

Configuration Manual

MSc Research Project
Data Analytics

Aditya Raj
Student ID: x20143311

School of Computing
National College of Ireland

Supervisor: Dr. Martin Alain

National College of Ireland
Project Submission Sheet
School of Computing



Student Name:	Aditya Raj
Student ID:	x20143311
Programme:	Data Analytics
Year:	2021
Module:	MSc Research Project
Supervisor:	Dr. Martin Alain
Submission Due Date:	31/01/2021
Project Title:	Configuration Manual
Word Count:	766
Page Count:	11

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.

ALL 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:	Aditya Raj
Date:	31st January 2022

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).	<input type="checkbox"/>
Attach a Moodle submission receipt of the online project submission , to each project (including multiple copies).	<input type="checkbox"/>
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.	<input type="checkbox"/>

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

Aditya Raj
x20143311

1 Introduction

In this configuration handbook, the detailed instructions of implementation are elucidated clearly. The steps pertaining to Tier - I and Tier - II of the KDD data mining process such as Data Selection, Data Preprocessing, Exploratory Data Analysis, Data Transformation, Modelling and Evaluation are supported with Code snippets, Screenshots, and Instructions for execution.

2 Hardware and Software Requirements

Table 1: Hardware Specifications

Device Name/OS	Macbook Pro/macOS Big Sur Version 11.0.1
RAM/CPU	16 GB 1867 MHz DDR3/2.7 GHz Dual-Core Intel Core i5
Hard Disk	256 GB SSD
GPU	Intel Iris Graphics 6100 1536 MB

Table 2: Software Specifications

Programming Language	Python Version 3.9
IDE	Jupyter Notebook
Browser	Google Chrome

3 Data Selection

3.1 Download the IBM HR Analytics Employee Attrition Dataset

- Open the URL <https://www.kaggle.com/pavansubhasht/ibm-hr-analytics-attrition-dataset>, Click on Download to download the dataset on to your system. Once downloaded, place the file in new folder which can be referenced as a file path for data source.

3.2 Installing and Importing Libraries and Data on Jupyter Notebook File

1. Select New and Python 3 from Jupyter Notebook Home to create a new Jupyter notebook file.
2. Install and Import the prerequisite libraries for the research. Use pip install <package-name> to install any libraries.

```
In [ ]: 1 pip install pandas
        2 pip install numpy
```

```
1. Importing Libraries
In [109]: 1 import missingno as msno
          2 from pandas_profiling import ProfileReport
          3
          4 # Numerical transformations
          5 import numpy as np
          6 # Data Processing
          7 import pandas as pd
          8
          9 # Visualisation
         10 import matplotlib.pyplot as plt
         11 from matplotlib import rc
         12 import seaborn as sns
         13 %matplotlib inline
         14 import plotly.offline as py
         15 py.init_notebook_mode(connected=True)
         16 import plotly.graph_objs as go
         17 import plotly.tools as tls
         18 import plotly.figure_factory as ff
         19 pd.options.display.max_columns = None
         20 import warnings
         21 warnings.filterwarnings('ignore')
         22 from matplotlib.colors import ListedColormap
         23
         24 from datetime import datetime
         25 from sklearn.preprocessing import StandardScaler, LabelEncoder
         26 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV, cross_val_score, learning_curve, train_test_s
         27 from sklearn.metrics import precision_score, roc_auc_score, recall_score, confusion_matrix, roc_curve, precision_re
         28 from lightgbm import LGBMClassifier
         29
         30 # IMBlearn
         31 #imblearn
         32 from imblearn.over_sampling import RandomOverSampler
         33 from imblearn.under_sampling import RandomUnderSampler
         34 from imblearn.over_sampling import SMOTE
         35 from imblearn.over_sampling import ADASYN
         36
         37
```

Figure 1: Import Libraries

3.3 Data Understanding

1. The data is imported using pandas - read_csv function.

```
2. Data Understanding
In [110]: 1 df = pd.read_csv("../data/WA_Fn-UseC_HR-Employee-Attrition.csv")
          2 df.head(10)
Out[110]:
```

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EmployeeCount	EmployeeNumber	EnvironmentSatisf
0	41	Yes	Travel_Rarely	1102	Sales	1	2	Life Sciences	1	1	
1	49	No	Travel_Frequently	279	Research & Development	8	1	Life Sciences	1	2	
2	37	Yes	Travel_Rarely	1373	Research & Development	2	2	Other	1	4	
3	33	No	Travel_Frequently	1392	Research & Development	3	4	Life Sciences	1	5	
4	27	No	Travel_Rarely	591	Research & Development	2	1	Medical	1	7	
5	32	No	Travel_Frequently	1005	Research & Development	2	2	Life Sciences	1	8	
6	59	No	Travel_Rarely	1324	Research & Development	3	3	Medical	1	10	
7	30	No	Travel_Rarely	1358	Research & Development	24	1	Life Sciences	1	11	
8	38	No	Travel_Frequently	216	Research & Development	23	3	Life Sciences	1	12	
9	36	No	Travel_Rarely	1299	Research & Development	27	3	Medical	1	13	

Figure 2: Data Import

4 Data Preprocessing

- The data is further checked for validations such as missing values, profiling, outliers detection etc. The missing values were evaluated on the dataset by importing `missingno` library.

