

```
In [39]: 1  ## Import Library ##
2  import pandas as pd
3  import numpy as np
4  import matplotlib.pyplot as plt
5  import seaborn as sns
6
7  ## For Data Encoding ##
8  from sklearn.preprocessing import LabelEncoder
9
10 ## For Model Evaluation ##
11 from sklearn.model_selection import KFold
12
13 ## Machine Learning Model ##
14 from sklearn.ensemble import RandomForestClassifier
15 from sklearn.tree import DecisionTreeClassifier
16 from sklearn.linear_model import LogisticRegression
17 from sklearn.linear_model import LinearRegression
18
19 ## For Model Performance ##
20 from sklearn.metrics import accuracy_score, confusion_matrix, precision_s
```

```
In [62]: 1  ## Model Performance
2  def evaluation(gt, pred):
3      acc = accuracy_score(gt, pred)
4      precision = precision_score(gt, pred)
5      recall = recall_score(gt, pred)
6      f1 = f1_score(gt, pred)
7      matrix = confusion_matrix(gt, pred)
8
9      return acc, precision, recall, f1, matrix
```

```
In [41]: 1  raw_data = pd.read_csv("./EmployeeAttrition.csv", index_col=0)
```

In [77]: 1 raw_data

Out[77]:

	Age	Attrition	BusinessTravel	Department	DistanceFromHome	Education	EducationField	En
0	41	1	Travel_Rarely	Sales	1	2	Life Sciences	
1	49	0	Travel_Frequently	Research & Development	8	1	Life Sciences	
2	37	1	Travel_Rarely	Research & Development	2	2	Other	
3	33	0	Travel_Frequently	Research & Development	3	4	Life Sciences	
4	27	0	Travel_Rarely	Research & Development	2	1	Medical	
...
1465	36	0	Travel_Frequently	Research & Development	23	2	Medical	
1466	39	0	Travel_Rarely	Research & Development	6	1	Medical	
1467	27	0	Travel_Rarely	Research & Development	4	3	Life Sciences	
1468	49	0	Travel_Frequently	Sales	2	3	Medical	
1469	34	0	Travel_Rarely	Research & Development	8	3	Medical	

1470 rows × 27 columns

In [43]: 1 raw_data.columns

Out[43]: Index(['Age', 'Attrition', 'BusinessTravel', 'Department', 'DistanceFromHome',
 'Education', 'EducationField', 'EmployeeNumber',
 'EnvironmentSatisfaction', 'Gender', 'JobInvolvement', 'JobLevel',
 'JobRole', 'JobSatisfaction', 'MaritalStatus', 'MonthlyIncome',
 'NumCompaniesWorked', 'PerformanceRating', 'RelationshipSatisfaction',
 'StockOptionLevel', 'TotalWorkingYears', 'TrainingTimesLastYear',
 'WorkLifeBalance', 'YearsAtCompany', 'YearsInCurrentRole',
 'YearsSinceLastPromotion', 'YearsWithCurrManager'],
 dtype='object')

