

EMPLOYEE CHURN PREDICTION

**A project report submitted to Department of Computer Science and Engineering for
partial fulfilments of the requirements for the award of the degree**

of

BACHELOR OF

TECHNOLOGY IN

COMPUTER SCIENCE AND ENGINEERING

Submitted by

B.Teja Sri Venkat, 121910317023

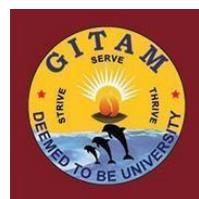
Jana Venkata Chaitanya, 121910317060

M.Subhani, 121910317038

S.Sai Teja, 121910317048

Under the esteemed guidance of

Smt N. Madhuri



DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING

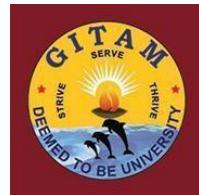
GITAM

(Deemed to be University)

VISAKHAPATNAM

OCTOBER 2022

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
GITAM SCHOOL OF TECHNOLOGY
GITAM
(Deemed to be University)



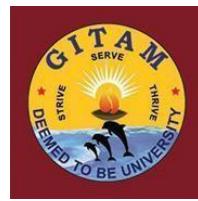
DECLARATION

We, hereby declare that the project report entitled "**Employee Churn Prediction**" is an original work done in the Department of Computer Science and Engineering, GITAM School of Technology, GITAM (Deemed to be University) submitted in partial fulfillment of the requirements for the award of the degree of B.Tech. in Computer Science and Engineering. The work has not been submitted to any other college or University for the award of any degree or diploma.

Date:17-10-2022

Registration No(s).	Name(s)	Signature(s)
121910317023	B.Teja Sri Venkat	
121910317060	Jana Venkata Chaitanya	
121910317038	D.Subhani	
121910317048	S.Sai Teja	

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
GITAM SCHOOL OF TECHNOLOGY
GITAM
(Deemed to be University)



CERTIFICATE

This is to certify that the project report entitled "**Employee Churn Prediction**" is a bonafide record of work carried out by **B.Teja Sri Venkat(121910317023), D.Subhani(121910317038), J.V.Chaitanya(121910317060), S.Sai Teja(121910317048)** students submitted in partial fulfillment of requirement for the award of degree of Bachelors of Technology in Computer Science and Engineering.

Project Guide

Smt N. Madhuri

Assistant Professor

**Head of the
Department**

Dr. R. Sireesha

Professor

TABLE OF CONTENTS

1.	Abstract	1
2.	Introduction	2
3.	Literature Review	3
4.	Problem Identification & Objectives	4
5.	System Methodology	5
6.	Overview of Technologies	8
7.	Implementation	10
8.	Results	12
9.	Conclusion	13
10.	References	14

TABLE OF DIAGRAMS

Fig 1	Use Case Diagram	5
Fig 2	Sequence Diagram	5
Fig 3	Class Diagram	6
Fig 4	State Chart Diagram	6
Fig 5	Component Diagram	7
Fig 6	Communication Diagram	7

ABSTRACT

An employee is a valuable resource who works for a company. If one of them abruptly departs the company, it may have a significant impact and expense to the specific company. And, in addition to time and money, the newly hired employee requires some time to make the particular firm cost-effective. This model will assist in predicting the rate at which people quit occupations based on available analytic data and will apply various machine learning algorithms to reduce prediction error. The prediction of a personalized or individual employee varies depending on the setting in which they work. While it has become clear that employee churn forecasts differ depending on geography, lifestyle, and environment, the associated knowledge and understanding remains fragmented.

INTRODUCTION

Employee churn refers to the departure or retirement of a knowledgeable individual or asset from a business. Instead, we can state that churn occurs when an employee departs an organization. This churn could be forced or voluntary on the part of the employee. Employee churn refers to any employee who wishes to leave a company or whose contract duration has expired, resulting in an unfavorable reduction in the organization. When former employees are replaced by newly hired employees, this is referred to as employee turnover. Both turnover and attrition are linked. This raises the entire cost of acquiring new personnel, managing the recruitment process, and training newly hired employees. All minor issues that arise following an employee's departure from the organization contribute to the integration of subsequent issues that may have an indirect impact on the company. Employee turnover is unavoidable in any sector. Avoiding such a massive disaster is a difficult task for everyone. Companies can avoid this problem by knowing the factor generating abrupt attrition. They can also inspire employees to stay with the organization. The employment and termination conditions determine the frequency with which employees leave an organization. Each employee may have the same or a different reason for leaving a company. To reduce turnover, the human resource team and management must carefully select the parameters.

Employee turnover occurs for a variety of reasons. If the company is losing money, some employees may be let go, or they may leave for a better future and job satisfaction. This prompted the search for a new job replacement. In order for the firm to work effectively and profitably, we must implement a good churn prediction system that will not only reduce the possibility of major financial losses but will also support overall growth of each and every employee.

LITERATURE REVIEW

The Techniques used in [1] are Area Under Curve, Standard mechanism and Social network analysis and the model takes full account of tiny features in the dataset while training, thus detecting smaller objects more efficiently. The AUD & SNA is a lightweight model, multiple.

The techniques used in [2] are Stratified K-fold cross validation technique, the model consists of Different algorithms were investigated to get good precision and recall. Through a series of experiments, It is concluded that Random Forest model Performs better than other algorithms.

The techniques used in [3] are Artificial neural networks, logistic regression, classification and regression trees (CART), classification trees (C5.0), and discriminant analysis) .This technique is capable of handling multiple objects using ANN, regression and classification techniques. It is capable of dealing large data using the simplest algorithms present.

The techniques used in [4] are Naïve Bayes, Support Vector Machines, Logistic Regression, Decision Trees and Random Forests. Demonstration of SVM technique can be used to build reliable and accurate predictive models for employee churn

The techniques used in [5] are Multiple analysis of variance (MANOVA), independent t-test, factor analysis, multiple regression, and binomial logistic regression analysis. The hypothesis was not completely supported when all of the factors required to continue with the prediction were considered. It is limited to a single company and cannot be used for all convenience.

PROBLEM IDENTIFICATION

Employee Churn is an expensive issue for businesses. The true cost of replacing an employee is frequently fairly high. According to a research conducted by the Center for American Progress, corporations normally pay around one-fifth of an employee's compensation to replace that individual, and the cost can skyrocket if CEOs or the highest-paid staff have to be replaced. In other words, the cost of replacing staff remains high for the majority of firms. This is due to the time spent interviewing and finding a successor, sign-on bonuses, and the loss of production for several months as the new employee adjusts to the new role. Previously, the majority of the focus was on "rates" such as attrition and retention rates. Using data warehousing techniques, HR managers compute prior rates and attempt to anticipate future rates. These rates represent the collective impact of churn, but they only provide half of the picture. Another strategy is to concentrate on individual records rather than aggregates. There are numerous case studies available on client churn. Employee turnover is analogous to consumer turnover. It is primarily concerned with the employee rather than the customer. You can predict who and when an employee will leave the company. Employee attrition is costly, and minor improvements will yield huge results.

OBJECTIVES

The main Objective of the Project is to solve the issue of employee attrition and retention. So the HR department can recruit the employee as they desire. The algorithm works of the previous and past source data of the employee, the data collected is in the form of a dataset and undergoes various process as discussed below and generates the attrition rate

1. DATA COLLECTION: A prediction system's first phase entails gathering data and dividing it into training and testing datasets. In this project, 30% of the dataset was used for testing and 70% was used for training. Employee datasets include all the information needed to predict employee attrition, including employee id, satisfaction level, most recent evaluation, number of projects, average monthly hours worked, time spent at the company (years), departments, and salary level.

2. ATTRIBUTE ANALYSIS: Radar charts, correlation matrices, and pair plots are used as data visualization techniques to examine the dataset's properties. The characteristics that significantly influence employee attrition are examined. During the procedure, irrelevant attributes like employee ID are removed.

3. DATA PREPROCESSING: Data preprocessing is crucial to machine learning models since it converts the data to the necessary format. This keeps the results from being misinterpreted. The issues of redundant data, unbalanced data, missing values, and categorical values are dealt with during data pre-processing. Dummy variables are dealt with using label encoding.

4. EMPLOYEE CHURN: Different machine learning classification models are used for classification, including Decision Tree, KNN, SVC, and Light GBM. Comparative analysis is used to compare the algorithms. In the suggested approach, the algorithm that has the highest accuracy rating is used to forecast Employee churn.