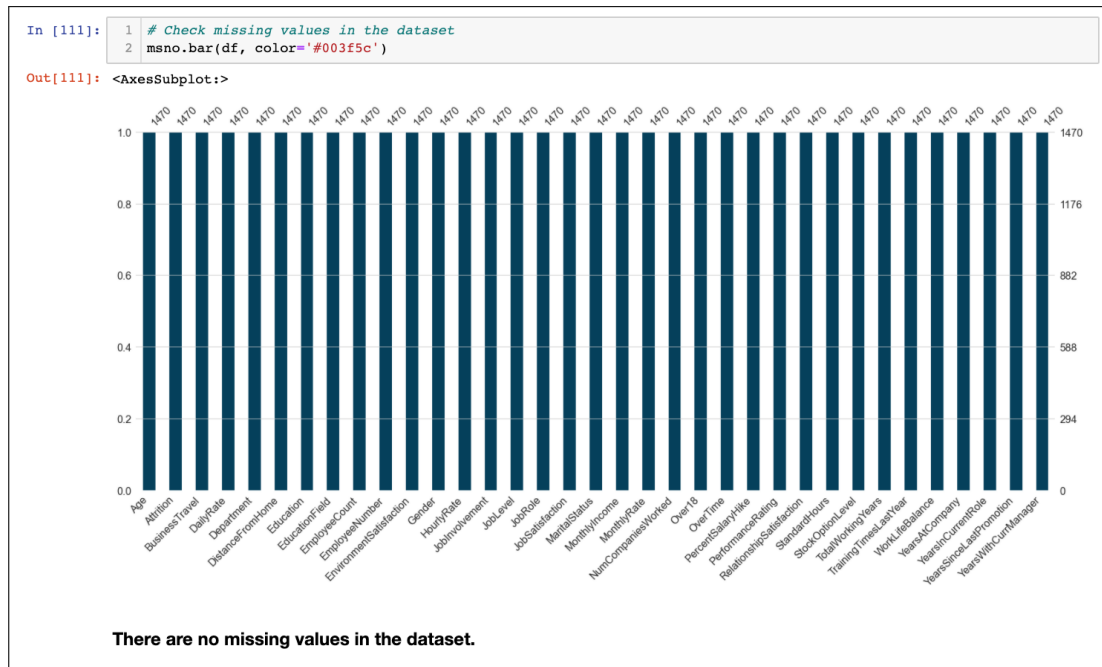


Figure 3: Missing Values

- The data profiling is powered by `pandas-profiling` library that prepared a **ProfileReport** as illustrated in figure fig:Profiling. Based on the insights gained from profiling report, 4 variables are dropped from the dataset, bringing down the number of variables in the dataset to 31.

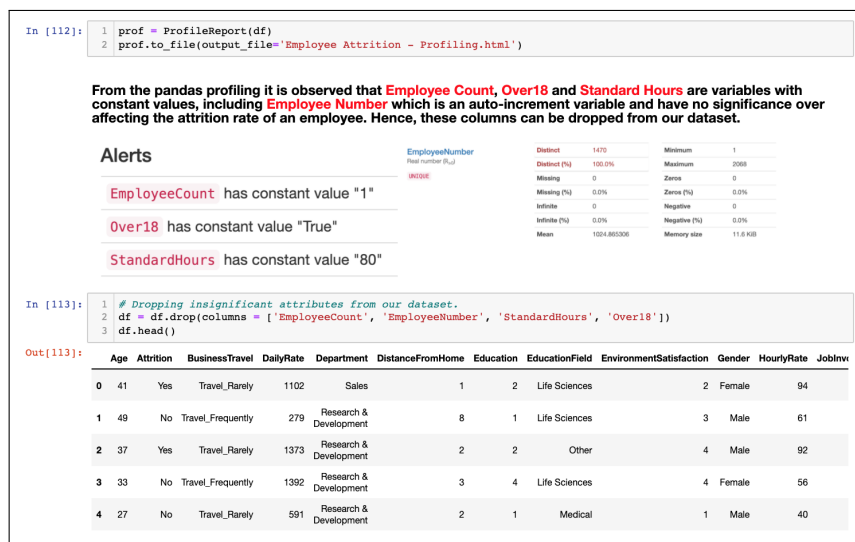


Figure 4: Profiling

3. Outlier Analysis performed using **boxplot** powered by **seaborn** library. There are significant outliers observed in the data that needs to be treated.