In [44]: 1 raw_data.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 1470 entries, 0 to 1469
Data columns (total 27 columns):
Age                                1470 non-null int64
Attrition                         1470 non-null object
BusinessTravel                    1470 non-null object
Department                        1470 non-null object
DistanceFromHome                  1470 non-null int64
Education                         1470 non-null int64
EducationField                    1470 non-null object
EmployeeNumber                    1470 non-null int64
EnvironmentSatisfaction            1470 non-null int64
Gender                            1470 non-null object
JobInvolvement                    1470 non-null int64
JobLevel                          1470 non-null int64
JobRole                           1470 non-null object
JobSatisfaction                   1470 non-null int64
MaritalStatus                     1470 non-null object
MonthlyIncome                     1470 non-null int64
NumCompaniesWorked                1470 non-null int64
PerformanceRating                 1470 non-null int64
RelationshipSatisfaction           1470 non-null int64
StockOptionLevel                  1470 non-null int64
TotalWorkingYears                 1470 non-null int64
TrainingTimesLastYear             1470 non-null int64
WorkLifeBalance                   1470 non-null int64
YearsAtCompany                    1470 non-null int64
YearsInCurrentRole                1470 non-null int64
YearsSinceLastPromotion           1470 non-null int64
YearsWithCurrManager              1470 non-null int64
dtypes: int64(20), object(7)
memory usage: 321.6+ KB
```

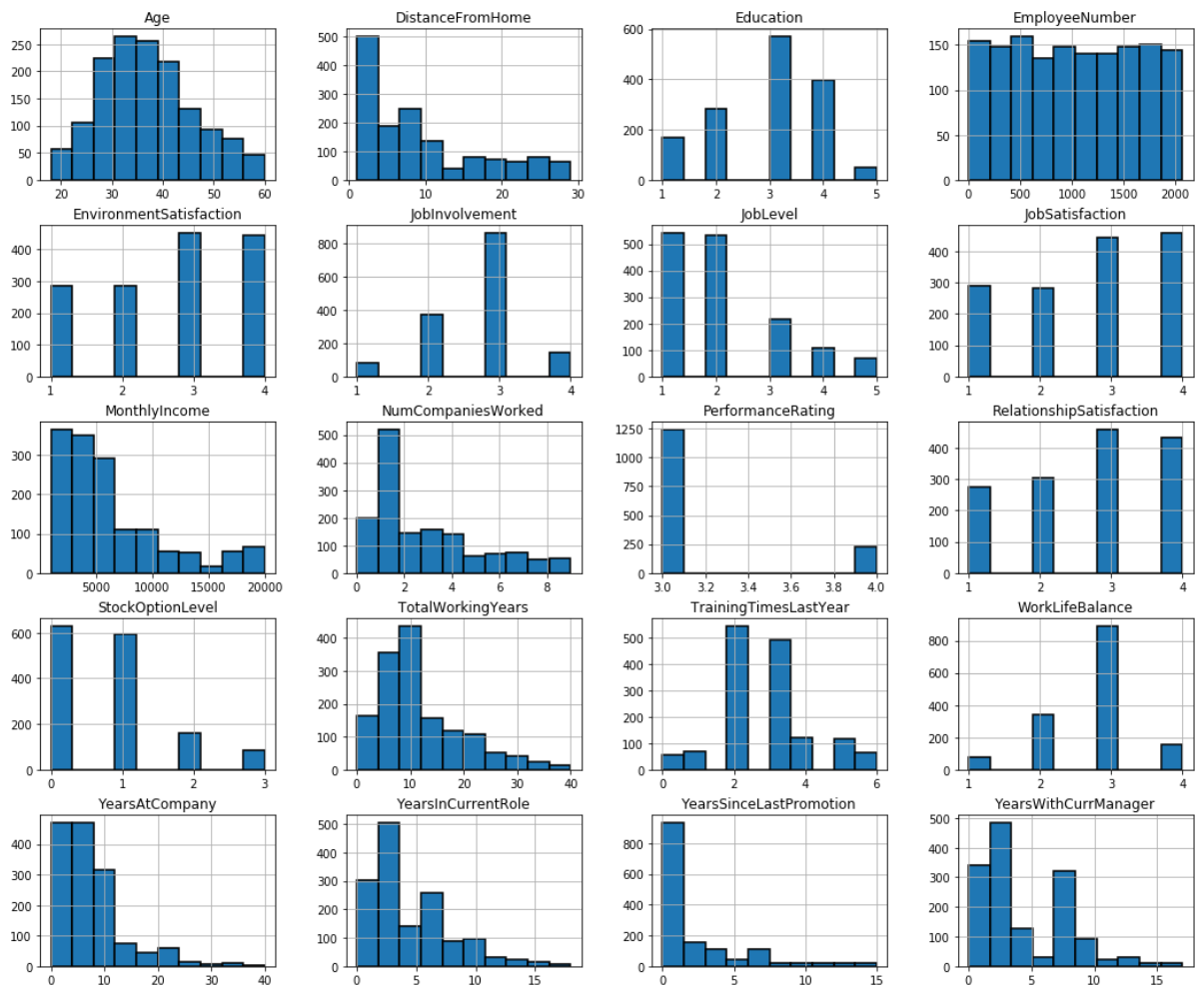
```
In [45]: 1 categorical_col = []
2 numeric_col = []
3 for col in raw_data.columns:
4     if raw_data[col].dtype == object and col != "Attrition":
5         categorical_col.append(col)
6         print(col, raw_data[col].unique())
7         print("=====")
8     elif raw_data[col].dtype == int and col != "Attrition":
9         numeric_col.append(col)
```

```
BusinessTravel ['Travel_Rarely' 'Travel_Frequently' 'Non-Travel']
=====
Department ['Sales' 'Research & Development' 'Human Resources']
=====
EducationField ['Life Sciences' 'Other' 'Medical' 'Marketing' 'Technical De
gree'
'Human Resources']
=====
Gender ['Female' 'Male']
=====
JobRole ['Sales Executive' 'Research Scientist' 'Laboratory Technician'
'Manufacturing Director' 'Healthcare Representative' 'Manager'
'Sales Representative' 'Research Director' 'Human Resources']
=====
MaritalStatus ['Single' 'Married' 'Divorced']
=====
```

```
In [46]: 1  ## Encode Label
2  raw_data["Attrition"] = raw_data["Attrition"].astype("category").cat.codes
3  raw_data["Attrition"].value_counts()
```