SYSTEM METHODOLOGY

The project's suggested system architecture provides a summary of how the system functions. The project begins with data gathering, followed by attribute analysis and data pre-processing. The data is preprocessed into the necessary format so that machine learning classification techniques can be used. The data must be divided into training and testing groups in the following phase. On the training dataset, machine learning techniques are used, and on the testing dataset, performance metric analysis is performed later. The following modules are included in this project:

1. Dataset collection
2. Attribute analysis
3. Data Preprocessing
4. Employee churn

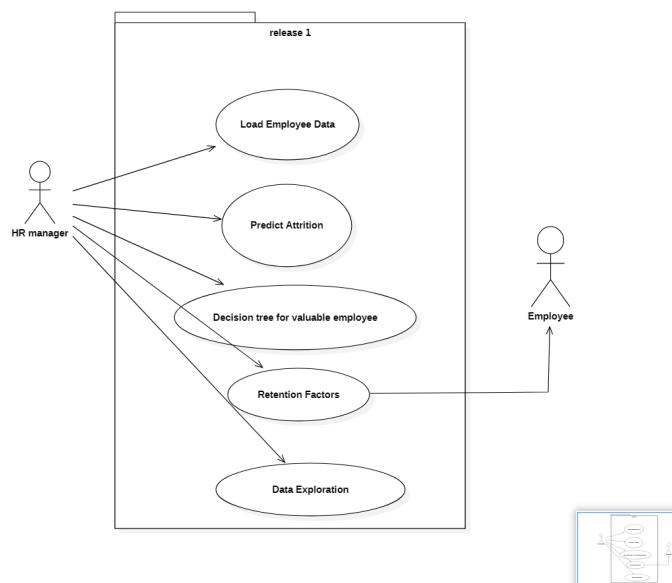


Figure 1: Use Case Diagram

As the diagram shows both the HR manager and employee have their roles to play .And their Operations are performed in a sequential manner .First the HR loads the data and basing on the data the attrition prediction will be done and basing on the attrition rate retention factors of the employee are observed.

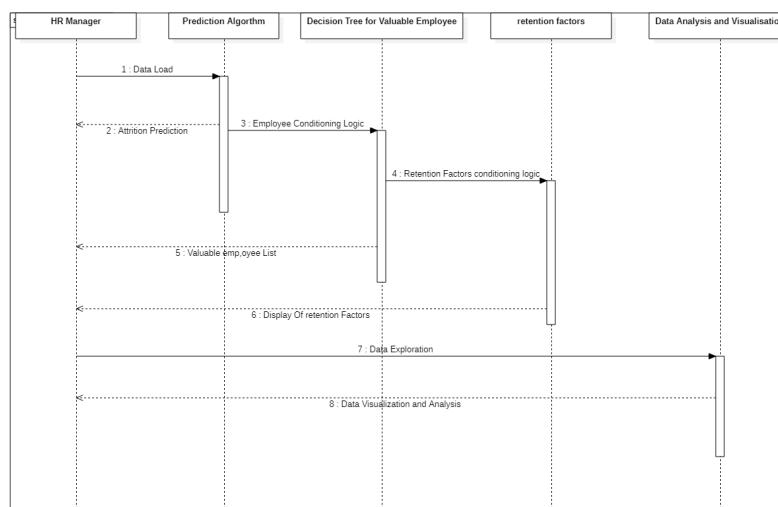


Figure 2: Sequence Diagram

This Diagram depicts the overall operations occurring in a correct sequential manner performed from time to time by various users

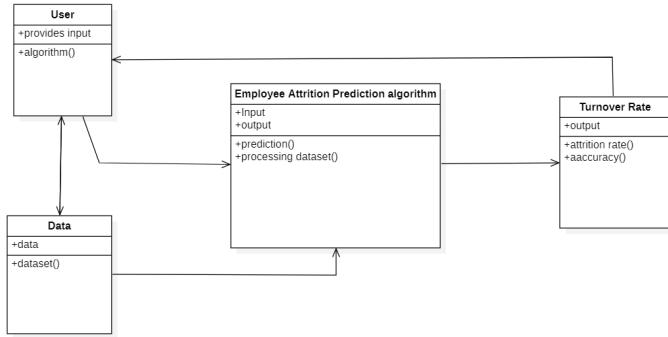


Figure 3: Class Diagram

The attributes and operations performed by various Class entities are given .And this also shows how various entities are correlated to each other

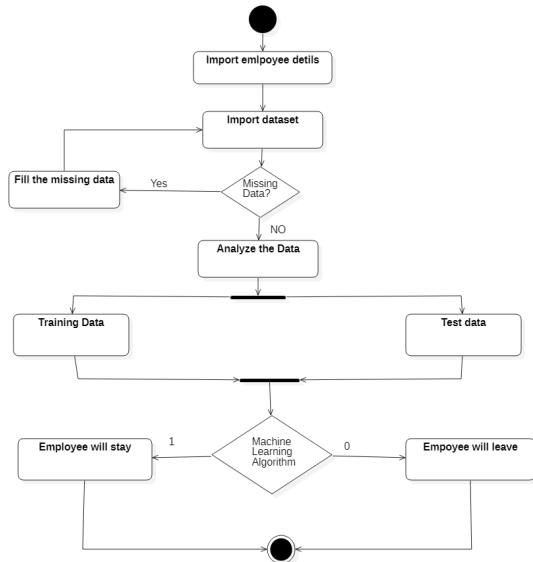


Figure 4: State Chart Diagram

In Activity and State Chart diagrams ,it shows various transactions between objects and other related data in a step by step order

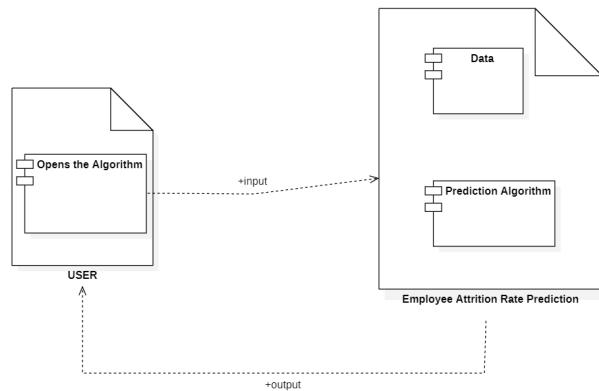


Figure 5: Component Diagram

Shows all the Components and the internal processes occurring in the algorithm .This practically displays all the internal operations from Input to the Output

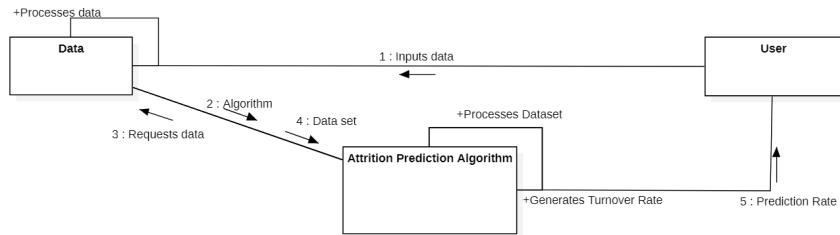


Figure 6: Communication Diagram

Provides all the information of how the internal communication is going between all the working components in the algorithm

OVERVIEW OF TECHNOLOGIES

GRADIENT BOOSTING ALGORITHM:

One of the most potent algorithms used in machine learning is the gradient boosting technique. As is well known, bias and variance mistakes can be used to broadly categorize errors in machine learning systems. As one of the boosting strategies, gradient boosting is used to reduce the model's bias error.

We are unable to discuss the base estimator in the gradient boosting process, unlike the Adaboosting algorithm. The Gradient Boost algorithm's default estimator, Decision Stump, is fixed. The gradient boosting algorithm's n estimator can be adjusted, just like AdaBoost. The default value of n estimator for this algorithm, however, is 100 if we do not specify a number for it.

Both continuous and categorical target variables can be predicted using the gradient boosting approach (as a Regressor) (as a Classifier). The cost function is Mean Square Error (MSE) when it is used as a regressor, while it is Log loss when it is used as a classifier.

GRADIENT DESCENT CLASSIFIER:

In order to train machine learning models, gradient descent is one of the most often utilized optimization techniques. It does this by reducing the discrepancy between the actual and expected outputs. Furthermore, neural networks are trained using gradient descent.