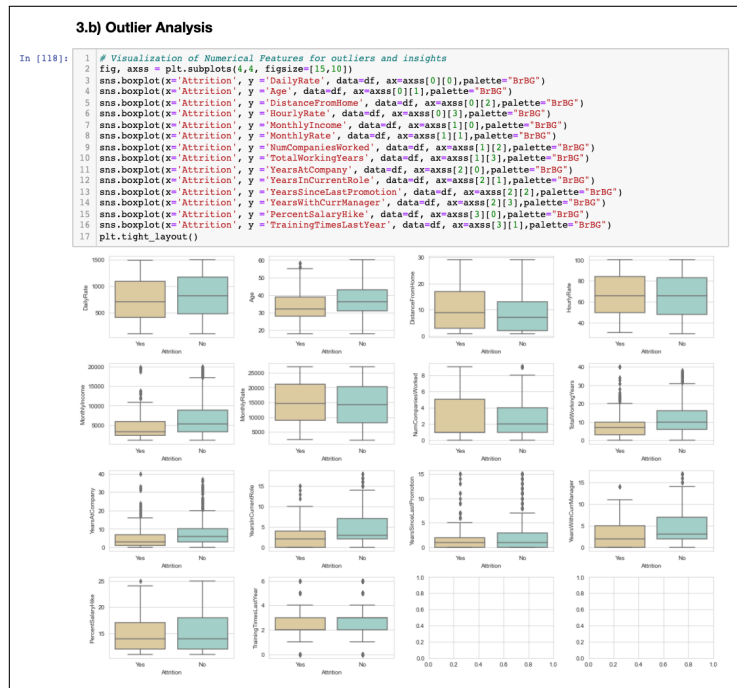


Figure 5: Outlier Analysis

4. Chi-squared test is performed on categorical variable to check their statistical significance. In order to achieve this, **chi2_contingency** was inherited from **scipy.stats** module



Figure 6: Chi-squared Test

5 Exploratory Data Analysis

The data is visualized against the target variable "Attrition" with respect to other features to extract insights. **Pie charts**, **Bar charts** and **Pie-donut** were built using **matplotlib.pyplot**, **countplot** from **seaborn** library, various plots (**bar**, **scatter**, **pie**) from **plotly.graph_objs**.

```

3. Exploratory Data Analysis

3.a) Attrition Rate Analysis

In [116]: 1 # Creating a new dataframe to analyze the Attrition Rate Analysis for different Age Groups
2 def agg(x):
3     names = {
4         'Attritted': x['Age'][x['Attrition']=='Yes'].count(),
5         'Retained': x['Age'][x['Attrition']=='No'].count()
6     }
7     return pd.Series(names)
8 df_AR=pd.DataFrame(df.groupby('Age').apply(agg))
9 df_AR['Attrition Rate by Age']=round(df_AR['Attritted']/df_AR['Retained'],2)
10 df_AR.reset_index(level=0, inplace=True)
11

In [117]: 1 fig, ax = plt.subplots(4,1, figsize=[20,20])
2
3 p1=sns.countplot(df.Age, palette="icefire", ax=ax[0])
4 # p1=plt.title("Distribution of Age Across Dataset")
5 p1.title.set_text("Distribution of Age Across dataset")
6 for i, bar in enumerate(p1.patches):
7     h = bar.get_height()
8     p1.text(
9         i, # bar index (x coordinate of text)
10        h*0.5, # y coordinate of text
11        '{}'.format(int(h)), # y label
12        ha='center',
13        va='center',
14        fontweight='light',
15        size=12)
16
17 p2=sns.countplot(df.Age[df['Attrition']=='Yes'], palette="flare", ax=ax[1])
18 p2.title.set_text("Distribution of Age Across Attrition=Yes")
19
20 for i, bar in enumerate(p2.patches):
21     h = bar.get_height()
22     p2.text(
23         i, # bar index (x coordinate of text)
24        h*0.45, # y coordinate of text
25        '{}'.format(int(h)), # y label
26        ha='center',
27        va='center',
28        fontweight='light',
29        size=12)
30
31
32 p3=sns.countplot(df.Age[df['Attrition']=='No'], palette="light:#5A9", ax=ax[2])
33 p3.title.set_text("Distribution of Age Across Attrition=No")
34

