# CompSci 671 In-Class Kaggle Competition Report

Yutong Shao

UID: Claireyt

December 9, 2022
Kaggle Competition Link

# 1 Exploratory Data Analysis

## 1.1 Data Description

This training data set contains 34 features describing a worker's attributes, and one label concerning attrition status of this worker. We were to train models on training set and predict the attrition status of test set.

To check these statistical properties, I first drew a table to check whether the data needs preprocessing, and observe each feature's statistical properties. It turns out that this data set is clean and tidy - no missing values, data type of each columns is consistent and well-defined. So, it saves me from many data-cleaning work - I only need to encode categorical features to numerical variables, then do feature engineering, which will be mentioned later.

## 1.2 Data visualization

However, the methods above are not enough to get the most out of the data. I still need to do some visualization to check the distributions of each feature in case there are some 'redundant' feature that does not contribute much, or contains outliers that need to be eliminated. Therefore, for categorical variable, I drew bar plots. For numerical features, I drew histograms and box plots, because histograms are good at showing distributions, and box plots add things by showing 0.5 and 0.75 percentile, meanwhile identify outliers. Diagrams are show below. (only include the graphs of one feature for brevity)
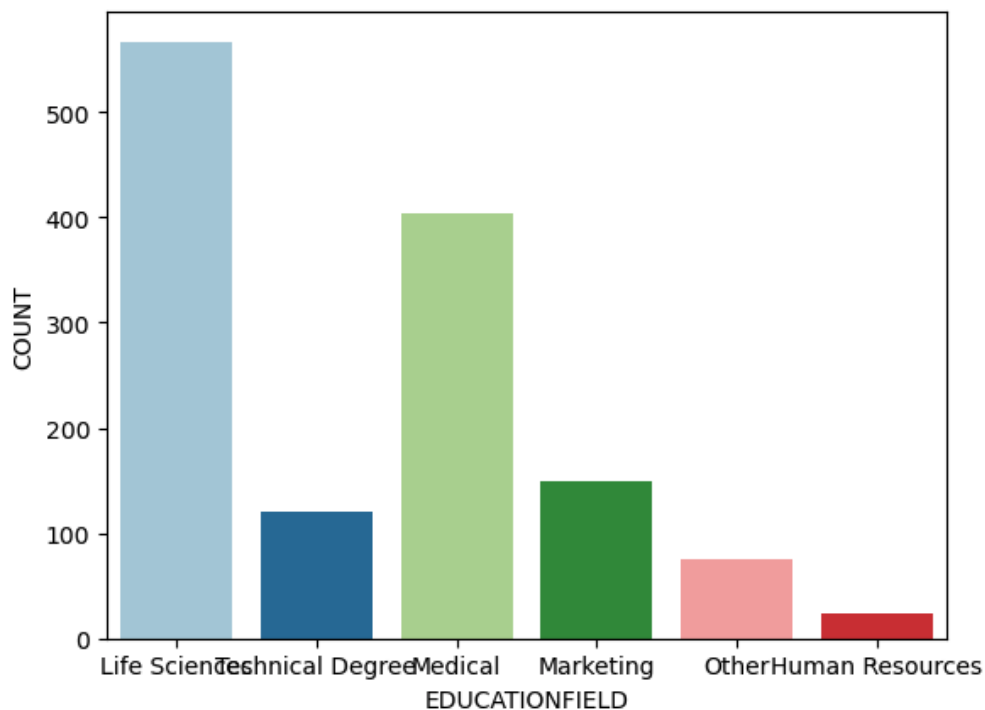
Figure 1: Bar plot for EducationField

The bar plot of "Attrition" shows there are 88% of people did not attrit, only around 11% of people attritted. So, we can see the ratio of attrition is not high, and this can serve as a criterion for checking whether our prediction is reasonable. Other bar plots are quite usual except for one - "Over18", which has only one bar on the plot. It means all of the people in the training set is over 18, and therefore this feature does not contribute much to our model. So, we can safely drop this column before model training. Another thing worth noticing is that we also need to drop "EmployeeID", which also has no actual meaning to our algorithms.
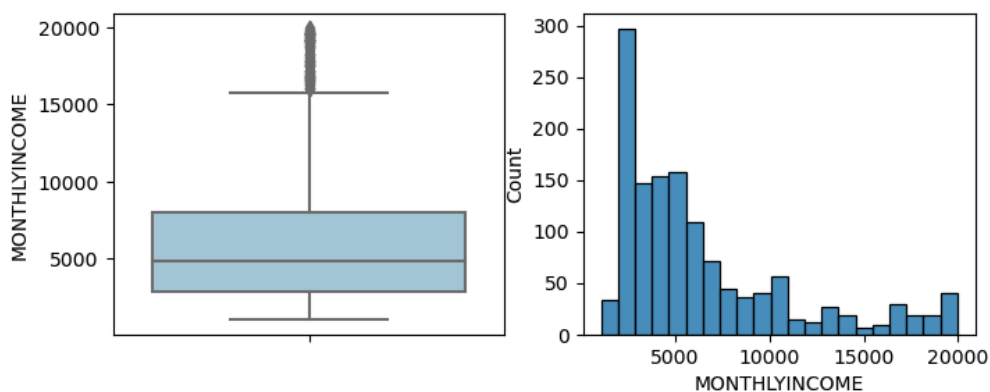


