

IBM MACHINE LEARNING

SUPERVISED MACHINE LEARNING: CLASSIFICATION MODELS FOR EMPLOYEE ATTRITION



Fatima, Sayeda

8/17/2022

Table of Contents

1) Project Overview	2
2) About the Dataset	3
2a) Brief description of the data set you chose:.....	3
2b) Summary of Data Attributes	3
2c) Main Objectives of Analysis	3
3) Data Exploration, Data Cleansing and Features Engineering	4
3a) Data Exploration:	4
3b) Data Cleansing Actions:.....	6
3c) Features Engineering.....	8
4) Summary of Training Different Classifier Models	12
4a) Machine Learning Algorithm Approaches	12
I) Data Level Approaches:	15
II) Algorithm Ensemble Approach:	15
4b) Summarizing Employed Models	16
1) Logistic Regression (LR) Models:	16
2) Random Forest Models:	18
3) XGB Model.....	24
5) Recommended Model	25
5a) Result Summary	25
5b) Overall Visual Summary	25
5c) Individual Model Visual Summary	26
5d) Model Choice and Justification.....	26
6) Summary Key Findings and Insights	27
6a) Summarizing Model Drivers:	27
6b) Enlisting Top Contributory Factors	27
6c) Visualizing Top Contributory Factors to Employee Attrition	28
7) Link to Other Useful Models	29
8) Github Link to Assignment Notebook	29

1) Project Overview

A fundamental issue facing organisations is attraction and retention of best talent. Given the cost of retraining new employees, it is important for a business to prevent loss of good talent. Hence, identification of key factors driving employee churning or turnover is important for the organization's Human Resource (HR) Department.

It is here that machine Learning models can be very useful to gain deeper insight into underlying factors and their relationship in driving employee turnover.

Hence, the main aim of the following machine learning modelling and analysis is to enable the business to:

- * To identify different factors predict employee churn
- * To gain insight into factors contributing to employee churning
- * To enable the business maximize employee attrition

2) About the Dataset

2a) Brief description of the data set you chose:

This project uses a hypothetical dataset 'IBM HR Analytics Employee Attrition & Performance' which was downloaded from the following link:

<https://www.kaggle.com/datasets/pavansubhasht/ibm-hr-analytics-attrition-dataset?resource=download>

2b) Summary of Data Attributes

The dataset exhibits 1,470 data points (rows) and 35 features (columns) reflecting on employees' background and characteristics and can be downloaded from the following link:

The data also comes with 'Attrition' Column to show current employees and leavers which represents the Class we are trying to predict.

2c) Main Objectives of Analysis

Organizational performance is largely dependent on its employees, their quality and experience. Hence, organizations are continuously faced with the challenge to reduce employee attrition and increase retention. Consequently, this analysis is targeted towards answering the following queries

- What are the various factors contributory to employee attrition?
- Which business units face higher employee attrition rate?

As a consequence, implementation of the model will enable the organization to:

- devise suitable measures to increase employee retention
- to save valuable resources in retraining new employees hired in place of leavers

3) Data Exploration, Data Cleansing and Features Engineering

Since the quality of any machine learning model highly depends on quality of data, hence, this stage is not only most important but is also time consuming. Hence, it was conducted in a step-by-step process.

3a) Data Exploration:

- Data was first loaded into pandas dataframe

Load & Read Dataset

```

1 # Load the dataset
2 url = ("C:/Users/fatima.s/Documents/PythonScripts/DATA SCIENCE/IBM Machine Learning Intermediate/MODULE 3 SUPERVISED MACHINE
3 df = pd.read_csv(url, index_col=False) # keep_default_na = False # na_filter=False,
4 df.head(100)
5

```

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EmployeeCount	EmployeeNumber	...	Relationship
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	...	

- Column types were explored

Check data set column types

```

1 df.info()

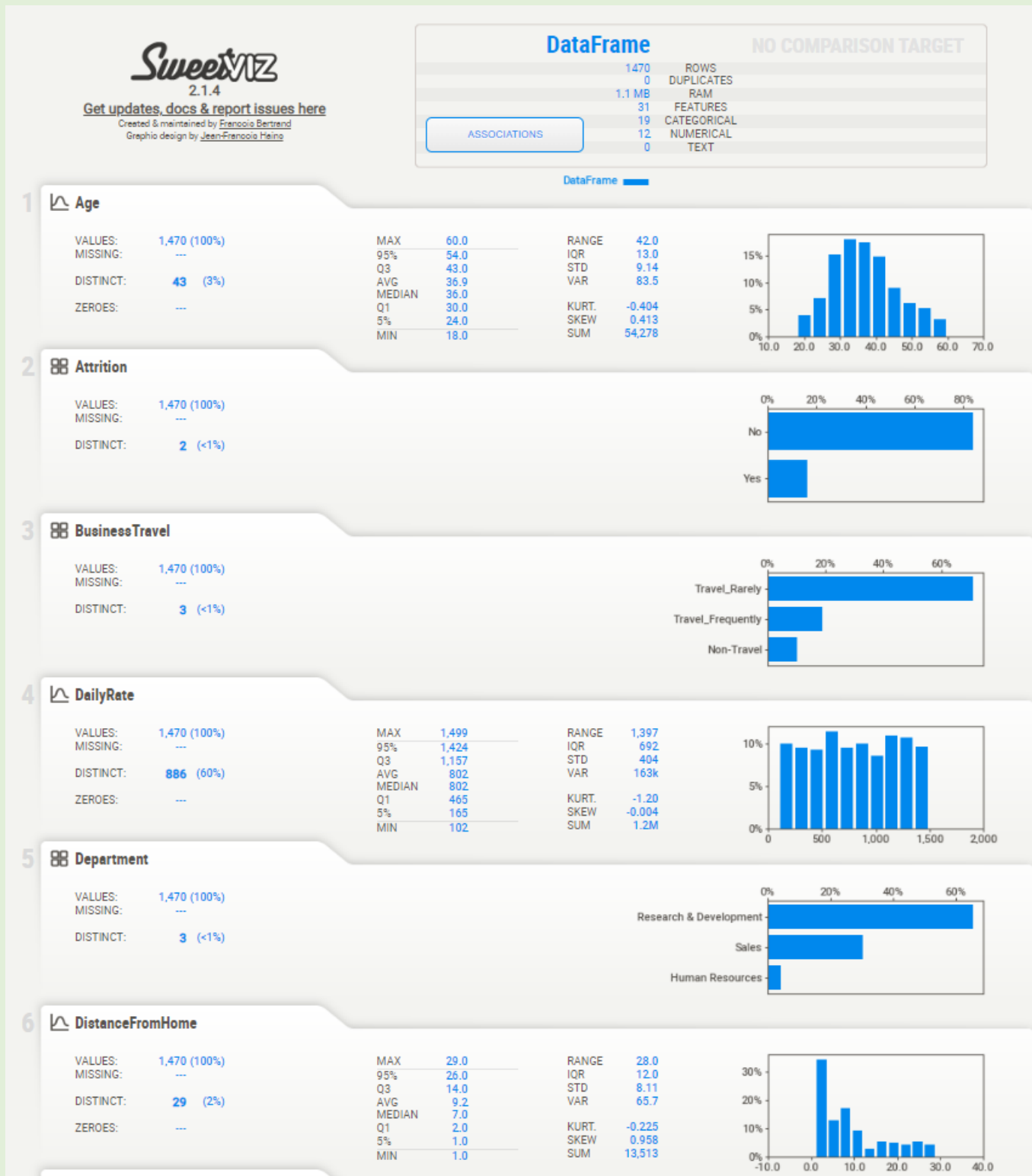
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1470 entries, 0 to 1469
Data columns (total 35 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Age                                   1470 non-null   int64
1   Attrition                             1470 non-null   object
2   BusinessTravel                         1470 non-null   object
3   DailyRate                             1470 non-null   int64
4   Department                             1470 non-null   object
5   DistanceFromHome                       1470 non-null   int64
6   Education                             1470 non-null   int64
7   EducationField                         1470 non-null   object
8   EmployeeCount                         1470 non-null   int64
9   EmployeeNumber                         1470 non-null   int64
10  EnvironmentSatisfaction                1470 non-null   int64
11  Gender                                 1470 non-null   object
12  HourlyRate                             1470 non-null   int64
13  JobInvolvement                         1470 non-null   int64

```

- Automated Exploratory Data Analysis was performed using Sweetviz to check



- Descriptive statistics were computed to summarize shape of a dataset's distribution, its dispersion and central tendency

Compute Descriptive Statistics: To summarize shape of a dataset's distribution, its dispersion and central tendency.

```
1 #To get description of all columns
2 df.describe(include = 'all')
```

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EmployeeCount	EmployeeNumber
count	1470.000000	1470	1470	1470.000000	1470	1470.000000	1470.000000	1470	1470.0	1470.000000
unique	NaN	2	3	NaN	3	NaN	NaN	6	NaN	NaN
top	NaN	No	Travel_Rarely	NaN	Research & Development	NaN	NaN	Life Sciences	NaN	NaN
freq	NaN	1233	1043	NaN	961	NaN	NaN	606	NaN	NaN
mean	36.923810	NaN	NaN	802.485714	NaN	9.192517	2.912925	NaN	1.0	1024.865306
std	9.135373	NaN	NaN	403.509100	NaN	8.106864	1.024165	NaN	0.0	602.024335
min	18.000000	NaN	NaN	102.000000	NaN	1.000000	1.000000	NaN	1.0	1.000000
25%	30.000000	NaN	NaN	465.000000	NaN	2.000000	2.000000	NaN	1.0	491.250000
50%	36.000000	NaN	NaN	802.000000	NaN	7.000000	3.000000	NaN	1.0	1020.500000
75%	43.000000	NaN	NaN	1157.000000	NaN	14.000000	4.000000	NaN	1.0	1555.750000
max	60.000000	NaN	NaN	1499.000000	NaN	29.000000	5.000000	NaN	1.0	2068.000000

3b) Data Cleansing Actions:

- Empty or nearly empty columns were removed using "drop_thresh" to drop columns if 90% of data was empty

Drop Columns if 90% data is empty

```
drop_thresh = df.shape[0]*.10
df = df.loc[:, df.isin([' ', 'NULL', 'NaN', 0]).mean() < drop_thresh]
df = df.dropna(thresh=drop_thresh, how='all', axis='columns').copy()
df.info()
```

```
1 print(df.isin([' ', 'NULL', 'NaN', 0]).mean())
2 drop_thresh = .90
3 df = df.loc[:, df.isin([' ', 'NULL', 'NaN', 0]).mean() < drop_thresh]
4 print(df.isin([' ', 'NULL', 'NaN', 0]).mean())
```

```
Age                0.000000
Attrition           0.000000
BusinessTravel      0.000000
DailyRate           0.000000
Department          0.000000
DistanceFromHome    0.000000
Education           0.000000
EducationField       0.000000
EmployeeCount        0.000000
EmployeeNumber       0.000000
EnvironmentSatisfaction 0.000000
Gender              0.000000
HourlyRate          0.000000
JobInvolvement       0.000000
JobLevel            0.000000
JobRole             0.000000
JobSatisfaction      0.000000
MaritalStatus        0.000000
MonthlyIncome        0.000000
MonthlyRate         0.000000
NumCompaniesWorked   0.134014
Over18              0.000000
OverTime            0.000000
OvertimeHours        0.000000
```