```

Figure 7: Attrition Rate Analysis Code

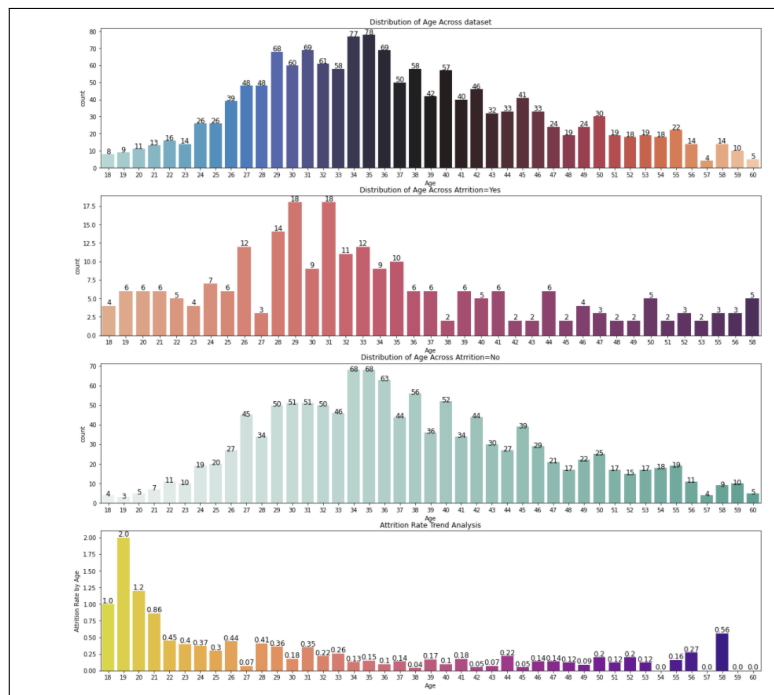


Figure 8: Attrition Rate Analysis

6 Data Preparation

6.1 Feature Engineering

In this research, 11 new features were engineered from the existing variables and their relationship with the target variable were analyzed. Based on the analysis, the following new features created were SalesDpt, RDDpt ModJobInv, ModTraining, MeanSatisfaction, OverSatRating, LongDis, Hrate Mrate, Stability, TotalCompWorked and Loyalty.

```
In [151]:
1 def SalesDpt(df):
2     if df['Department'] == 'Sales':
3         return 1
4     else:
5         return 0
6 df['SalesDpt'] = df.apply(lambda df:SalesDpt(df),axis = 1)
7
8 def RDDpt(df):
9     if df['Department'] == 'Research & Development':
10        return 1
11    else:
12        return 0
13 df['RDDpt'] = df.apply(lambda df:RDDpt(df),axis = 1)
14
15 def ModJobInv(df):
16     if df['JobInvolvement'] < 3 :
17         return 1
18     else:
19         return 0
20 df['ModJobInv'] = df.apply(lambda df:ModJobInv(df),axis = 1)
21
22 def ModTraining(df):
23     if df['TrainingTimesLastYear'] >= 2 and df['TrainingTimesLastYear'] <= 4:
24         return 1
25     else:
26         return 0
27 df['ModTraining'] = df.apply(lambda df:ModTraining(df),axis = 1)
28
29 df['MeanSatisfaction'] = (df['RelationshipSatisfaction'] + df['EnvironmentSatisfaction'] + df['JobSatisfaction'])/3
30
31 def OverSatRating(df):
32     if df['MeanSatisfaction'] > 2.7 :
33         return 1
34     else:
35         return 0
36 df['OverSatRating'] = df.apply(lambda df:OverSatRating(df),axis = 1)
37
38 def LongDis(df):
39     if df['DistanceFromHome'] > 9:
40         return 1
41     else:
42         return 0
43 df['LongDis'] = df.apply(lambda df:LongDis(df),axis = 1)
44 df['Hrate_Mrate'] = df['HourlyRate']/df['MonthlyRate']
45 df['Stability'] = df['YearsInCurrentRole']/df['YearsAtCompany']
46 df['Stability'] = df['Stability'].fillna(df['Stability'].mean(), inplace=True)
47 def TotalCompWorked(df):
48     if df['NumCompaniesWorked'] <= 0:
49         return df['NumCompaniesWorked']
50     else:
51         return 1
52 df['TotalCompWorked'] = df.apply(lambda df:TotalCompWorked(df),axis = 1)
53 df['Loyalty'] = (df['TotalCompWorked'])/df['TotalWorkingYears']
54 df['Loyalty'] = df['Loyalty'].replace(np.nan, 0)
55
```

Figure 9: Feature Engineering

6.2 Feature Encoding and Scaling

All the binary columns (Features with 2 unique values) were converted into numerical values using Label Encoding. All the other category columns which had less than 10 unique values within the entire dataset were subjected to dummy encoding. The numerical columns were identified from the data and standard scaling was performed. **Feature Encoding** used `LabelEncoder` and `get_dummies` from `sklearn.preprocessing` and `pandas` package respectively. The normalization of values was performed using `StandardScaler` from `sklearn.preprocessing`.

```
5.a) Features Encoding and Scaling]]
In [153]:
1 #Target column
2 target_col = ['Attrition']
3 #Categorical columns
4 cat_cols = df.unique()[df.unique() == 1].keys().tolist()
5 cat_cols = [x for x in cat_cols if x not in target_col]
6 #Numerical columns
7 num_cols = [x for x in df.columns if x not in cat_cols + target_col]
8 #Binary columns with 2 values
9 bin_cols = df.unique()[df.unique() == 2].keys().tolist()
10 bin_cols = [x for x in bin_cols if x not in target_col]
11 #Columns more than 2 values
12 multi_cols = [x for x in cat_cols if x not in bin_cols]
13
14
In [154]:
1 print("\033[0mMultiple Category Columns:\033[0m",multi_cols)
2 print("\033[0mBinary Columns:\033[0m",bin_cols)
3 print("\033[0mNumeric Columns:\033[0m",num_cols)
4
Multiple Category Columns: ['BusinessTravel', 'EducationField', 'JobLevel', 'JobRole', 'MaritalStatus', 'StockOptionLevel']
Binary Columns: ['Overtime', 'SalesDpt', 'RDDpt', 'ModJobInv', 'ModTraining', 'OverSatRating', 'LongDis']
Numeric Columns: ['Age', 'DailyRate', 'MonthlyIncome', 'PercentSalaryHike', 'YearsSinceLastPromotion', 'YearsWithCurrManager', 'MeanSatisfaction', 'Hrate_Mrate', 'Stability', 'Loyalty']
5
In [155]:
1 # Using Label Encoder for Binary Variables
2 le = LabelEncoder()
3 for i in bin_cols:
4     df[i] = le.fit_transform(df[i])
5
6 # Duplicating columns for multi value columns
7 df = pd.get_dummies(data = df,columns = multi_cols)
8
9 df.head()
10 # Scaling Numerical columns
11 std = StandardScaler()
12 scaled = std.fit_transform(df[num_cols])
13 scaled = pd.DataFrame(scaled,columns=num_cols)
14
15 # dropping original values merging scaled values for numerical columns
16 df.drop(columns = num_cols,axis = 1)
17 df = df.drop(columns = num_cols,axis = 1)
18 df = df.merge(scaled,left_index=True,right_index=True,how = 'left')
19
```