The term "optimization algorithm" in mathematics describes the process of minimizing or maximizing an objective function $f(x)$ parameterized by x . Similar to this, the goal of optimization in machine learning is to reduce the cost function, which is parameterized by the model's parameters. Utilizing parameter updates in iterations, the primary goal of gradient descent is to minimize the convex function. These machine learning models can be utilized as effective tools for artificial intelligence and a variety of computer science applications once they have been tuned.

KNN ALGORITHM:

With k-NN, all computation is postponed until after the function has been evaluated and the function is only locally approximated. Since this technique relies on distance for classification, normalizing the training data can significantly increase accuracy if the features reflect several physical units or have distinct sizes.

Assigning weights to neighbor contributions can be a helpful strategy for both classification and regression, making the closer neighbor contribute more to the average than the farther neighbor. As an illustration, a typical weighting method assigns each neighbor a weight of $1/d$, where d is the distance between the neighbors.

When using k-NN classification or regression, the neighbors are chosen from a set of objects for which the class or object property value is known. Although there is no need for an explicit training step, this can be considered of as the algorithm's training set.

The k-NN algorithm has the feature of being sensitive to the local structure of the data. The existence of noisy or irrelevant features, or if the feature scales are inconsistent with their relevance, can significantly reduce the accuracy of the k-NN method. A lot of study has gone into choosing or scaling characteristics to enhance classification. The application of evolutionary algorithms to enhance feature scaling is a particularly well-liked method. Scaling features based on the mutual information of training data and training classes is a popular alternative.

RANDOM FOREST CLASSIFIER:

A large number of decision trees are built during the training phase of the random forests or random decision forests ensemble learning approach, which is used for classification, regression, and other tasks. The class that the majority of the trees chose is the output of the random forest for classification problems. The mean or average prediction of each individual tree is returned for regression tasks. The tendency of decision trees to overfit their training set is corrected by random decision forests. Although they frequently outperform decision trees, gradient boosted trees are more accurate than random forests. [Reference needed] However, their effectiveness may be impacted by data peculiarities.

Random forest predictors automatically produce a dissimilarity metric between the observations as part of their creation. The objective is to build a random forest predictor that separates the "observed" data from adequately generated synthetic data. This is also known as a random forest dissimilarity measure between unlabeled data. The synthetic data are taken from a reference distribution, whereas the observed data are the original unlabeled data.

CATBOOST CLASSIFIER:

CatBoost is a technique for decision trees that uses gradient boosting. It is created by Yandex researchers and engineers and is used for a variety of jobs at Yandex and in other businesses, such as CERN, Cloudflare, and Careem taxi, including search, recommendation systems, personal assistants, self-driving cars, weather prediction, and many more. Everyone is welcome to use it because it is open-source.

There are numerous methods for handling categorical features in boosted trees, which are frequently present in datasets. CatBoost automatically handles categorical features in contrast to other gradient boosting techniques (which need numeric input). One-hot encoding is one of the most popular methods for dealing with categorical data, however it is impractical for many characteristics. Target statistics are used to categorize features in order to address this (estimate target value for each category). There are several techniques to calculate the target statistics: greedy, hold out, leave one out, and ordered. Target statistics are ordered in CatBoost.

IMPLEMENTATION

1. Coding

```
[1] pip install pyforest  
pip install --user pycaret  
pip install pandas_profiling  
pip install squarify  
pip install termcolor  
pip install catboost  
pip install pyclustertend  
  
[2] > import pandas_profiling  
import pyforest  
  
import ipywidgets  
from ipywidgets import interact  
  
import numpy as np  
import pandas as pd  
  
import seaborn as sns  
import matplotlib.pyplot as plt  
%matplotlib inline  
import matplotlib.ticker as mticker  
import squarify as sq  
  
[3] import plotly.express as px  
import cufflinks as cf  
import plotly.offline  
cf.go_offline()  
cf.set_config_file(offline=False, world_readable=True)  
  
...  
  
[4] pip install yellowbrick  
  
[5] pip install --upgrade scikit-learn  
  
[6] import colorama  
from colorama import Fore, Style  
from termcolor import colored  
from termcolor import cprint  
  
import scipy.stats as stats  
from scipy.cluster.hierarchy import linkage, dendrogram  
import statsmodels.api as sm  
import statsmodels.formula.api as smf  
import missingno as msno  
  
import datetime as dt  
from datetime import datetime  
  
from pyclustertend import hopkins  
from sklearn.cluster import KMeans, AgglomerativeClustering  
from sklearn.compose import make_column_transformer  
from sklearn.decomposition import PCA  
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis  
from sklearn.ensemble import RandomForestRegressor, RandomForestClassifier, GradientBoostingRegressor  
from sklearn.ensemble import ExtraTreesRegressor, AdaBoostClassifier  
from sklearn.feature_selection import SelectKBest, SelectPercentile, f_classif, f_regression, mutual_info_regression
```

```

from scipy.cluster.hierarchy import linkage, dendrogram
import statsmodels.api as sm
import statsmodels.formula.api as smf
import missingno as mno

import datetime as dt
from datetime import datetime

from pyclustertend import hopkins
from sklearn.cluster import KMeans, AgglomerativeClustering
from sklearn.compose import make_column_transformer
from sklearn.decomposition import PCA
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
from sklearn.ensemble import RandomForestRegressor, RandomForestClassifier, GradientBoostingRegressor
from sklearn.ensemble import ExtraTreesRegressor, AdaBoostClassifier
from sklearn.feature_selection import SelectKBest, SelectPercentile, f_classif, f_regression, mutual_info_regression
from sklearn.impute import SimpleImputer, KNNImputer
from sklearn.linear_model import LinearRegression, Lasso, Ridge, ElasticNet

from sklearn.metrics import silhouette_samples, silhouette_score
from sklearn.metrics.cluster import adjusted_rand_score
from sklearn.model_selection import RandomizedSearchCV
from sklearn.model_selection import RepeatedStratifiedKFold, KFold, cross_val_predict, train_test_split
from sklearn.model_selection import StratifiedKFold, GridSearchCV, cross_val_score, cross_validate
from sklearn.naive_bayes import GaussianNB
from sklearn.neighbors import KNeighborsRegressor, KNeighborsClassifier
from sklearn.pipeline import make_pipeline, Pipeline
from sklearn.preprocessing import scale, StandardScaler, MinMaxScaler, RobustScaler
from sklearn.preprocessing import PolynomialFeatures, OneHotEncoder, PowerTransformer, LabelEncoder
from sklearn.svm import SVR, SVC
from sklearn.tree import plot_tree, DecisionTreeClassifier
from yellowbrick.classifier import ClassificationReport

from xgboost import XGBRegressor, XGBClassifier, plot_importance

def missing_values(df):
    missing_number = df.isnull().sum().sort_values(ascending = False)
    missing_percent = (df.isnull().sum() / df.isnull().count()).sort_values(ascending = False)
    missing_values = pd.concat([missing_number, missing_percent], axis = 1, keys = ['Missing_Number', 'Missing_Percent'])
    return missing_values[missing_values['Missing_Number'] > 0]

def first_looking(df):
    print(colored("Shape:", attrs=['bold']), df.shape, '\n',
          colored("****100", 'red', attrs = ['bold']),
          colored("\nInfo:\n", attrs = ['bold']), sep = ''))
    print(df.info(), '\n',
          colored("****100", 'red', attrs = ['bold']), sep = ''))
    print(colored("Number of Unique:\n", attrs = ['bold']), df.nunique(), '\n',
          colored("****100", 'red', attrs = ['bold']), sep = ''))
    print(colored("Missing Values:\n", attrs = ['bold']), missing_values(df), '\n',
          colored("****100", 'red', attrs = ['bold']), sep = ''))
    print(colored("All columns:", attrs = ['bold']), list(df.columns), '\n',
          colored("****100", 'red', attrs = ['bold']), sep = ''))
    df.columns = df.columns.str.lower().str.replace('&', '_').str.replace(' ', '_')
    print(colored("Columns after rename:", attrs = ['bold']), list(df.columns), '\n',
          colored("****100", 'red', attrs = ['bold']), sep = ''))
    print(colored("Columns after renames:", attrs = ['bold']), list(df.columns), '\n',
          colored("****100", 'red', attrs = ['bold']), sep = ''))
    print(colored("Descriptive Statistics \n", attrs = ['bold']), df.describe().round(2), '\n',
          colored("****100", 'red', attrs = ['bold']), sep = '') # gives a statistical breakdown of the data.
    print(colored("Descriptive Statistics (Categorical columns) \n", attrs = ['bold']), df.describe(include = object).T, '\n',
          colored("****100", 'red', attrs = ['bold']), sep = '') # gives a statistical breakdown of the data.

def multicolinearity_control(df):
    feature = []
    collinear = []
    for col in df.corr().columns:
        for i in df.corr().index:
            if (abs(df.corr()[col][i]) > .9 and abs(df.corr()[col][i]) < 1):
                feature.append(col)
                collinear.append(i)
                print(colored(F" Multicolinearity alert in between:{col} - {i}",
                             "red", attrs = ['bold']), df.shape, '\n',
                      colored("****100", 'red', attrs = ['bold']), sep = ''))

def multicolinearity_control(df):
    feature = []
    collinear = []
    for col in df.corr().columns:
        for i in df.corr().index:
            if (abs(df.corr()[col][i]) > .9 and abs(df.corr()[col][i]) < 1):
                feature.append(col)
                collinear.append(i)
                print(colored(F" Multicolinearity alert in between:{col} - {i}",
                             "red", attrs = ['bold']), df.shape, '\n',
                      colored("****100", 'red', attrs = ['bold']), sep = ''))

def duplicate_values(df):
    print(colored("Duplicate check...", attrs = ['bold']), sep = '')
    print("There are", df.duplicated(subset = None, keep = 'first').sum(), "duplicated observations in the dataset.")
    duplicate_values = df.duplicated(subset = None, keep = 'first').sum()
    if duplicate_values > 0:
        df.drop_duplicates(keep = 'first', inplace = True)
        print(duplicate_values, colored(" duplicates were dropped!"), '\n',
              colored("****100", 'red', attrs = ['bold']), sep = ''))

def drop_null(df, limit):
    print('Shape:', df.shape)
    for i in df.isnull().sum().index:
        if (df.isnull().sum()[i] / df.shape[0]*100) > limit:
            print(df.isnull().sum()[i], 'percent of', i, 'null and were dropped')
            df.drop(i, axis = 1, inplace = True)
            print('new shape:', df.shape)

    print('New shape after missing value control:', df.shape)

def first_look(col):
    print("column name : ", col)
    print("-----")
    print("No. of Nulls : ", " %", round((df[col].isnull().sum() / df.shape[0])*100, 2))

