features'].apply(lambda x: df[x].mean())</pre>	
In [142]:	1 2 3	<pre># dropping multiple encounters while keeping either first or last encounter of these patients df2 = df.drop_duplicates(subset= ['patient_nbr'], keep = 'first') df2.shape</pre>	
Out[142]:	(67	580, 55)	
In [143]:	1 2 3	<pre># standardize function def standardize(raw_data): return ((raw_data - np.mean(raw_data, axis = 0)) / np.std(raw_data, axis = 0))</pre>	
In [144]:	1 2 3	<pre>df2[numerics] = standardize(df2[numerics]) import scipy as sp df2 = df2[(np.abs(sp.stats.zscore(df2[numerics])) < 3).all(axis=1)]</pre>	

Feature Scaling is done to normalize or standardize the feature values in the attributes, so the results are optimized.



Outliers are detected and plotted in the form of a box plot and are removed as it influences the results giving biased predictions.

6.5 Data Modelling

The most important step of the data mining process is building a model. For building model, Scikit-learn machine learning library is used which have packages for data pre-processing, transformation, building model and evaluation metrics.

The base classifiers build for this project are – Random Forest, Decision Tree, k Nearest Neighbors, Support Vector Machines, Logistic Regression, Gradient Boosting and Extreme

Gradient Boosting. Then based on the performance of each classifier, the best models are selected for building a stacked ensemble model.



This Figure shows the implementation of the Random Forest Classifier and the results obtained by using the testing data.



This figure shows the implementation of the Decision Tree Classifier and the results obtained by using the testing data.



This figure shows the implementation of k Nearest Neighbors, where the value of 'k' is selected optimally by implementing the values from 1 to 15 as shown in the plot and the results obtained by using the testing data while building a classifier for k=2.



This figure shows the implementation of the Support Vector Machine file and the results obtained by using the testing data, there were other experiments also done for other kernels like polynomial, sigmoid which is there in the code.

In [186]:	<pre>1 train_input_new = pd.DataFrame(train_input_new, columns = list(X.columns)) 2 from sklearn.model_selection_import_train_test_split</pre>									
	3 from sklearn.linear_model import LogisticRegression									
	4 from sklearn.model_selection import cross_val_score									
	5 X train, X test, y train, y test = train test_split(train_input_new, train_output_new, test_size=0.50, random_state=0)									
	6 # Logistic Regression with the L1 penalty 7 logit = LogisticRegression(fit intercept=True, penalty='11', solver='liblinear')									
	7 logit = LogisticRegression(fit_intercept=True, penalty='11', solver='liblinear')									
	8 logit.fit(X_train, y_train)									
Out[186]:	LogisticRegression(penalty='11', solver='1iblinear')									
In [187]:	1 # Logistic Regression prediction									
	2 logit_pred2 = logit.predict(X_test)									
	3 pd.crosstab(pd.Series(y_test, name = 'Actual'), pd.Series(logit_pred2, name = 'Predict'), margins = True)									
Out[187]:										
	Predict 0 1 All									
	Actual									
	0 9866 10004 19870									
	1 856 964 1820									
	All 10722 10968 21690									
In [188]:	<pre>1 print_report(y_test, logit_pred2, thresh)</pre>									
	AUC:0.745									
	accuracy:0.745									
	recall:0.747									
	precision:0.746									
	fscore:0.747									
	specificity:0.744									
Out[188]:	(0.7453666196782243,									
	0.7453730182333966,									
	0.7469467736055414,									
	0.7461307356154406,									
	0.7465385316686719,									
	0.7437804057509073)									

This figure shows the implementation of Logistic Regression and the results obtained by using the testing data.



This figure shows the implementation of Extreme Gradient Boosting and the results obtained by using the testing data.

Gradient Boost

```
In [44]: 1 # Gradient Boost mode
                2 gradient_model = GradientBoostingClassifier(random_state=42)
3 gradient_model.fit(X_train, y_train)
Out[44]: GradientBoostingClassifier(random state=42)
In [45]: 1 # Training and validation prediction from Gradient Boosting
2 gb_train_pred = gradient_model.predict_proba(X_train)[i,1]
3 gb_vala_pred = gradient_model.predict_proba(X_val)[i,1]
                5 print("Gradient Boosing")
                    print(
                              Training:')
                   gbc_train_auc, gbc_train_accuracy, gbc_train_recall, gbc_train_precision,gbc_train_fscore, gbc_train_specificity = print_rep
print('Validation:')
                gbc_val_auc, gbc_val_accuracy, gbc_val_recall, gbc_val_precision, gbc_val_fscore, gbc_val_specificity = print_report(y_val,g
             Gradient Boosing
              Training:
AUC:0.697
             AUC:0.697
accuracy:0.641
recall:0.604
precision:0.649
fscore:0.626
specificity:0.678
             Validation:
              AUC:0.663
             accuracy:0.626
recall:0.595
             precision:0.644
fscore:0.619
              specificity:0.658
```

This figure shows the implementation of Gradient Boosting and the results obtained by using the testing data.