```
Out[46]: 0    1233
1      237
Name: Attrition, dtype: int64
```

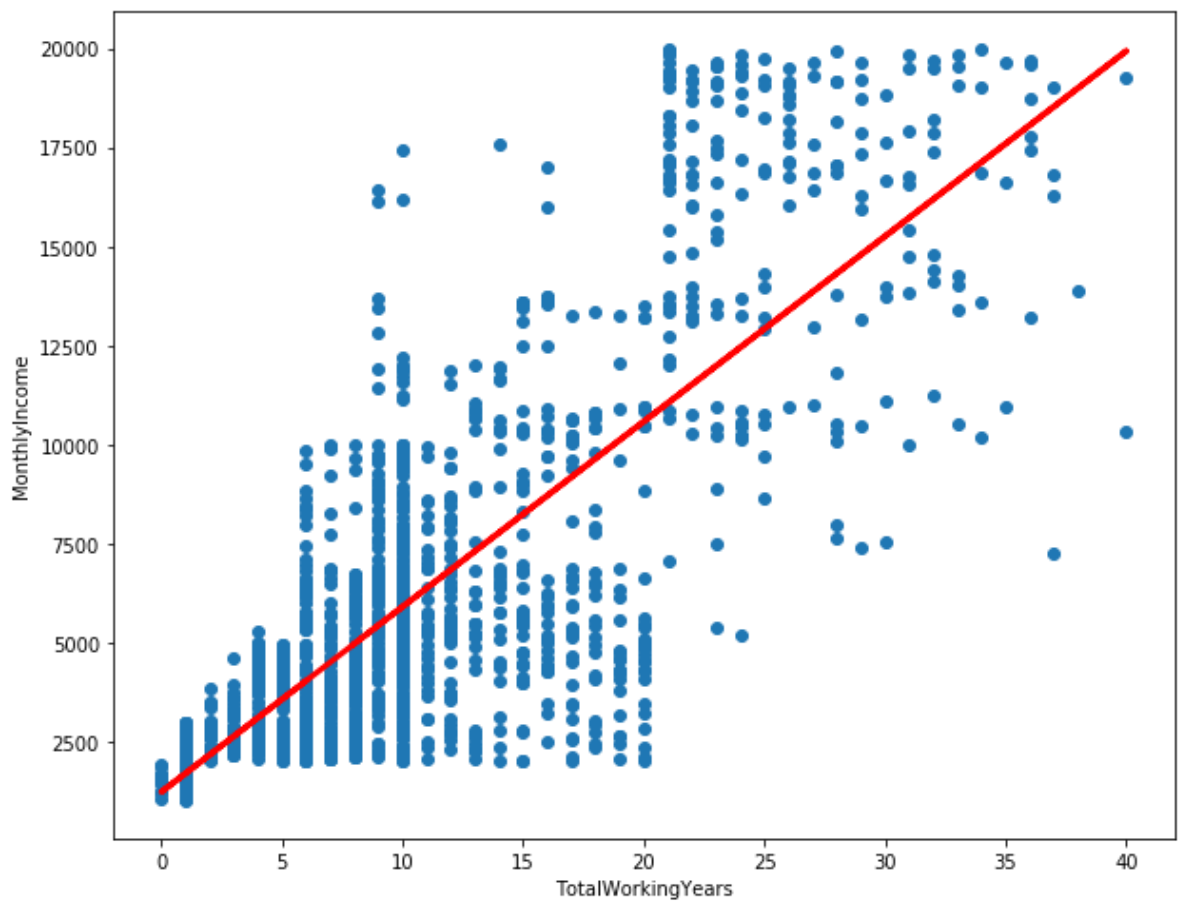
```
In [52]: 1  ## Data Visualization for numerical data
2  raw_data[numeric_col].hist(edgecolor='black', linewidth=1.5, figsize=(18,
3  plt.show())
```



```
In [55]: 1  X = np.array(raw_data["TotalWorkingYears"]).reshape(len(raw_data), 1)
2  y = np.array(raw_data["MonthlyIncome"])
3  model = LinearRegression()
4  model = model.fit(X, y)
5  pred = model.predict(X)
6  print("y = %f * X + %f" % (model.coef_[0], model.intercept_))
```