Figure 10: Feature Encoding and Scaling

6.3 Correlation Matrix

It is a visual representation of correlation coefficients between independent variables in the form of a matrix. Based on the threshold, 3 variables were found to be multicollinear and hence, dropped from the dataset. The correlation between the variables was produced using `corr` function from `pandas.DataFrame` package and illustrated using `create_distplot_fromplotly.figure` factory.

6.4 Train and Test Split, Handling Class Imbalance

The data was then split into train and test set using `train_test_split` imported from `sklearn.model_selection` library. Once the training and test samples were created, the class balancing techniques were incorporated using `RandomOversampler` and `SMOTE`, functions imported from `imblearn.over sampling`.

```
5.c) Train and Test Split
In [59]: 1 X = df.loc[:, df.columns == "Attrition"]
          2 y = df.loc[:, df.columns == "Attrition"]

In [60]: 1 #Split the data 80 - 20 train/test using stratifying Split
          2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, stratify=y, random_state=1)

In [61]: 1 #Save Split Data
          2 with open('data/X_train.pickle', 'wb') as f:
          3     pickle.dump(X_train,f)
          4 with open('data/y_train.pickle', 'wb') as f:
          5     pickle.dump(y_train,f)
          6 with open('data/X_test.pickle', 'wb') as f:
          7     pickle.dump(X_test,f)
          8 with open('data/y_test.pickle', 'wb') as f:
          9     pickle.dump(y_test,f)

5.d) Random Oversampling and SMOTE Oversampling
In [62]: 1 # Show add some random oversampling of the minority classes using imblearn
          2 ros = RandomOversampler(random_state=0)
          3 X_resampled_up, y_resampled_up = ros.fit_resample(X_train,y_train)

In [63]: 1 #Save Oversampled Minority Class
          2 with open('data/X_resampled_up.pickle', 'wb') as f:
          3     pickle.dump(X_resampled_up,f)
          4 with open('data/y_resampled_up.pickle', 'wb') as f:
          5     pickle.dump(y_resampled_up,f)

In [64]: 1 oversampler=SMOTE(random_state=0)
          2 X_train_smote, y_train_smote = oversampler.fit_resample(X_train,y_train)

In [65]: 1 #Save SMOTE Oversampled Minority Class
          2 with open('data/X_smote_up.pickle', 'wb') as f:
          3     pickle.dump(X_train_smote,f)
          4 with open('data/y_smote_up.pickle', 'wb') as f:
          5     pickle.dump(y_train_smote,f)

In [66]: 1 print("Original shape", X_train.shape, y_train.shape)
          2 print("Random Upsampled Shape", X_resampled_up.shape, y_resampled_up.shape)
          3 print("SMOTE Upsampled Shape", X_train_smote.shape, y_train_smote.shape)

Original shape: (1176, 44) (1176, 1)
Random Upsampled Shape: (1972, 44) (1972, 1)
SMOTE sampled Shape: (1972, 44) (1972, 1)
```

Figure 12: Train and Test Split Oversampling

7 Modelling

Once the data was ready for model training, the baseline models `RandomForestClassifier` was imported from `sklearn.ensemble` and `LGBMClassifier` was imported from `lightgbm`. Random Forest Classifier and LightGBM were the proposed models for the research in combination with ensemble methods for hyperparameter tuning such as `GridSearchCV` and `Recursive Feature Elimination`. The approaches defined within the Modelling section different libraries adhering to the proposed research were imported such as `make_pipeline` from `sklearn.pipeline` to be used an estimator for `GridSearchCV` imported from `sklearn.model_selection` to identify the best parameters and best Cross validation score. Another approach pertaining to `Recursive Feature Elimination` required to import `RFECV` from `sklearn.feature selection`.