Then finally the stacked ensemble model is built by selecting only the best performing models and then the results of those base models are used to train the meta classifier Random Forest and final prediction results are obtained. The 'vecstack' package is used to stack together with the base models.

Stacked Ensemble Model - Final Model



In [207]: 1 !pip install vecstack

This figure shows the stacking of models for hospital readmission prediction.



Stacked Ensemble Model - Final Model

This figure shows the stacking of models for the length of stay prediction.

In [209]:	1 #	Trainin	g the Stacked Ense	emble model with 5 cross vaidations for the list of models selected						
	2 5	_train,	5_test = stacking	models,						
	3			X_train, y_train, X_test,						
	4			regression=False,						
	5			<pre>mode='oof_pred_bag', needs_proba=False,</pre>						
	6									
	/									
	8									
	9			and the law second						
	10			save_dir=None,						
	11									
	12			metric=accuracy_score,						
	15									
	14			Π_ΤΟΙΔ5=>,						
	15			stastifi adaTuus						
	17			stratified=frue,						
	18			chufflo-True						
	10			Shuffie-Inde,						
	20			pandom stato=0						
	21									
	22			verbose=2)						
	task:		[classification]							
	n clas	sses:	[2]							
	metric	= :	[accuracy score]							
	mode:		[oof pred bag]							
	n mode	els:	[7]							
	model	0:	[GradientBoosting	Classifier]						
	fo	old 0:	[0.88184834]							
	fo	old 1:	[0.88458019]							
	fo	old 2:	[0.87714482]							
	fo	old 3:	[0.89212995]							
	fo	old 4:	[0.88469458]							
	ME	EAN:	[0.88407957] + [0	.00486840]						

Then the base models are trained using 5-fold cross-validation where the first four folds are used to train the data and the last fold is used to test the data and get the final predictions which are then augmented to the dataset. This augmented dataset is then used to train the meta classifier Random Forest which is shown in the below figure.

```
In [210]: 1 # initializing generalizer model i.e., RF classifier in our case
model = RandomForestClassifier(n_estimators = 10, max_depth=25, criterion = "gini", min_samples_split=10)
3 model = model.fit(S_train, y_train)
5 y_pred = model.predict(S_test)
6 print('Final prediction score: [%.8f]' % accuracy_score(y_test, y_pred))
Final prediction score: [0.94397291]
```



This figure shows the prediction results of the hospital readmission prediction using the Stacked Ensemble model.

In [139]:	1	# Results and	confusio	on matrix	for the	tacked er	semble m	odel				
	2	CM=confusion_	matrix(y	_test,y_p	red) =="VlGnBu!	fmt-'a'	5					
	4	sits.neachap(c	n, annoc	- True, cilla	- 110100	Junc- 8	/					
	5	TN = CM[0][0]										
	6	FN = CM[1][0]										
	7	TP = CM[1][1]										
	9	specificity =	TN/(TN+	FP)								
	10	loss log = lo	g loss(y	test, y p	pred)							
	11	<pre>11 acc=accuracy_score(y_test, y_pred) 12 roc=roc_auc_score(y_test, y_pred) 13 prec = precision_score(y_test, y_pred) 14 rec = recall_score(y_test, y_pred) 15 f1 = f1 score(y_test, y_pred)</pre>										
	12											
	13											
	15											
	16		())	1-6/								
	17	mathew = matt	hews_cori	<pre>rcoef(y_t</pre>	est, y_pre	ed)						
	18	model_results	=pd.Data	aFrame([[STacked (lassifier	2',acc,	prec,rec	,specific	ity, f1,roc, lo	ss_log,mathew]],	
	20		COLUMN	s = [mode	er, Acci	iracy , Pr	ecision	, Sensi	tivity,	specificity,	ri score, kuc, Log_	Loss , mathew_
	21	model results										
		-										
Out[139]:		Model	Ассигасу	Precision	Sensitivity	Specificity	E1 Score	ROC		mathew corrcoef		
	-											
	0	STacked Classifier2	0.901444	0.900973	0.902794	0.900084	0.901883	0.901439	3.404057	0.802885		
					- 1	6000						
						4000						
	0	17242		1914	- 1	4000						
					- 1	2000						
					- 1	0000						
					- 8	000						
	-	1875		17414	- e	000						
		1875		17414	- 6	000						
	-	1875		17414	- e - 4	000						

This figure shows the prediction results of the length of stay prediction using the Stacked Ensemble model.

7 Conclusion

All the pre-requisites – hardware and software configuration and requirement with the library and packages used for building models have been described in this document. This report explains the complete project development process in a structured, concise, and detailed manner which will help in understanding the flow of implementation.

References

Strack, B. *et al.* (2014) 'Impact of HbA1c Measurement on Hospital Readmission Rates: Analysis of 70,000 Clinical Database Patient Records', *BioMed Research International*. Edited by A. Rizvi, 2014, p. 781670. doi: 10.1155/2014/781670.