

IBM MACHINE LEARNING

IBM EXPLORATORY DATA ANALYSIS FOR EMPLOYEE ATTRITION MODEL



Fatima, Sayeda

8/17/2022

Table of Contents

1) Project Overview	2
2) About the Dataset	3
2a) Brief description of the data set you chose:.....	3
2b) Summary of Data Attributes	3
2c) Main Objectives of Analysis	3
3) Data Exploration, Data Cleansing and Features Engineering	7
3a) Data Exploration:	4
3b) Data Cleansing Actions:.....	7
3c) Features Engineering.....	9
4) Summary of Training Different Classifier Models	Error! Bookmark not defined.
4a) Machine Learning Algorithm Approaches	Error! Bookmark not defined.
I) Data Level Approaches:	Error! Bookmark not defined.
II) Algorithm Ensemble Approach:	Error! Bookmark not defined.
4b) Summarizing Employed Models	Error! Bookmark not defined.
1) Logistic Regression (LR) Models:	Error! Bookmark not defined.
2) Random Forest Models:.....	Error! Bookmark not defined.
3) XGB Model.....	Error! Bookmark not defined.
5) Recommended Model.....	Error! Bookmark not defined.
5a) Result Summary	Error! Bookmark not defined.
5b) Overall Visual Summary	Error! Bookmark not defined.
5c) Individual Model Visual Summary	Error! Bookmark not defined.
5d) Model Choice and Justification.....	Error! Bookmark not defined.
6) Summary Key Findings and Insights	Error! Bookmark not defined.
6a) Summarizing Model Drivers:	Error! Bookmark not defined.
6b) Enlisting Top Contributory Factors	Error! Bookmark not defined.
6c) Visualizing Top Contributory Factors to Employee Attrition	Error! Bookmark not defined.
7) Link to Other Useful Models	22
8) Github Link to Assignment Notebook	22

1) Project Overview

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

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

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

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

2) About the Dataset

2a) Brief description of chosen data set:

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

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

2b) Summary of Data Attributes

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

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

2c) Main Objectives of Analysis

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

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

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

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

3) Initial Plan for Data Exploration

3a) Data Exploration

- Data was first loaded into pandas dataframe

Load & Read Dataset

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

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EmployeeCount	EmployeeNumber	...	Relationship
0	41	Yes	Travel_Rarely	1102	Sales	1	2	Life Sciences	1	1	...	
1	49	No	Travel_Frequently	279	Research & Development	8	1	Life Sciences	1	2	...	
2	37	Yes	Travel_Rarely	1373	Research & Development	2	2	Other	1	4	...	