```

```

def first_look(col):
    print("column name : ", col)
    print("-----")
    print("Per_of_Nulls : ", "%", round(df[col].isnull().sum() / df.shape[0]*100, 2))
    print("Num_of_Nulls : ", df[col].isnull().sum())
    print("Num_of_Unique : ", df[col].nunique())
    print("Duplicates : ", df.duplicated(subset = None, keep = 'first').sum())
    print(df[col].value_counts(dropna = False))

def fill_most(df, group_col, col_name):
    '''Fills the missing values with the most existing value (mode) in the relevant column according to single-stage grouping'''
    for group in list(df[group_col].unique()):
        cond = df[group_col] == group
        mode = list(df[cond][col_name].mode())
        if mode != []:
            df.loc[cond, col_name] = df.loc[cond, col_name].fillna(df[cond][col_name].mode()[0])
        else:
            df.loc[cond, col_name] = df.loc[cond, col_name].fillna(df[col_name].mode()[0])
    print("Number of Null : ", df[col_name].isnull().sum())
    print("-----")
    print(df[col_name].value_counts(dropna = False))

def show_values_on_bars(axs):
    def _show_on_single_plot(ax):
        for p in ax.patches:
            _x = p.get_x() + p.get_width() / 2
            _y = p.get_y() + p.get_height()
            value = '{:.2f}'.format(p.get_height())
            ax.text(_x, _y, value, ha="center")
        if isinstance(axs, np.ndarray):
            for idx, ax in np.ndenumerate(axs):
                _show_on_single_plot(ax)
        else:
            _show_on_single_plot(axs)

df0 = pd.read_csv('HR_Dataset.csv')
df = df0.copy()
df.head(3)

```

Python

```

df.profile_report()

```

Python

```

first_looking(df)
duplicate_values(df)
print(colored("Shape:", attrs = ['bold']), df.shape, '\n', colored('*'*100, 'red', attrs = ['bold']))

```

Python

```

df.columns

```

Python

```

df.rename({'departments_': 'department'}, axis=1, inplace=True)
df.head(1)

```

Python

```

df = df[['satisfaction_level', 'last_evaluation', 'number_project',
         'average_monthly_hours', 'time_spend_company', 'work_accident',
         'promotion_last_5years', 'department', 'salary', 'left']]
df.head(1)

```

Python

```

print("Have a First Look to 'left' Column",'green')
first_look('left')

```

Python

```

fig = px.pie(df, values = df['left'].value_counts(),
              names = (df['left'].value_counts()).index,
              title = "'left' Column Distribution")
fig.show()

```

Python

```

y = df['left']
print(f'Percentage of left-1: % {round(y.value_counts(normalize=True)[1]*100,2)} --> \
({y.value_counts()[1]} observations for left-1)\nPercentage of left-0: % {round(y.value_counts(normalize=True)[0]*100,2)} --> ({y.value_counts()[0]} observations for left-0)')

```

Python

```

df.groupby('left').mean()

```

Python

```

cprint('Dataset describe results according to the "left==1" condition','green')
df[df['left'] == 1].describe().T.style.background_gradient(subset = ['mean', 'min', '50%', 'max'], cmap = 'RdPu')

```

Python

```

cprint('Dataset describe results according to the "left==0" condition','green')
df[df['left'] == 0].describe().T.style.background_gradient(subset = ['mean', 'min', '50%', 'max'], cmap = 'RdPu')

```

Python

```

cprint("Have a First Look to 'left' Column",'green')
first_look('satisfaction_level')

```

Python

```

df['satisfaction_level'].value_counts().iplot(kind="bar", title = "'satisfaction_level' Column Distribution")

```

Python

```

px.histogram(df, x = df['satisfaction_level'], color='left', marginal = "box", hover_data = df.columns,
             title = 'satisfaction level and left')

```

Python

```

pd.crosstab(df['satisfaction_level'], df['left']).iplot(kind="bar", title = 'satisfaction_level and left')

```

Python

```

[23]   sns.scatterplot(data=df, x="satisfaction_level", y="left", hue="left");
[24] df.corr(method='pearson')[7:]
[25] cprint("Have a First Look to 'last_evaluation' Column", 'green')
      first_look('last_evaluation')

[26] df['last_evaluation'].value_counts().iplot(kind="bar", title = "'last_evaluation' Column Distribution")
[27] px.histogram(df, x = df['last_evaluation'], color='left', marginal = "box", hover_data = df.columns,
[28] pd.crosstab(df['last_evaluation'], df['left']).iplot(kind="bar", title = 'last_evaluation and left')
[29] cprint("Have a First look to 'number_project' Column", 'green')
      first_look('number_project')

[30] df['number_project'].value_counts().iplot(kind="bar", title = "'number_project' Column Distribution")
[31] fig = px.pie(df, values = df['number_project'].value_counts(),
[32]                 names = (df['number_project'].value_counts()).index,
[33]                 title = "'number_project' Column Distribution")
      fig.show()

[34] px.histogram(df, x = df['number_project'], color='left', marginal = "box", hover_data = df.columns,
[35]                 title = 'number_project and left')
[36] pd.crosstab(df['number_project'], df['left']).iplot(kind='bar', title = 'number_project and left')

[37] cprint("Have a First Look to 'average_montly_hours' Column",'green')
      first_look('average_montly_hours')

[38] plt.figure(figsize=(15, 8))
fig = sns.histplot(
    df,
    x="average_montly_hours", hue="left",
    multiple="stack",
    palette="light:m_r",
    edgecolor=".3",
    linewidth=.5
)
fig2 = fig.twinx()

fig.yaxis.set_label_position('left')
fig2.yaxis.set_label_position('right')
fig2.set_ylabel('Frequency [%]')

for p in fig.patches:
    x=p.get_bbox().get_points()[0,0]
    y=p.get_bbox().get_points()[1,1]
    fig.annotate(('.1f%' .format(100.*y/len(df)), (x.mean(), y),
                 ha="center", va="bottom") # set the alignment of the text
    fig2.grid(None)
fig2.set_ylim(0,100)

[39] numerical= df.drop(['left'], axis=1).select_dtypes('number').columns
categorical = df.select_dtypes('object').columns

print(colored("Numerical Columns:", attrs=['bold']), list(df[numerical].columns),'\n',
      colored("-"*124, 'red', attrs=['bold']), sep='')
print(colored("Categorical Columns:", attrs=['bold']), list(df[categorical].columns),'\n',
      colored("-"*124, 'red', attrs=['bold']), sep='')

[40] df[categorical].head().T
[41] df[categorical].describe().T.style.background_gradient(subset=['unique','freq','count'], cmap='RdPu')

[42] df[categorical].nunique()
[43] for col in enumerate(df[categorical].columns):
    xtab = pd.crosstab(df[col], df['left'], normalize=True)
    print(colored('*'*40, 'red', attrs=['bold']), sep='')
    print(xtab*100)

[44] for i, col in enumerate(df[categorical].columns):
    xtab = pd.crosstab(df[col], df['left'], normalize=True)
    print(colored('*'*40, 'red', attrs=['bold']), sep='')
    print(xtab*100)