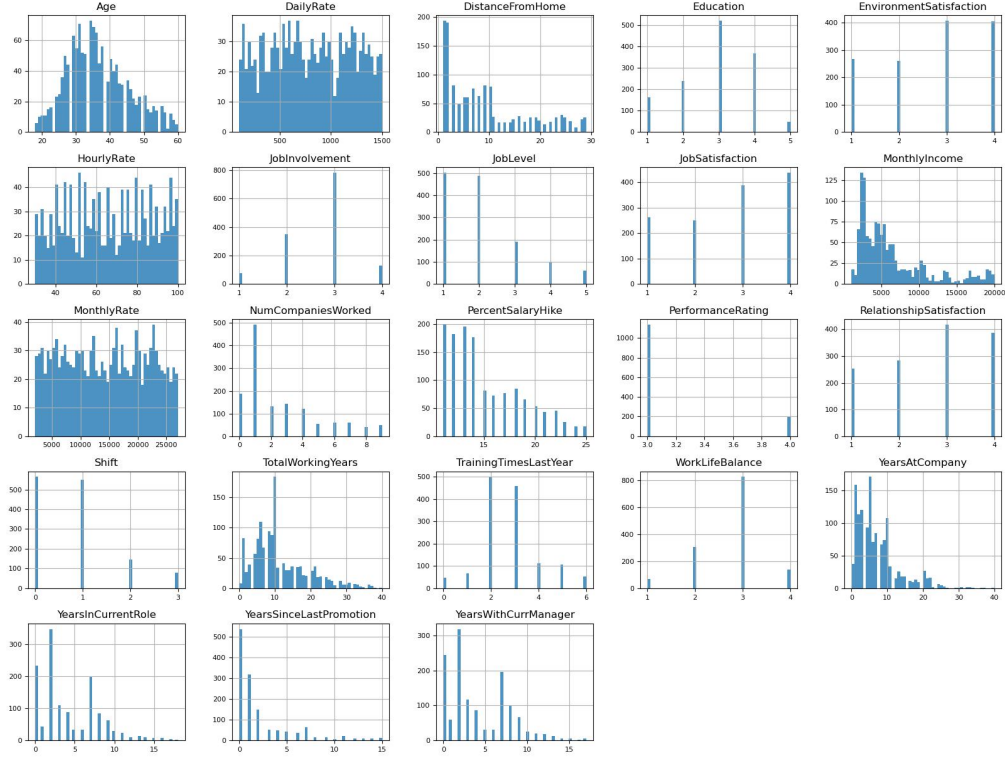Figure 2: Box plot and histogram for MonthlyIncome

Figure 3: Histogram

For numerical data, I draw both box plots and histograms. After browsing all the plots we can see there are a few abnormal plots. One is "EmployeeCount", because every data entry counts one, so we can safely drop this column. And "StandardHours" of all individuals are the same. So, the features we need to drop are "EmployeeID","StandardHours", "EmployeeCount","Over18".

Moreover, the box plot also show another valuable information - some features have outliers that may have bad impact for our model. Since outliers have many sources, including error input, extreme values, they will influence the overall distribution of the data set and cause inaccurate model estimation. Thus, we need to remove outliers of these features.

# 2 Feature Engineering

In this part, I need to do some feature engineering to select the most useful features. Because data and features define the upper limit of machine learning models, while we are just trying to converge to this upper limit using models.

## 2.1 Correlation Analysis

Correlation analysis is an important step of feature engineering. To ensure each feature is valuable and contribute to models, we need to remove the highly-correlated ones. And the most straightforward approach is using correlation matrix, setting a threshold and delete features that has correlation above that threshold. I draw a heat map to show the correlations, and the closer the color is to blue, the higher the correlations of the two features are.
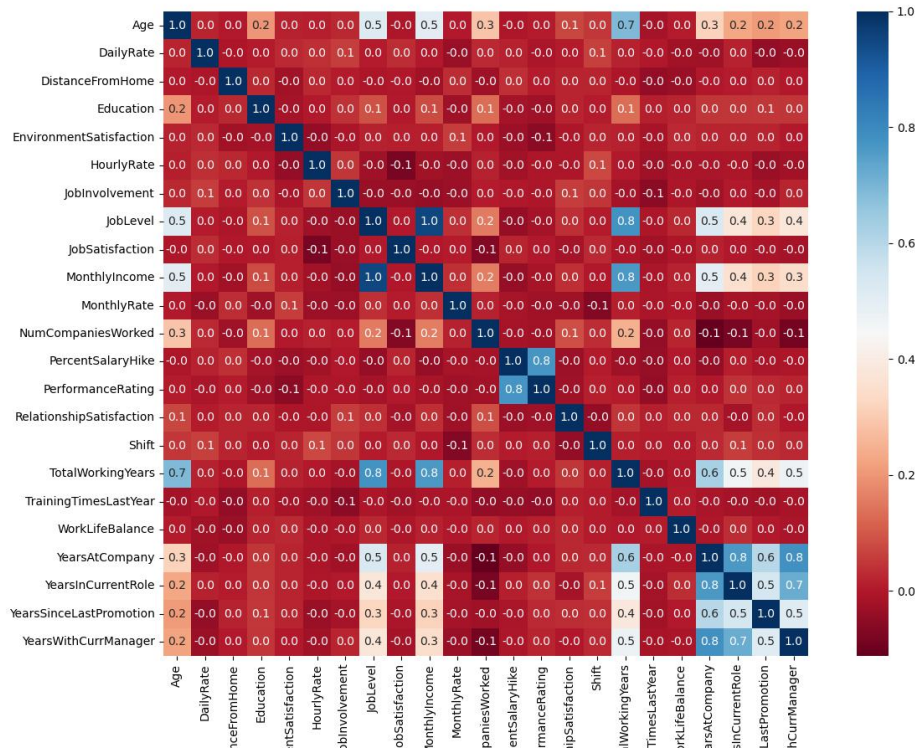