- Column types were explored

Check data set column types

```
1 df.info()
```

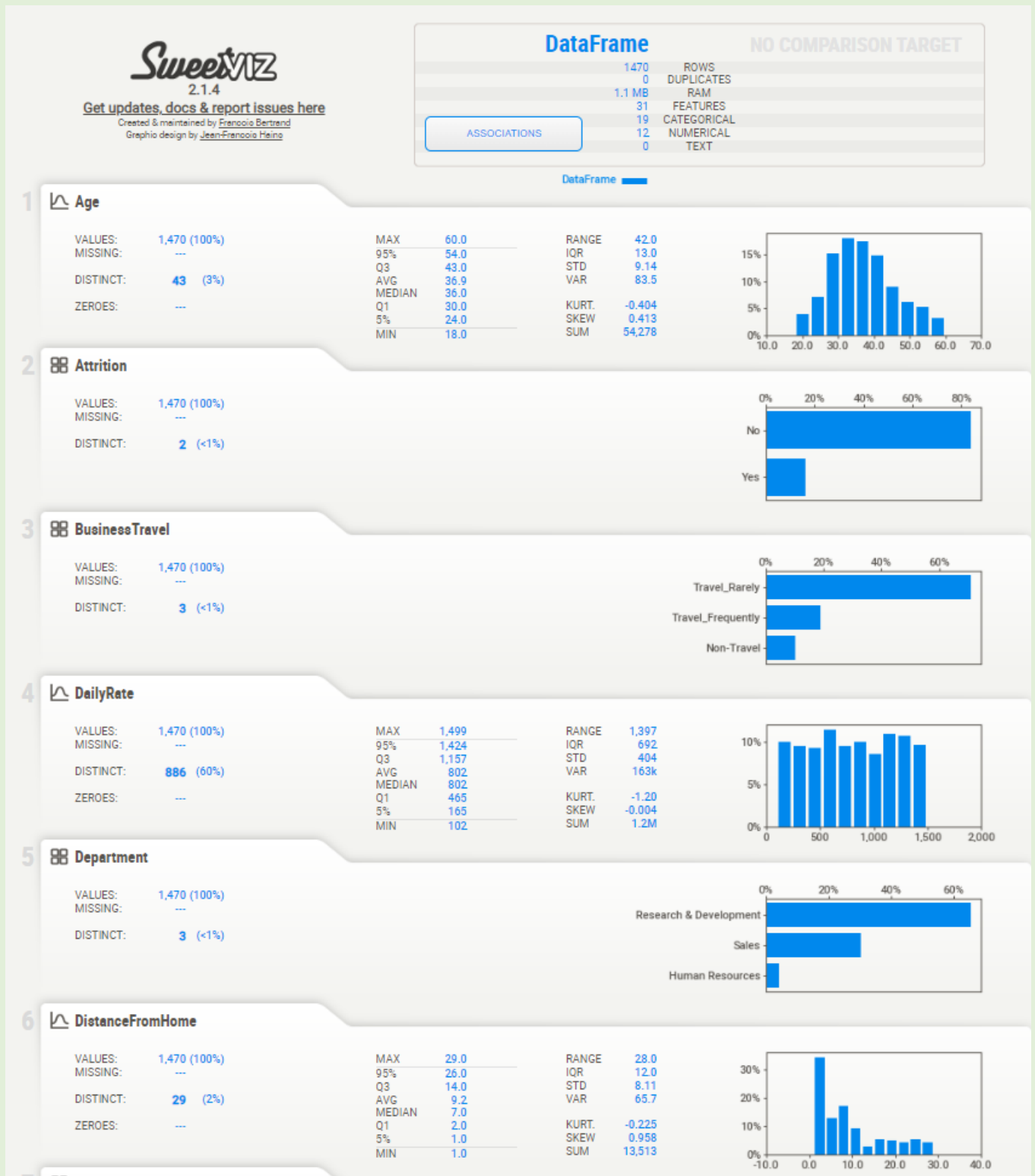
```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 1470 entries, 0 to 1469
```

```
Data columns (total 35 columns):
```

#	Column	Non-Null Count	Dtype
0	Age	1470 non-null	int64
1	Attrition	1470 non-null	object
2	BusinessTravel	1470 non-null	object
3	DailyRate	1470 non-null	int64
4	Department	1470 non-null	object
5	DistanceFromHome	1470 non-null	int64
6	Education	1470 non-null	int64
7	EducationField	1470 non-null	object
8	EmployeeCount	1470 non-null	int64
9	EmployeeNumber	1470 non-null	int64
10	EnvironmentSatisfaction	1470 non-null	int64
11	Gender	1470 non-null	object
12	HourlyRate	1470 non-null	int64
13	JobInvolvement	1470 non-null	int64

- Automated Exploratory Data Analysis was performed using Sweetviz to check



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

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

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

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

4) Actions taken for Data Cleansing and Features Engineering

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

4a) Data Cleansing Actions

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

```
Drop Columns if 90% data is empty

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

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

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

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

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

```

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

```

Data has 0 null values.

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

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

4b) Features Engineering

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

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

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

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

```

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

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

```

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