```

```

[42]     numerical= df.drop(['left'], axis = 1).select_dtypes('number').columns
    categorical = df.select_dtypes('object').columns
    print('-----')
    print("Numerical Columns: " + str(df[numerical].columns))
    print("Categorical Columns: " + str(df[categorical].columns))
    print('-----')

[43]
    df[numerical].describe().T.style.background_gradient(subset = ['mean','std','50%','count'], cmap = 'RdPu')

[44]
    df[numerical].iplot(kind = 'histogram', subplots = True, bins = 50)

[45]
    for i in numerical:
        df[i].iplot(kind = 'box', title = i, boxpoints = 'all')

[46]
    sns.pairplot(df, hue = 'left', corner = True);

[47]

[48]    df.corr()['left'].sort_values().drop('left').iplot(kind = 'barh', colors='purple');

[49]    df_temp = df.corr()
    count = 'Done'
    feature = []
    collinear= []
    for col in df_temp.columns:
        for i in df_temp.index:
            if (df_temp[col][i] > .9 and df_temp[col][i] < 1) or (df_temp[col][i] < -.9 and df_temp[col][i] > -1):
                feature.append(col)
                collinear.append(i)
                print(Fore.RED + f'\033[1mMulticollinearity alert in between\033[0m {col} - {i}')
            else:
                print(f'For {col} and {i}, there is NO multicollinearity problem')
    print(f'\033[1mThe number of strong correlated features:\033[0m {count}')

[50]

[51]    plt.figure(figsize = (7,5))
    sns.countplot(data = df, x = 'left');
    for index,value in enumerate(df.left.value_counts()):
        plt.text(index, value, f'{value}', ha = 'center', va = 'bottom', fontsize = 13);

[52]

[53]    print("left" Column Distribution,'green')
    fig = plt.figure(figsize = (11, 6))
    ax = fig.add_axes([0, 0, 1, 1])
    ax.bar(df.left.value_counts().index, df.left.value_counts().values, color = 'green')
    plt.title("left" Column Distribution)
    plt.xlabel("left")
    plt.ylabel("Number of Employees")
    for index,value in enumerate(df.left.value_counts()):
        plt.text(index, value, f'{value}', ha = 'center', va = 'bottom', fontsize = 13)
    plt.show()

[54]

[55]    cprint("number_project" Column Distribution,'green')
    df.number_project.value_counts()

[56]

[57]    df['number_project'].value_counts().iplot(kind="bar", title = '"number_project" Column Distribution')

[58]    plt.rcParams['figure.constrained_layout.use'] = False
    for index,value in enumerate(df.number_project.value_counts().sort_values(ascending=False)):
        plt.text(index, value, f'{value}', ha = 'center', va = 'bottom', fontsize = 13)
    plt.show()

[59]

[60]    cprint("time_spend_company" Column Distribution,'green')
    df.time_spend_company.value_counts()

[61]

[62]    df['time_spend_company'].value_counts().iplot(kind="bar", title = '"time_spend_company" Column Distribution')

[63]    plt.xlabel("Number of Employees")
    plt.xticks(rotation = 0)
    for index,value in enumerate(df.time_spend_company.value_counts().sort_values(ascending=False)):
        plt.text(index, value, f'{value}', ha = 'center', va = 'bottom', fontsize = 13)
    plt.show()

[64]

[65]    for i in df:
        df[i].iplot(kind = 'histogram', subplots = True, bins = 50)

[66]

[67]    df1 = df.drop('left', axis = 1)
    df1.head(1)

[68]

[69]    df1 = pd.get_dummies(df1, columns = ['department','salary'], drop_first = True)
    df1.head(1)

[70]

[71]    df1.head()

```

```

scaler = MinMaxScaler()
scaler.fit(df1)

df1_scaled= scaler.transform(df1)

```

[7] Python

```

df1_scaled

```

[8] Python

```

> v
df1_scaled

```

[96] Python

```

from pyclustertend import hopkins
hopkins(df1_scaled, df1.shape[0])

```

[97] Python

```

ks = range(1,10)
inertias=[]
for k in ks :
    kc = KMeans(n_clusters=k,random_state=1)
    kc.fit(df1_scaled)
    inertias.append(kc.inertia_)

f, ax = plt.subplots(figsize=(8, 6))
plt.plot(ks, inertias,'o')
plt.xlabel('Number of clusters, k')
plt.ylabel('Inertia')
plt.xticks(ks)
plt.style.use('ggplot')
plt.title('What is the Best Number for KMeans ?')
plt.show()

```

[100] Python

```

from yellowbrick.cluster import KElbowVisualizer

kmeans = KMeans()
visu = KElbowVisualizer(kmeans, k = (1,10))
visu.fit(df1_scaled)
visu.show();

```

[101] Python

```

from sklearn.metrics import silhouette_samples,silhouette_score
ssd =[]
K = range(2, 10)

for k in K:
    model = KMeans(n_clusters=k)
    model.fit(df1_scaled)
    ssd.append(model.inertia_)
    print("Silhouette Score for {} clusters: {}".format(k, silhouette_score(df1_scaled, model.labels_)))

```

[102] Python

```

from sklearn.cluster import KMeans
from yellowbrick.cluster import SilhouetteVisualizer

model_4 = KMeans(n_clusters=4, random_state=101)
visualizer = SilhouetteVisualizer(model_4)
visualizer.fit(df1_scaled)
visualizer.poof();

```

[103] Python

```

from sklearn.cluster import KMeans
from yellowbrick.cluster import SilhouetteVisualizer

model_3 = KMeans(n_clusters=3, random_state=101)
visualizer = SilhouetteVisualizer(model_3)
visualizer.fit(df1_scaled)
visualizer.poof();

```

[104] Python

```

k_means_model = KMeans(n_clusters = 3, random_state = 101)
k_means_model.fit(df1_scaled)
labels = k_means_model.labels_
labels

```

[105] Python

```

df['predicted_clusters'] = labels
df

```

[106] Python

```

df['predicted_clusters'].value_counts()

```

[107] Python

```

fig = px.pie(df, values = df['predicted_clusters'].value_counts(),
              names = (df['predicted_clusters'].value_counts().index,
                       title = 'Predicted_Clusters Distribution')
fig.show()

```

[108] Python

```

fig = px.pie(df, values = df[df['left']==0]['predicted_clusters'].value_counts(),
              names = df[df['left']==0]['predicted_clusters'].value_counts().index,
              title = 'Predicted_Clusters_Almost_Lost Distribution')
fig.show()

```

[109] Python

```

fig = px.pie(df, values = df[df['left']==1]['predicted_clusters'].value_counts(),
              names = df[df['left']==1]['predicted_clusters'].value_counts().index,
              title = 'Predicted_Clusters_Almost_Lost Distribution')
fig.show()

```

[110] Python

```

pd.crosstab(df['left'],
            df['predicted_clusters']).plot(kind="bar", title = 'Compare (left vs predicted-clusters)',
                                             xTitle = 'left & clusters', yTitle = 'counts')

```

[111] Python

```

df.groupby('left').mean()