```
y = 467.658412 * X + 1227.935288
```

```
In [56]: 1 plt.figure(figsize=(10, 8))
2 plt.scatter(raw_data["TotalWorkingYears"], raw_data["MonthlyIncome"])
3 plt.plot(raw_data["TotalWorkingYears"], pred, color="red", linewidth = 3)
4 plt.xlabel("TotalWorkingYears")
5 plt.ylabel("MonthlyIncome")
6 plt.show()
```



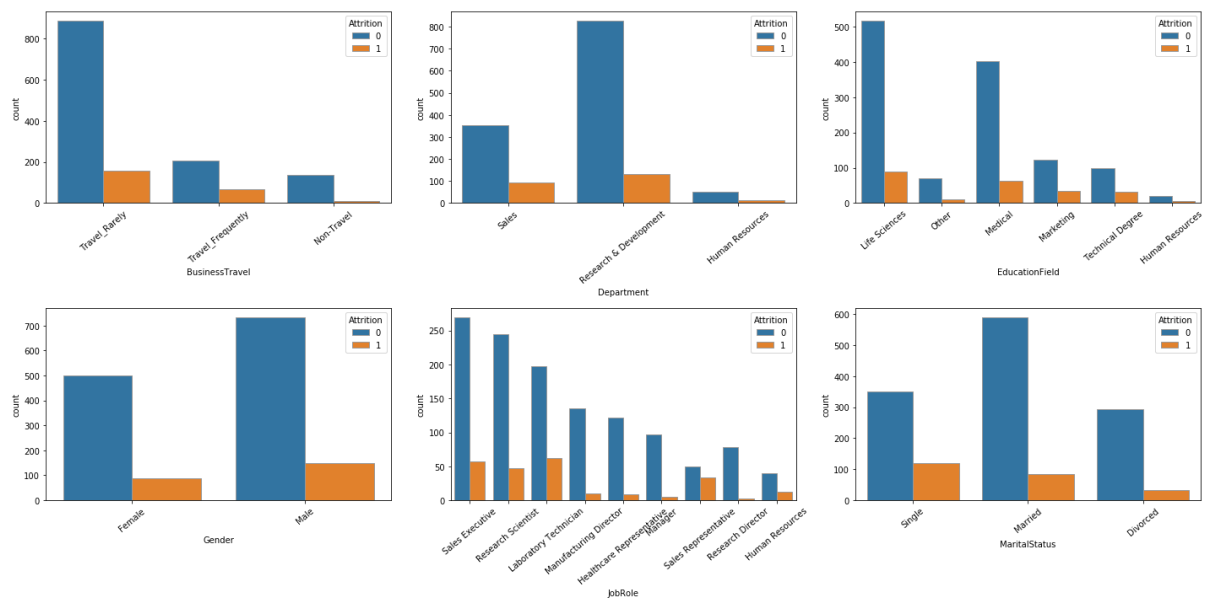
```
In [50]: 1 categorical_col
```

```
Out[50]: ['BusinessTravel',
'Department',
'EducationField',
'Gender',
'JobRole',
'MaritalStatus']
```

```

In [57]: 1  ## Data Visualization for categorical data
          2  fig, axes = plt.subplots(2, 3, figsize=(20, 10))
          3  for index, col in enumerate(categorical_col):
          4      row_num = int(index / 3)
          5      col_num = (index % 3)
          6      label = list(raw_data[col].unique())
          7      sns.countplot(x=str(col), hue="Attrition", edgecolor=".6", data=raw_data)
          8      axes[row_num, col_num].set_xticklabels(label, rotation=40)
          9  plt.tight_layout()
         10  plt.show()

```



```

In [58]: 1  ## Data Encoding (one-hot encoding)
          2  one_hot_encoding_df = pd.get_dummies(raw_data, columns=categorical_col)

```

```

In [16]: 1  categorical_col

```

```

Out[16]: ['BusinessTravel',
          'Department',
          'EducationField',
          'Gender',
          'JobRole',
          'MaritalStatus']

```

```
In [60]: 1 set(one_hot_encoding_df.columns) - set(numeric_col)
```

```
Out[60]: {'Attrition',  
          'BusinessTravel_Non-Travel',  
          'BusinessTravel_Travel_Frequently',  
          'BusinessTravel_Travel_Rarely',  
          'Department_Human Resources',  
          'Department_Research & Development',  
          'Department_Sales',  
          'EducationField_Human Resources',  
          'EducationField_Life Sciences',  
          'EducationField_Marketing',  
          'EducationField_Medical',  
          'EducationField_Other',  
          'EducationField_Technical Degree',  
          'Gender_Female',  
          'Gender_Male',  
          'JobRole_Healthcare Representative',  
          'JobRole_Human Resources',  
          'JobRole_Laboratory Technician',  
          'JobRole_Manager',  
          'JobRole_Manufacturing Director',  
          'JobRole_Research Director',  
          'JobRole_Research Scientist',  
          'JobRole_Sales Executive',  
          'JobRole_Sales Representative',  
          'MaritalStatus_Divorced',  
          'MaritalStatus_Married',  
          'MaritalStatus_Single'}
```