```

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

```

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

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

2) Data Encoding of Object/String Columns:

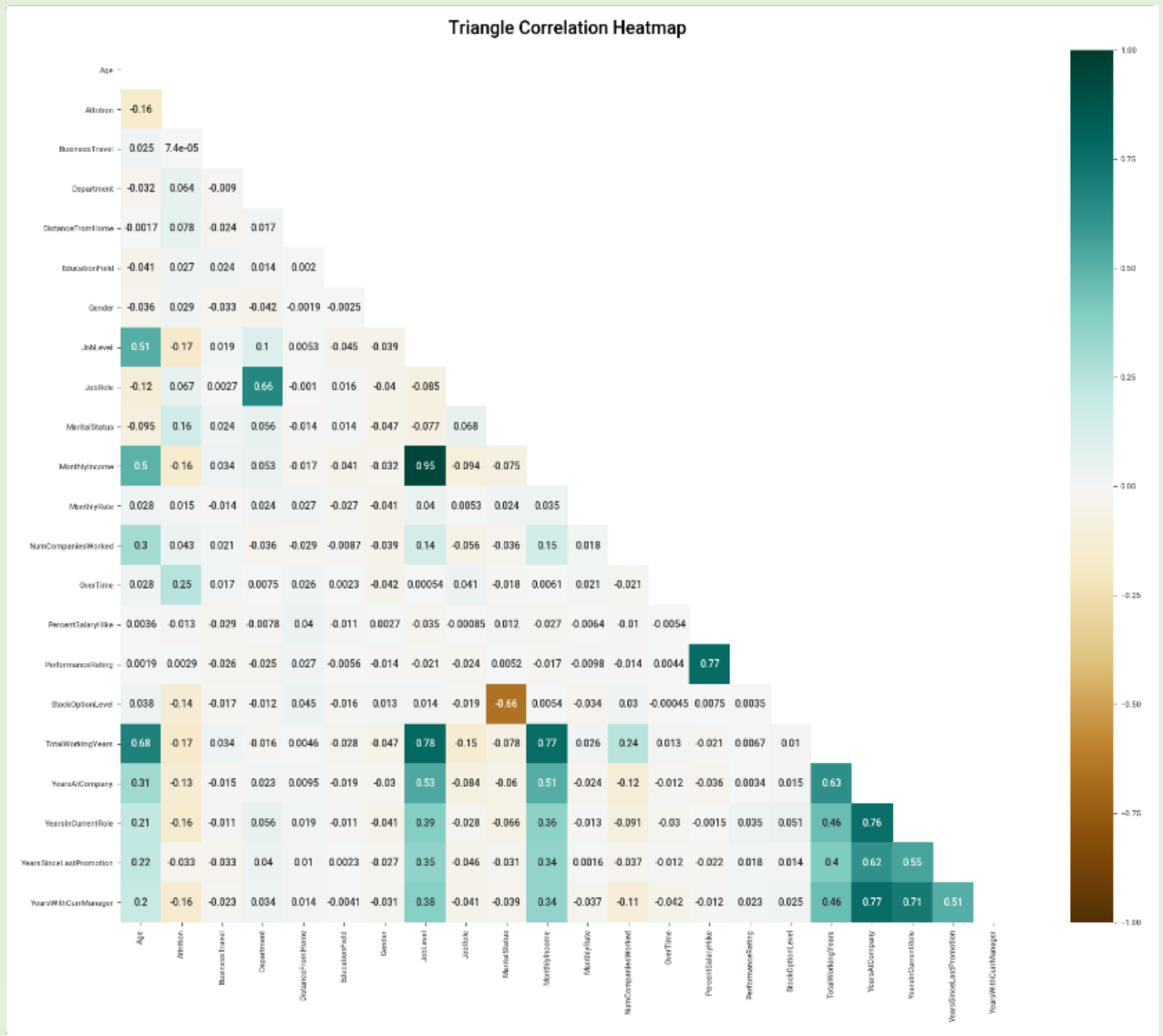
- * List all Object/String Columns
- * Deep copy the original data
- * Create Function to employ Scikit-learn label encoding to encode object data
- * Create a new dataframe with encoded data description to attach to model outcomes

```

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

```

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



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

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

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