```

[112] Python

```
[81] df.groupby(['left', 'predicted_clusters']).mean()
[82] pd.crosstab(df['predicted_clusters'],
[83]     df['left']).plot(kind='bar', title = 'Compare (predicted-clusters vs left)',
[84]     xTitle = 'clusters & left', yTitle = 'counts')
[85] df.left.value_counts()
[86] df.groupby('predicted_clusters').mean()
[87] df.groupby(['predicted_clusters', 'left']).mean()
[88] df2 = df.drop('predicted_clusters', axis = 1)
[89] df2.head(1)
[90] df2 = pd.get_dummies(df2, columns = ['department','salary'], drop_first = True)
[91] df2.head(1)
[92] X = df2.drop('left', axis = 1)
[93] y = df2['left']
[94] X_train, X_test, y_train, y_test = train_test_split(X, y, stratify = y, test_size = 0.3, random_state = 101)
[95] 
[96] scaler = MinMaxScaler()
[97] X_train = scaler.fit_transform(X_train)
[98] X_test = scaler.transform(X_test)
[99] 
[100] def eval(model, X_train, X_test):
[101]     y_pred = model.predict(X_test)
[102]     y_pred_train = model.predict(X_train)
[103] 
[104]     print("Confusion matrix(y_test, y_pred)")
[105]     print("Test Set")
[106]     print(classification_report(y_test,y_pred))
[107]     print("Train Set")
[108]     print(classification_report(y_train,y_pred_train))
[109]     cm = confusion_matrix(y_test,y_pred)
[110]     disp = ConfusionMatrixDisplay(confusion_matrix=cm)
[111]     disp.plot(model, X_test, y_test, cmap="plasma")
[112] 
[113] def train_val(y_train, y_train_pred, y_test, y_pred):
[114] 
[115]     scores = {"train_set": {"Accuracy" : accuracy_score(y_train, y_train_pred),
[116]                         "Precision" : precision_score(y_train, y_train_pred),
[117]                         "Recall" : recall_score(y_train, y_train_pred),
[118]                         "f1" : f1_score(y_train, y_train_pred)},
[119] 
[120]             "test_set": {"Accuracy" : accuracy_score(y_test, y_pred),
[121]                         "Precision" : precision_score(y_test, y_pred),
[122]                         "Recall" : recall_score(y_test, y_pred),
[123]                         "f1" : f1_score(y_test, y_pred)}}
[124] 
[125]     return pd.DataFrame(scores)
[126] 
[127] 
[128] from sklearn.metrics import RocCurveDisplay, roc_auc_score, roc_curve, f1_score, accuracy_score, recall_score
[129] 
[130] 
[131] GB_model = GradientBoostingClassifier(random_state = 101)
[132] GB_model.fit(X_train, y_train)
[133] y_pred = GB_model.predict(X_test)
[134] y_train_pred = GB_model.predict(X_train)
[135] 
[136] GB_model_f1 = f1_score(y_test, y_pred)
[137] GB_model_acc = accuracy_score(y_test, y_pred)
[138] GB_model_recall = recall_score(y_test, y_pred)
[139] GB_model_auc = roc_auc_score(y_test, y_pred)
[140] ...
[141] 
[142] 
[143] from sklearn.metrics import classification_report
[144] from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
[145] 
[146] 
[147] print("GB Model")
[148] print("-----")
[149] eval(GB_model, X_train, X_test)
[150] 
[151] 
[152] cf_matrix = confusion_matrix(y_test, y_pred)
[153] group_names = ["True Negatives (TN)", "False Positives (FP)\n(Type I Error)", "False Negatives (FN)\n(Type II Error)", "True Positives (TP)"]
[154] group_counts = ["%(0:0.0f)".format(value) for value in
[155]                 cf_matrix.flatten()]
[156] group_percentages = ["%(0:.2%)".format(value) for value in
[157]                      cf_matrix.flatten()/np.sum(cf_matrix)]
[158] labels = ["%(V1)\n(V2)\n(V3)" % for V1, V2, V3 in
[159]           zip(group_names, group_counts, group_percentages)]
[160] labels = np.asarray(labels).reshape(2, 2)
[161] 
[162] ax = sns.heatmap(cf_matrix, annot=labels, fmt="", cmap="Blues")
[163] ax.set(xlabel="Predicted Class", ylabel = "Actual Class")
```

```

[119]
plt.figure(figsize=(12, 8))
visualizer = ClassificationReport(GB_model, support=True)
visualizer.fit(X_train, y_train)
visualizer.score(X_test, y_test)
visualizer.show()
[119] Python

[120]
>>> from yellowbrick.classifier import ClassPredictionError
visualizer = ClassPredictionError(GB_model)
visualizer.fit(X_train, y_train)
visualizer.score(X_test, y_test)
visualizer.poof();
[120] Python

[121]
GB_feature_imp = pd.DataFrame(index=X.columns, data = GB_model.feature_importances_, columns = ['Importance']).sort_values('Importance', ascending = False)
GB_feature_imp
[121] Python

[122]
plt.figure(figsize = (12,8))
sns.barplot(data = GB_feature_imp.sort_values('Importance', ascending = False), x = GB_feature_imp.sort_values('Importance', ascending = False).index, y = 'Importance')
plt.title('Feature Importance')
[122] Python

[123]
GB_cv = GradientBoostingClassifier(random_state = 101)

GB_cv_scores = cross_validate(GB_cv, X_train, y_train,
                             scoring = ('accuracy', 'precision','recall', 'f1', 'roc_auc'), cv = 10)
GB_cv_scores = pd.DataFrame(GB_cv_scores, index = range(1, 11))

GB_cv_scores.mean()[2:]
[123] Python

[124]
param_grid = {'n_estimators':[100, 200, 300],
              'subsample':[0.5, 1],
              'learning_rate': [0.01, 0.05, 0.1],
              'max_depth':[3, 4, 5, 6]}
[124] Python

[125]
GB_grid = GradientBoostingClassifier(random_state = 101)
GB_grid_model = GridSearchCV(GB_grid, param_grid, scoring = "f1", verbose = 0, n_jobs = -1).fit(X_train, y_train)
[125] Python

...
[126]
GB_grid_model.best_estimator_
[126] Python

[127]
GB_tuned = GradientBoostingClassifier(learning_rate = 0.01,
                                      max_depth = 6,
                                      n_estimators = 200,
                                      subsample = 0.5,
                                      random_state = 101).fit(X_train, y_train)
[127] Python

...
[128]
y_pred = GB_tuned.predict(X_test)
y_train_pred = GB_tuned.predict(X_train)

GB_tuned_f1 = f1_score(y_test, y_pred)
GB_tuned_acc = accuracy_score(y_test, y_pred)
GB_tuned_recall = recall_score(y_test, y_pred)
GB_tuned_auc = roc_auc_score(y_test, y_pred)
[128] Python

[129]
y_pred = GB_tuned.predict(X_test)
y_train_pred = GB_tuned.predict(X_train)

GB_tuned_f1 = f1_score(y_test, y_pred)
GB_tuned_acc = accuracy_score(y_test, y_pred)
GB_tuned_recall = recall_score(y_test, y_pred)
GB_tuned_auc = roc_auc_score(y_test, y_pred)
[129] Python

...
[130]
print("GB_tuned")
print ("-----")
eval(GB_tuned, X_train, X_test)
[130] Python

[131]
cf_matrix = confusion_matrix(y_test, y_pred)

group_names = ["True Negatives (TN)", "False Positives (FP)\n(Type I Error)", "False Negatives (FN)\n(Type II Error)", "True Positives (TP)"]
group_counts = ["{:0.0f}".format(value) for value in cf_matrix.flatten()]
group_percentages = ["{:0.0%}".format(value) for value in cf_matrix.flatten()/np.sum(cf_matrix)]
labels = [f'{v1}\n{v2}\n{v3}' for v1, v2, v3 in zip(group_names, group_counts, group_percentages)]
labels = np.asarray(labels).reshape(2, 2)

ax = sns.heatmap(cf_matrix, annot=labels, fmt="", cmap="Blues")
ax.set(xlabel="Predicted Class", ylabel = "Actual Class")
[131] Python

[132]
plt.figure(figsize=(12, 8))
visualizer = ClassificationReport(GB_tuned, support=True)
visualizer.fit(X_train, y_train)
visualizer.score(X_test, y_test)
visualizer.show()
[132] Python

[133]
from sklearn.metrics import mean_squared_error
from sklearn.metrics import roc_curve
from sklearn.metrics import RocCurveDisplay
[133] Python

[134]
from yellowbrick.classifier import ClassPredictionError
visualizer = ClassPredictionError(GB_tuned)

visualizer.fit(X_train, y_train)

visualizer.score(X_test, y_test)

visualizer.poof();
[134] Python

[135]
RocCurveDisplay.from_estimator(GB_model, X_test, y_test);
[135] Python

