# Classification Analysis Report

This report outlines the process of performing classification analysis on a dataset to predict employee attrition. The analysis includes data preprocessing, feature exploration, and the implementation of multiple classification algorithms. Evaluation metrics such as accuracy, F1 score, precision, recall, and ROC-AUC are computed to compare the performance of the models. The following steps are implemented:

1. Problem Statement Predict whether an employee will leave the organization (attrition) or stay, based on their attributes such as demographics, job satisfaction, and performance metrics.

2. Data Preprocessing

   - Class Imbalance Handling: Addressed using class weights in models to account for the imbalance between attrition and non-attrition cases.
   - Feature Scaling: Features were normalized to improve model performance.
   - Dataset Splitting: The dataset was split into training and testing sets for model evaluation.

3. Exploratory Data Analysis (EDA)

   - Correlation Analysis: A heatmap of the correlation matrix was created to explore relationships between features and the target variable.
   - Outlier Detection: Visualizations such as boxplots and histograms were used to identify and address outliers.
   - Feature Insights: Key features influencing employee attrition were identified through statistical analysis and visual exploration.

4. Model Implementation

   - Support Vector Machine (SVM): A balanced SVM with a linear kernel was implemented to handle class imbalance effectively.
   - Logistic Regression: Implemented as a baseline model for comparison.
   - Decision Trees: Used to capture non-linear relationships between features.
   - Random Forest: An ensemble learning method was employed to improve classification accuracy and capture feature importance.

5. Model Evaluation

- Evaluation Metrics: Accuracy, F1 score, precision, recall, and ROC-AUC were computed to assess model performance.

# Importing, Previewing, and Exploring Data

We first imported and previewed the data.

In [270…

```python
#importing necessary libraries, etc
import pandas as pd
import numpy as np
import seaborn as sns
import os
np.random.seed(0)
import matplotlib.pyplot as plt
import statsmodels.api as sm
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
from sklearn.model_selection import GridSearchCV
from sklearn.preprocessing import OneHotEncoder
import warnings
warnings.filterwarnings('ignore')
%matplotlib inline
#importing datafile
df = pd.read_csv('employee_attrition.csv')
#previewing data
df.head()
```

Out[270...

| | Age | Attrition | BusinessTravel | DailyRate | Department | DistanceFromHome | Education | EducationField | EmployeeCount | Em |
|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 41 | Yes | Travel_Rarely | 1102 | Sales | 1 | 2 | Life Sciences | 1 | |
| **1** | 49 | No | Travel_Frequently | 279 | Research & Development | 8 | 1 | Life Sciences | 1 | |
| **2** | 37 | Yes | Travel_Rarely | 1373 | Research & Development | 2 | 2 | Other | 1 | |
| **3** | 33 | No | Travel_Frequently | 1392 | Research & Development | 3 | 4 | Life Sciences | 1 | |
| **4** | 27 | No | Travel_Rarely | 591 | Research & Development | 2 | 1 | Medical | 1 | |

5 rows × 35 columns

## Comments:

It looks like we have numerical and categorical data with mostly self-explanatory columns. A few features are not clear how the survey questions translate to actual answers. For example, the feature Education contains integers. This project assumes that level of Education increases as the numbers increase, but it is unclear with a dictionary key. Similarly, Distance From Home contains a set of integers. Again, this project assumes a higher number equates to a home farther away from the workplace.

In [271...
```python
#Verifying there are no duplicate rows.
df.duplicated().value_counts()
```

Out[271...
```
False    1470
Name: count, dtype: int64
```

## Comment on Findings:

There are no duplicate rows in this dataset.

In [272...
```python
#verifying all features have all rows and no null values.
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1470 entries, 0 to 1469
Data columns (total 35 columns):
 #   Column                    Non-Null Count  Dtype
---  ------                    --------------  -----
 0   Age                       1470 non-null   int64
 1   Attrition                 1470 non-null   object
 2   BusinessTravel            1470 non-null   object
 3   DailyRate                 1470 non-null   int64
 4   Department                1470 non-null   object
 5   DistanceFromHome          1470 non-null   int64
 6   Education                 1470 non-null   int64
 7   EducationField            1470 non-null   object
 8   EmployeeCount             1470 non-null   int64
 9   EmployeeNumber            1470 non-null   int64
 10  EnvironmentSatisfaction   1470 non-null   int64
 11  Gender                    1470 non-null   object
 12  HourlyRate                1470 non-null   int64
 13  JobInvolvement            1470 non-null   int64
 14  JobLevel                  1470 non-null   int64
 15  JobRole                   1470 non-null   object
 16  JobSatisfaction           1470 non-null   int64
 17  MaritalStatus             1470 non-null   object
 18  MonthlyIncome             1470 non-null   int64
 19  MonthlyRate               1470 non-null   int64
 20  NumCompaniesWorked        1470 non-null   int64
 21  Over18                    1470 non-null   object
 22  OverTime                  1470 non-null   object
 23  PercentSalaryHike         1470 non-null   int64
 24  PerformanceRating         1470 non-null   int64
 25  RelationshipSatisfaction  1470 non-null   int64
 26  StandardHours             1470 non-null   int64
 27  StockOptionLevel          1470 non-null   int64
 28  TotalWorkingYears         1470 non-null   int64
 29  TrainingTimesLastYear     1470 non-null   int64
 30  WorkLifeBalance           1470 non-null   int64
 31  YearsAtCompany            1470 non-null   int64
 32  YearsInCurrentRole        1470 non-null   int64
 33  YearsSinceLastPromotion   1470 non-null   int64
 34  YearsWithCurrManager      1470 non-null   int64
dtypes: int64(26), object(9)
memory usage: 402.1+ KB
```