```

In [66]: 1  ## Data Splitting and Model Learning (Decision Tree)
2  avg_acc = 0
3  avg_precision = 0
4  avg_recall = 0
5  avg_f1 = 0
6  avg_confusion_matrix = []
7  avg_feature_importance = []
8
9  kf = KFold(n_splits=5)
10 fold_count = 0
11 for train_index, test_index in kf.split(one_hot_encoding_df):
12     print("Training Data: %d, Testing Data: %d" % (len(train_index), len(
13         train_X = one_hot_encoding_df.loc[train_index, one_hot_encoding_df.co
14         train_y = one_hot_encoding_df.loc[train_index]["Attrition"]
15         test_X = one_hot_encoding_df.loc[test_index, one_hot_encoding_df.colu
16         test_y = one_hot_encoding_df.loc[test_index]["Attrition"]
17
18         model = DecisionTreeClassifier(random_state=200)
19         model = model.fit(train_X, train_y)
20         test_predict = model.predict(test_X)
21         avg_feature_importance.append(model.feature_importances_)
22
23         acc, precision, recall, f1, matrix = evaluation(test_y, test_predict)
24
25         print("Fold: %d, Accuracy: %f, Precision: %f, Recall: %f, F1: %f" % (
26         avg_acc += acc
27         avg_precision += precision
28         avg_recall += recall
29         avg_f1 += f1
30         avg_confusion_matrix.append(matrix)
31         fold_count += 1
32
33 print("=====
34 print("Avg Accuracy: %f, Avg Precision: %f, Avg Recall: %f, Avg F1: %f" %
35
36
37

```

Training Data: 1176, Testing Data: 294

Fold: 1, Accuracy: 0.752000, Precision: 0.268000, Recall: 0.319000, F1: 0.291000

Training Data: 1176, Testing Data: 294

Fold: 2, Accuracy: 0.789000, Precision: 0.333000, Recall: 0.442000, F1: 0.380000

Training Data: 1176, Testing Data: 294

Fold: 3, Accuracy: 0.741000, Precision: 0.311000, Recall: 0.237000, F1: 0.269000

Training Data: 1176, Testing Data: 294

Fold: 4, Accuracy: 0.738000, Precision: 0.233000, Recall: 0.311000, F1: 0.267000

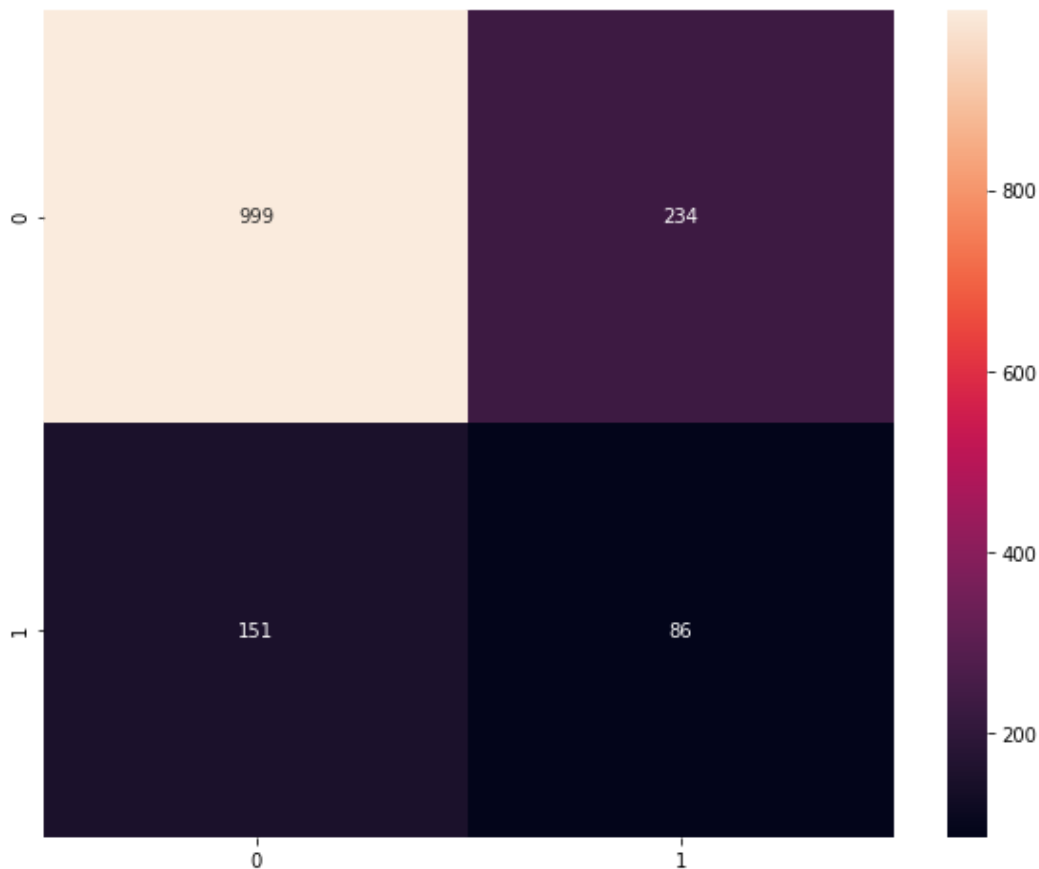
Training Data: 1176, Testing Data: 294

Fold: 5, Accuracy: 0.670000, Precision: 0.235000, Recall: 0.558000, F1: 0.331000

=====

Avg Accuracy: 0.738000, Avg Precision: 0.276000, Avg Recall: 0.374000, Avg F1: 0.308000