- Duplicates were dropped using pandas "df.drop_duplicates()" method

Handle Missing Values: Replace remaining ["None","nan", "NaN", ""] values with Zero

```
1 df = df.replace(["None","nan", "NaN", ""], "0") # Replace all Nan Values with Zero
2 null = (df.isin(["None","nan", "NaN", ""]).sum()) # Sum as series
3 null_df=pd.DataFrame({'cols':null.index, 'sum':null.values}).sort_values(by=['sum'],ascending=False)
4
5 print(colored("Data has ", 'green', attrs=['bold']))
6     +colored((null_df.at[0,'sum']), 'red', attrs=['bold'])
7     +colored(" null values.\n ", 'green', attrs=['bold'])
8     +colored(null_df.tail(35), 'red', attrs=['bold'])) # print first two rows
```

Data has 0 null values.

	cols	sum
0	Age	0
26	StandardHours	0
20	NumCompaniesWorked	0
21	Over18	0
22	OverTime	0
23	PercentSalaryHike	0
24	PerformanceRating	0
25	RelationshipSatisfaction	0
27	StockOptionLevel	0
18	MonthlyIncome	0
28	TotalWorkingYears	0
29	TrainingTimesLastYear	0
30	WorkLifeBalance	0
31	YearsAtCompany	0
32	YearsInCurrentRole	0
33	YearsSinceLastPromotion	0
19	MonthlyRate	0
17	MaritalStatus	0
1	Attrition	0
8	EmployeeCount	0
2	BusinessTravel	0
3	DailyRate	0
4	Department	0
5	DistanceFromHome	0
6	Education	0
7	EducationField	0
9	EmployeeNumber	0
16	JobSatisfaction	0
10	EnvironmentSatisfaction	0
11	Gender	0
12	HourlyRate	0
13	JobInvolvement	0
14	JobLevel	0
15	JobRole	0
34	YearsWithCurrManager	0

- Null values were summed and Data was found to exhibit zero null values. Thus, no filling of null values was required

3c) Features Engineering

In machine learning, feature selection is the method to reduce the number of input variables during developing predictive modelling. This reduction in input variables is necessary not only to minimize computational cost of modeling but also to achieve performance improvement of the model.

Among widely practices feature selection approaches include statistical-based feature selection methods which use statistical measures to evaluate relationship between each input variable and the target variable and then select those exhibiting strongest relationship with the latter. While these methods can be both speedy and effective, however, the ultimate choice of statistical measure is largely dependant on data types of both of these variables.

Irrespective of the statistical measure being employed, two dominant feature selection techniques, that is supervised and unsupervised, exist where the former can be further categorized into wrapper, filter and intrinsic techniques. Filter-based feature selection methods employs statistical measures to evaluate correlation between input and output variables so that those exhibiting highest correlations are selected. Statistical measures employed in filter-based feature selection are normally univariate in nature since they evaluate relationship of single input variables one by one with target variable, disregarding their interaction with each other.

Consequently, adopting filter-based feature selection methods, the employee attrition model approached filter engineering in three steps. Firstly, unique values for all columns were computed after which columns with unique values less than 2 were dropped.

```

1) Assessing Columns for Feature Selection:
Get unique counts to determine threshold for dropping columns
Drop Columns from dataframe if uniqueness is less than threshold (eg. 2)

1 unique_counts = pd.DataFrame.from_records([(col, df[col].nunique()) for col in df.columns], # get unique counts
2         columns=['Column_Name', 'Unique']).sort_values(by=['Unique'])
3 print(colored("\nThis can help us determine threshold for which columns to exclude from Features.\n", 'blue', attrs=['bold']
4         + colored(type(unique_counts), 'green', attrs=['bold'])
5         + colored("\n\n", 'green', attrs=['bold'])
6         + colored(unique_counts, 'red', attrs=['bold'])
7     )
8
9 unique = unique_counts[(unique_counts['Unique'] < 2)] #If threshold is less than 2 then
10 drop_unique = (unique['Column_Name'].tolist()) # List of columns to drop
11
12 cols_to_exclude = ['EmployeeNumber']
13 cols_to_exclude = ['EmployeeNumber'] + drop_unique
14
15 print(colored("\n\n", 'blue', attrs=['bold'])
16         + colored(type(unique), 'green', attrs=['bold'])
17         + colored("\n", 'green', attrs=['bold'])
18         + colored(unique, 'red', attrs=['bold'])
19
20         + colored("\n\nList of columns to drop\n", 'blue', attrs=['bold'])
21         + colored(cols_to_exclude, 'red', attrs=['bold'])
22     )
23
24 #Function to Drop Columns & Convert to Categories
25 for col in df.columns:
26     if col in cols_to_exclude:
27         df = df.loc[:, ~df.columns.isin(cols_to_exclude)]
28 df.info()

```

This can help us determine threshold for which columns to exclude from Features.

```

<class 'pandas.core.frame.DataFrame'>

```

	Column_Name	Unique
21	Over18	1
26	StandardHours	1
8	EmployeeCount	1
11	Gender	2
1	Attrition	2
24	PerformanceRating	2
22	OverTime	2
17	MaritalStatus	3

Prior to final features selection Data Encoding of Object or String Columns was carried out to facilitate any statistical computation during features selection process. Hence, after deep copying of original dataset, a function was created and employed to encode object data using Scikit-learn label encoder.

2) Data Encoding of Object/String Columns:

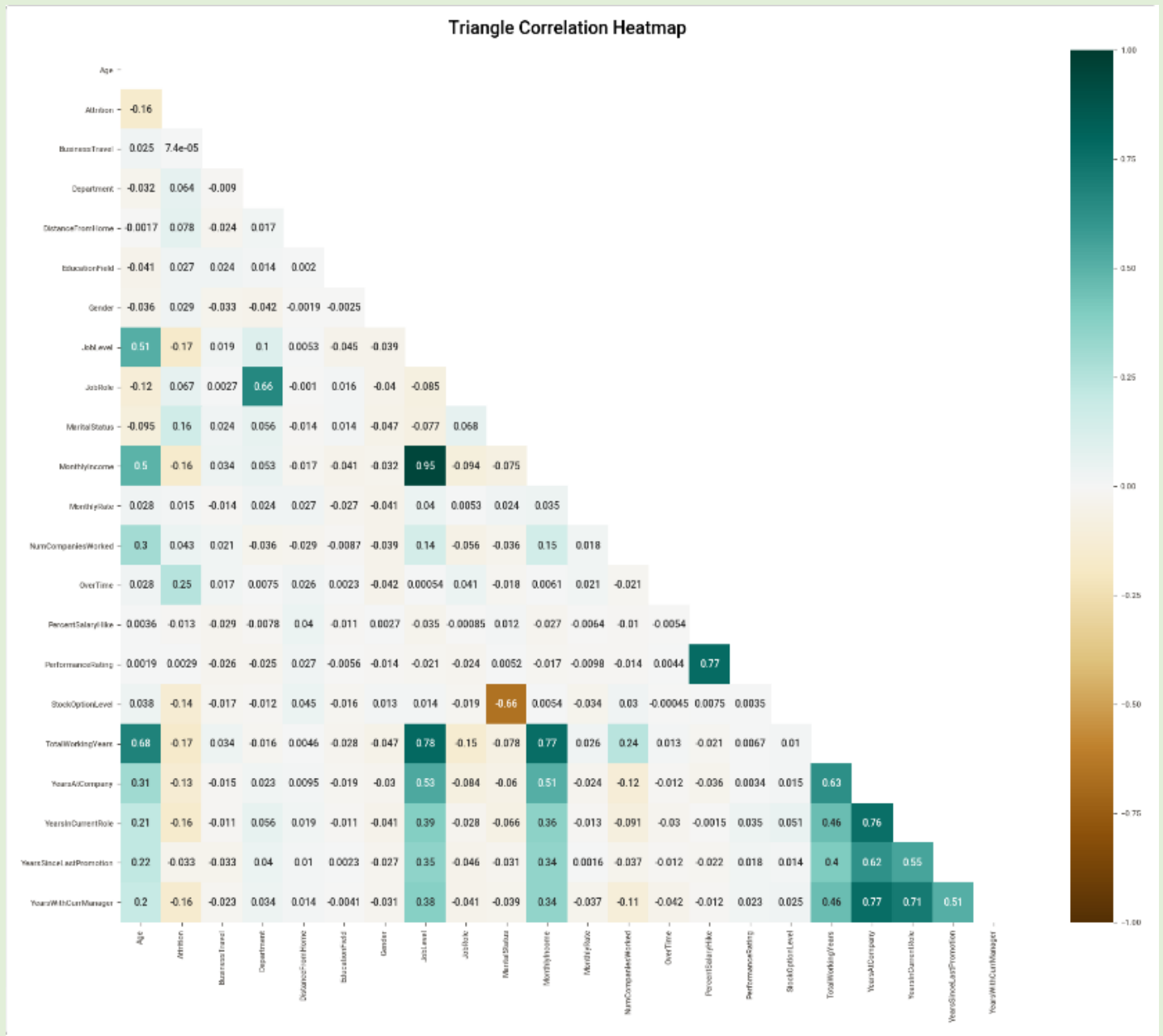
- * List all Object/String Columns
- * Deep copy the original data
- * Create Function to employ Scikit-learn label encoding to encode object data
- * Create a new dataframe with encoded data description to attach to model outcomes

```

1 #Function to encode object/string columns
2
3 #List all Object/String Columns
4 from sklearn import preprocessing
5 cat_columns = df.select_dtypes(include=[object]) # Get Object Type Columns to Convert to Encoded Categories
6 cat_columns.info()
7
8 categorical_column = list(cat_columns.columns)# List of columns to for label encoding
9
10 print(colored("\n\nColumns Requiring Encoding: \n", 'blue', attrs=['bold'])
11       + colored(categorical_column, 'green', attrs=['bold']))
12
13 #Deep copy the original data
14 df_encoded = df.copy(deep=True)
15
16 # Make Empty Dataframe to decode encoded data Later
17 decode_features = pd.DataFrame()
18
19 ##### Employ Scikit-Learn Label encoding to encode object data #####
20 lab_enc = preprocessing.LabelEncoder()
21 for col in categorical_column:
22     df_encoded[col] = lab_enc.fit_transform(df[col])
23     le_name_mapping = dict(zip(lab_enc.classes_, lab_enc.transform(lab_enc.classes_)))
24
25     ##### Decode Encoded Data #####
26     feature_df = pd.DataFrame([le_name_mapping])
27     feature_df = feature_df.astype(str)
28     print(feature_df)
29     feature_df = (col + "_" + feature_df.iloc[0:])
30     feature_df["Feature"] = col
31     print(feature_df)
32     decode_features = decode_features.append(feature_df)# Append Dictionaries to Empty Dataframe for Later Decoding
33
34     ##### Print Encoded Data #####
35     print(colored("Feature: \n", 'blue', attrs=['bold'])
36           + colored(col, 'red', attrs=['bold'])
37           + colored("\nMapping: \n", 'blue', attrs=['bold'])
38           + colored(le_name_mapping, 'green', attrs=['bold'])
39           + colored("\n\n", 'blue', attrs=['bold'])
40           )
41 df_encoded.head(3)
42
43 ##### Make Decoded Factor Dataframe with Description #####
44 #print(decode_features)
45 factor_list = decode_features.T # Transpose Dataframe and place in new dataframe
46 factor_list = factor_list.replace(np.nan, "/") # nan values with forward slash
47 factor_list["Factors"] = factor_list.astype(str).agg("".join,axis=1).replace(r'[^\\w\\s]|/', '', regex=True) # Aggregate ALL (
48 factor_list.reset_index() # Reset index before copying/assigning it to a new column
49 factor_list["Description"] = factor_list.index # Assign index to column
50

```

Statistical measures were then employed with supervised filter-based feature selection technique. Using Pearson's Correlation, the first set of features are selected based on the strength of positive correlation with target variable 'Attrition'. Additionally, Pearson's Correlation Matrix was also computed to select feature pairs exhibiting positive correlations with each other.