## Comment on Findings:

Nothing that needs cleaning, and we have no null values. Our target variable, Attrition, currently contains the strings "Yes" or "No" so we'll have to change that for our analysis.

```
In [273...   #checking percentage of employees in data with attrition
             df.Attrition.value_counts(normalize=True)
             #
```

```
Out[273...   Attrition
             No     0.838776
             Yes    0.161224
             Name: proportion, dtype: float64
```

## Comment on Findings:

Random chance would be right as a baseline 83.88% of the time. The following models will aim for higher than that. An unbalanced dataset with many more "No" than "Yes" for the Attrition column.

# Manipulating the Data, Logistic Regression Models

```
In [274...   #Changing target feature Attrition to numerical values instead of categorical
             df.Attrition = df.Attrition.map({'Yes': 1, 'No': 0})
```

```
In [275...   df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1470 entries, 0 to 1469
Data columns (total 35 columns):
 #   Column                    Non-Null Count  Dtype
---  ------                    --------------  -----
 0   Age                       1470 non-null   int64
 1   Attrition                 1470 non-null   int64
 2   BusinessTravel            1470 non-null   object
 3   DailyRate                 1470 non-null   int64
 4   Department                1470 non-null   object
 5   DistanceFromHome          1470 non-null   int64
 6   Education                 1470 non-null   int64
 7   EducationField            1470 non-null   object
 8   EmployeeCount             1470 non-null   int64
 9   EmployeeNumber            1470 non-null   int64
 10  EnvironmentSatisfaction   1470 non-null   int64
 11  Gender                    1470 non-null   object
 12  HourlyRate                1470 non-null   int64
 13  JobInvolvement            1470 non-null   int64
 14  JobLevel                  1470 non-null   int64
 15  JobRole                   1470 non-null   object
 16  JobSatisfaction           1470 non-null   int64
 17  MaritalStatus             1470 non-null   object
 18  MonthlyIncome             1470 non-null   int64
 19  MonthlyRate               1470 non-null   int64
 20  NumCompaniesWorked        1470 non-null   int64
 21  Over18                    1470 non-null   object
 22  OverTime                  1470 non-null   object
 23  PercentSalaryHike         1470 non-null   int64
 24  PerformanceRating         1470 non-null   int64
 25  RelationshipSatisfaction  1470 non-null   int64
 26  StandardHours             1470 non-null   int64
 27  StockOptionLevel          1470 non-null   int64
 28  TotalWorkingYears         1470 non-null   int64
 29  TrainingTimesLastYear     1470 non-null   int64
 30  WorkLifeBalance           1470 non-null   int64
 31  YearsAtCompany            1470 non-null   int64
 32  YearsInCurrentRole        1470 non-null   int64
 33  YearsSinceLastPromotion   1470 non-null   int64
 34  YearsWithCurrManager      1470 non-null   int64
dtypes: int64(27), object(8)
memory usage: 402.1+ KB
```

## Comments:

.map() worked, and our Attrition column now contains integers in place of the strings, "Yes" and "No."

```python
#previewing the change to numerical values for Attrition feature
df.head()
```

| | Age | Attrition | BusinessTravel | DailyRate | Department | DistanceFromHome | Education | EducationField | EmployeeCount | En |
|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 41 | 1 | Travel_Rarely | 1102 | Sales | 1 | 2 | Life Sciences | 1 | |
| **1** | 49 | 0 | Travel_Frequently | 279 | Research & Development | 8 | 1 | Life Sciences | 1 | |
| **2** | 37 | 1 | Travel_Rarely | 1373 | Research & Development | 2 | 2 | Other | 1 | |
| **3** | 33 | 0 | Travel_Frequently | 1392 | Research & Development | 3 | 4 | Life Sciences | 1 | |
| **4** | 27 | 0 | Travel_Rarely | 591 | Research & Development | 2 | 1 | Medical | 1 | |

5 rows × 35 columns

## Dividing the Training and Testing Data

```python
y = df.Attrition
X = df.drop('Attrition', axis = 1)
#normalizing data
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=42)
```

## Visualizing Most Correlated Features

We are curious which features may closely correlate with Attrition of the get-go.

```python
In [278...    # Create a new 'images' folder in directory
             if not os.path.exists('images'):
                 os.makedirs('images')

             # Convert categorical data to numerical
             X_train_encoded = pd.get_dummies(X_train, drop_first=True)
             y_train_encoded = y_train.copy()

             # Combine target and features
             heatmap_data = pd.concat([pd.Series(y_train_encoded, name='Attrition'), X_train_encoded], axis=1)

             # Compute correlation matrix
             corr = heatmap_data.corr()

             fig, ax = plt.subplots(figsize=(20, 20))

             # Plots a heatmap of the correlation matrix, with both
             # numbers and colors indicating the correlations
             sns.heatmap(
                 # Specifies the data to be plotted
                 data=corr,
                 # The mask means we only show half the values,
                 # instead of showing duplicates. It's optional.
                 mask=np.triu(np.ones_like(corr, dtype=bool)),
                 # Specifies that we should use the existing axes
                 ax=ax,
                 # Specifies that we want labels, not just colors
                 annot=True,
                 # Customizes colorbar appearance
                 cbar_kws={"label": "Correlation", "orientation": "horizontal", "pad": .2, "extend": "both"}
             )

             # Customize the plot appearance
             ax.set_title("Heatmap of Correlation Between Attributes (Including Target)");

             # Save the plot as a JPG image with the title as the filename
             filename = f"images/Correlation Between Attributes.jpg"
             plt.savefig(filename)
```