Figure 4: Correlation heat map

We can see a few squares have very blueish color with correlation over 0.6. Because I want to maintain as much information as possible while getting rid of redundant features,

I set the threshold to be 0.7 - any features has a correlation above 0.7 with another one will be removed. Finally, we have a list of features to drop: ["PerformanceRating", "TotalWorkingYears", "YearsInCurrentRole", "YearsWithCurrManager"] Dropping these are also reasonable from intuition. Because the total working years should be highly correlated with the years spent with current role or manager, and employee's monthly income could also be hooked with his or her performance.

## 2.2   Encoding Categorical data

After removing some highly correlated features, the next step I need to do is converting categorical data to numerical data. And I chose `Sklearn.LabelEncoder` to encode these features.

## 2.3   Outlier Detection

As I mentioned before, we need to remove outliers from numerical data. Therefore, I wrote a function that can detect outliers and return there indices. I chose the quantile method - first compute interquantile range (IQR), i.e. the difference between 0.75 and 0.25 quantile, then set the upper bounds to be 0.75 quantile + 1.5 * IQR, and lower bounds are similarly defined, so that any data lies outside of the interval - [lower bound, upper bound] - will be considered as outliers.

## 2.4   Feature Selection

I use `Sklearn.SelectKbest` to select the best features and set k to be 20. This method is basically computing the dependence of each feature to the label and retain those that have higher dependence.

## 2.5   Scaling Numerical Data

Since we have over 20 features and each numerical feature has different meaning and unit. For example, the 'MonthlyIncome' should be in dollars, while 'DistanceFromHome' should be in miles. The primary reason for scaling these data is that we want the algorithm to treat the *ranking* or *relative magnitude* of each entry seriously instead of the *absolute values*. Because in the latter case, the algorithm would think a monthly income of 16183 is more

important than 10 total working years. Fortunately, rescale features will largely solve this problem.

Therefore, I use `Sklearn.Preprocessing.StandardScaler` to transform data to be in the same range for better training and prediction.

# 3  Models

I chose tree-based model because they are non-parameterized and are easier to visualize, I used `AdaBoost` and `RandomForest`.

## 3.1  Model 1

Reasons why I chose `AdaBoost` are as follows.

1. robustness - it is relatively robust to overfitting in low noise datasets, because it ensembles many weak tree classifers that underfit (high bias) and combine those predictions into a stronger model, and reduce the overfitting (high variance);

2. ease of hyperparameter tuning - it has only a few hyperparameters that need to be tuned;

## 3.2  Model 2

Random forests are bagged decision tree models that split on randomly selected subset of features on each split. Reasons why I chose `RandomForest` are as follows.

1. it can handle all kinds of data, including numerical, categorical, binary features, little preprocessing need to be done before model training;

2. we also do not need to rescale data if using Random Forest;

3. efficiency of training - it is faster than other boosted decision trees because we are only training on subsets of features, which greatly improves training speed and efficiency.

# 4  Training

## 4.1  Data splits

I used `Sklearn.model_selection.train_test_split` to split training set and validation set - 80% are train set and 20% are validation set. When doing cross validation, I chose k to be 5.

## 4.2  Algorithms

### 4.2.1  AdaBoost

1. Algorithm

   (a) Assign each data point the same weight;

   (b) Train weak classifiers, and identify the classification result;

   (c) Update the weights of each data points. increase the weights of wrongly classified data and decrease the weights of those that are correctly classified;

   (d) Iterate the above steps for T times;

   (e) Finally, output the final classifier, which is a strong classifer.

2. Run time

   Adaboost does not have very apparent computation efficiency improvement, so the run time is slower than random forest, around 0.3 seconds.

### 4.2.2  Random Forest

1. Algorithm

   (a) Bootstrapping a sample of size $n$ from the training set.