All feature lists were then combined to filter out dataframe columns not included in the ‘final_features’ list.

```
1 df_encoded = df_encoded.filter(final_features)
2 df_encoded.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 1470 entries, 0 to 1469
Data columns (total 22 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Age                                  1470 non-null   int64
1   Attrition                           1470 non-null   int32
2   BusinessTravel                       1470 non-null   int32
3   Department                           1470 non-null   int32
4   DistanceFromHome                     1470 non-null   int64
5   EducationField                       1470 non-null   int32
6   Gender                               1470 non-null   int32
7   JobLevel                             1470 non-null   int64
8   JobRole                              1470 non-null   int32
9   MaritalStatus                       1470 non-null   int32
10  MonthlyIncome                        1470 non-null   int64
11  MonthlyRate                          1470 non-null   int64
12  NumCompaniesWorked                   1470 non-null   int64
13  OverTime                             1470 non-null   int32
14  PercentSalaryHike                    1470 non-null   int64
15  PerformanceRating                    1470 non-null   int64
16  StockOptionLevel                     1470 non-null   int64
17  TotalWorkingYears                    1470 non-null   int64
18  YearsAtCompany                       1470 non-null   int64
19  YearsInCurrentRole                   1470 non-null   int64
20  YearsSinceLastPromotion               1470 non-null   int64
21  YearsWithCurrManager                  1470 non-null   int64
dtypes: int32(8), int64(14)
memory usage: 218.2 KB
```

4) Summary of Training Different Classifier Models

4a) Machine Learning Algorithm Approaches

Since EDA revealed a highly imbalanced class distribution, this necessitated achieving appropriate class balancing. While data level and algorithm ensemble approaches do exist for dealing with imbalanced datasets, nevertheless, an automated optimal parameter search method was created to achieve best class reweighting along with isolating other optimal model parameters. This approach was employed because best hyper-parameters are not automatically learnt within estimators and its manual search not only slows down model development but may also lead to ineffective model construction. Hence, exhaustive cv grid search approach was used to pass parameter arguments to the constructor in order to find optimal parameters for each model.

Logistic Regression (LR) Models

Grid Search Method to Find 'Best Parameters' to 'Build Logistic Regression WITH Best Class Weights'

```

1 # Grid Search Method to find Best Hyperparameters for a Logistic Regression Model
2 def grid_search_lr(X_train, y_train):
3     # Parameters
4     params_grid = {
5         'class_weight': [{0:0.1, 1:0.9}, {0:0.2, 1:0.8}, {0:0.3, 1:0.7}],
6         'solver': ['lbfgs', 'saga', 'liblinear', 'newton-cg', 'sag']
7     }
8     # LR Model
9     lr_model = LogisticRegression(random_state=rs, max_iter=1000)
10
11     # Search Best Parameters
12     grid_search = GridSearchCV(estimator = lr_model,
13                               param_grid = params_grid,
14                               scoring='f1',
15                               cv = 5, verbose = 1)
16     # Train Model with Best Parameters
17     grid_search.fit(X_train, y_train)
18
19     # Get Best/optimal parameters
20     best_lrparams = grid_search.best_params_
21     return best_lrparams

```

Get Optimal Parameters for LR Model using Grid Search LR Method above

```

1 final_lrparams = grid_search_lr(X_train, y_train) # From the cell above, Call grid_search_rf(X_train, y_train)
2
3 final_lrparams_df = pd.DataFrame([final_lrparams]) # Dictionary To dataframe
4 print(final_lrparams_df)
5
6 # Make Optimal Variables
7 optimal_lr_class_weight = (final_lrparams_df.at[0, 'class_weight'])
8 optimal_solver = (final_lrparams_df.at[0, 'solver'])
9 print('Optimal LR Class Weights: ', optimal_lr_class_weight)
10 print('Optimal Solver: ', optimal_solver)
11
12 # Define Optimal Parameters
13 optimal_lr_params = {'class_weight': optimal_lr_class_weight, 'solver': optimal_solver}
14 print(optimal_lr_params)

```

```

Fitting 5 folds for each of 15 candidates, totalling 75 fits
  class_weight  solver
0 {0: 0.2, 1: 0.8} newton-cg
Optimal LR Class Weights: {0: 0.2, 1: 0.8}
Optimal Solver: newton-cg
{'class_weight': {0: 0.2, 1: 0.8}, 'solver': 'newton-cg'}

```

Random Forest (RF) Models

Grid Search Method to Find 'Best Parameters' to 'Build Random Forest WITH Class Weights'

```

1 # Method for Grid Search Hyperparameters for a Random Forest Model
2 def grid_search_rf(X_train, y_train):
3     # Parameters
4     params_grid = {
5         'max_depth': [2*n+1 for n in range(10)], #[5, 10, 15, 20],
6         'n_estimators': [2*n+1 for n in range(20)], #[25, 50, 100],
7         'min_samples_split': [2, 5],
8         'class_weight': [{0:0.1, 1:0.9}, {0:0.2, 1:0.8}, {0:0.3, 1:0.7}]
9     }
10    # RF Model
11    rf_model = RandomForestClassifier(random_state=rs)
12
13    # Search Best Parameters
14    grid_search = GridSearchCV(estimator = rf_model,
15                               param_grid = params_grid,
16                               scoring='f1',
17                               cv = 5, verbose = 1)
18    # Train Model with Best Parameters
19    grid_search.fit(X_train, y_train)
20
21    # Get Best/optimal parameters
22    best_params = grid_search.best_params_
23    #Best_Score = grid_search.best_score_
24    #accuracy = get_accuracy(X_train, X_test, y_train, y_test, grid_search.best_estimator_)
25    return best_params

```

Get 'Optimal Parameters' for RF Model using Grid Search RF Method above ¶

```

1 #Calculate StartTime to Measure Script Execution Time at the End of Script
2 start_time = datetime.now()
3
4 #Get Optimal Parameters for RF Model using "grid_search_rf" Method to Find 'Best/Optimal Parameters'
5 best_params = grid_search_rf(X_train, y_train) # From the cell above, Call grid_search_rf(X_train, y_train)
6
7 best_params_df = pd.DataFrame([best_params]) # Dictionary To dataframe
8 print(best_params_df)
9
10 # Make Optimal Parameter Variables
11 optimal_class_weight = (best_params_df.at[0,'class_weight'])
12 print(optimal_class_weight)
13 optimal_max_depth = (best_params_df.at[0,'max_depth'])
14 print(optimal_max_depth)
15 optimal_min_samples_split = (best_params_df.at[0,'min_samples_split'])
16 print(optimal_min_samples_split)
17 optimal_n_estimators = (best_params_df.at[0,'n_estimators'])
18 print(optimal_n_estimators)
19
20 # Define Optimal Parameters
21 optimal_rf_params = {'bootstrap': True,
22                      'class_weight': optimal_class_weight,
23                      'max_depth': optimal_max_depth,
24                      'min_samples_split': optimal_min_samples_split,
25                      'n_estimators': optimal_n_estimators}
26 print(optimal_rf_params)
27
28 #Print Total Execution Time
29 print('Model Execution Time: ', datetime.now() - start_time)

```

Fitting 5 folds for each of 1200 candidates, totalling 6000 fits

```

class_weight max_depth min_samples_split n_estimators
0 {0: 0.1, 1: 0.9}          7              5          23
{0: 0.1, 1: 0.9}
7
5
23
{'bootstrap': True, 'class_weight': {0: 0.1, 1: 0.9}, 'max_depth': 7, 'min_samples_split': 5, 'n_estimators': 23}

```

eXtreme Gradient Boosting (XGB) Model

Grid Search Method to Find 'Best Parameters' to 'Build XGB WITH Best Class Weights'

```

1 #Calculate StartTime to Measure Model Execution Time at the End
2 xgb_start_time = datetime.now()
3
4 # Method for Grid Search Hyperparameters for a Random Forest Model
5 def grid_search_xgb(X_train, y_train):
6     # Parameters
7     params_grid = {
8         'max_depth': [5, 10, 15, 20], #[2*n+1 for n in range(10)],
9         'n_estimators': [100, 300, 500],#[2*n+1 for n in range(20)],
10        'min_samples_split': [2, 5, 8],
11    }
12    # RF Model
13    xgb_model = GradientBoostingClassifier(random_state=rs)
14
15    # Search Best Parameters
16    grid_search = GridSearchCV(estimator = xgb_model,
17                               param_grid = params_grid,
18                               scoring='f1',
19                               cv = 5, verbose = 1)
20    # Train Model with Best Parameters
21    grid_search.fit(X_train, y_train)
22
23    # Get Best/optimal parameters
24    best_params = grid_search.best_params_
25    #Best_Score = grid_search.best_score_
26    #accuracy = get_accuracy(X_train, X_test, y_train, y_test, grid_search.best_estimator_)
27    return best_params

```

Get 'Optimal Parameters' for XGB Model using Grid Search XGB Method above

```

1 xgb_params = grid_search_xgb(X_train, y_train) # From the cell above, Call grid_search_xgb(X_train, y_train)

```

Fitting 5 folds for each of 36 candidates, totalling 180 fits

```

1 xgb_params_df = pd.DataFrame([xgb_params]) # Dictionary To dataframe
2 print(xgb_params_df)
3
4 # Make Optimal Parameter Variables|
5 xgb_opt_max_depth = (xgb_params_df.at[0, 'max_depth'])
6 print(xgb_opt_max_depth)
7 xgb_opt_min_samples_split = (xgb_params_df.at[0, 'min_samples_split'])
8 print(xgb_opt_min_samples_split)
9 xgb_opt_n_estimators = (xgb_params_df.at[0, 'n_estimators'])
10 print(xgb_opt_n_estimators)
11
12 # Define Optimal Parameters
13 xgb_optimal_params = {'bootstrap': True,
14                       #'class_weight': xgb_opt_class_weight,
15                       'max_depth': xgb_opt_max_depth,
16                       'min_samples_split': xgb_opt_min_samples_split,
17                       'n_estimators': xgb_opt_n_estimators}
18 print(xgb_optimal_params)

```

	max_depth	min_samples_split	n_estimators
0	5	2	500

```

5
2
500
{'max_depth': 5, 'min_samples_split': 2, 'n_estimators': 500}

```


Additionally, the following approaches were also combined with optimal parameters to find a model with best scores:

I) Data Level Approaches:

i) Synthetic Minority Over-sampling Technique (SMOTE):

Due to highly imbalance class distribution, employee data contains very few instances of minority class for any classification model to explicitly learn decision boundary. A popular approach to tackle this problem is oversampling minority class examples which are close in the feature space using SMOTE. This approach allowed us to achieve a balanced class distribution.

ii) Random under-sampling:

Using random under-sampling examples from majority class were deleted to achieve a balanced class distribution.

iii) Random over-sampling:

Random oversampling was employed to duplicate examples from minority class to achieve a balanced class distribution.