Heatmap of Correlation Between Attributes (Including Target)

# Comment on Findings:

It looks like none of our features are very highly correlated. The features below are the least correlated, and therefore correspond to the least attrition are:

- Total Working Years
- Years in Current Role
- Job Level As these features increase, the employee's risk of attrition descreases

Conversely, the features that were the most highly correlated, and therefore correspond to the most attrition are:

- Number of Companies Worked
- Distance from Home
- Monthly Rate As these features increase, the employee's risk of attrition increases.

we are curious as to the distribution of the least correlated features (the ones that correlate with the lowest attrition), so we take a look at those below.

In [279...
```
#Describing top 3 most correlating features
df[['YearsInCurrentRole', 'TotalWorkingYears', 'JobLevel']].describe()
```

Out[279...

|       | YearsInCurrentRole | TotalWorkingYears | JobLevel    |
|-------|--------------------|-------------------|-------------|
| count | 1470.000000        | 1470.000000       | 1470.000000 |
| mean  | 4.229252           | 11.279592         | 2.063946    |
| std   | 3.623137           | 7.780782          | 1.106940    |
| min   | 0.000000           | 0.000000          | 1.000000    |
| 25%   | 2.000000           | 6.000000          | 1.000000    |
| 50%   | 3.000000           | 10.000000         | 2.000000    |
| 75%   | 7.000000           | 15.000000         | 3.000000    |
| max   | 18.000000          | 40.000000         | 5.000000    |

## Comments:

The range is interesting, and the majority of the data falls in between 2 and 7 years for the Years In Current Role, 6-15 for total working years, and 1-3 for Job level. It would be interesting to investigate how high the attrition is for the lowest quartile, and how low the attrition is for the highest quartile.
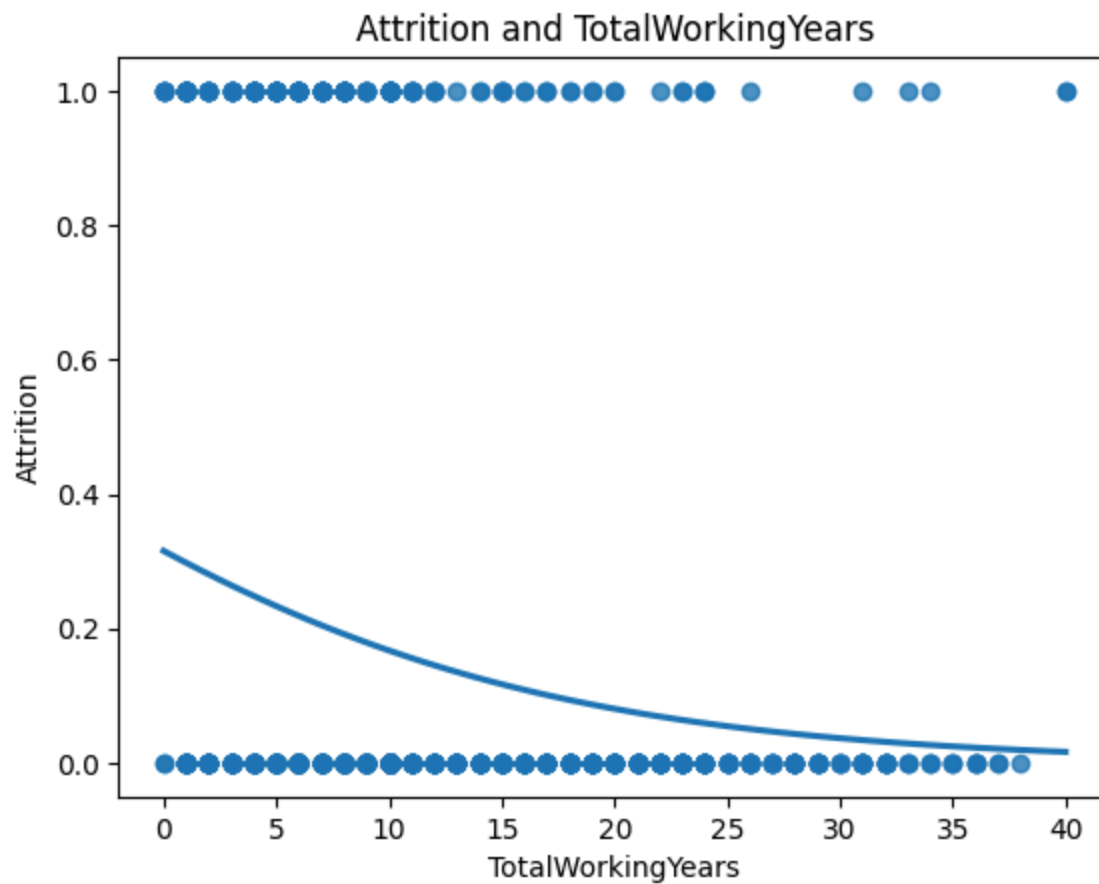
## Visualizing Most and Least Correlated Features