```

```

[285]: GB_Pred = ("Actual": y_test, "GB_Pred":y_pred)
GB_Pred = pd.DataFrame.from_dict(GB_Pred)
GB_Pred.head(20)

[286]: Python

[287]: GB_Pred = ("Actual": y_test, "GB_Pred":y_pred)
GB_Pred = pd.DataFrame.from_dict(GB_Pred)
GB_Pred.head(20)

[288]: Python

[289]: Model_Preds = GB_Pred
Model_Preds.head(1)

[290]: Python

[291]: KNN_model = KNeighborsClassifier(n_neighbors = 5)
KNN_model.fit(X_train, y_train)
y_pred = KNN_model.predict(X_test)
y_train_pred = KNN_model.predict(X_train)

KNN_model_f1 = f1_score(y_test, y_pred)
KNN_model_acc = accuracy_score(y_test, y_pred)
KNN_model_recall = recall_score(y_test, y_pred)
KNN_model_auc = roc_auc_score(y_test, y_pred)

[292]: Python

[293]: Model_Preds = GB_Pred
Model_Preds.head(1)

[294]: Python

[295]: KNN_model = KNeighborsClassifier(n_neighbors = 5)
KNN_model.fit(X_train, y_train)
y_pred = KNN_model.predict(X_test)
y_train_pred = KNN_model.predict(X_train)

KNN_model_f1 = f1_score(y_test, y_pred)
KNN_model_acc = accuracy_score(y_test, y_pred)
KNN_model_recall = recall_score(y_test, y_pred)
KNN_model_auc = roc_auc_score(y_test, y_pred)

[296]: Python

[297]: print("KNN_Model")
print("-----")
eval(KNN_model, X_train, X_test)

[298]: Python

[299]: cf_matrix = confusion_matrix(y_test, y_pred)

group_names = ["True Negatives (TN)", "False Positives (FP)\n(Type I Error)", "False Negatives (FN)\n(Type II Error)", "True Positives (TP)"]
group_counts = [(0.0,0.0)]*4
for value in cf_matrix:
    group_counts += [f'{(0.0,0.0)}'.format(value) for value in value]
    group_percentages = [(0.0,0.0)]*4
    for value in cf_matrix:
        group_percentages += [(0.0,0.0)].format(value) for value in value
        cf_matrix = cf_matrix/np.sum(cf_matrix)

labels = [(f'({v1},{v2},{v3})') for v1, v2, v3 in zip(group_names, group_counts, group_percentages)]
labels = np.asarray(labels).reshape(2, 2)

ax = sns.heatmap(cf_matrix, annot=labels, fmt="", cmap="Blues")
ax.set(xlabel="Predicted Class", ylabel = "Actual Class");

[300]: Python

[301]: plt.figure(figsize=(12, 8))
visualizer = ClassificationReport(KNN_model, support=True)
visualizer.fit(X_train, y_train)
visualizer.score(X_test, y_test)

[302]: Python

[303]: from yellowbrick.classifier import ClassPredictionError
visualizer = ClassPredictionError(KNN_model)

visualizer.fit(X_train, y_train)

[304]: Python

[305]: visualizer.poof();

[306]: Python

[307]: KNN_cv = KNeighborsClassifier(n_neighbors = 5)
KNN_cv_scores = cross_validate(KNN_cv, X_train, y_train,
scoring = ['accuracy', 'precision','recall', 'f1', 'roc_auc'], cv = 10)
KNN_cv_scores = pd.DataFrame(KNN_cv_scores, index = range(1, 11))

KNN_cv_scores.mean()[:2]

[308]: Python

[309]: plt.figure(figsize = (15, 8))
plt.plot(range(1, 30), test_error_rates, color = 'blue', linestyle = '--', marker = 'o',
markerfacecolor = 'red', markersize = 10)
plt.title('Test Error Rates vs. K Value')
plt.xlabel('Value')
plt.ylabel('Error Rate')

plt.hlines(y = 0.045974456078039, xmin = 0, xmax = 30, colors = 'r', linestyles = '--', label = 'K-Value = 2 Line')
plt.hlines(y = 0.0453008720579312, xmin = 0, xmax = 30, colors = 'r', linestyles = '--', label = 'K-Value = 4 Line')
plt.hlines(y = 0.0558663600050404, xmin = 0, xmax = 30, colors = 'blue', linestyles = '--', label = 'Default K-Value = 5 Line')
plt.legend(prop = {'size':15})

[310]: Python

[311]: Knn5 = KNeighborsClassifier(n_neighbors = 5)
knn5.fit(X_train,y_train)
pred = knn5.predict(X_test)

print('WITH K=5')
print("-----")
print(confusion_matrix(y_test, pred))
print("-----")
print(classification_report(y_test, pred))

[312]: Python

[313]: knn2 = KNeighborsClassifier(n_neighbors = 2)
knn2.fit(X_train,y_train)
pred = knn2.predict(X_test)

print('WITH K=2')
print("-----")
print(confusion_matrix(y_test, pred))
print("-----")
print(classification_report(y_test, pred))

[314]: Python

[315]: k_values = range(1, 30)
param_grid = {"n_neighbors": k_values, "p": [1, 2], "weights": ['uniform', 'distance']}

[316]: Python

[317]: KNN_grid = KNeighborsClassifier()
KNN_grid_model = GridSearchCV(KNN_grid, param_grid, cv = 10, scoring = 'recall')
KNN_grid_model.fit(X_train, y_train)

[318]: Python

```

```

[224] print(colored("\033[1mBest Parameters of GridSearchCV for KNN Model:\033[0m", 'blue'), colored(KNN_grid_model.best_params_, 'cyan'))
print("-----")
print(colored("\033[1mBest Estimator of GridSearchCV for KNN Model:\033[0m", 'blue'), colored(KNN_grid_model.best_estimator_, 'cyan'))

[225]

KNN_tuned3 = KNeighborsClassifier(n_neighbors = 3, p = 1)
KNN_tuned3.fit(X_train, y_train)
y_pred = KNN_tuned3.predict(X_test)
y_train_pred = KNN_tuned3.predict(X_train)

KNN_tuned3_f1 = f1_score(y_test, y_pred)
KNN_tuned3_acc = accuracy_score(y_test, y_pred)
KNN_tuned3_recall = recall_score(y_test, y_pred)
KNN_tuned3_auc = roc_auc_score(y_test, y_pred)

print("KNN_tuned (k=3)")
print("-----")
eval(KNN_tuned3, X_train, X_test)
train_val(y_train, y_train_pred, y_test, y_pred)

[226] Python

[227]

KNN_tuned1 = KNeighborsClassifier(n_neighbors = 1, p = 1)
KNN_tuned1.fit(X_train, y_train)
y_pred = KNN_tuned1.predict(X_test)
y_train_pred = KNN_tuned1.predict(X_train)

KNN_tuned1_f1 = f1_score(y_test, y_pred)
KNN_tuned1_acc = accuracy_score(y_test, y_pred)
KNN_tuned1_recall = recall_score(y_test, y_pred)
KNN_tuned1_auc = roc_auc_score(y_test, y_pred)

print("KNN_tuned (k=1)")
print("-----")
eval(KNN_tuned1, X_train, X_test)

[228] RocCurveDisplay.from_estimator(KNN_model, X_test, y_test);

[229] ...

```

```

KNN_Pred = ("Actual": y_test, "KNN_Pred":y_pred)
KNN_Pred = pd.DataFrame.from_dict(KNN_Pred)
KNN_Pred.head(20)