6.a) Random Forest Classifier using Grid Search Cross Validation and Selection of Best Hyperparameters - Model 1

Random Forest Classifier is implemented using 10-fold Grid Search Cross Validation on our original and artificially oversampled datasets to train our model. We will assess the scores to select the optimum dataset for model training and choose best parameters based on the results.

```
In [70]: 1 # Build random forest classifier
2 def random_forest_classifier(X_train, y_train):
3     methods_data = {"Original": (X_train, y_train),
4                     "Up-Resampled": (X_resampled_up, y_resampled_up),
5                     "SMOTE Upsampled": (X_train_smote, y_train_smote)}
6
7     for method in methods_data.keys():
8         pip_rf = make_pipeline(StandardScaler(),
9                                RandomForestClassifier(n_estimators=500,
10                                                       class_weight="balanced",
11                                                       random_state=123))
12
13         hyperparam_grid = {
14             "randomforestclassifier_n_estimators": [10, 30, 50, 80, 100, 150, 200, 500],
15             "randomforestclassifier_max_features": ["sqrt", "log2", 0.4, 0.5],
16             "randomforestclassifier_min_samples_leaf": [1, 3, 5],
17             "randomforestclassifier_criterion": ["gini", "entropy"]}
18         gs_rf = GridSearchCV(pip_rf,
19                              hyperparam_grid,
20                              scoring="f1",
21                              cv=10,
22                              n_jobs=-1)
23
24         # Training the model on oversampled folds of training set
25         gs_rf.fit(methods_data[method][0], methods_data[method][1])
26
27         print(f"\033[1m\033[90mThe best hyperparameters for {method} data:")
28         for hyperparam in gs_rf.best_params_.keys():
29             print(hyperparam[hyperparam.find("__") + 2:], ", ", gs_rf.best_params_[hyperparam])
30
31         print(f"\033[1m\033[90mBest {gs_rf.cv}-folds CV F1 score: {gs_rf.best_score_ * 100:.2f}%.")
32
```

```
In [71]: 1 # Running Random Forest for 10-fold Cross Validation
2 random_forest_classifier(X_train, y_train)
```

```
The best hyperparameters for Original data:
criterion : entropy
max_features : sqrt
min_samples_leaf : 5
n_estimators : 80
Best 10-folds CV F1 score: 48.00%.
The best hyperparameters for Up-Resampled data:
criterion : gini
max_features : log2
min_samples_leaf : 1
n_estimators : 80
Best 10-folds CV F1 score: 98.65%.
The best hyperparameters for SMOTE Upsampled data:
criterion : gini
max_features : log2
min_samples_leaf : 1
n_estimators : 100
Best 10-folds CV F1 score: 91.82%.
```

Figure 13: Model 1 RF Classifier GridSearchCV

6.b) Random Forest Classifier using Recursive Feature Elimination with 10-fold Cross Validation - Model 2

Random Forest Classifier is implemented using Recursive Feature Elimination with 10-fold Cross Validation on randomly upsampled data. The main objective of conducting this experiment is to check if the performance can be enhanced using less number of features.

The total number of existing features in the dataset to train the model is 44.

```
In [81]: 1 min_features_to_select=1
2 rf=RandomForestClassifier(n_estimators = 80,criterion = 'gini',max_features = 'log2',min_samples_leaf = 1)
3 rfecv = RFECV(estimator=rf, step=1, cv=10, scoring='f1', verbose=2)
4 rfecv.fit(X_train, y_train)
5 print("Optimal Number : %d" % rfecv.n_features_)
6
7 # Plot Number of Features VS. Cross-validation scores on Train
8 plt.figure(figsize=(10,8))
9 plt.xlabel("Total Features")
10 plt.ylabel("Cross Validation Score (F1 Score)")
11 plt.plot(
12     range(min_features_to_select, len(rfecv.grid_scores_) + min_features_to_select),
13     rfecv.grid_scores_,
14 )
15 plt.show()
```

Fitting estimator with 44 features.
Fitting estimator with 43 features.

Figure 14: Model 2 RF RFECV

6.c) Random Forest Classifier using Manual Hyperparameter Tuning - Model 3

Random Forest Classifier is implemented by performing Hyperparameter Tuning. The parameters of the model are manually tweaked through iterations over array of values for *n_estimators*, *criterion*, *max_depth*, *min_samples_leaf*, *min_samples_split* and *class_weight*. The objective of this experiment is to achieve better F1 score, Accuracy and ROC AUC for Test data.

```
In [92]: 1 # Random Forest Classifier - Model 3 manual tuning hyperparameters
2 estimators = [10,80,100,110,120,130,140,150,200,210,220,230,250,280,300,320,350, 400,430,450,500,530,550,600,630]
3 criterion = ['gini', 'entropy']
4 max_depth = [1,5,10,12,14,15,20,25,30,35,40,45,50,55,60,65]
5 max_features = ['sqrt', 'log2', 'auto']
6 min_samples_leaf = [1,3,5,7,9,10,11,12,13,15,17,19,21]
7 min_samples_split = [2,5,10,15,20,25,30,35]
8 class_weight = ['balanced', 'balanced_subsample', None]
9 train_scores = []
10 test_scores = []
11 for i in estimators:
12     clf_RF = RandomForestClassifier(n_estimators=i,
13                                 criterion='gini',
14                                 max_depth=14,
15                                 max_features='sqrt',
16                                 min_samples_leaf=11,
17                                 min_samples_split=2,
18                                 class_weight='balanced',
19                                 random_state=25)
20     clf_RF.fit(X_train,y_train)
21     train_sc = f1_score(y_train,clf_RF.predict(X_train))
22     test_sc = f1_score(y_test,clf_RF.predict(X_test))
23     test_scores.append(test_sc)
24     train_scores.append(train_sc)
25     print('estimators = ',i, 'Train Score', train_sc, 'test Score', test_sc)
26 plt.plot(estimators,train_scores,label='Train Score',color='ffcc2b',linewidth=2)
27 plt.plot(estimators,test_scores,label='Test Score',color='005777',linewidth=2)
28 plt.xlabel('estimators')
29 plt.ylabel('Score')
30 plt.legend()
31 plt.title('estimators vs score')
32 plt.grid()
```