### Comments:

We thought it important to visualize these significant features for ourselves. A heat map is useful, but it doesn't necessarily convey the type of relationship between the two variables in a way that would be easily understood by everyone.

```
least_correlated_feature1 = 'TotalWorkingYears'
least_correlated_feature2 = 'YearsInCurrentRole'
least_correlated_feature3 = 'JobLevel'
most_correlated_feature1 = 'NumCompaniesWorked'
most_correlated_feature2 = 'DistanceFromHome'
most_correlated_feature3 = 'MonthlyRate'

sns.regplot(x=X_train[least_correlated_feature1], y=y_train, data=df, logistic=True, ci=None).set(title=
    'Attrition and {}'.format(least_correlated_feature1));

# Save the plot as a JPG image with the title as the filename
filename = f"images/Attrition and Total Working Years.jpg"
plt.savefig(filename)
```
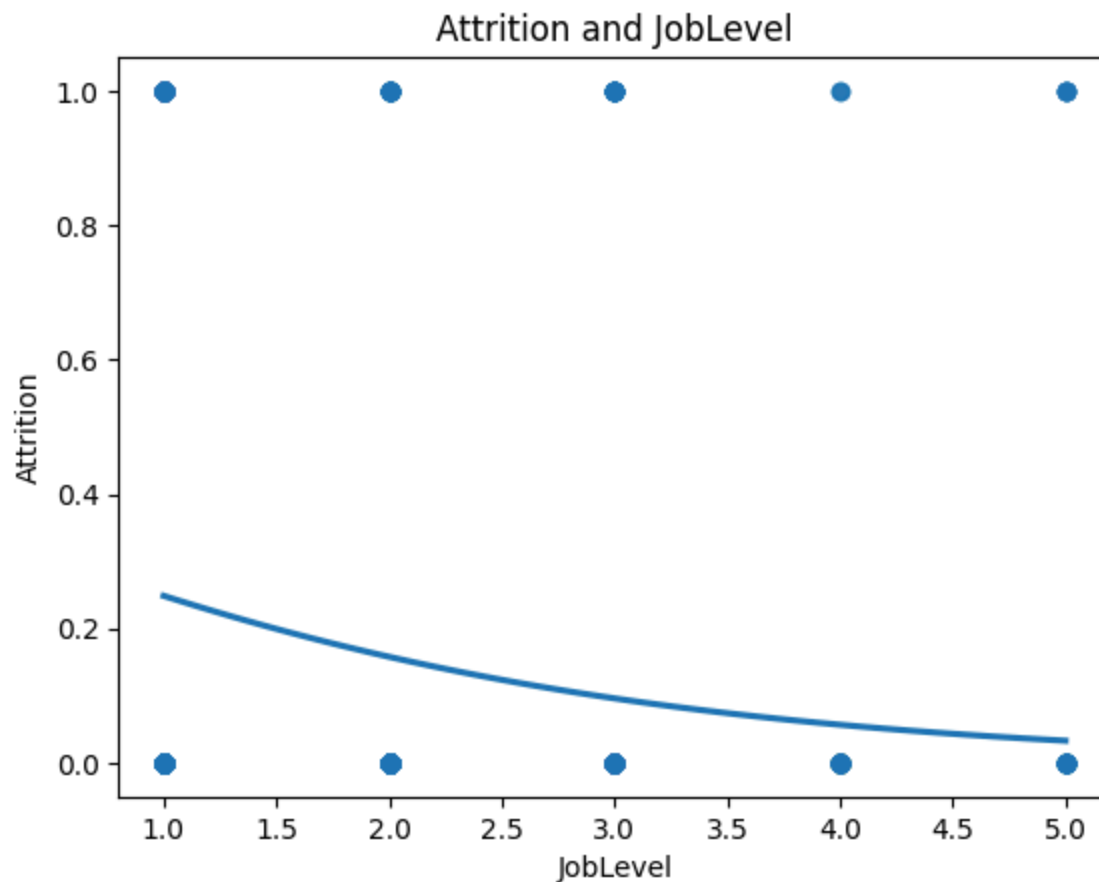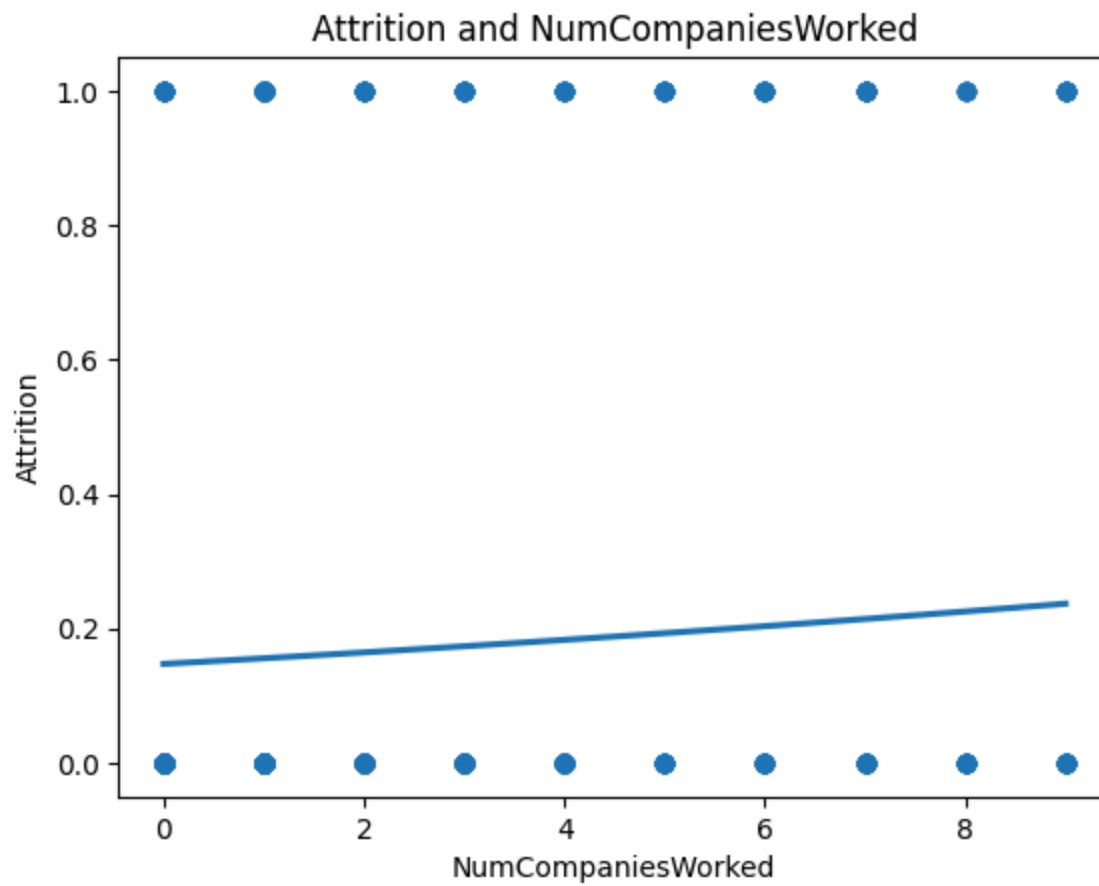
In [280...