   (b) Draw a single decision tree follow the steps below:

   - Randomly choose $m$ features;
   - Evaluate the splitting criteria on these features, and split on the best feature;
   - Repeat until the observations of this node reach a certain minimum threshold;

   (c) Iterate the above steps for T times.

   (d) Output all trees, then use the majority-vote rule to decide the best tree and use this tree to make predictions.

2. Run time

Since random forest only trains on a subset of features every time, I assume the training time won't be longer than that of Adaboost. Around 0.2 seconds shall be reasonable.

# 5   Hyperparameter Tuning

In each of the tuning process, I used `sklearn.model_selection.GridSearchCV`.

## 5.1   Hyperparameter Tuning for AdaBoost

The tuning parameters of Adaboost is quite lightweighted - there are only two hyperparameters need to be tuned: `n_estimators` and `learning_rate`. I set the parameter grid to be

```
1  params = {'n_estimators': [10, 25, 50, 75, 100, 200, 300, 500],
2  'learning_rate': [0.001, 0.01, 0.1, 0.5, 1., 5, 10, 20, 50, 100]}
```

The first parameter is `n_estimators`, for which I have tried a wider range of numbers, from 10 to 500. The second parameter is `learning_rate`, and I have tried 10 different values ranging from 0.01 to 100. And the functional relationship between predictive accuracy and learning rate is
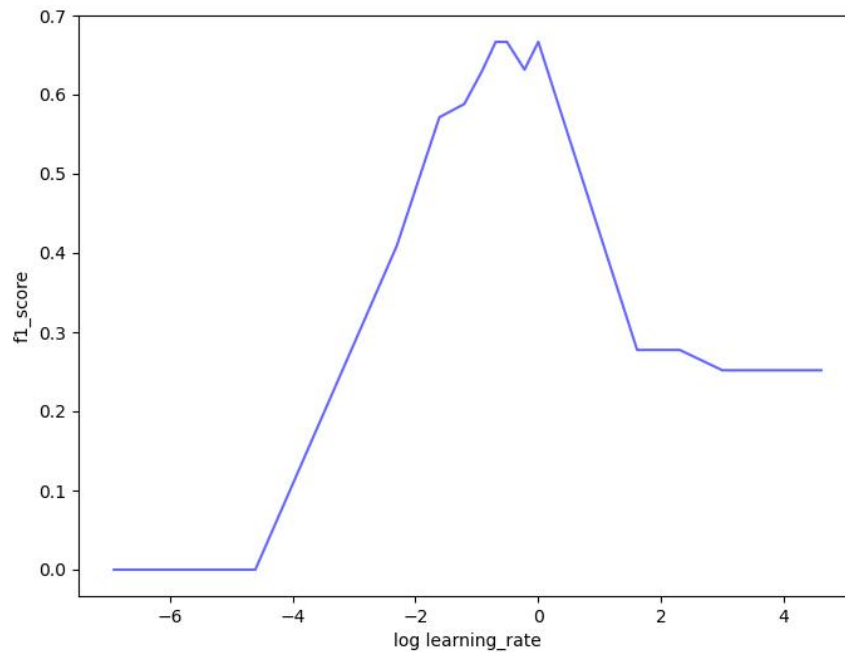
Figure 5: Learning rate VS f1 score

We can see that when `learning_rate` is around 0.5, the f1_score is the highest.

The best parameters are

```
{'learning_rate': 0.5, 'n_estimators': 200}
```

## 5.2  Hyperparameter Tuning for Random Forest

Hyperparameter tuning for random forest involves more parameter.

```
paramgrid_rf = {'max_depth': [5, 10, 15, 20, 25],
'min_samples_leaf': [i for i in np.linspace(1, 10, num = 5)],
'min_samples_split': [i for i in np.linspace(1, 10, num = 5)],
'n_estimators': [10,50,100],
'criterion': ['gini', 'entropy', 'log_loss']}
```

Note that random forest algorithm has more parameters than what I have tuned, but the change of other parameters did not yield significant improvement on the prediction result and accuracy of the model.

And the functional relationship between predictive accuracy and `min_samples_leaf` is
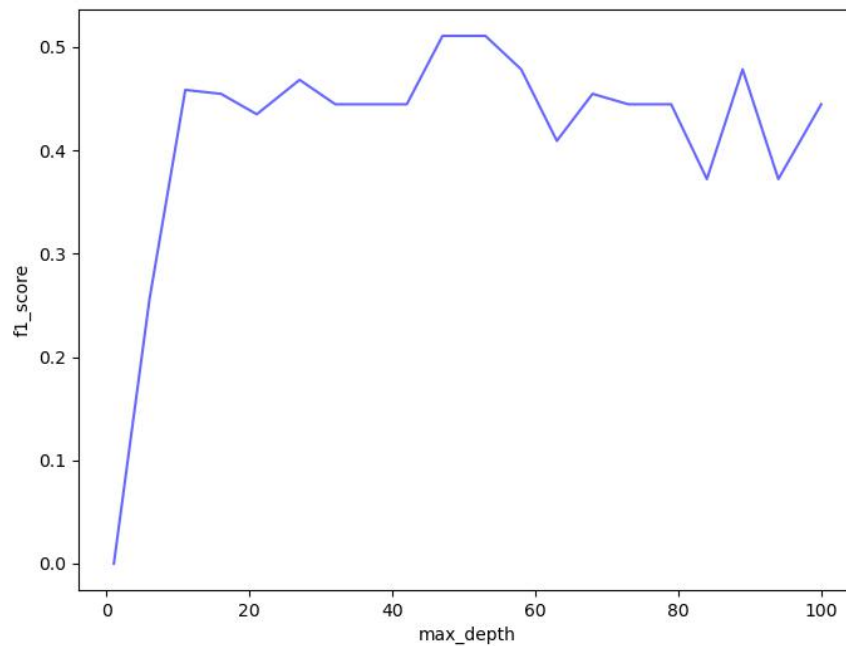
Figure 6: max depths VS f1 score