estimators = 50 Train Score 0.9404106159238858 test Score 0.5901639344262295
estimators = 80 Train Score 0.9450000000000001 test Score 0.6238487384957983
estimators = 100 Train Score 0.9411764705882352 test Score 0.6178861788617885
estimators = 110 Train Score 0.941422118172659 test Score 0.6016260162601627
estimators = 120 Train Score 0.943528235882059 test Score 0.6065573770491803
estimators = 130 Train Score 0.9399399393939399 test Score 0.5853658536585366
estimators = 140 Train Score 0.9399399393939399 test Score 0.5867741935483811
estimators = 150 Train Score 0.9411764705882352 test Score 0.595041223140495
estimators = 200 Train Score 0.939530248825587 test Score 0.576271184440678
estimators = 210 Train Score 0.939530248825587 test Score 0.5714285714285714
estimators = 220 Train Score 0.942528735632184 test Score 0.5714285714285714
estimators = 225 Train Score 0.9420594205942 test Score 0.5714285714285714
estimators = 230 Train Score 0.9415876185721418 test Score 0.5714285714285714
estimators = 250 Train Score 0.9378757515030061 test Score 0.5714285714285714
estimators = 280 Train Score 0.938164493480442 test Score 0.5641025641025642
estimators = 300 Train Score 0.9404702351175587 test Score 0.5666666666666667
estimators = 320 Train Score 0.940949493480442 test Score 0.559322033898305
estimators = 350 Train Score 0.9383458646616543 test Score 0.559322033898305
estimators = 400 Train Score 0.9392875062719518 test Score 0.559322033898305
estimators = 430 Train Score 0.9392875062719518 test Score 0.559322033898305
estimators = 450 Train Score 0.9397590361445785 test Score 0.559322033898305
estimators = 500 Train Score 0.9397590361445785 test Score 0.5546218487394958
estimators = 530 Train Score 0.9392875062719518 test Score 0.5641025641025642
estimators = 550 Train Score 0.9388164493480442 test Score 0.559322033898305
estimators = 600 Train Score 0.9388164493480442 test Score 0.5641025641025642
estimators = 650 Train Score 0.9388164493480442 test Score 0.559322033898305
estimators = 700 Train Score 0.9388164493480442 test Score 0.559322033898305
estimators = 800 Train Score 0.9388164493480442 test Score 0.5546218487394958
estimators = 900 Train Score 0.9388164493480442 test Score 0.5546218487394958
estimators = 1000 Train Score 0.9392875062719518 test Score 0.5546218487394958
estimators = 1050 Train Score 0.9402310396785535 test Score 0.5546218487394958
estimators = 1100 Train Score 0.9397590361445785 test Score 0.5546218487394958
estimators = 1200 Train Score 0.9397590361445785 test Score 0.5546218487394958
estimators = 1300 Train Score 0.9392875062719518 test Score 0.5546218487394958

Figure 15: Model 3 RF Man Hyp



Figure 16: Model 4 LGBM

8 Evaluation

The performance of each classifier is assessed on the basis of F1-score, Accuracy and ROC AUC score. The required libraries for `f1_score`, `accuracy_score` and `roc_auc_score`, `confusion_matrix` and `classification_report` are imported from `sklearn.metrics`.