```
Int64Index: 1470 entries, 0 to 1469
```

```
Data columns (total 22 columns):
```

#	Column	Non-Null Count	Dtype
0	Age	1470 non-null	int64
1	Attrition	1470 non-null	int32
2	BusinessTravel	1470 non-null	int32
3	Department	1470 non-null	int32
4	DistanceFromHome	1470 non-null	int64
5	EducationField	1470 non-null	int32
6	Gender	1470 non-null	int32
7	JobLevel	1470 non-null	int64
8	JobRole	1470 non-null	int32
9	MaritalStatus	1470 non-null	int32
10	MonthlyIncome	1470 non-null	int64
11	MonthlyRate	1470 non-null	int64
12	NumCompaniesWorked	1470 non-null	int64
13	Overtime	1470 non-null	int32
14	PercentSalaryHike	1470 non-null	int64
15	PerformanceRating	1470 non-null	int64
16	StockOptionLevel	1470 non-null	int64
17	TotalWorkingYears	1470 non-null	int64
18	YearsAtCompany	1470 non-null	int64
19	YearsInCurrentRole	1470 non-null	int64
20	YearsSinceLastPromotion	1470 non-null	int64
21	YearsWithCurrManager	1470 non-null	int64

```
dtypes: int32(8), int64(14)
```

```
memory usage: 218.2 KB
```

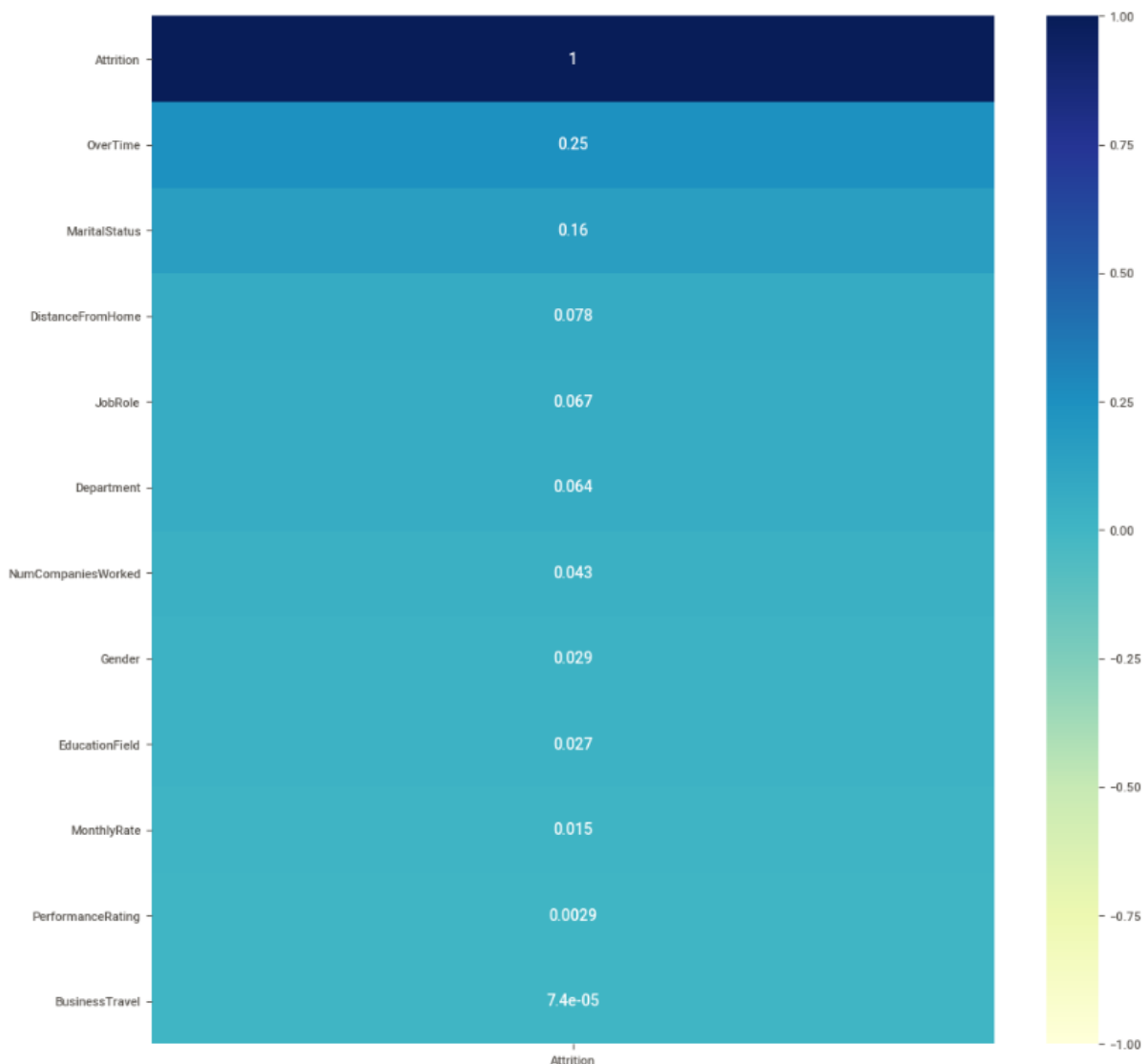
5) Key Findings and Insights, which synthesizes the results of Exploratory Data Analysis in an insightful and actionable manner

5a) Summary of Features Exhibiting Positive Correlation with Target

The following features exhibited positive correlations with target variable:

Positive Correlations between Individual Features and Target Variable

	Features1	Corr
18	OverTime	0.246118
14	MaritalStatus	0.162070
4	DistanceFromHome	0.077924
12	JobRole	0.067151
3	Department	0.063991
17	NumCompaniesWorked	0.043494
8	Gender	0.029453
6	EducationField	0.026846
16	MonthlyRate	0.015170
20	PerformanceRating	0.002889
1	BusinessTravel	0.000074



5b) Summary of Positive Correlations for Feature Pairs

In addition, the following feature pairs displayed high correlation