```
In [67]: 1 plt.figure(figsize=(10, 8))
2 sns.heatmap(np.sum(np.array(avg_confusion_matrix), axis=0), annot=True, f
3 plt.show()
```



```
In [68]: 1 importance_dict = {}
2 for col, importance in zip(train_X.columns, np.mean(np.array(avg_feature_
3 importance_dict[col] = importance
4
5 sorted(importance_dict.items(), key=lambda x: -x[1])[:10])
```

```
Out[68]: [('MonthlyIncome', 0.10020746536225093),
('TotalWorkingYears', 0.08286263996666518),
('Age', 0.07866270282338719),
('DistanceFromHome', 0.07651699500507207),
('EmployeeNumber', 0.06251880620094238),
('NumCompaniesWorked', 0.054148502062084915),
('TrainingTimesLastYear', 0.04497546911268388),
('JobInvolvement', 0.03828324416180727),
('StockOptionLevel', 0.035453500524868876),
('JobSatisfaction', 0.03481479445074119)]
```

```

In [69]: 1  ## Data Splitting and Model Learning (Random Forest)
2  avg_acc = 0
3  avg_precision = 0
4  avg_recall = 0
5  avg_f1 = 0
6  avg_confusion_matrix = []
7  avg_feature_importance = []
8
9  kf = KFold(n_splits=5)
10 fold_count = 0
11 for train_index, test_index in kf.split(one_hot_encoding_df):
12     print("Training Data: %d, Testing Data: %d" % (len(train_index), len(
13         train_X = one_hot_encoding_df.loc[train_index, one_hot_encoding_df.co
14         train_y = one_hot_encoding_df.loc[train_index]["Attrition"]
15         test_X = one_hot_encoding_df.loc[test_index, one_hot_encoding_df.colu
16         test_y = one_hot_encoding_df.loc[test_index]["Attrition"]
17
18         model = RandomForestClassifier(n_estimators=300)
19         model = model.fit(train_X, train_y)
20         test_predict = model.predict(test_X)
21
22         avg_feature_importance.append(model.feature_importances_)
23
24         acc, precision, recall, f1, matrix = evaluation(test_y, test_predict)
25         print("Fold: %d, Accuracy: %f, Precision: %f, Recall: %f, F1: %f" % (
26         avg_acc += acc
27         avg_precision += precision
28         avg_recall += recall
29         avg_f1 += f1
30         avg_confusion_matrix.append(matrix)
31         fold_count += 1
32
33 print("=====
34 print("Avg Accuracy: %f, Avg Precision: %f, Avg Recall: %f, Avg F1: %f" %
35
36
37

```

Training Data: 1176, Testing Data: 294

Fold: 1, Accuracy: 0.847000, Precision: 0.625000, Recall: 0.106000, F1: 0.182000

Training Data: 1176, Testing Data: 294

Fold: 2, Accuracy: 0.871000, Precision: 0.727000, Recall: 0.186000, F1: 0.296000

Training Data: 1176, Testing Data: 294

Fold: 3, Accuracy: 0.806000, Precision: 0.667000, Recall: 0.068000, F1: 0.123000

Training Data: 1176, Testing Data: 294

Fold: 4, Accuracy: 0.861000, Precision: 0.667000, Recall: 0.178000, F1: 0.281000

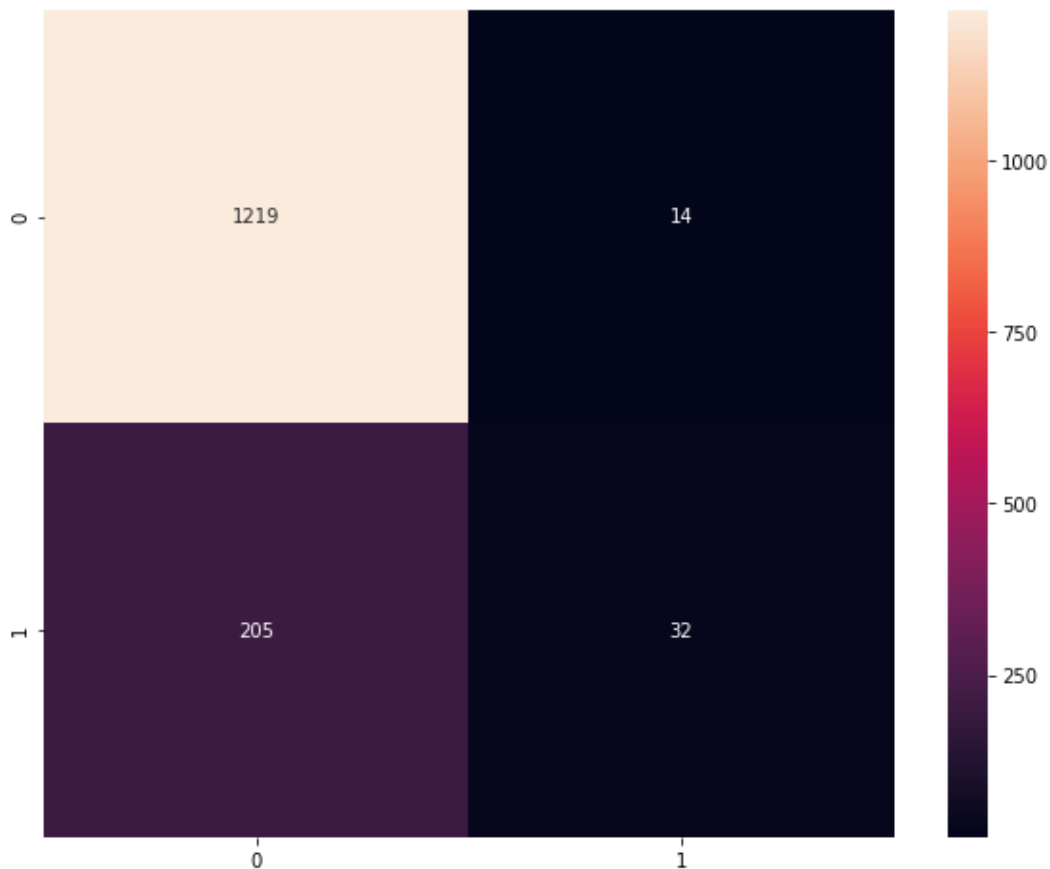
Training Data: 1176, Testing Data: 294

Fold: 5, Accuracy: 0.871000, Precision: 0.778000, Recall: 0.163000, F1: 0.269000

=====

Avg Accuracy: 0.851000, Avg Precision: 0.693000, Avg Recall: 0.140000, Avg F1: 0.230000