II) Algorithm Ensemble Approach:

i) Boosting:

Using boosting, a sequential aggregate of base classifier was created on weighted versions of training data set which focused on misclassified samples at every stage of creating new classifiers based on sample weights that were altered as per classifier's performance. Boosting was achieved using XGB Classifier.

4b) Summarizing Employed Models

Following three main classifier models have been used to predict employee attrition.

1) Logistic Regression (LR) Models:

Due to the binary nature of predictive variable 'y' where the output can only result in whether the employees will fall in attrition class or not, logistic regression classification algorithm was employed. A single method was created to measure predictive capability of the following two LR models:

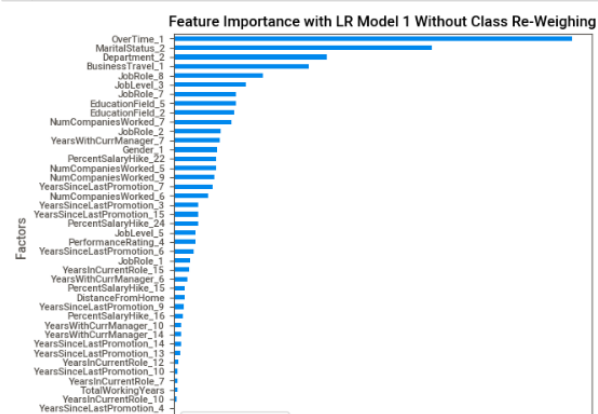
Method to 'Build Logistic Regression WITH & Without Optimal Class Weights'

```
1 # Build a Logistic regression model with Optimal Class Weights
2 def build_op_lr(X_train, y_train, X_test, threshold=0.5, best_params=None):
3
4     model = LogisticRegression(random_state=rs,
5                               max_iter = 1000)
6     # If best parameters are provided
7     if best_params:
8         model = LogisticRegression(penalty = 'l2',
9                                   random_state=rs,
10                                  max_iter = 1000,
11                                  class_weight = optimal_lr_class_weight,
12                                  solver = optimal_solver
13                                  )
14     # Train the model
15     model.fit(X_train, y_train)
16     # If predicted probability is larger than threshold (default value is 0.5), generate a positive label
17     predicted_proba = model.predict_proba(X_test)
18     yp = (predicted_proba[:,1] >= threshold).astype('int')
19     return yp, model
20
```

1a) LR Model 1 Without Class Re-Weighing: A simple algorithm was created without achieving any class balance or hyperparameter tuning.

LR Model 1 Without 'Class Re-Weighing'

```
1 #Running the Model...Call Method to 'Build LR Without Adjusted Class Weights'
2 lr_preds, lr_model = build_op_lr(X_train, y_train, X_test, best_params=None) # Call Method to 'Build Logistic Regression WIT
3 lr_results = evaluate(y_test, lr_preds, val_type="Log Reg")
4 lr_results["Factors"] = 'LR1_features_df'
5 #print(lr_results)
6
7 # Get Factors Contributing to Attrition
8 lr1_features = lr_model.coef_[0]##Logistic regression does not have an attribute for ranking feature so we cannot use featur
9
10 # Convert to Series and Plot Features
11 LR1_features = pd.Series(lr1_features, index=X.columns) # Attach Importance to data before split
12
13 # Get Features into Dataframe
14 LR1_features_df = LR1_features.to_frame()
15 LR1_features_df.reset_index(inplace=True)# Reset index before copying index to new column
16 LR1_features_df = LR1_features_df.rename(columns = {'index':'Factors', 0: 'Importance'}).sort_values(by=['Importance'],ascen
17 LR1_features_df = LR1_features_df[LR1_features_df['Importance'] > 0] # Filter Features with Negative Correlation
18
19 # Plot Features
20 LR1_features_df.set_index('Factors', inplace=True)
21 LR1_features_df.plot(kind='barh', figsize=(8, 6))
22 plt.title('Feature Importance with LR Model 1 Without Class Re-Weighing')
23
24 # Merge Description to Encoded Features/Factors
25 LR1_features_df = LR1_features_df.rename_axis('Factors').reset_index() # Reset index before copying it to column
26 LR1_features_df = pd.merge(LR1_features_df,factor_list,on='Factors',how='left',indicator=True).replace(np.nan, "") # Merge D
27 LR1_features_df['Description'] = np.where(LR1_features_df['Description'] == '/', LR1_features_df['Factors'], LR1_features_df
28 LR1_features_df['Description'] = LR1_features_df['Description'].replace(r'_+', '', regex=True) # Replace Underscore with sin
29 final_columns = ['Factors', 'Importance', 'Description'] # List of Columns to keep
30 LR1_features_df = LR1_features_df.drop(columns=[col for col in LR1_features_df if col not in final_columns]) # Drop Columns
```



1b) LR Model 2 WITH Auto-Hypertuned Class Reweighting: A modified algorithm with auto-hypertuned parameters was created using CV grid search method described above.

LR Model 2 WITH Auto-Hypertuned Class Reweighting

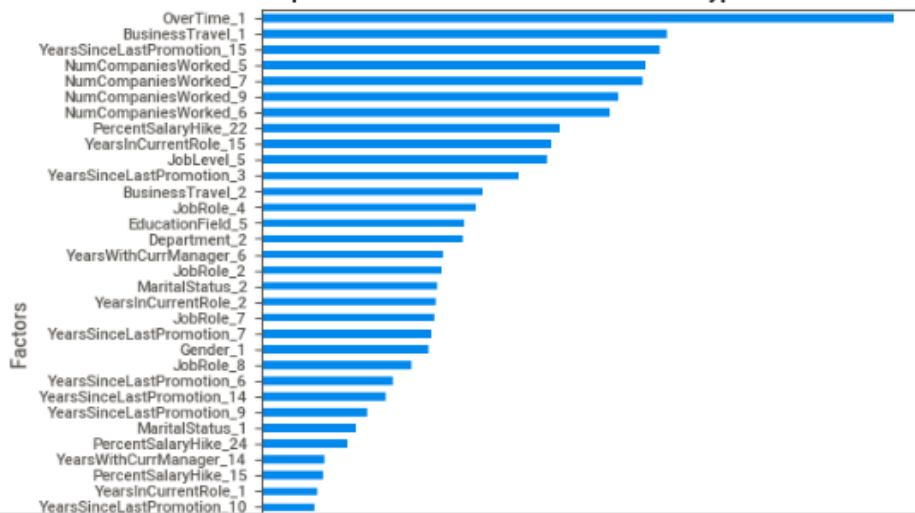
```

1 #Running the Model...Call Method to 'Build LR With Optimal Class Weights'
2 op_lr_preds, op_lr_model = build_op_lr(X_train, y_train, X_test, best_params=optimal_lr_params) # Call Method to 'Build Logi
3 op_lr_results = evaluate(y_test, op_lr_preds, eval_type="LogReg Opt")
4 op_lr_results["Factors"] = 'LR2_features_df'
5 print(op_lr_results)
6
7 # Get Factors Contributing to Attrition
8 lr2_features = op_lr_model.coef_[0]##Logistic regression does not have an attribute for ranking feature so we cannot use fea
9
10 # Convert to Series and Plot Features
11 LR2_features = pd.Series(lr2_features, index=X.columns) # Attach Importance to data before split
12
13 # Get Features into Dataframe
14 LR2_features_df = LR2_features.to_frame()
15 LR2_features_df.reset_index(inplace=True)# Reset index before copying index to new column
16 LR2_features_df = LR2_features_df.rename(columns = {'index':'Factors', 0: 'Importance'}).sort_values(by=['Importance'],ascen
17 LR2_features_df = LR2_features_df[LR2_features_df['Importance'] > 0] # Filter Features with Negative Correlation
18
19 # Plot Features
20 LR2_features_df.set_index('Factors', inplace=True)
21 LR2_features_df.plot(kind='barh', figsize=(6, 6))
22 plt.title('Feature Importance with LR Model 2 WITH Auto-Hypertuned Class Reweighting')
23
24 # Merge Description to Encoded Features/Factors
25 LR2_features_df = LR2_features_df.rename_axis('Factors').reset_index() # Reset index before copying it to column
26 LR2_features_df = pd.merge(LR2_features_df,factor_list,on='Factors',how='left',indicator=True).replace(np.nan, "") # Merge D
27 LR2_features_df['Description'] = np.where(LR2_features_df['Description'] == '/', LR2_features_df['Factors'], LR2_features_df
28 LR2_features_df['Description'] = LR2_features_df['Description'].replace(r'_',' ', regex=True) # Replace Underscore with sin
29 final_columns = ['Factors', 'Importance', 'Description'] # List of Columns to keep
30 LR2_features_df = LR2_features_df.drop(columns=[col for col in LR2_features_df if col not in final_columns]) # Drop Columns

```

{'type': 'LogReg Opt', 'accuracy': 0.7891156462585034, 'recall': 0.6808510638297872, 'precision': 0.4050632911392405, 'fscore': 0.6634768740031898, 'auc': 0.7452838315100354, 'Factors': 'LR2_features_df'}

Feature Importance with LR Model 2 WITH Auto-Hypertuned Class Reweighting



2) Random Forest Models:

A single method was employed to run random forest models with and without optimal parameters and class weighing:

Method to 'Build Random Forest WITH or WITHOUT Class Weights'

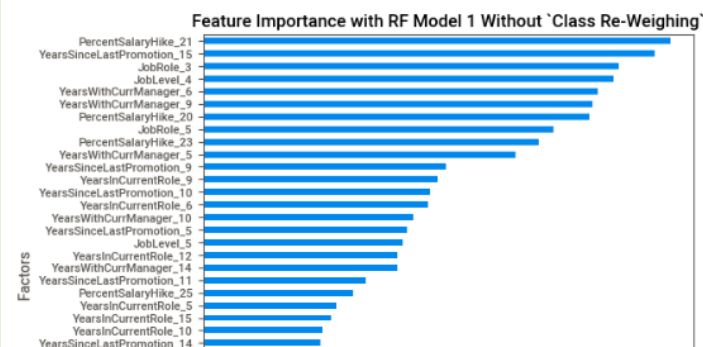
```
1 def build_op_rf(X_train, y_train, X_test, threshold=0.5, best_params=None):
2     model = RandomForestClassifier(random_state = rs)
3     # If best parameters are provided
4     if best_params:
5         model = RandomForestClassifier(random_state = rs,
6                                       # If bootstrap sampling is used
7                                       bootstrap = True, # False if Bootstaping is being kept as None
8                                       # Max depth of each tree
9                                       max_depth = optimal_rf_params['max_depth'],
10                                      # Class weight parameters
11                                      class_weight=optimal_rf_params['class_weight'],
12                                      # Number of trees
13                                      n_estimators=optimal_rf_params['n_estimators'],
14                                      # Minimal samples to split
15                                      min_samples_split=optimal_rf_params['min_samples_split'])
16
17     # Train the model
18     model.fit(X_train, y_train)
19     # If predicted probability is largr than threshold (default value is 0.5), generate a positive Label
20     predicted_proba = model.predict_proba(X_test)
21     yp = (predicted_proba[:,1] >= threshold).astype('int')
22     return yp, model
```

2a) RF Model 1 Without Class Re-Weighing: A simple algorithm was created without achieving any class balance or hyperparameter tuning.