Correlation Matrix Results for Feature Pairs

	Features2	Features3	Correlation_abs
31	MonthlyIncome	JobLevel	0.950300
33	TotalWorkingYears	JobLevel	0.782208
35	PerformanceRating	PercentSalaryHike	0.773550
37	TotalWorkingYears	MonthlyIncome	0.772893
39	YearsWithCurrManager	YearsAtCompany	0.769212
41	YearsAtCompany	YearsInCurrentRole	0.758754
43	YearsWithCurrManager	YearsInCurrentRole	0.714365
45	Age	TotalWorkingYears	0.680381
47	StockOptionLevel	MaritalStatus	0.662577
49	JobRole	Department	0.662431
51	TotalWorkingYears	YearsAtCompany	0.628133
53	YearsAtCompany	YearsSinceLastPromotion	0.618409

Based on the above two correlation findings, the following final features were selected while others were dropped as they were considered to bear no or less impacts towards employee attrition.

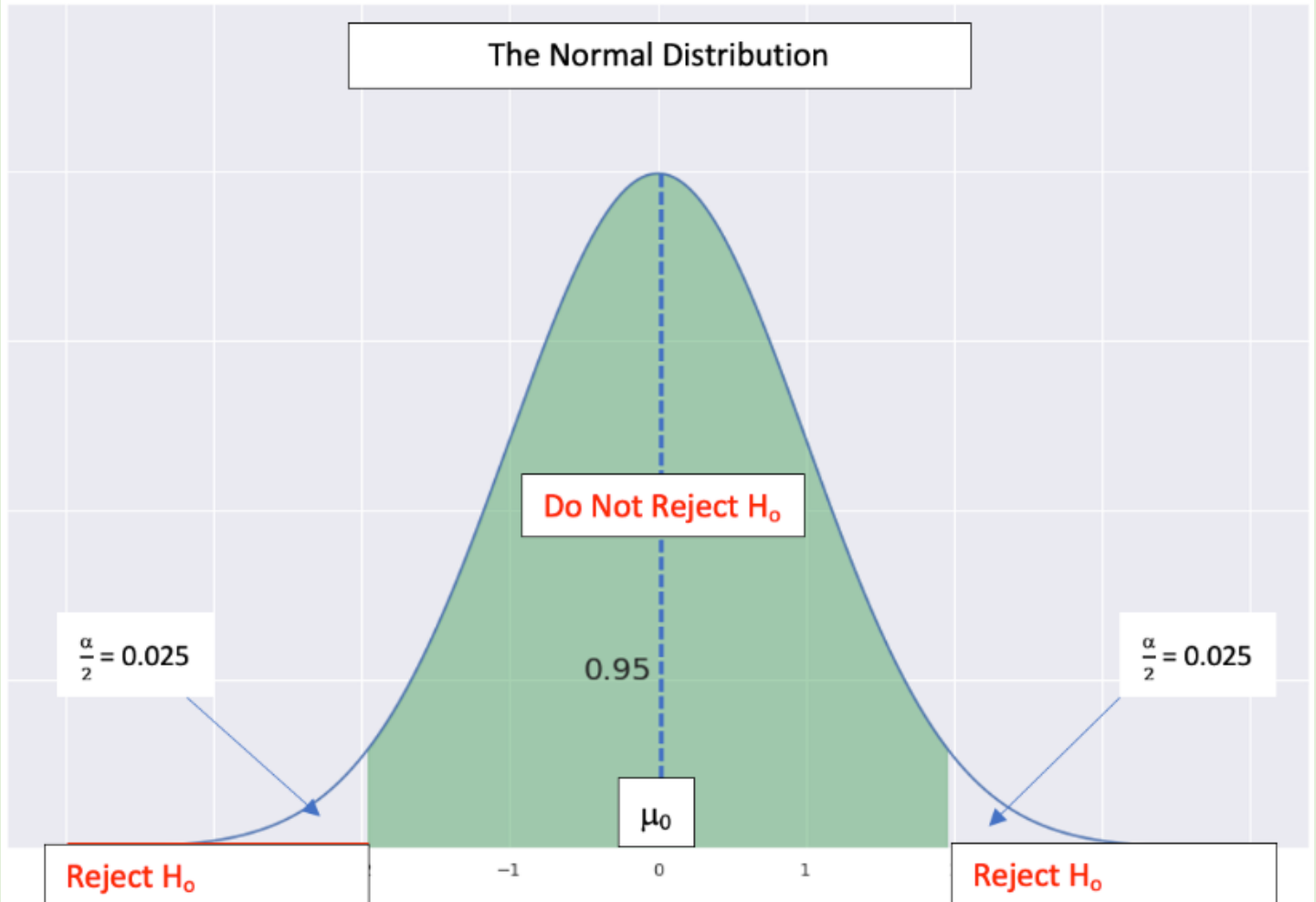
List of Final Features