# 6 Predictive Accuracy

I finally got a f1 score of 0.606. The graphs below show the AUC curve and confusion matrix.

# 7 Errors and Mistakes

I think the most tricky part is data preprocessing and feature engineering. The first thing I ignored was that I should remove the outliers before model training. Because outliers will cause high bias. And I also did not rescale the data to the same range, which is very important because data with different magnitudes will be treated differently while in fact only the relative ranking matters. After implementing these two steps, f1_score improved from 0.5045 to 0.54.

And another thing is I should select 20 features instead of using all of the 26 features to train the model. Because some of the features are not relevant to the label, so applying `SelectKbest`, we can pick out the most relavant features.
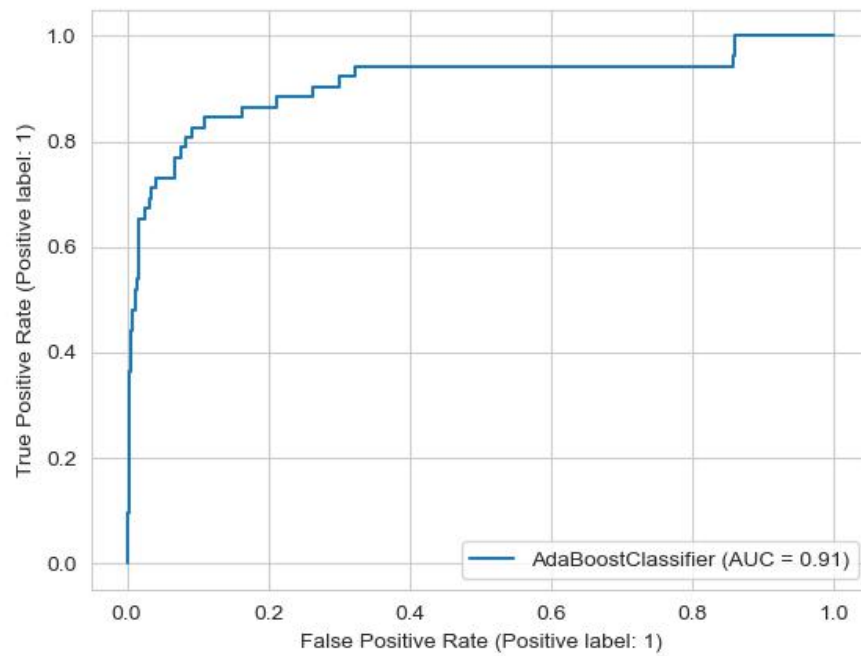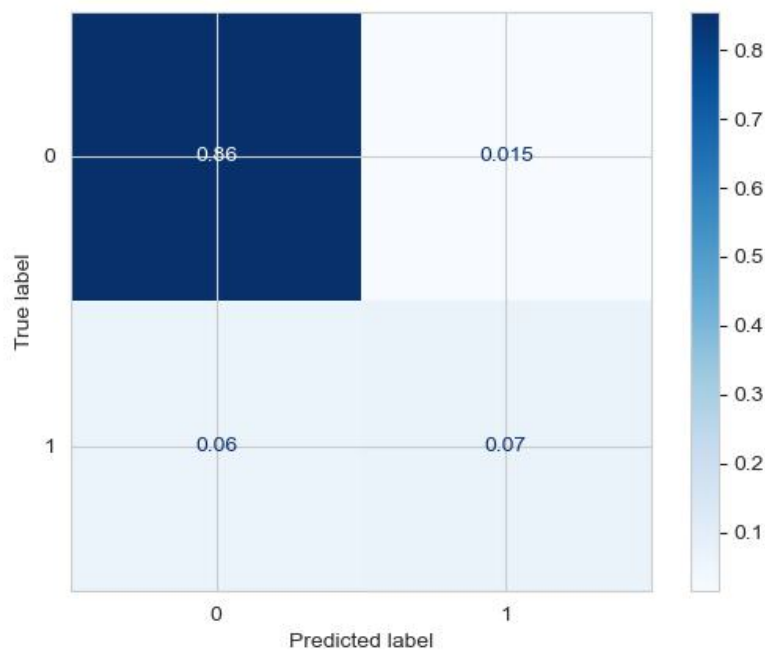
Figure 7: ROC curve



Figure 8: Confusion matrix

# Code

Copy and paste your code into your write-up document. Also, attach all the code needed for your competition. The code should be commented so that the grader can understand what is going on. Points will be taken off if if the code or the comments do not explain what is taking place. Your code should be executable with the installed libraries and only minor modifications.