RF Model 1 Without 'Class Re-Weighing'

```
1 #Running the Model...Call Method to 'Build Simple Random Forest Without Class Weights'
2 preds, model = build_op_rf(X_train, y_train, X_test, best_params=None) # Call Method to 'Build Simple Random Forest Without
3
4 #Collect Results
5 original_results = evaluate(y_test, preds, eval_type="RForest")
6 original_results["Factors"] = 'RF1_features_df'
7 print(original_results)
8
9 # Get Factors Contributing to Attrition in Dataframe
10 RF1_features_df = feature_importance(X, model)[:40]
11 RF1_features_df = RF1_features_df[RF1_features_df['Importance'] > 0] # Filter Features with Negative Correlation
12 RF1_features_df = RF1_features_df.sort_values(by=['Importance'],ascending=True) # Attach Importance to data before split
13
14 # Capture Model Drivers in Dataframe
15 RF1_drivers = [['Uses train test set to achieve best hypertuning to resample data from train set to achieve balanced class.
16                ['No optimal parameters were used.'],
17                ['Runs random forest model with above drivers to get top contributory features.']] # initialize list of lists
18 RF1_drivers = pd.DataFrame(RF1_drivers, columns=['Main Drivers behind RF1 model are as follows:']) # Create the pandas DataF
19 RF1_drivers.index = np.arange(1,len(RF1_drivers)+1)
20
21 # Plot Features
22 RF1_features_df.set_index('Factors', inplace=True)
23 RF1_features_df.plot(kind='barh', figsize=(6, 6))
24 plt.title('Feature Importance with RF Model 1 Without 'Class Re-Weighing')
25
26 # Merge Description to Encoded Features/Factors
27 RF1_features_df = RF1_features_df.rename_axis('Factors').reset_index() # Reset index before copying it to column
28 RF1_features_df = pd.merge(RF1_features_df, factor_list, on='Factors', how='left', indicator=True).replace(np.nan, "") # Merge D
29 RF1_features_df['Description'] = np.where(RF1_features_df['Description'] == '/', RF1_features_df['Factors'], RF1_features_df
30 RF1_features_df['Description'] = RF1_features_df['Description'].replace(r'_', ' ', regex=True) # Replace Underscore with sin
31 final_columns = ['Factors', 'Importance', 'Description'] # List of Columns to keep
32 RF1_features_df = RF1_features_df.drop(columns=[col for col in RF1_features_df if col not in final_columns]) # Drop Columns
33
```

{'type': 'RForest', 'accuracy': 0.8537414965986394, 'recall': 0.1276595744680851, 'precision': 0.75, 'fscore': 0.13186813186813184, 'auc': 0.5597812042380912, 'Factors': 'RF1_features_df'}



2b) RF Model 2 WITH Auto-Hypertuned Class Reweighing: A modified algorithm with auto-hypertuned parameters was created using CV grid search method described above.

RF Model 2 WITH Auto-Hypertuned Class Reweighing

```

1 #Running the Model...Call Method to 'Build Simple Random Forest With Class Weights'
2 op_preds, op_model = build_op_rf(X_train, y_train, X_test, best_params=optimal_rf_params) # Call Method to 'Build Simple Ran
3
4 #Collect Results
5 print(original_results)
6 adjusted_results = evaluate(y_test, op_preds, eval_type="RForest Opt")
7 adjusted_results["Factors"] = 'RF2_features_df'
8 print(adjusted_results)
9
10 # Get Factors Contributing to Attrition in Dataframe
11 RF2_features_df = feature_importance(X, op_model)[:40]
12 RF2_features_df = RF2_features_df[RF2_features_df['Importance'] > 0] # Filter Features with Negative Correlation
13 RF2_features_df = RF2_features_df.sort_values(by=['Importance'],ascending=True) # Attach Importance to data before split
14
15 # Capture Model Drivers in Dataframe
16 RF2_drivers = [['Uses train test set to achieve best hypertuning to resample data from train set to achieve balanced class'
17                ['Uses optimal parameters to achieve best hypertuning'],
18                ['Runs random forest model with above drivers to get top contributory features']] # initialize List of Lists
19 RF2_drivers = pd.DataFrame(RF2_drivers, columns=['Main Drivers behind RF2 model are as follows:']) # Create the pandas DataF
20 RF2_drivers.index = np.arange(1,len(RF2_drivers)+1)
21
22 # Plot Features
23 RF2_features_df.set_index('Factors', inplace=True)
24 RF2_features_df.plot(kind='barh', figsize=(6, 6))
25 plt.title('Feature Importance with RF Model 2 WITH Auto-Hypertuned Class Reweighing')
26
27 # Merge Description to Encoded Features/Factors
28 RF2_features_df = RF2_features_df.rename_axis('Factors').reset_index() # Reset index before copying it to column
29 RF2_features_df = pd.merge(RF2_features_df,factor_list,on='Factors',how='left',indicator=True).replace(np.nan, "") # Merge D
30 RF2_features_df['Description'] = np.where(RF2_features_df['Description'] == '/', RF2_features_df['Factors'], RF2_features_df
31 RF2_features_df['Description'] = RF2_features_df['Description'].replace(r'_',' ', regex=True) # Replace Underscore with sp
32 final_columns = ['Factors', 'Importance', 'Description'] # List of Columns to keep
33 RF2_features_df = RF2_features_df.drop(columns=[col for col in RF2_features_df if col not in final_columns]) # Drop Columns

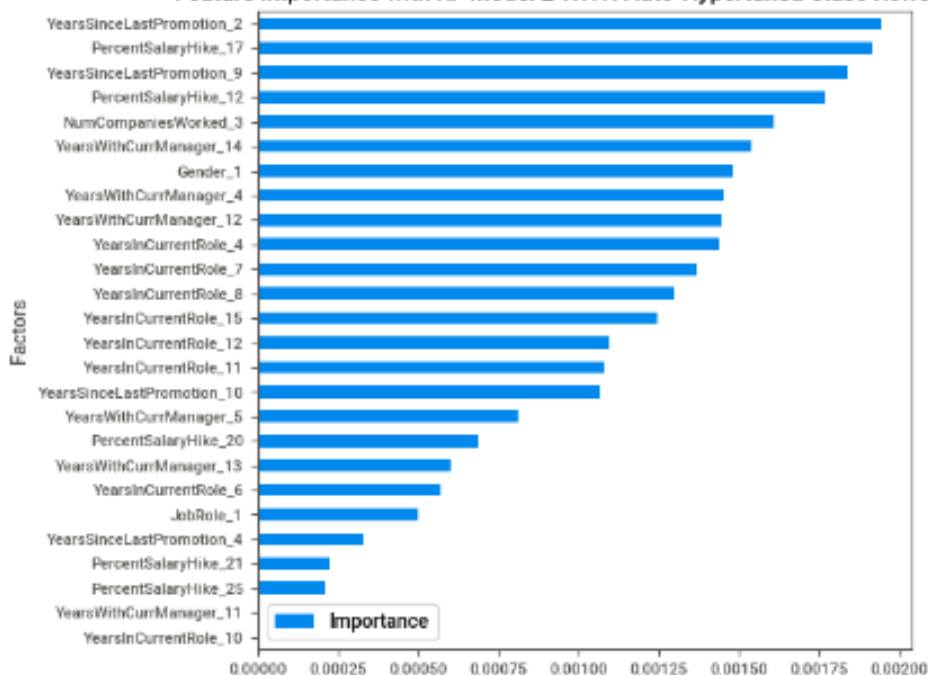
```

```

{'type': 'RForest', 'accuracy': 0.8537414965986394, 'recall': 0.1276595744680851, 'precision': 0.75, 'fscore': 0.13186813186813
184, 'auc': 0.5597812042380912, 'Factors': 'RF1_features_df'}
{'type': 'RForest Opt', 'accuracy': 0.6972789115646258, 'recall': 0.574468085106383, 'precision': 0.28125, 'fscore': 0.55232100
70810387, 'auc': 0.6475579291928676, 'Factors': 'RF2_features_df'}

```

Feature Importance with RF Model 2 WITH Auto-Hypertuned Class Reweighing



2c) RF Model 3 SMOTE Without Auto-Hypertuned Class Reweighing: A simple algorithm was created without achieving any class balance or hyperparameter tuning which employed smote sampling to achieve class balance.

RF Model 3 SMOTE Without Auto-Hypertuned Class Reweighing

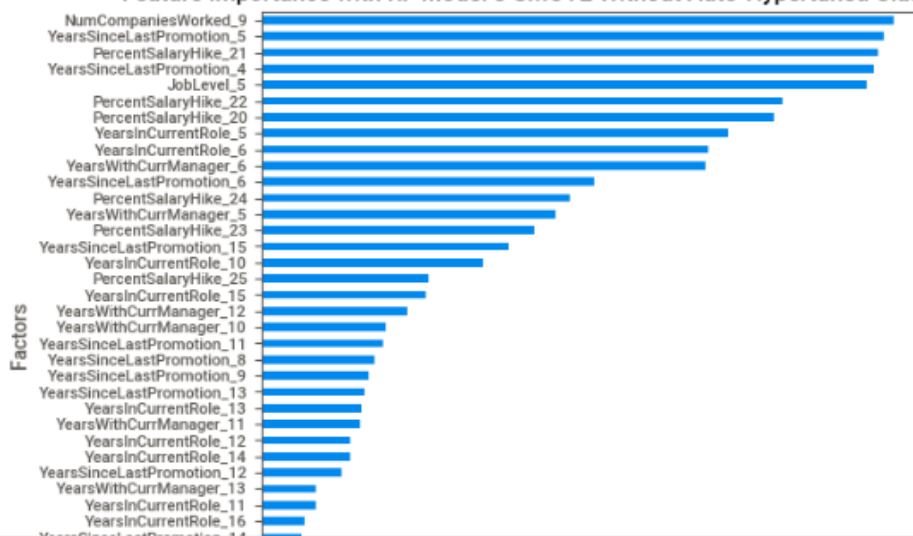
```

1 smote_preds, smo_model = build_op_rf(X_smo, y_smo, X_test, best_params = None)
2 smote_results = evaluate(y_test, smote_preds, eval_type="RForest Smo")
3 smote_results["Factors"] = 'RF3_features_df' # Add Important Factors to Dictionary
4 print(smote_results)
5
6 # Get Factors Contributing to Attrition in Dataframe
7 RF3_features_df = feature_importance(X, smo_model)[:40]
8 RF3_features_df = RF3_features_df[RF3_features_df['Importance'] > 0] # Filter Features with Negative Correlation
9 RF3_features_df = RF3_features_df.sort_values(by=['Importance'], ascending=True) # Attach Importance to data before split
10
11 # Capture Model Drivers in Dataframe
12 RF3_drivers = [['Uses smot parameters to resample data from train set to achieve balanced class'],
13               ['Runs random forest model without optimal parameters with above drivers to get top contributory features']]
14 RF3_drivers = pd.DataFrame(RF3_drivers, columns='Main Drivers behind RF4 model are as follows:') # Create the pandas DataF
15 RF3_drivers.index = np.arange(1, len(RF3_drivers)+1)
16
17 # Plot Features
18 RF3_features_df.set_index('Factors', inplace=True)
19 RF3_features_df.plot(kind='barh', figsize=(6, 6))
20 plt.title('Feature Importance with RF Model 3 SMOTE Without Auto-Hypertuned Class Reweighing')
21
22 # Merge Description to Encoded Features/Factors
23 RF3_features_df = RF3_features_df.rename_axis('Factors').reset_index() # Reset index before copying it to column
24 RF3_features_df = pd.merge(RF3_features_df, factor_list, on='Factors', how='left', indicator=True).replace(np.nan, "") # Merge D
25 RF3_features_df['Description'] = np.where(RF3_features_df['Description'] == '/', RF3_features_df['Factors'], RF3_features_df
26 RF3_features_df['Description'] = RF3_features_df['Description'].replace(r'_ ', ' ', regex=True) # Replace Underscore with sin
27 final_columns = ['Factors', 'Importance', 'Description'] # List of Columns to keep
28 RF3_features_df = RF3_features_df.drop(columns=[col for col in RF3_features_df if col not in final_columns]) # Drop Columns

```

```
{'type': 'RForest Smo', 'accuracy': 0.8469387755102041, 'recall': 0.23404255319148937, 'precision': 0.55, 'fscore': 0.239330543
93305436, 'auc': 0.5988026531139633, 'Factors': 'RF3_features_df'}
```

Feature Importance with RF Model 3 SMOTE Without Auto-Hypertuned Class Reweighing



2d) RF Model 4 SMOTE WITH Auto-Hypertuned Class Reweighing: A modified algorithm with auto-hypertuned parameters was created using CV grid search method described above which employed smote sampling to achieve class balance.