```
['Age' 'Attrition' 'BusinessTravel' 'Department' 'DistanceFromHome'
 'EducationField' 'Gender' 'JobLevel' 'JobRole' 'MaritalStatus'
 'MonthlyIncome' 'MonthlyRate' 'NumCompaniesWorked' 'OverTime'
 'PercentSalaryHike' 'PerformanceRating' 'StockOptionLevel'
 'TotalWorkingYears' 'YearsAtCompany' 'YearsInCurrentRole'
 'YearsSinceLastPromotion' 'YearsWithCurrManager']
```

6) Three Major Hypothesis

Null hypothesis (H_0) is a statistical hypothesis which postulates random factors causing difference in observations.

Alternative hypothesis (H_A) is a statistical hypothesis which postulates real impacts causing difference in observations.



Based on the above diagram, the significance level is the decision point for null hypothesis acceptance or vice versa.

This significance level is normally set at 5%, 1% or 0.5%, depending on business requirements.

Hence, given a 5% significance level, $\alpha(\alpha) = 0.05$.

Consequently, for 2-tailed test, alpha should be divided by 2, which will yield 0.025 for alpha set at 0.05.

Hence, Null Hypothesis will be Rejected if computed p-value is less than alpha and vice versa.

6a) Hypothesizing Gender Differences in Attrition Rates