## Attrition and TotalWorkingYears



In [281... 
```python
sns.regplot(x=X_train[least_correlated_feature2], y=y_train, data=df, logistic=True, ci=None).set(title=
    'Attrition and {}'.format(least_correlated_feature2));

# Save the plot as a JPG image with the title as the filename
filename = f"images/Attrition and Years in Current Role.jpg"
plt.savefig(filename)
```
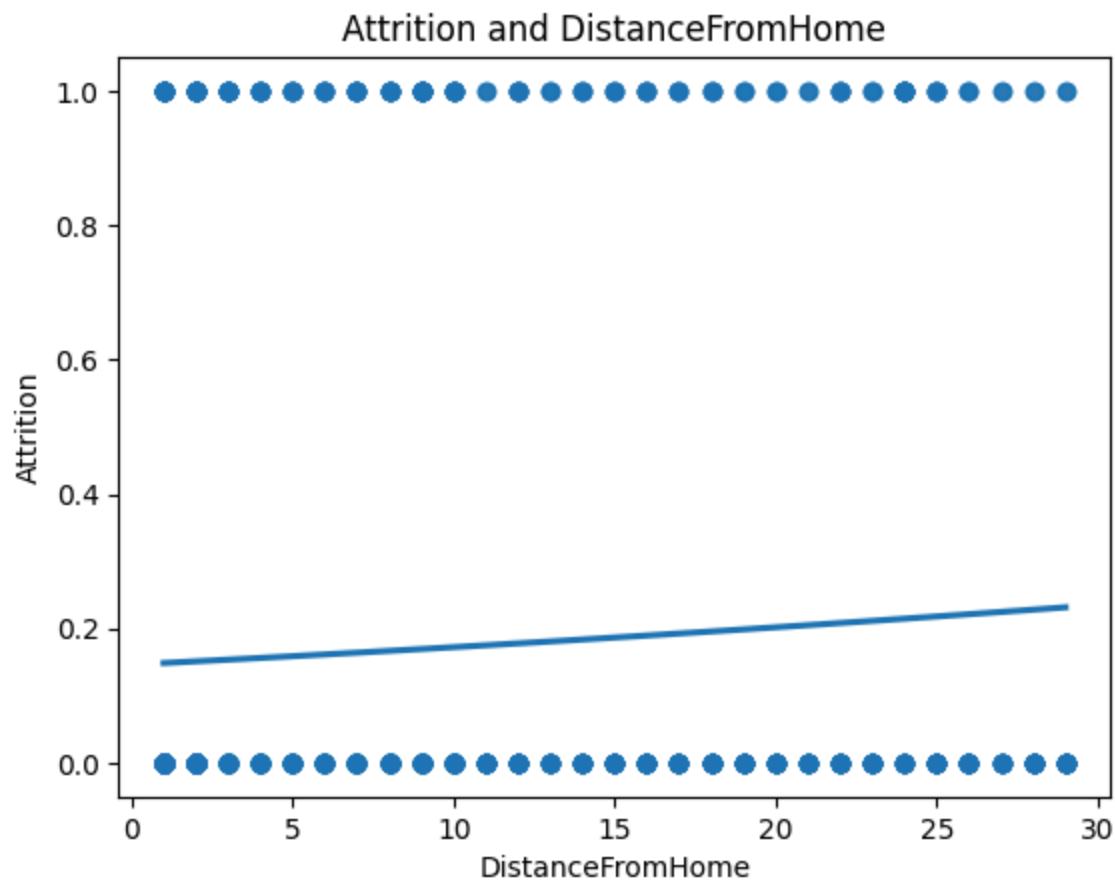
## Attrition and YearsInCurrentRole



```python
sns.regplot(x=X_train[least_correlated_feature3], y=y_train, data=df, logistic=True, ci=None).set(title=
    'Attrition and {}'.format(least_correlated_feature3));

# Save the plot as a JPG image with the title as the filename
filename = f"images/Attrition and Job Level.jpg"
plt.savefig(filename)
```