RF Model 4 SMOTE WITH Auto-Hypertuned Class Reweighing

```

1 op_smote_preds, op_smo_model = build_op_rf(X_smo, y_smo, X_test, best_params=optimal_rf_params)
2 op_smote_results = evaluate(y_test, op_smote_preds, eval_type="RForest Opt Smote")
3 op_smote_results["Factors"] = 'RF4_features_df' # Add Important Factors to Dictionary
4 print(op_smote_results)
5
6 # Get Factors Contributing to Attrition in Dataframe
7 RF4_features_df = feature_importance(X, op_smo_model)[:40]
8 RF4_features_df = RF4_features_df[RF4_features_df['Importance'] > 0] # Filter Features with Negative Correlation
9 RF4_features_df = RF4_features_df.sort_values(by=['Importance'], ascending=True) # Attach Importance to data before split
10
11 # Capture Model Drivers in Dataframe
12 RF4_drivers = [['Uses smot parameters to achieve best hypertuning to resample data from train set to achieve balanced class
13                'Uses optimal parameters to achieve best hypertuning'],
14                ['Runs random forest model with above drivers to get top contributory features']] # initialize List of Lists
15 RF4_drivers = pd.DataFrame(RF4_drivers, columns=['Main Drivers behind RF4 model are as follows:']) # Create the pandas DataF
16 RF4_drivers.index = np.arange(1, len(RF4_drivers)+1)
17
18 # Plot Features
19 RF4_features_df.set_index('Factors', inplace=True)
20 RF4_features_df.plot(kind='barh', figsize=(6, 6))
21 plt.title('Feature Importance with RF Model 4 SMOTE WITH Auto-Hypertuned Class Reweighing')
22
23 # Merge Description to Encoded Features/Factors
24 RF4_features_df = RF4_features_df.rename_axis('Factors').reset_index() # Reset index before copying it to column
25 RF4_features_df = pd.merge(RF4_features_df, factor_list, on='Factors', how='left', indicator=True).replace(np.nan, "") # Merge D
26 RF4_features_df['Description'] = np.where(RF4_features_df['Description'] == '/', RF4_features_df['Factors'], RF4_features_df
27 RF4_features_df['Description'] = RF4_features_df['Description'].replace(r'_', ' ', regex=True) # Replace Underscore with sin
28 final_columns = ['Factors', 'Importance', 'Description'] # List of Columns to keep
29 RF4_features_df = RF4_features_df.drop(columns=[col for col in RF4_features_df if col not in final_columns]) # Drop Columns

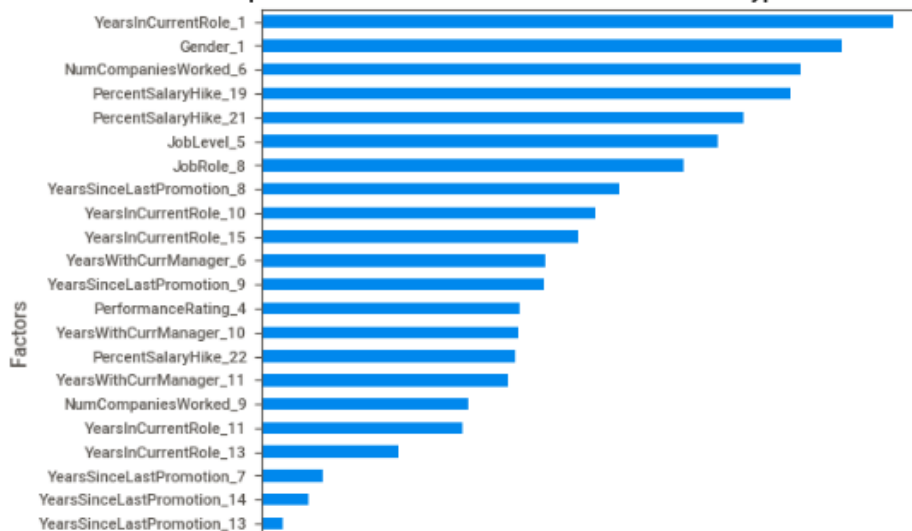
```

```

{'type': 'RForest Opt Smote', 'accuracy': 0.45918367346938777, 'recall': 0.8085106382978723, 'precision': 0.20212765957446807,
'fscore': 0.7248716067498167, 'auc': 0.6006115944525799, 'Factors': 'RF4_features_df'}

```

Feature Importance with RF Model 4 SMOTE WITH Auto-Hypertuned Class Reweighing



2e) RF Model 5 UNDER SAMPLING WITH Auto-Hypertuned Class Reweighing:

A modified algorithm with auto-hypertuned parameters was created using CV grid search method described above which employed random under sampling to achieve class balance.

RF Model 5 UNDER SAMPLING WITH Auto-Hypertuned Class Reweighing

```

1 under_preds, under_model = build_op_rf(X_under, y_under, X_test, best_params=optimal_rf_params)
2 under_results = evaluate(y_test, under_preds, eval_type="RForest Opt Under")
3 under_results["Factors"] = 'RF5_features_df' # Add Important Factors to Dictionary
4 print(under_results)
5
6 # Get Factors Contributing to Attrition in Dataframe
7 RF5_features_df = feature_importance(X, under_model)[:40]
8 RF5_features_df = RF5_features_df[RF5_features_df['Importance'] > 0] # Filter Features with Negative Correlation
9 RF5_features_df = RF5_features_df.sort_values(by=['Importance'], ascending=True) # Attach Importance to data before split
10
11 # Capture Model Drivers in Dataframe
12 RF5_drivers = [['Uses undersampling to delete examples from majority class in order to achieve balance among classes'],
13               ['Uses cv grid method to find optimal parameters to achieve best hypertuning'],
14               ['Runs random forest model with above drivers to get top contributory features']] # initialize list of lists
15 RF5_drivers = pd.DataFrame(RF5_drivers, columns=['Main Drivers behind RF5 model are as follows:']) # Create the pandas DataF
16 RF5_drivers.index = np.arange(1, len(RF5_drivers)+1)
17
18 # Plot Features
19 RF5_features_df.set_index('Factors', inplace=True)
20 RF5_features_df.plot(kind='barh', figsize=(6, 6))
21 plt.title('Feature Importance with RF Model 5 UNDER SAMPLING WITH Auto-Hypertuned Class Reweighing')
22
23 # Merge Description to Encoded Features/Factors
24 RF5_features_df = RF5_features_df.rename_axis('Factors').reset_index() # Reset index before copying it to column
25 RF5_features_df = pd.merge(RF5_features_df, factor_list, on='Factors', how='left', indicator=True).replace(np.nan, "/") # Merge
26 RF5_features_df['Description'] = np.where(RF5_features_df['Description'] == '/', RF5_features_df['Factors'], RF5_features_df
27 RF5_features_df['Description'] = RF5_features_df['Description'].replace(r'_ ', ' ', regex=True) # Replace Underscore with sin
28 final_columns = ['Factors', 'Importance', 'Description'] # List of Columns to keep
29 RF5_features_df = RF5_features_df.drop(columns=[col for col in RF5_features_df if col not in final_columns]) # Drop Columns

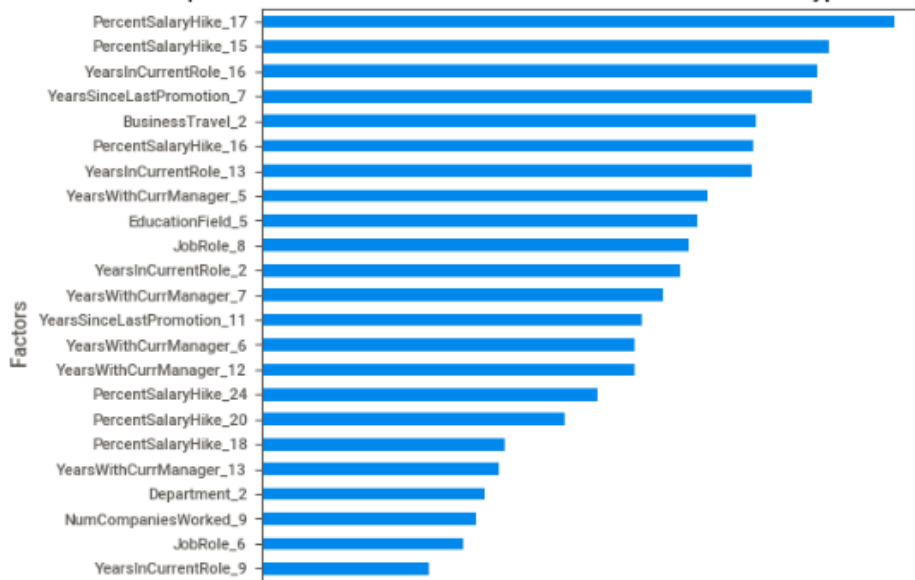
```

```

{'type': 'RForest Opt Under', 'accuracy': 0.29931972789115646, 'recall': 0.9148936170212766, 'precision': 0.17551020408163265,
'fscore': 0.7873239436619718, 'auc': 0.5485399259195451, 'Factors': 'RF5_features_df'}

```

Feature Importance with RF Model 5 UNDER SAMPLING WITH Auto-Hypertuned Class Reweighing



2f) RF Model 6 OVER SAMPLING WITHOUT Auto-Hypertuned Class Reweighing: A modified algorithm was created without achieving any class balance or hyperparameter tuning which employed random over sampling to achieve class balance.