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import random
import seaborn as sns
from sklearn.svm import SVC
from sklearn.preprocessing import OneHotEncoder
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression

from sklearn.ensemble import RandomForestRegressor
from sklearn.kernel_ridge import KernelRidge
from sklearn.linear_model import Ridge, Lasso
from sklearn.metrics import plot_roc_curve, roc_auc_score, f1_score,
    confusion_matrix
import time

import warnings
warnings.filterwarnings('ignore')

train = pd.read_csv('train.csv')
test = pd.read_csv('test.csv')
```

Data decription

```python
def DescribeData(dataframe):
    print('\n', '-' * 20, 'shape'.upper(), 20 * '-')
    print(dataframe.shape)
    print('\n', '-' * 20, 'dtypes'.upper(), 20 * '-')
    print(dataframe.dtypes)
    print('\n', '-' * 20, 'columns'.upper(), 20 * '-')
    print(dataframe.columns)
    print('\n', '-' * 20, 'info'.upper(), 20 * '-')
    print(dataframe.info())
    print('\n', '-' * 20, 'any null values'.upper(), 20 * '-')
    print(dataframe.isnull().values.any())
```

```python
    print('\n', '-' * 20, 'null values'.upper(), 20 * '-')
    print(dataframe.isnull().sum().sort_values(ascending=False))
    print('\n', '-' * 20, 'descriptive statistics'.upper(), 20 * '-')
    print(dataframe.describe())
DescribeData(train)


# find categorical and numerical features
def FindNumOrCategFeature(dataframe, dtype):
    '''''
    Report the numerical or categorical features
    for the dataframe.
    dtype = 'numerical' or 'categorical' or 'both'
    '''''
    # categorical features
    cate = (dataframe.dtypes == 'object')
    categorical = list(cate[cate].index)


    # numerical features
    num = (dataframe.dtypes == 'int64')
    nums = list(num[num].index)


    if dtype == 'numerical':
        print('Categorical featrues:')
        print(categorical)

    elif dtype == 'categorical':
        print('Numerical featrues:')
        print(nums)

    elif dtype == 'both':
        print('Categorical featrues:')
        print(categorical)
        print('Numerical featrues:')
        print(nums)

    return categorical, nums
cates, nums = FindNumOrCategFeature(train, 'both')

def DescribeCateFeature(dataframe, categorical, plot=False):

    print('\n', '-' * 10, categorical.upper(), 10 * '-')
    print(pd.DataFrame({
        categorical: dataframe[categorical].value_counts(),
        'RATIO (%)': round(100 * (dataframe[categorical].value_counts() /
```

```
          len(dataframe)), 2)
56    }))

57

58    if plot == True:
59        sns.countplot(x=dataframe[categorical], palette='Paired')
60        plt.ylabel('COUNT')
61        plt.xlabel(categorical.upper())
62        plt.show()

63

64 for c in cates:
65     DescribeCateFeature(dataframe=train, categorical=c, plot=True)
```

drop useless columns

```
1 train1 = train.copy()
2 test1 = test.copy()
3 # drop 'EmployeeID','STANDARDHOURS', 'EmployeeCount','Over18'
4 # all are over 18
5 # employee count is always 1
6 # standard hours is always 80
7 # employee ID does not conribute to model estimation
8 train2 = train1.drop(columns=['EmployeeID','StandardHours', 'EmployeeCount
      ','Over18'])

9

10 train2.hist(figsize=(20, 15), bins=50, xlabelsize=8, ylabelsize=8, alpha
      =0.8);
11 plt.savefig('histogram.jpg')
```

Feature engineering for training set

```
1 corr = train2.corr()
2 corr

3

4 def CorrAnalysis(dataframe, corr_threshold=0.7, plot=False):
5     corr = dataframe.corr()
6     cor_matrix = corr.abs()
7     upper_triangle_matrix = cor_matrix.where(
8         np.triu(np.ones(cor_matrix.shape), k=1).astype(np.bool_))
9     drop_list = [col for col in upper_triangle_matrix.columns
10                 if any(upper_triangle_matrix[col] > corr_threshold)]

11

12    if plot == True:
13        plt.figure(figsize=(15, 10))
14        sns.heatmap(corr, cmap='RdBu', square=True, annot=True, fmt='.1f')
15 #         plt.show()
```

```
16          plt.savefig('corr.jpg')
17
18      return drop_list
19
20  CorrAnalysis(train2, plot=True)
21
22  # drop highly correlated features
23  features_to_drop = ['MonthlyIncome',
24   'PerformanceRating',
25   'TotalWorkingYears',
26   'YearsInCurrentRole',
27   'YearsWithCurrManager']
28  train3 = train2.drop(columns=features_to_drop)
29  train3.shape
30  train3.columns
```

Encoding cateorical data