```
In [72]: 1 plt.figure(figsize=(10, 8))
2          sns.heatmap(np.sum(np.array(avg_confusion_matrix), axis=0), annot=True, f
3          plt.show()
```



```
In [73]: 1 importance_dict = {}
2          for col, importance in zip(train_X.columns, np.mean(np.array(avg_feature_
3              importance_dict[col] = importance
4
5          sorted(importance_dict.items(), key=lambda x: -x[1])[:10])
```

```
Out[73]: [('MonthlyIncome', 0.08101629049886984),
('Age', 0.0716518888272812),
('EmployeeNumber', 0.06303118631947577),
('TotalWorkingYears', 0.057268848061669896),
('DistanceFromHome', 0.0561689074827289),
('YearsAtCompany', 0.0471668933929604),
('NumCompaniesWorked', 0.03970843810345892),
('YearsWithCurrManager', 0.03646828928980059),
('EnvironmentSatisfaction', 0.03451136743068421),
('YearsInCurrentRole', 0.03261817566506732)]
```

```

In [74]: 1  ## Data Splitting and Model Learning (Logistic Regression)
2  avg_acc = 0
3  avg_precision = 0
4  avg_recall = 0
5  avg_f1 = 0
6  avg_feature_importance = []
7  avg_confusion_matrix = []
8
9  kf = KFold(n_splits=5)
10 fold_count = 0
11 for train_index, test_index in kf.split(one_hot_encoding_df):
12     print("Training Data: %d, Testing Data: %d" % (len(train_index), len(
13         train_X = one_hot_encoding_df.loc[train_index, one_hot_encoding_df.co
14         train_y = one_hot_encoding_df.loc[train_index]["Attrition"]
15         test_X = one_hot_encoding_df.loc[test_index, one_hot_encoding_df.colu
16         test_y = one_hot_encoding_df.loc[test_index]["Attrition"]
17