RF Model 6 OVER SAMPLING WITHOUT Auto-Hypertuned Class Reweighing

```

1 over_preds, over_model = build_op_rf(X_over, y_over, X_test, best_params=None)
2 over_results = evaluate(y_test, over_preds, eval_type="RForest Over")
3 over_results["Factors"] = 'RF6_features_df' # Add Important Factors to Dictionary
4 print(over_results)
5
6 # Get Factors Contributing to Attrition in Dataframe
7 RF6_features_df = feature_importance(X, over_model)[:40]
8 RF6_features_df = RF6_features_df[RF6_features_df['Importance'] > 0] # Filter Features with Negative Correlation
9 RF6_features_df = RF6_features_df.sort_values(by=['Importance'], ascending=True) # Attach Importance to data before split
10
11 # Capture Model Drivers in Dataframe
12 RF6_drivers = [['Uses oversampling to randomly selected examples from majority class in order to achieve balance among class',
13               'Runs random forest model with above drivers to get top contributory features']] # initialize list of lists
14 RF6_drivers = pd.DataFrame(RF6_drivers, columns=['Main Drivers behind RF6 model are as follows:']) # Create the pandas DataFrame
15 RF6_drivers.index = np.arange(1, len(RF6_drivers)+1)
16
17 # Plot Features
18 RF6_features_df.set_index('Factors', inplace=True)
19 RF6_features_df.plot(kind='barh', figsize=(6, 6))
20 plt.title('Feature Importance with RF Model 6 OVER SAMPLING WITHOUT Auto-Hypertuned Class Reweighing')
21
22 # Merge Description to Encoded Features/Factors
23 RF6_features_df = RF6_features_df.rename_axis('Factors').reset_index() # Reset index before copying it to column
24 RF6_features_df = pd.merge(RF6_features_df, factor_list, on='Factors', how='left', indicator=True).replace(np.nan, "") # Merge D
25 RF6_features_df['Description'] = np.where(RF6_features_df['Description'] == '/', RF6_features_df['Factors'], RF6_features_df
26 RF6_features_df['Description'] = RF6_features_df['Description'].replace(r'_ ', ' ', regex=True) # Replace Underscore with sin
27 final_columns = ['Factors', 'Importance', 'Description'] # List of Columns to keep
28 RF6_features_df = RF6_features_df.drop(columns=[col for col in RF6_features_df if col not in final_columns]) # Drop Columns

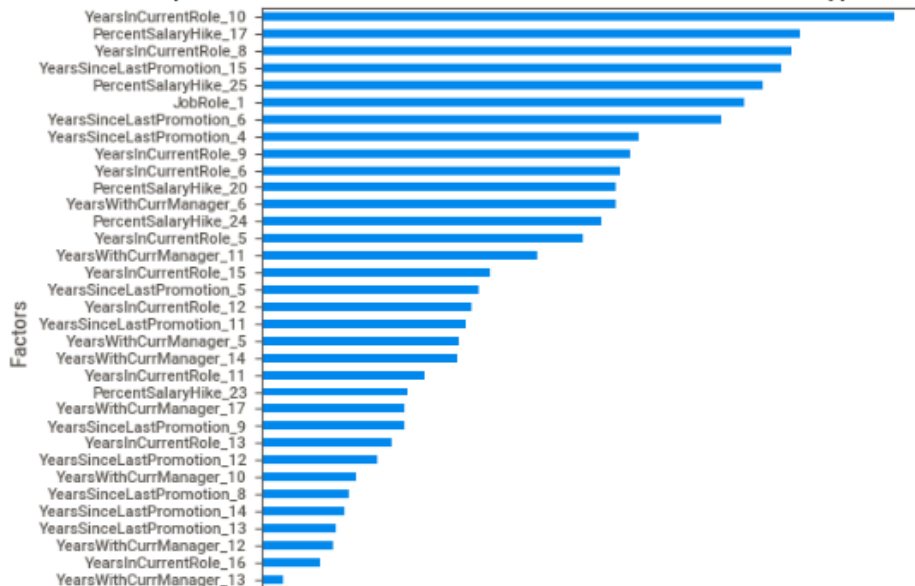
```

```

{'type': 'RForest Over', 'accuracy': 0.6768707482993197, 'recall': 0.723404255319149, 'precision': 0.29310344827586204, 'fscore': 0.6847405112316035, 'auc': 0.6957102248255663, 'Factors': 'RF6_features_df'}

```

Feature Importance with RF Model 6 OVER SAMPLING WITHOUT Auto-Hypertuned Class Reweighing



3) XGB Model

A single method was employed to run XGB models with and without optimal parameters and class weighing:

Method to 'Build XGB WITH or WITHOUT Class Weights'

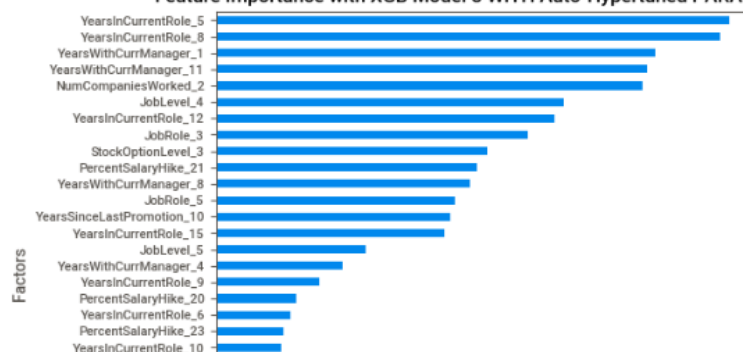
```
1 def build_xgb(X_train, y_train, X_test, threshold=0.5, best_params=None):
2     xgb_model = GradientBoostingClassifier(random_state = rs)
3     # If best parameters are provided
4     if best_params:
5         xgb_model = GradientBoostingClassifier(random_state = rs,
6         # Max depth of each tree
7         max_depth = xgb_optimal_params['max_depth'],
8         # Class weight parameters
9         #class_weight=optimal_params['class_weight'],
10        # Number of trees
11        n_estimators=xgb_optimal_params['n_estimators'],
12        # Minimal samples to split
13        min_samples_split=xgb_optimal_params['min_samples_split'])
14
15    # Train the model
16    xgb_model.fit(X_train, y_train)
17    # If predicted probability is larger than threshold (default value is 0.5), generate a positive label
18    predicted_proba = xgb_model.predict_proba(X_test)
19    yp = (predicted_proba[:,1] >= threshold).astype('int')
20    return yp, xgb_model
```

XGB Model 8 WITH Auto-Hyptuned PARAMETERS

```
1 xgb_preds, xgb_model = build_xgb(X_train,y_train, X_test, best_params=xgb_optimal_params)
2 xgb_results = evaluate(y_test, xgb_preds, eval_type="XGB Opt")
3 xgb_results["Factors"] = 'XGB_features_df' # Add Important Factors to Dictionary
4 print(xgb_results)
5
6 # Get Factors Contributing to Attrition in Dataframe
7 XGB_features_df = feature_importance(X, xgb_model)[:40]
8 XGB_features_df = XGB_features_df[XGB_features_df['Importance'] > 0] # Filter Features with Negative Correlation
9 XGB_features_df = XGB_features_df.sort_values(by=['Importance'],ascending=True) # Attach Importance to data before split
10
11 # Capture Model Drivers in Dataframe
12 XGB_drivers = [['Uses cv grid method to find optimal parameters to achieve best hypertuning'],
13                ['Uses train test and optimal parameters to achieve right balance among classes'],
14                ['Runs XGB model with above drivers to get top contributory features']] # initialize List of Lists
15 XGB_drivers = pd.DataFrame(XGB_drivers, columns=['Main Drivers behind RF6 model are as follows:']) # Create the pandas DataF
16 XGB_drivers.index = np.arange(1,len(XGB_drivers)+1)
17
18 # Plot Features
19 XGB_features_df.set_index('Factors', inplace=True)
20 XGB_features_df.plot(kind='barh', figsize=(6, 6))
21 plt.title('Feature Importance with XGB Model 8 WITH Auto-Hyptuned PARAMETERS')
22
23 # Merge Description to Encoded Features/Factors
24 XGB_features_df = XGB_features_df.rename_axis('Factors').reset_index() # Reset index before copying it to column
25 XGB_features_df = pd.merge(XGB_features_df,factor_list,on='Factors',how='left',indicator=True).replace(np.nan, "") # Merge D
26 XGB_features_df['Description'] = np.where(XGB_features_df['Description'] == '/', XGB_features_df['Factors'], XGB_features_df
27 XGB_features_df['Description'] = XGB_features_df['Description'].replace(r'_',' ', regex=True) # Replace Underscore with sin
28 final_columns = ['Factors', 'Importance', 'Description'] # List of Columns to keep
29 XGB_features_df = XGB_features_df.drop(columns=[col for col in XGB_features_df if col not in final_columns]) # Drop Columns
30
31 # Get Model Execution Time
32 print('XGB Model Execution Time: ', datetime.now() - xgb_start_time)
```

```
{'type': 'XGB Opt', 'accuracy': 0.8469387755102041, 'recall': 0.2553191489361702, 'precision': 0.5454545454545454, 'fscore': 0.
2606516290726817, 'auc': 0.6074166594883279, 'Factors': 'XGB_features_df'}
XGB Model Execution Time: 4:02:32.350523
```

Feature Importance with XGB Model 8 WITH Auto-Hyptuned PARAMETERS



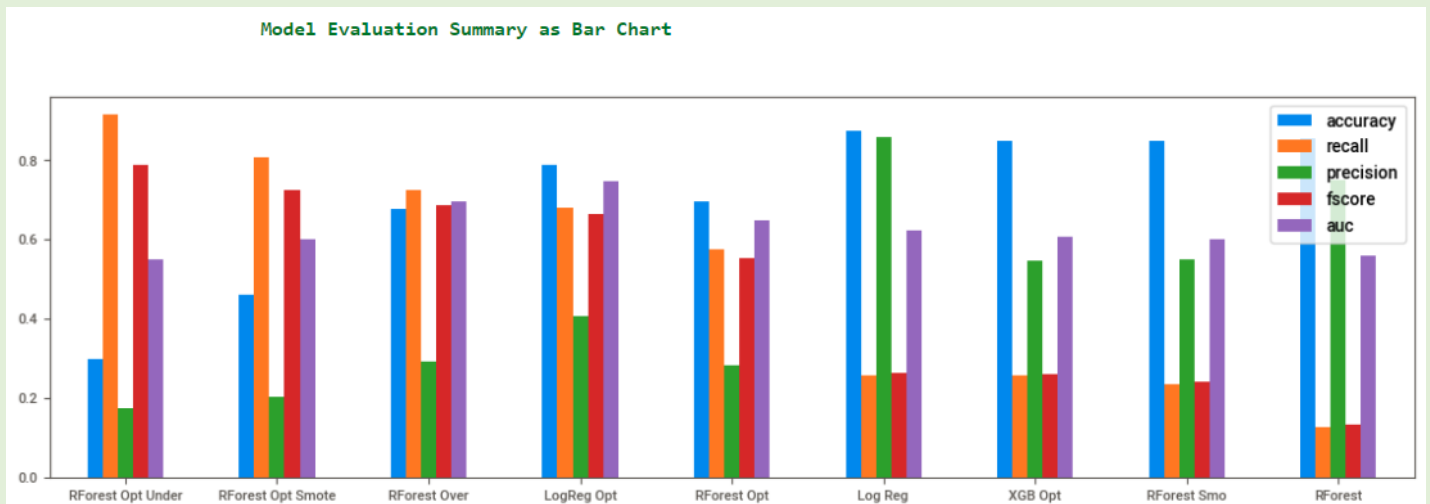
5) Recommended Model

5a) Result Summary

Result Summary is given below:

	type	accuracy	recall	precision	fscore
0	RForest Opt Under	0.299320	0.914894	0.175510	0.787324
1	RForest Opt Smote	0.459184	0.808511	0.202128	0.724872
2	RForest Over	0.676871	0.723404	0.293103	0.684741
3	LogReg Opt	0.789116	0.680851	0.405063	0.663477
4	RForest Opt	0.697279	0.574468	0.281250	0.552321
5	Log Reg	0.874150	0.255319	0.857143	0.262405
6	XGB Opt	0.846939	0.255319	0.545455	0.260652
7	RForest Smo	0.846939	0.234043	0.550000	0.239331
8	RForest	0.853741	0.127660	0.750000	0.131868