```

```

[230] Python

D: ~
KNN_Pred.drop("Actual", axis = 1, inplace = True)
Model_Preds = pd.merge(Model_Preds, KNN_Pred, left_index = True, right_index = True)
Model_Preds.head(20)

[231] Python

[232]

RF_model = RandomForestClassifier(class_weight = "balanced", random_state = 101)
RF_model.fit(X_train, y_train)
y_pred = RF_model.predict(X_test)
y_train_pred = RF_model.predict(X_train)

RF_model_f1 = f1_score(y_test, y_pred)
RF_model_acc = accuracy_score(y_test, y_pred)
RF_model_recall = recall_score(y_test, y_pred)
RF_model_auc = roc_auc_score(y_test, y_pred)

[233] Python

[234]
print("RF Model")
print("-----")
eval(RF_model, X_train, X_test)

[235] Python

[236]

cf_matrix = confusion_matrix(y_test, y_pred)

group_names = ["True Negatives (TN)", "False Positives (FP)\n(Type I Error)", "False Negatives (FN)\n(Type II Error)", "True Positives (TP)"]
group_counts = ["(0:0.01)".format(value) for value in cf_matrix.flatten()]
group_percentages = ["{0:(0.001)}".format(value) for value in cf_matrix.flatten() / np.sum(cf_matrix)]
labels = [(v1)\|(v2)\|(v3)" for v1, v2, v3 in zip(group_names, group_counts, group_percentages)]
labels = np.asarray(labels).reshape(2, 2)

ax = sns.heatmap(cf_matrix, annot=labels, fmt="", cmap="Blues")
ax.set(xlabel="Predicted Class", ylabel = "Actual Class");

[237] Python

[238]

plt.figure(figsize=(12, 8))
visualizer = ClassificationReport(RF_model, support=True)
visualizer.fit(X_train, y_train)
visualizer.score(X_test, y_test)
visualizer.show()

[239] from yellowbrick.classifier import ClassPredictionError
visualizer = ClassPredictionError(RF_model)

visualizer.fit(X_train, y_train)

[240]

visualizer.score(X_test, y_test)
visualizer.poof();

[241] Python

[242]

RF_feature_imp = pd.DataFrame(index=X.columns, data = RF_model.feature_importances_, columns = ["Importance"]).sort_values("Importance", ascending = False)
RF_feature_imp

[243] Python

[244]

plt.figure(figsize = (12,6))
sns.barplot(data = RF_feature_imp.sort_values('Importance', ascending = False), x = RF_feature_imp.sort_values('Importance', ascending = False).index, y = 'Importance')
plt.xticks(rotation = 75);

[245] Python

[246]

RF_cv = RandomForestClassifier(class_weight = "balanced", random_state = 101)

RF_cv_scores = cross_validate(RF_cv, X_train, y_train,
                             scoring = ['accuracy', 'precision','recall', 'f1', 'roc_auc'], cv = 10)
RF_cv_scores = pd.DataFrame(RF_cv_scores, index = range(1, 11))

RF_cv_scores.mean()[2:1]

[247] Python

[248]

RF_cv = RandomForestClassifier(class_weight = "balanced", random_state = 101)

RF_cv_scores = cross_validate(RF_cv, X_train, y_train,
                             scoring = ['accuracy', 'precision','recall', 'f1', 'roc_auc'], cv = 10)
RF_cv_scores = pd.DataFrame(RF_cv_scores, index = range(1, 11))

RF_cv_scores.mean()[2:1]

[249] Python

```

```

param_grid = {"n_estimators": [50, 100, 300],
              'max_features': [2, 3, 4],
              'max_depth': [3, 5, 7, 9],
              'min_samples_split': [2, 5, 8]}

[240] RF_grid = RandomForestClassifier(class_weight = 'balanced', random_state = 101)
RF_grid.grid = GridSearchCV( estimator = RF_grid,
                            param_grid = param_grid,
                            scoring = "recall",
                            n_jobs = -1, verbose = 0)
RF_grid.grid.fit(X_train, y_train)

[241] RF_grid.grid.best_estimator_
[242]

[243] RF_tuned = RandomForestClassifier(class_weight = 'balanced',
                                      max_depth = 3,
                                      min_features = 4,
                                      n_estimators = 300,
                                      random_state = 101).fit(X_train, y_train)

[244] y_pred = RF_tuned.predict(X_test)
y_train_pred = RF_tuned.predict(X_train)

RF_tuned.F1 = f1_score(y_test, y_pred)
RF_tuned.acc = accuracy_score(y_test, y_pred)
RF_tuned.recall = recall_score(y_test, y_pred)
RF_tuned.auc = roc_auc_score(y_test, y_pred)

[245] print("RF_tuned")
print ("-----")
eval(RF_tuned, X_train, X_test)

[246] cf_matrix = confusion_matrix(y_test, y_pred)

group_names = ["True Negatives (TN)", "False Positives (FP)\n(Type I Error)", "False Negatives (FN)\n(Type II Error)", "True Positives (TP)"]

group_counts = ["{:0.0f}".format(value) for value in
cf_matrix.flatten()]

group_percentages = ["{:.0%}".format(value) for value in
cf_matrix.flatten()/np.sum(cf_matrix)]

labels = [(v1,v2,v3) for v1, v2, v3 in
zip(group_names, group_counts, group_percentages)]
labels = np.asarray(labels).reshape(2, 2)

ax = sns.heatmap(cf_matrix, annot=labels, fmt="", cmap="Blues")
ax.set(xlabel="Predicted Class", ylabel = "Actual Class");

[247] plt.figure(figsize=(12, 8))
visualizer = ClassificationReport(RF_tuned, support=True)
visualizer.fit(X_train, y_train)
visualizer.score(X_test, y_test)
visualizer.poof()

[248] from yellowbrick.classifier import ClassPredictionError
visualizer = ClassPredictionError(RF_tuned)

visualizer.fit(X_train, y_train)

[249] from yellowbrick.classifier import ClassPredictionError
visualizer = ClassPredictionError(RF_tuned)

visualizer.fit(X_train, y_train)

visualizer.score(X_test, y_test)

visualizer.poof();

[250] RocCurveDisplay.from_estimator(RF_model, X_test, y_test);

[251]

RF_Pred = ("Actual": y_test, "RF_Pred":y_pred)
RF_Pred = pd.DataFrame.from_dict(RF_Pred)
RF_Pred.head(10)

[252] Model_Preds.sample(10)

[253]

compare = pd.DataFrame([{"Model": ["GB_model", "GB_tuned", "KNN_Model", "KNN_tuned3", "RF_model", "RF_tuned"],
                         "F1_Score": [GB_model.F1, GB_tuned.F1, KNN_model.F1, KNN_tuned3.F1, RF_model.F1, RF_tuned.F1],
                         "Accuracy_Score": [GB_model.acc, GB_tuned.acc, KNN_model.acc, KNN_tuned3.acc, RF_model.acc, RF_tuned.acc],
                         "Recall_Score": [GB_model.recall, GB_tuned.recall, KNN_model.recall, KNN_tuned3.recall, RF_model.recall, RF_tuned.recall],
                         "ROC_AUC_Score": [GB_model.auc, GB_tuned.auc, KNN_model.auc, KNN_tuned3.auc, RF_model.auc, RF_tuned.auc]}])

compare = compare.sort_values(by="Recall_Score", ascending=True)
fig = px.bar(compare, x = "Recall_Score", y = "Model", title = "Recall_Score")
fig.show()

compare = compare.sort_values(by="F1_Score", ascending=True)
fig = px.bar(compare, x = "F1_Score", y = "Model", title = "F1_Score")
fig.show()

compare = compare.sort_values(by="Accuracy_Score", ascending=True)
fig = px.bar(compare, x = "Accuracy_Score", y = "Model", title = "Accuracy_Score")
fig.show()

compare = compare.sort_values(by="ROC_AUC_Score", ascending=True)
fig = px.bar(compare, x = "ROC_AUC_Score", y = "Model", title = "ROC_AUC_Score")
fig.show()

```

2. Testing

The goal of system testing, which consists of a variety of tests, is to fully indoctrinate the computer-based system. Even though each test has a specific goal, they are all created to check that every system component has been properly integrated and is carrying out its designated functions. To make sure the product performs exactly as intended, testing procedures are frequently used. The following goals are to be attained during the trial phase:

- To confirm the project's level of excellence.
- To locate and correct any lingering errors from previous steps.
- To guarantee the system's dependability in functioning
- To verify that the program is a viable fix for the initial problem.

RESULTS



CONCLUSION

An employee who works for a firm is a significant resource. If one of them leaves the firm unexpectedly, it could have a big effect and cost the particular company a lot of money. Additionally, the newly employed employee needs some time to make the particular company cost-effective in addition to the time and money needed. Based on analytical data that is currently accessible, this model will help estimate the pace at which people leave their jobs and will use a variety of machine learning algorithms to reduce prediction error. So basing on the accuracy of the prediction rate the HR of an organization can decide the peoples intake and provide good valuable employees to the company or organization.

REFERENCES

1. Base Paper
<https://www.igi-global.com/article/intelligent-employee-retention-system-for-attrition-rate-analysis-and-churn-prediction/273145>
2. Churn Prediction
<https://github.com/Tarun-D/Bank-Customer-Churn-Prediction>
- 3.https://github.com/izkabrhm/Human-Resources-Employee-Churn-Prediction/blob/main/Employee_Churn_Prediction_Analysis.ipynb
- 4.<https://github.com/zunicd/Bank-Churn-Prediction/blob/master/Bank%20Customer%20Churn%20Prediction%20-%20Summary.ipynb>
- 5.https://github.com/wenkangwei/Bank-Churn-Prediction/blob/main/Bank_Customer_Churn_Prediction.ipynb