```
1  train3c = train3.copy()
2  cates1 = list(set(cates).intersection(set(train3c.columns)))
3  cates1
4
5  from sklearn.preprocessing import LabelEncoder
6
7  le = LabelEncoder()
8  le.fit_transform(train3c[cates1[-3]])
9
10  def LabelEncoder1(dataframe, feature_list):
11
12      for f in feature_list:
13          le = LabelEncoder()
14          encoded = le.fit_transform(dataframe[f])
15          dataframe[f] = encoded
16
17      return dataframe
18
19  df = LabelEncoder1(dataframe=train3c, feature_list=cates1)
20  df.columns
```

Outlier detection

```
1  # define a function to detect and remove outliers
2
3  def DetectOutlier(dataframe, features, drop=False):
4      '''
```

```
 5      Detect outliers using quantile method
 6      Feature: numerical variables, can be passed as int, float or list
 7      '''
 8      for feature in features:
 9
10          IQR = dataframe[feature].quantile(0.75) - dataframe[feature].
    quantile(0.25)
11          lower = dataframe[feature].quantile(0.25) - (IQR * 2.5)
12          upper = dataframe[feature].quantile(0.75) + (IQR * 2.5)
13          out_of_bound = dataframe[feature].loc[(dataframe[feature] < lower)
14                                              | (dataframe[feature] >
    upper)]
15          oob_idx = list(out_of_bound.index)
16
17          if len(oob_idx) != 0:
18              print('Feature {feature} has outliers.'.format(feature=feature
    ))
19              print('Outliers index of feature {feature} are: \n'.format(
    feature=feature), oob_idx)
20              print('Outliers range of feature {feature} are < {lower} or >
    {upper}\n'.format(
21                  feature=feature, lower=lower, upper=upper))
22
23          if drop == True:
24              dataframe = dataframe.drop(index=oob_idx)
25
26      return dataframe
27
28 nums1 = list(set(nums).intersection(set(train3c.columns)))
29 df_o = DetectOutlier(dataframe=df, features=nums1, drop=True)
```

Feature selection

```
 1 from sklearn.feature_selection import SelectKBest, chi2
 2
 3 X_ = df_o.drop(columns='Attrition')
 4 y_ = df_o['Attrition']
 5 slkb = SelectKBest(score_func=chi2,k=20)
 6 X1 = slkb.fit_transform(X_, y_)
 7
 8 cols = slkb.get_support(indices=True)
 9 df_new = pd.DataFrame(X1)
10 df_new
```

Scaling numerical data

```python
from sklearn.preprocessing import MinMaxScaler, StandardScaler

nums2 = list(set(nums).intersection(set(df_new.columns)))

def ScaleData(data, feature, scaler, inplace=False):

    if scaler == 'MinMax':
        min_ = np.min(data[feature])
        max_ = np.max(data[feature])

        if inplace == True:
            data[feature] = (data[feature] - min_) / (max_ - min_)

    if scaler == 'Standard':
        mu = np.mean(data[feature])
        std = np.std(data[feature])

        if inplace == True:
            data[feature] = (data[feature] - mu) / std

    return data

for f in df_new.columns:
    ScaleData(df_new, f, 'Standard', inplace=True)

df_new
```

Model training

```python
# split data
X = df_new
y = y_
X_train, X_test, y_train, y_test = train_test_split(
                        X, y, test_size = 0.2, random_state = 2022)

from sklearn.metrics import roc_auc_score, plot_confusion_matrix,
    roc_curve, f1_score

def TrainModelAndEval(model, X_train, X_test, y_train, y_test,
                    plot_ROC=False, plot_confusion_mat=False):
    start = time.time()
    # train model
    model.fit(X_train, y_train)
```

```python
    y_train_pred = model.predict(X_train)
    # predict
    y_pred = model.predict(X_test)

    # training accuracy
    f1score = f1_score(y_train, y_train_pred)
    auc = roc_auc_score(y_train, y_train_pred)
    train_time = time.time() - start

    # test accuracy
    f1score_t = f1_score(y_test, y_pred)
    auc_t = roc_auc_score(y_test, y_pred)

    print('-------Training Accuracy--------')
    print('F1 score = {}'.format(f1score))
    print('AUC = {}'.format(auc))
    print('Training time = {} \n'.format(train_time))

    print("-------Testing Accuracy--------")
    print("F1 score of test set = {}".format(f1score_t))
    print("AUC of test set = {}".format(auc_t))

    if plot_ROC == True:
        probs = model.predict_proba(X_train)
        fpr, tpr, thresholds = roc_curve(y_test, probs)
        plot_roc_cur(tpr, fpr)

    if plot_confusion_mat == True:
        plot_confusion_matrix(model, X_test, y_test,cmap=plt.cm.Blues,
    normalize = 'all')

    return model, f1score_t, auc_t, train_time

# random forest
from sklearn.ensemble import RandomForestClassifier

params_rf = {'max_depth': 10,
             'min_samples_leaf': 2,
             'min_samples_split': 2,
             'n_estimators': 100,
             'random_state': 2022}

model_rf = RandomForestClassifier(**params_rf)
model_rf, f1_score_rf, auc_rf, train_time_rf = TrainModelAndEval(
```