5b) Overall Visual Summary

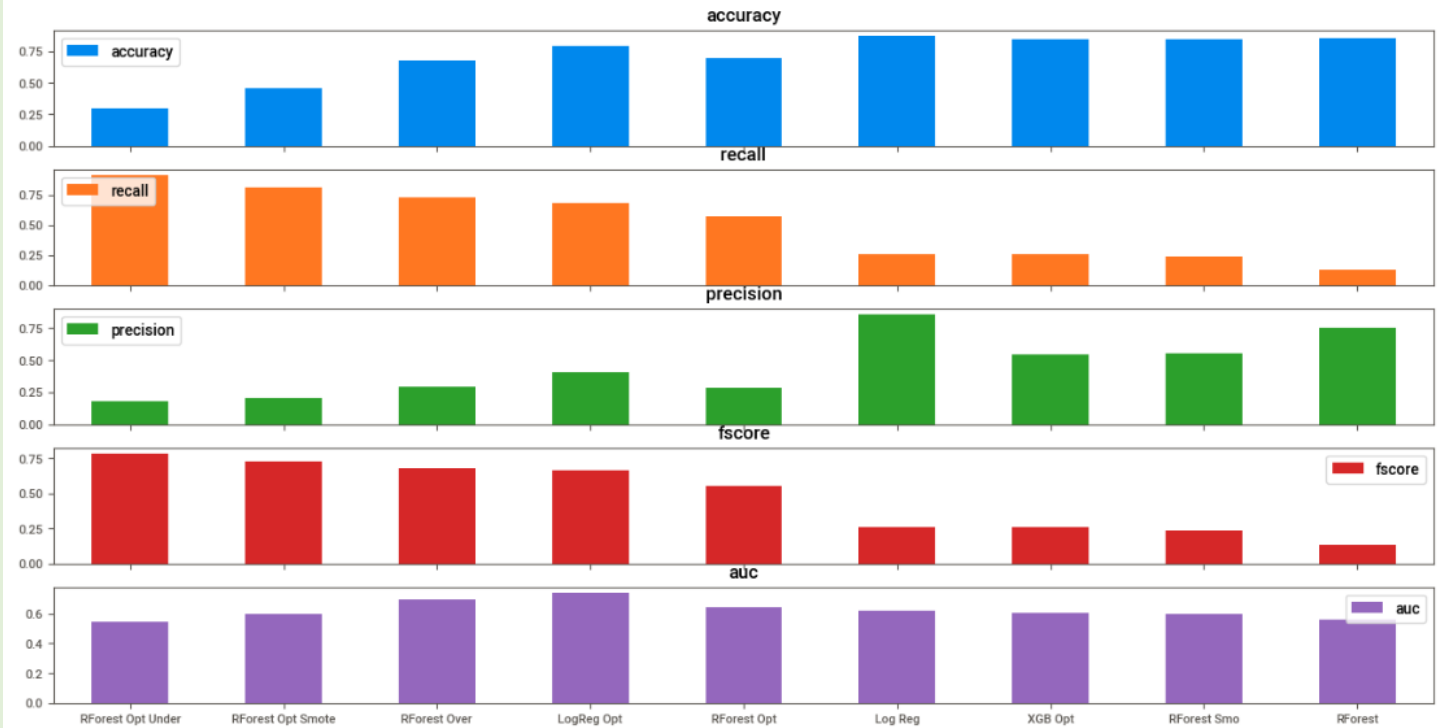


5c) Individual Model Visual Summary

3. Visualize Result Metrics Summary as Individual Bar Charts

```
1 axes = chart.plot.bar(rot=0, subplots=True, figsize=(16, 8))
2 axes[1].legend(loc=2)
```

<matplotlib.legend.Legend at 0x2e9ec702e80>



5d) Model Choice and Justification

For many machine learning tasks with imbalanced datasets, like Employee Attrition, we normally care more about Recall than precision. As a baseline, we want the model to be able to find all possible factors and so, we would allow the model to make false-positive errors because the cost of false positives is usually not very high (maybe it will just cost a false notification email or phone call to confirm with employee).

On the other hand, failing to recognize positive examples (such as employee wanting to leave) can be too costly for the organization. **As such, our first priority is to improve model's recall; then we will also want to keep precision as high as possible.**

In this case, the Model with Best Recall and F-Score is RForest Opt Under. Hence, we will select this model for Employee Attrition Prediction.

6) Summary Key Findings and Insights

6a) Summarizing Model Drivers:

Main Drivers behind top performing RF5 model are as follows:

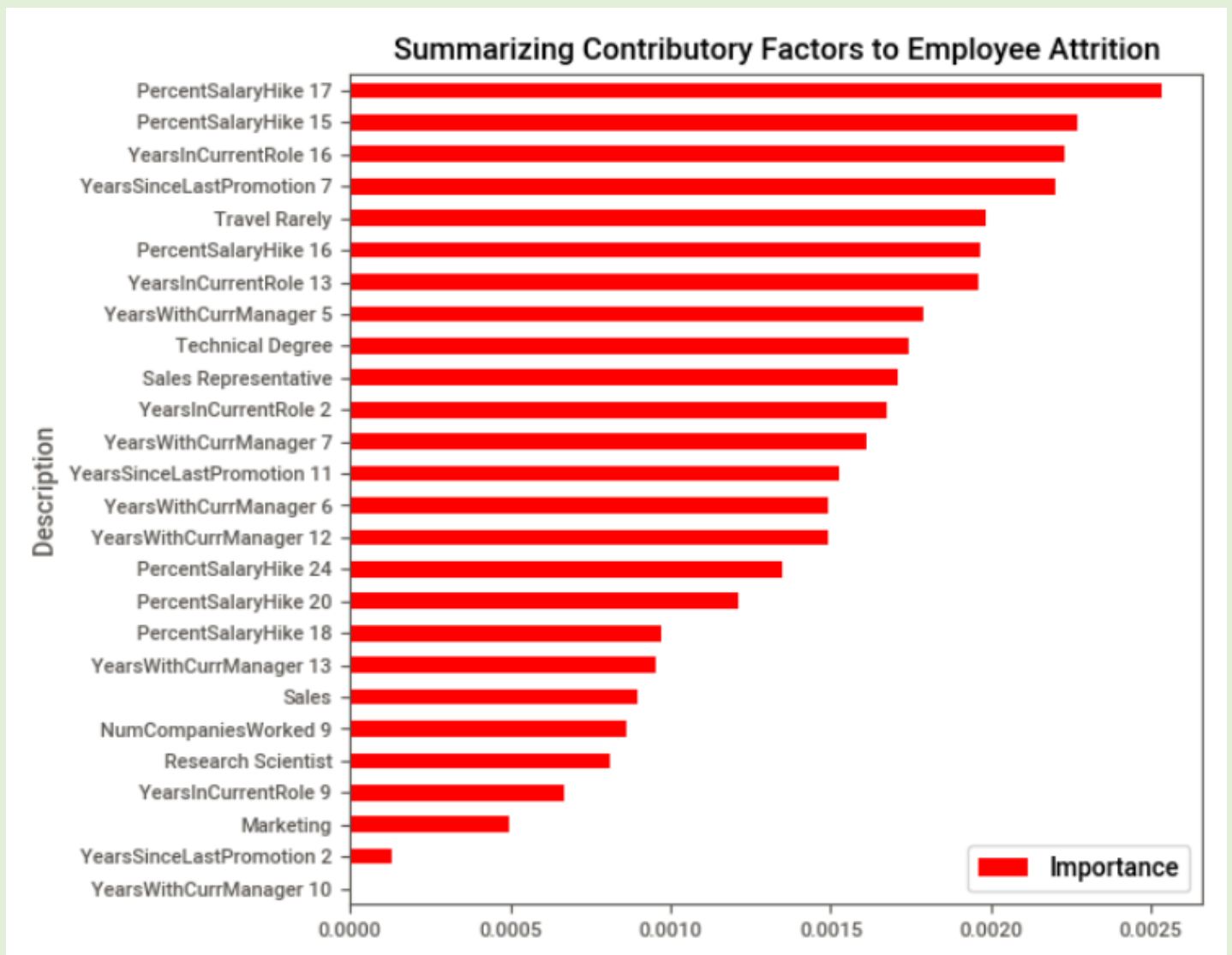
- Uses under-sampling to delete examples from majority class in order to achieve balance among classes
- Uses cv grid method to find optimal parameters to achieve best hyper-tuning
- Runs random forest model with above drivers to get top contributory features

6b) Enlisting Top Contributory Factors

Top Factors Contributing to Employee Attrition in ascending order of importance are as follows:

Top Contributory Factors:			
	Factors	Importance	Description
0	YearsWithCurrManager_10	2.707586e-18	YearsWithCurrManager 10
1	YearsSinceLastPromotion_2	1.274117e-04	YearsSinceLastPromotion 2
2	EducationField_2	4.953221e-04	Marketing
3	YearsInCurrentRole_9	6.702776e-04	YearsInCurrentRole 9
4	JobRole_6	8.088280e-04	Research Scientist
5	NumCompaniesWorked_9	8.604488e-04	NumCompaniesWorked 9
6	Department_2	8.957229e-04	Sales
7	YearsWithCurrManager_13	9.515419e-04	YearsWithCurrManager 13
8	PercentSalaryHike_18	9.725873e-04	PercentSalaryHike 18
9	PercentSalaryHike_20	1.214221e-03	PercentSalaryHike 20
10	PercentSalaryHike_24	1.346836e-03	PercentSalaryHike 24
11	YearsWithCurrManager_12	1.492372e-03	YearsWithCurrManager 12
12	YearsWithCurrManager_6	1.492726e-03	YearsWithCurrManager 6
13	YearsSinceLastPromotion_11	1.524371e-03	YearsSinceLastPromotion 11
14	YearsWithCurrManager_7	1.609366e-03	YearsWithCurrManager 7
15	YearsInCurrentRole_2	1.675148e-03	YearsInCurrentRole 2
16	JobRole_8	1.711574e-03	Sales Representative
17	EducationField_5	1.744438e-03	Technical Degree
18	YearsWithCurrManager_5	1.788844e-03	YearsWithCurrManager 5
19	YearsInCurrentRole_13	1.962220e-03	YearsInCurrentRole 13
20	PercentSalaryHike_16	1.969210e-03	PercentSalaryHike 16
21	BusinessTravel_2	1.982832e-03	Travel Rarely
22	YearsSinceLastPromotion_7	2.201833e-03	YearsSinceLastPromotion 7
23	YearsInCurrentRole_16	2.228832e-03	YearsInCurrentRole 16
24	PercentSalaryHike_15	2.271719e-03	PercentSalaryHike 15
25	PercentSalaryHike_17	2.533883e-03	PercentSalaryHike 17

6c) Visualizing Top Contributory Factors to Employee Attrition



7) Link to Other Useful Models

- a) <https://github.com/IBM/employee-attrition-aif360/blob/master/notebooks/employee-attrition.ipynb>
- b) https://github.com/JNYH/employee_attrition/blob/master/employee_attrition.ipynb
- c) <https://github.com/elastic/examples/tree/master/Machine%20Learning/Analytics%20Jupyter%20Notebooks>
- d) https://github.com/ganesh10-india/HR_Analytics-Employee_Attrition-Classification-Models/blob/main/HR_Analytics_Employee_Attrition_Classification_Models.ipynb

8) Github Link to Assignment Notebook

<https://github.com/FATIMASP/IBM-MACHINE-LEARNING-CERTIFICATION/tree/main/Supervised%20Machine%20Learning:%20Classification>