## Attrition and JobLevel



## Comments on Findings:

We can see that as these features increase, attrition decreases. Attrition lowers with higher number of working years, the years in current role, and the job level
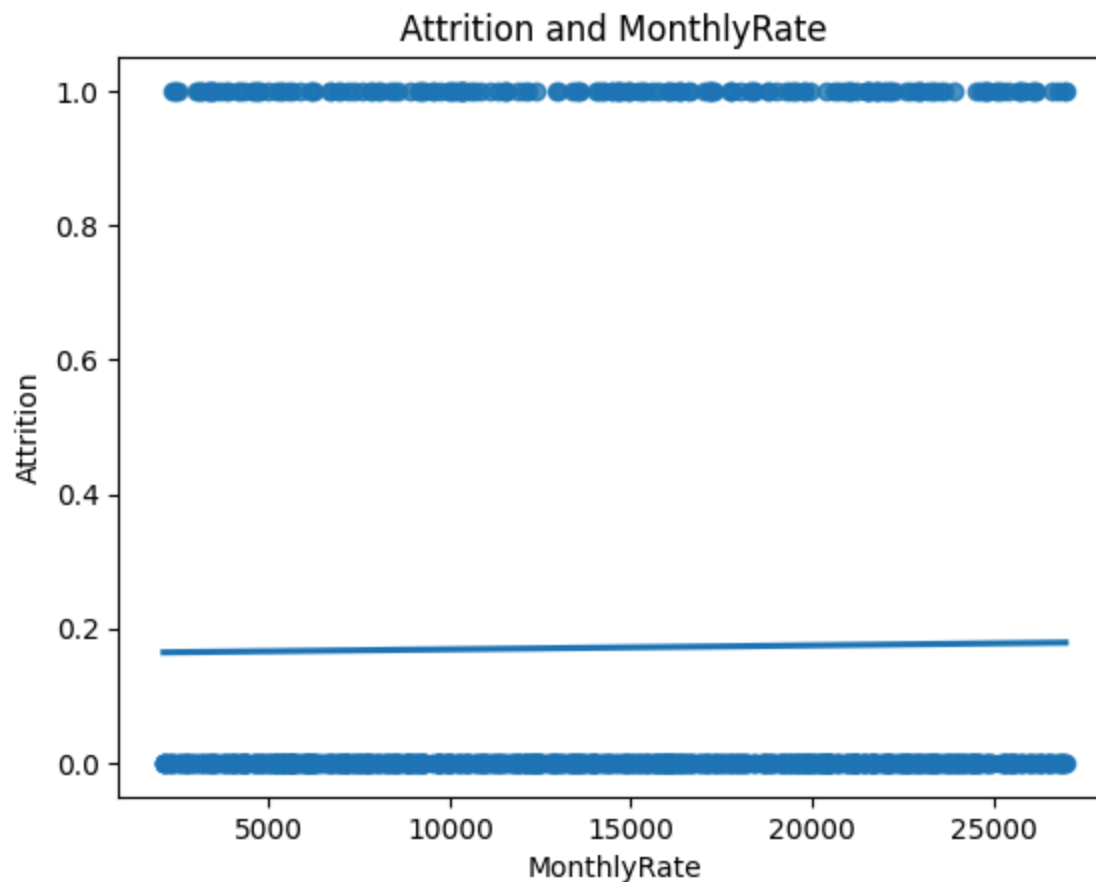
```
In [283...  sns.regplot(x=X_train[most_correlated_feature1], y=y_train, data=df, logistic=True, ci=None).set(title=
               'Attrition and {}'.format(most_correlated_feature1));
```

Attrition and NumCompaniesWorked

```
In [284...    sns.regplot(x=X_train[most_correlated_feature2], y=y_train, data=df, logistic=True, ci=None).set(title=
                  'Attrition and {}'.format(most_correlated_feature2));
```

Attrition and DistanceFromHome

```
sns.regplot(x=X_train[most_correlated_feature3], y=y_train, data=df, logistic=True, ci=None).set(title=
    'Attrition and {}'.format(most_correlated_feature3));
```

## Attrition and MonthlyRate



## Comments on Findings:

We can see that as these features increase, attrition increases. Attrition rises slightly with higher number of companies worked, the distance from home, and the monthly rate.

## Commented Out:

The below cell is an attempt at a model with only the highest correlating features, but this model was poorer than keeping all features.

```
In [286... #Removing high attrition features and keeping the low attrition features
         #running this cell and the following cells resulted in a less accurate model than keeping all features
```

```
# y2 = df.Attrition
# X2 = df[['Age', 'JobInvolvement', 'JobLevel', 'JobSatisfaction', 'MonthlyIncome', 'StockOptionLevel', 'TotalWorking
#         'YearsAtCompany', 'YearsInCurrentRole', 'YearsWithCurrManager']]
# #normalizing data
# X_train, X_test, y_train, y_test = train_test_split(X2, y2, random_state=42)
```

In [287…
```
#Verifying we have the current distribution of attrition/no attrition in training and test data
print("Train percent attrition/no attrition:", y_train.value_counts(normalize=True)[1])
print("Test percent attrition/no attrition: ", y_test.value_counts(normalize=True)[1])
```

```
Train percent attrition/no attrition: 0.17150635208711434
Test percent attrition/no attrition:  0.13043478260869565
```

### Comment on Findings

A reasonably close distribution of y values between train and test data.