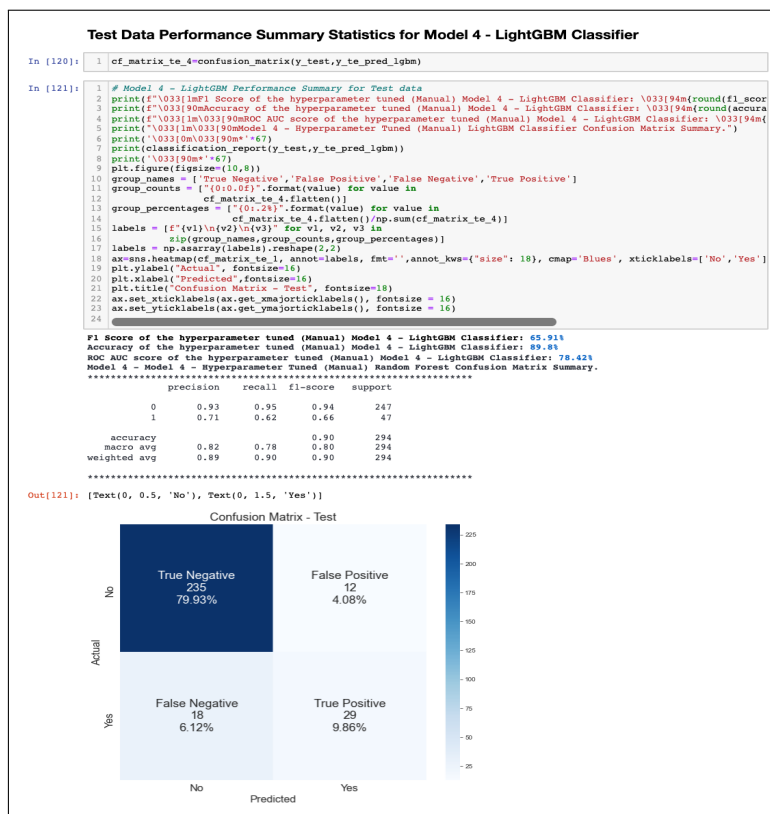


Figure 17: Model 4 LGBM F1

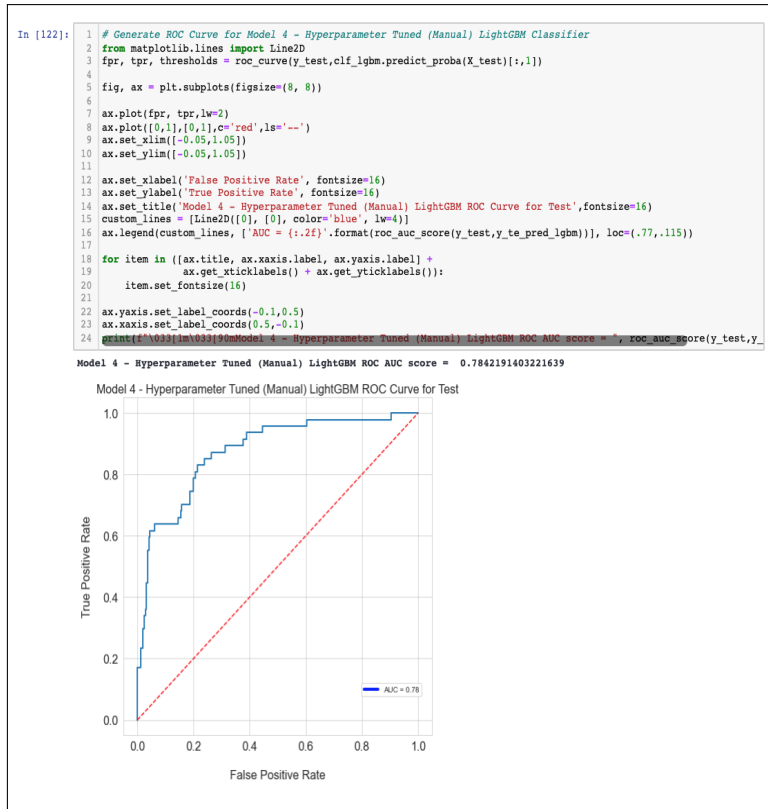


Figure 18: Model 4 LGBM AUC

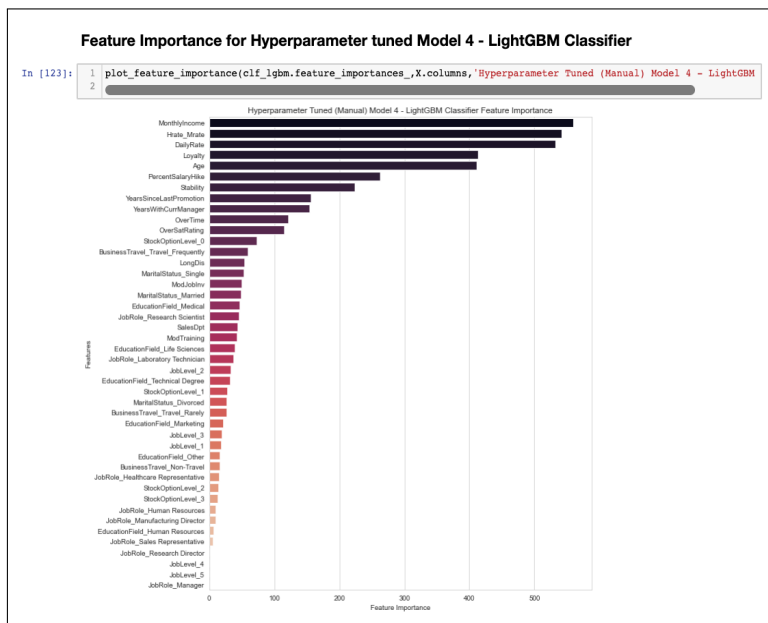


Figure 19: Model 4 LGBM Feature Importance