```
57                                              model_rf , X_train , X_test , y_train ,
     y_test )

58

59 # adaboost
60 from sklearn.ensemble import AdaBoostClassifier

61

62 params_ada = {'n_estimators': 100}
63 model_ada = AdaBoostClassifier(**params_ada)
64 model_ada , f1_score_ada , auc_ada , train_time_ada = TrainModelAndEval(
65                                              model_ada , X_train , X_test , y_train ,
     y_test )
```

Hyperparameter tuning

```
1 from sklearn.model_selection import GridSearchCV

2

3

4 params = {'n_estimators': [10, 25, 50, 75, 100, 200, 300, 500],
5           'learning_rate': [0.001, 0.01, 0.1, 0.5, 1., 5, 10, 20, 50, 100]}

6

7 adaboost = AdaBoostClassifier()
8 clf = GridSearchCV(adaboost , params)
9 clf.fit(X_train , y_train)
10 best_params = clf.best_params_
11 print(best_params)
12 model_ada_b = AdaBoostClassifier(**best_params)
13 model_ada_b , f1_score_ada_b , auc_ada_b , train_time_ada_b =
     TrainModelAndEval(
14                                              model_ada_b , X_train , X_test , y_train
     , y_test )

15

16 # plot f1 score VS learning rate

17

18 # lrate =   [0.001, 0.01, 0.1, 0.5, 1., 5, 10, 20, 50, 100]
19 lrate = [0.001, 0.01, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.8, 1, 5, 10, 20, 50,
     100]

20

21 f1s = []
22 for l in lrate:
23     ada_clf = AdaBoostClassifier(n_estimators=100, learning_rate =l)
24     ada_clf.fit(X_train , y_train)
25     y_pred = ada_clf.predict(X_test)
26     f1 = f1_score(y_test , y_pred)
27     f1s.append(f1)
```

```
28
29  fig1, ax1 = plt.subplots(figsize=(8,6))
30  log_lrate = np.log(lrate)
31  plt.plot(log_lrate, f1s, c='b',alpha=0.6)
32  plt.xlabel('log learning_rate')
33  plt.ylabel('f1_score')
34  plt.savefig('lrateADA.jpg')
```

```
1   paramgrid_rf = {'max_depth': [5, 10, 15, 20, 25],
2                'min_samples_leaf': [i for i in np.linspace(0.01, 1, num = 8)
        ],
3                'min_samples_split': [i for i in np.linspace(0.01, 1, num =
        8)],
4                'n_estimators': [10,50,100],
5                'criterion': ['gini', 'entropy', 'log_loss']
6                }
7   rf_clf = RandomForestClassifier()
8   clf = GridSearchCV(rf_clf, paramgrid_rf)
9   clf.fit(X_train, y_train)
10
11  depths = [int(i) for i in np.linspace(1, 100, num = 20)]
12
13  f1s_rf = []
14  for d in depths:
15      rf_clf = RandomForestClassifier(n_estimators=100, max_depth =d)
16      rf_clf.fit(X_train, y_train)
17      y_pred = rf_clf.predict(X_test)
18      f1_rf = f1_score(y_test, y_pred)
19      f1s_rf.append(f1_rf)
20
21  fig1, ax1 = plt.subplots(figsize=(8,6))
22  # log_split = np.log(lrate)
23  plt.plot(splits, f1s_rf, c='b',alpha=0.6)
24  plt.xlabel('max_depth')
25  plt.ylabel('f1_score')
26  plt.savefig('maxdepthRF.jpg')
```

Feature engineering for test set

```
1   # drop useless columns
2   test1 = test.drop(columns=['EmployeeID','StandardHours', 'EmployeeCount','
        Over18'])
3   # drop highly correlated features
4   test2 = test1.drop(columns=features_to_drop)
```

```
5  # select features
6  test3 = slkb.transform(test2)
7  # rescale data
8  scaler = StandardScaler()
9  test4 = scaler.fit_transform(test3)
10 test4
```

Predict and submit

```
1  y_pred_ada = model_ada.predict(test4)
2  y_pred_ada
3  y_pred_rf = model_rf_b.predict(test4)
4  y_pred_rf
5
6  ada = pd.DataFrame(y_pred_ada)
7  submit_ada = ada.reset_index()
8
9  submit_ada.columns=[['id','Predicted']]
10 submit_ada.to_csv('submission_ada_01.csv', index=False)
11
12 submit_rf = pd.DataFrame(y_pred_rf)
13 submit_rf = submit_rf.reset_index()
14
15 submit_rf.columns=[['id','Predicted']]
16 submit_rf.to_csv('submission_rf_01.csv', index=False)
```