## Dummy Variables for Categorical Features

Below We create dummy variables for the dataset's categorical features, and concatenate them together again.

In [288…
```
#Separating X data into continuous vs. categorical
X_train_cont = X_train.select_dtypes(include='number')
X_test_cont = X_test.select_dtypes(include='number')
X_train_cat = X_train.select_dtypes(exclude='number')
X_test_cat = X_test.select_dtypes(exclude='number')

#Dummy encoding categorical values using OneHotEncoder
ohe = OneHotEncoder(handle_unknown='ignore')
X_train_cat = ohe.fit_transform(X_train_cat)
X_test_cat = ohe.transform(X_test_cat)

X_train_preprocessed = np.concatenate([X_train_cont, X_train_cat.todense()], axis=1)
X_test_preprocessed = np.concatenate([X_test_cont, X_test_cat.todense()], axis=1)
```

In [289…
```
X_train = X_train_preprocessed
X_test = X_test_preprocessed
```

### Commented Out:

The below two cells is an attempt at a model with removal of the features with the highest P values, but this model was poorer than keeping all features. Here, We realized, that it makes sense why removing the highest or lowest correlating features might negatively impact the accuracy of the model. If the model doesn't have the lowest correlating features to work with, how could it determine the highest ones?

In [290...
```
# Building a logistic regression model using statsmodels
#log_reg = sm.Logit(y_train, X_train).fit()
#print(log_reg.summary())
```

In [291...
```
#Dropping features with high P Values
#This resulted in a model that was less accurate than including all features.

# df2 = df.drop(columns=['DailyRate', 'Education', 'EmployeeCount', "EmployeeNumber",
#                        'HourlyRate', 'JobLevel','MonthlyIncome', 'MonthlyRate',
#                        'PercentSalaryHike', 'PerformanceRating', 'StandardHours',
#                        'StockOptionLevel', 'BusinessTravel', 'EducationField',
#                        'MaritalStatus', 'OverTime'])

#dummy_dataframe = pd.get_dummies(df2)
# y3 = dummy_dataframe.Attrition
# X3 = dummy_dataframe.drop('Attrition', axis = 1)
# X_train, X_test, y_train, y_test = train_test_split(X3, y3, random_state=42)
```

## Initial Model

### Comments:

The initial model. We used the F1 Score and a confusion matrix to see the baseline for this first model, particularly because the accuracy score wouldn't be as helpful with such an unbalanced dataset.

In [292...
```
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import cross_val_score
from sklearn.metrics import accuracy_score, precision_score, recall_score, roc_auc_score

baseline_model = LogisticRegression(solver='newton-cg', class_weight='balanced')
X_train = np.asarray(X_train)
y_train = np.asarray(y_train)
```

```python
baseline_model.fit(X_train, y_train)

# Predicting on the test data
X_test = np.asarray(X_test)
pred_test = baseline_model.predict(X_test)

# Calculating and printing the f1 score
f1_test = f1_score(y_test, pred_test)
print('The f1 score for the testing data:', f1_test)

# Additional Metrics
accuracy_test = accuracy_score(y_test, pred_test)
precision_test = precision_score(y_test, pred_test)
recall_test = recall_score(y_test, pred_test)
roc_auc_test = roc_auc_score(y_test, baseline_model.predict_proba(X_test)[:, 1])

print(f'Accuracy: {accuracy_test:.3f}')
print(f'Precision: {precision_test:.3f}')
print(f'Recall: {recall_test:.3f}')
print(f'ROC-AUC Score: {roc_auc_test:.3f}')

# Creating a function for the confusion matrix
def conf_matrix(y_true, y_pred):
    cm = {'TP': 0, 'TN': 0, 'FP': 0, 'FN': 0}

    for ind, label in enumerate(y_true):
        pred = y_pred[ind]
        if label == 1:
            # CASE: TP
            if label == pred:
                cm['TP'] += 1
            # CASE: FN
            else:
                cm['FN'] += 1
        else:
            # CASE: TN
            if label == pred:
                cm['TN'] += 1
            # CASE: FP
            else:
                cm['FP'] += 1
    return cm
```

```
conf_matrix(y_test, pred_test)
```

The f1 score for the testing data: 0.3684210526315789
Accuracy: 0.739
Precision: 0.269
Recall: 0.583
ROC-AUC Score: 0.795

Out[292...    {'TP': 28, 'TN': 244, 'FP': 76, 'FN': 20}

SVM Implementation This cell trains a Support Vector Machine (SVM) classifier using a linear kernel and balanced class weights. It predicts test labels and computes evaluation metrics like F1 score, accuracy, precision, recall, and ROC-AUC.

In [293...
```python
# from sklearn.svm import SVC
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import f1_score, accuracy_score, precision_score, recall_score, roc_auc_score

# Feature Scaling
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)


# Initialize and fit the SVM model
svm_model = SVC(kernel='linear', probability=True, class_weight='balanced')
svm_model.fit(X_train, y_train)

# Predictions
svm_pred_test = svm_model.predict(X_test)

# Metrics
from sklearn.metrics import f1_score, accuracy_score, precision_score, recall_score, roc_auc_score
svm_f1_test = f1_score(y_test, svm_pred_test)
svm_accuracy_test = accuracy_score(y_test, svm_pred_test)
svm_precision_test = precision_score(y_test, svm_pred_test)
svm_recall_test = recall_score(y_test, svm_pred_test)
svm_roc_auc_test = roc_auc_score(y_test, svm_model.predict_proba(X_test)[:, 1])

print(f"SVM F1 Score: {svm_f1_test:.3f}")
print(f"SVM Accuracy: {svm_accuracy_test:.3f}")
print(f"SVM Precision: {svm_precision_test:.3f}")
```

```
print(f"SVM Recall: {svm_recall_test:.3f}")
print(f"SVM ROC-AUC Score: {svm_roc_auc_test:.3f}")
```

```
SVM F1 Score: 0.395
SVM Accuracy: 0.750
SVM Precision: 0.288
SVM Recall: 0.625
SVM ROC-AUC Score: 0.789
```

Decision Tree Implementation This cell trains a Decision Tree Classifier with balanced class weights. It predicts test labels and evaluates the model using the same metrics as the SVM.