18         model = LogisticRegression(solver='liblinear')
19         model = model.fit(train_X, train_y)
20         test_predict = model.predict(test_X)
21
22         acc, precision, recall, f1, matrix = evaluation(test_y, test_predict)
23         print("Fold: %d, Accuracy: %f, Precision: %f, Recall: %f, F1: %f" % (
24             avg_acc += acc
25             avg_precision += precision
26             avg_recall += recall
27             avg_f1 += f1
28             avg_confusion_matrix.append(matrix)
29             fold_count += 1
30
31 print("=====
32 print("Avg Accuracy: %f, Avg Precision: %f, Avg Recall: %f, Avg F1: %f" %
33
34
35

```

Training Data: 1176, Testing Data: 294

Fold: 1, Accuracy: 0.854000, Precision: 0.833000, Recall: 0.106000, F1: 0.189000

Training Data: 1176, Testing Data: 294

Fold: 2, Accuracy: 0.861000, Precision: 0.562000, Recall: 0.209000, F1: 0.305000

Training Data: 1176, Testing Data: 294

Fold: 3, Accuracy: 0.813000, Precision: 0.611000, Recall: 0.186000, F1: 0.286000

Training Data: 1176, Testing Data: 294

Fold: 4, Accuracy: 0.871000, Precision: 0.733000, Recall: 0.244000, F1: 0.367000

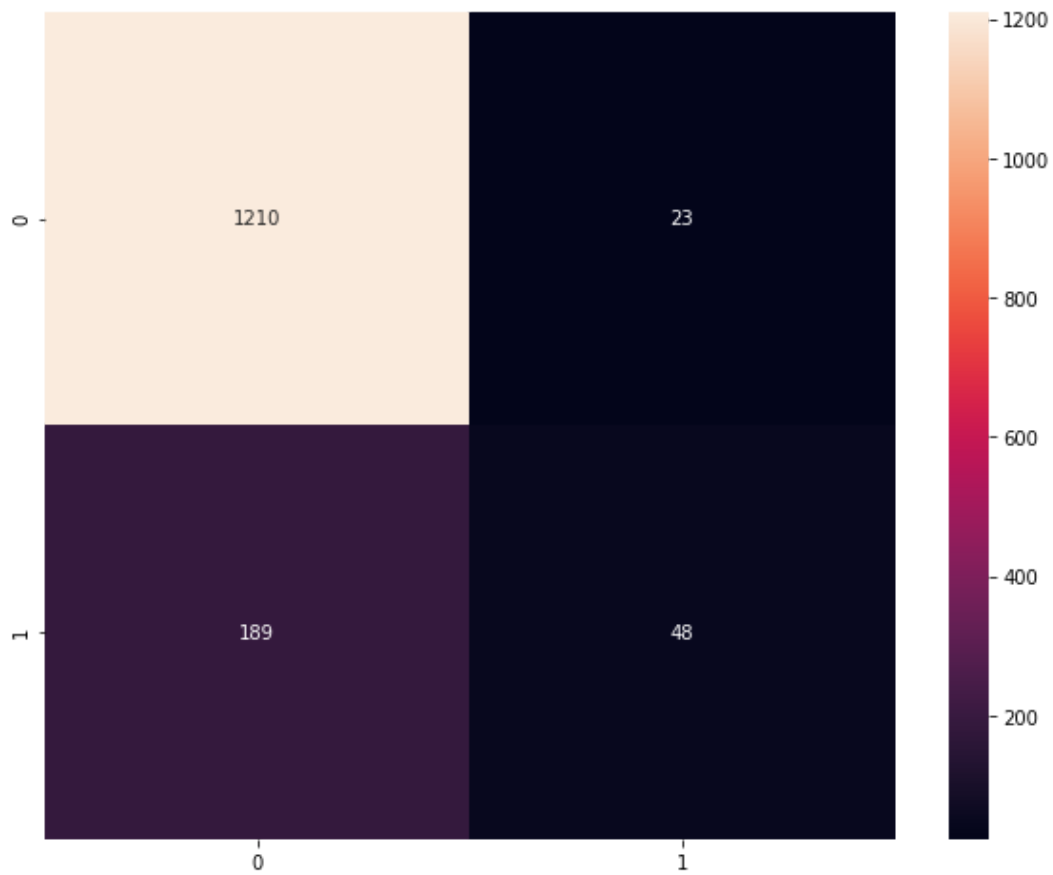
Training Data: 1176, Testing Data: 294

Fold: 5, Accuracy: 0.881000, Precision: 0.750000, Recall: 0.279000, F1: 0.407000

=====

=====
Avg Accuracy: 0.856000, Avg Precision: 0.698000, Avg Recall: 0.205000, Avg F1: 0.311000

```
In [76]: 1 plt.figure(figsize=(10, 8))  
2 sns.heatmap(np.sum(np.array(avg_confusion_matrix), axis=0), annot=True, f  
3 plt.show()
```



```
In [ ]:
```

```
1
```