In [294...
```python
from sklearn.tree import DecisionTreeClassifier

# Initialize and fit the Decision Tree model
dt_model = DecisionTreeClassifier(class_weight='balanced')
dt_model.fit(X_train, y_train)

# Predictions
dt_pred_test = dt_model.predict(X_test)

# Metrics
dt_f1_test = f1_score(y_test, dt_pred_test)
dt_accuracy_test = accuracy_score(y_test, dt_pred_test)
dt_precision_test = precision_score(y_test, dt_pred_test)
dt_recall_test = recall_score(y_test, dt_pred_test)
dt_roc_auc_test = roc_auc_score(y_test, dt_model.predict_proba(X_test)[:, 1])

print(f"Decision Tree F1 Score: {dt_f1_test:.3f}")
print(f"Decision Tree Accuracy: {dt_accuracy_test:.3f}")
print(f"Decision Tree Precision: {dt_precision_test:.3f}")
print(f"Decision Tree Recall: {dt_recall_test:.3f}")
print(f"Decision Tree ROC-AUC Score: {dt_roc_auc_test:.3f}")
```

```
Decision Tree F1 Score: 0.278
Decision Tree Accuracy: 0.788
Decision Tree Precision: 0.250
Decision Tree Recall: 0.312
Decision Tree ROC-AUC Score: 0.586
```

Random Forest Implementation This cell trains a Random Forest Classifier with balanced class weights and 100 decision trees. Predictions and evaluations are made using F1 score, accuracy, precision, recall, and ROC-AUC metrics.

```
In [295...   from sklearn.ensemble import RandomForestClassifier

             # Initialize and fit the Random Forest model
             rf_model = RandomForestClassifier(class_weight='balanced', n_estimators=100, random_state=42)
             rf_model.fit(X_train, y_train)

             # Predictions
             rf_pred_test = rf_model.predict(X_test)

             # Metrics
             rf_f1_test = f1_score(y_test, rf_pred_test)
             rf_accuracy_test = accuracy_score(y_test, rf_pred_test)
             rf_precision_test = precision_score(y_test, rf_pred_test)
             rf_recall_test = recall_score(y_test, rf_pred_test)
             rf_roc_auc_test = roc_auc_score(y_test, rf_model.predict_proba(X_test)[:, 1])

             print(f"Random Forest F1 Score: {rf_f1_test:.3f}")
             print(f"Random Forest Accuracy: {rf_accuracy_test:.3f}")
             print(f"Random Forest Precision: {rf_precision_test:.3f}")
             print(f"Random Forest Recall: {rf_recall_test:.3f}")
             print(f"Random Forest ROC-AUC Score: {rf_roc_auc_test:.3f}")
```

```
Random Forest F1 Score: 0.148
Random Forest Accuracy: 0.875
Random Forest Precision: 0.667
Random Forest Recall: 0.083
Random Forest ROC-AUC Score: 0.756
```

# 5. Model Evaluation

## Evaluation Metrics

To assess model performance, several evaluation metrics were computed, including accuracy, F1 score, precision, recall, and ROC-AUC. These metrics provide a comprehensive view of how well each model performs in predicting employee attrition, considering both overall performance and the ability to correctly identify positive and negative cases.

# Comparison

The performance of all models was compared to identify the most effective one for predicting employee attrition. Here's a summary of the results:

- **Logistic Regression:** The model demonstrated a relatively high accuracy (0.739) but had a lower F1 score (0.368), indicating that it struggles to balance precision and recall. The precision (0.269) is low, which suggests a higher rate of false positives, while recall (0.583) shows better performance in identifying actual positives.

- **SVM (Support Vector Machine):** This model showed competitive performance with an F1 score of 0.395 and a ROC-AUC score of 0.790, effectively balancing precision (0.288) and recall (0.625). It outperformed logistic regression in terms of F1 score and ROC-AUC, making it a strong candidate for further optimization.

- **Decision Tree:** The decision tree had a relatively high accuracy (0.788), but its F1 score (0.278) and recall (0.312) were lower, showing weaker performance in identifying true positives. Precision was also low (0.250), suggesting that it often misclassifies positive cases as negatives.

- **Random Forest:** Although the random forest had the highest accuracy (0.875), its F1 score (0.148) was the lowest, indicating poor performance in balancing precision and recall. The high precision (0.667) was offset by a very low recall (0.083), suggesting that it fails to identify many actual positive cases despite correctly classifying some negatives.

# Insights

Among all the models, **SVM** showed the most competitive performance, with an F1 score of 0.395 and a ROC-AUC score of 0.790. This model strikes a better balance between precision and recall compared to the other models, making it a strong contender for predicting employee attrition effectively. However, further fine-tuning could still improve these results.