$H_0: \mu_1 - \mu_2 = 0$ There is No Significant Gender Difference between Employee Attrition of Female and Male Workers.
 $H_A: \mu_1 - \mu_2 \neq 0$ There is Significant Gender Difference between Employee Attrition of Female and Male Workers.

%age Distribution of Gender

	%	Gender
1	60.0	1
0	40.0	0

Value Counts show an imbalanced Class Distribution with 60.0 % in class 1 and 40.0 % in class 0

Mean Attrition for Females is: 0.14795918367346939

Mean Attrition for Males is: 0.17006802721088435

$t_value1 = -1.1289761152328313$

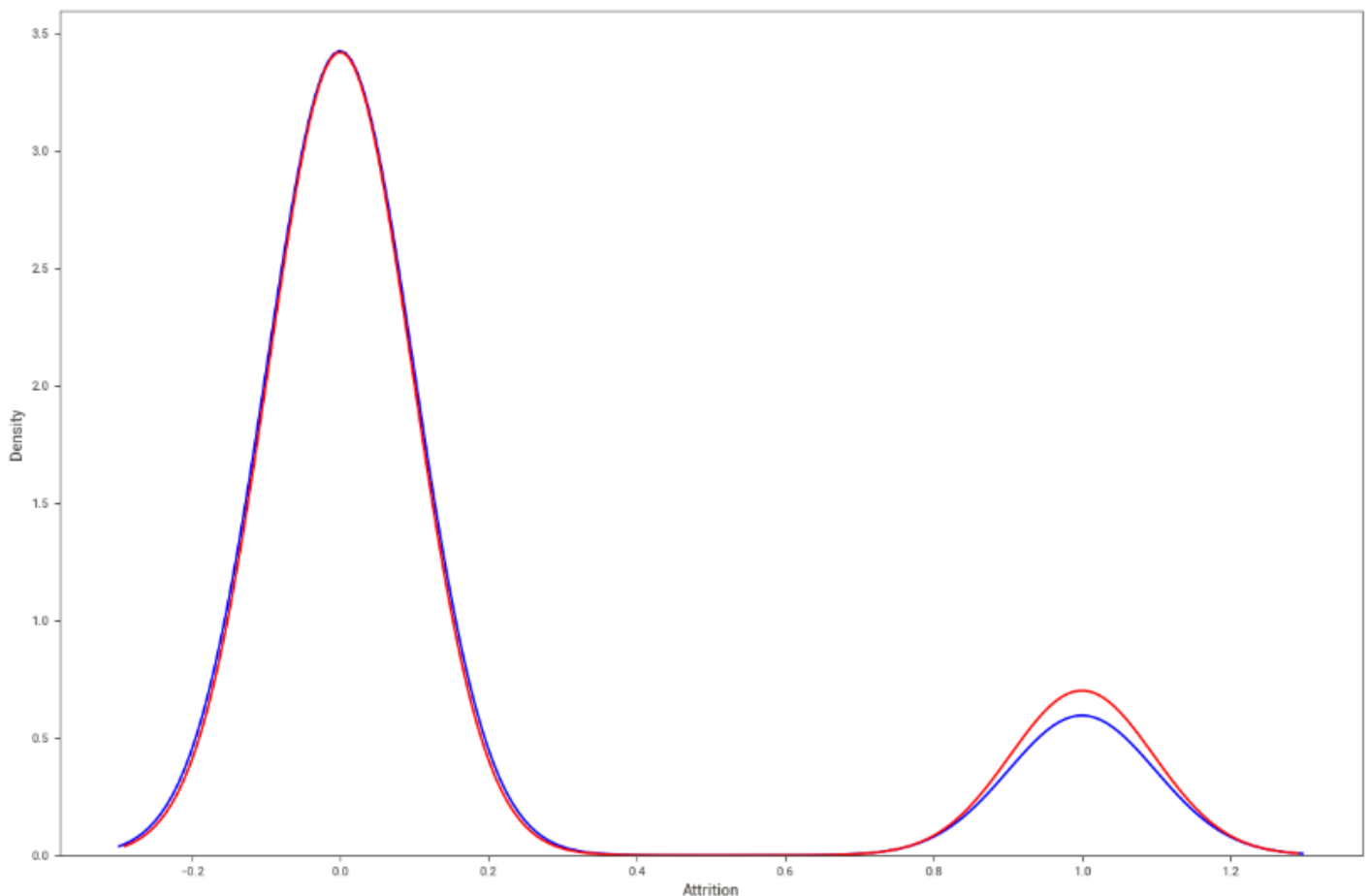
$p_value1 = 0.25909236414147996$

Conclusion:

Since p_value 0.25909236414147996 is greater than alpha 0.05

We Accept the Null Hypothesis that there is No Significant Gender Difference between Attrition Rates of Female and Male Workers.

<AxesSubplot:xlabel='Attrition', ylabel='Density'>



6b) Hypothesizing Differences in Salary Hike between Leavers and Stayers

$H_{02}: \mu_1 - \mu_2 = 0$ Average Salary Hike of Leavers is Less than or Equal to Stayers.

$H_{A2}: \mu_1 - \mu_2 \neq 0$ Average Salary Hike of Leavers is Greater than or Equal to Stayers.

Summary Statistics

AOV Table Results

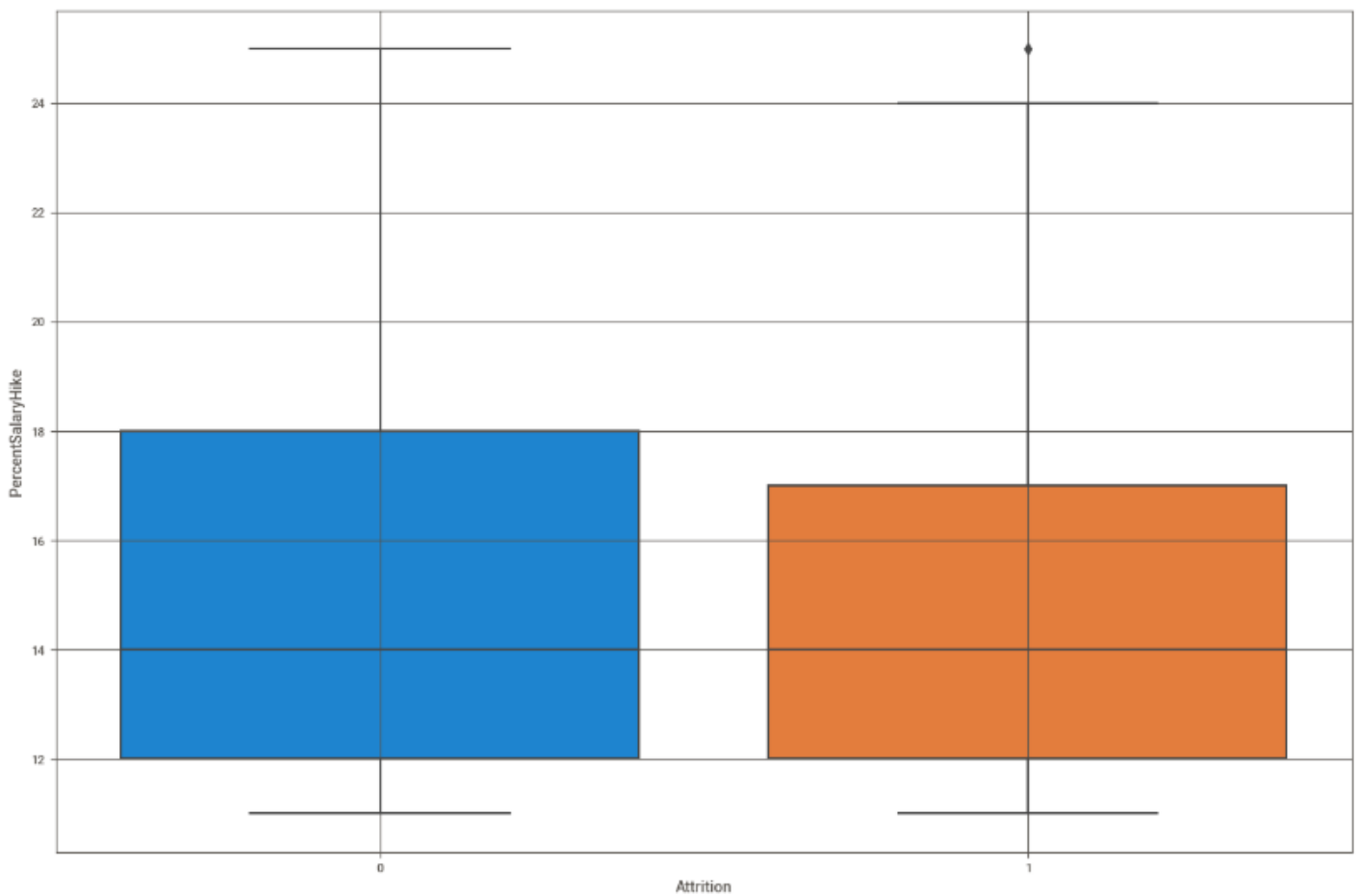
	df	sum_sq	mean_sq	F	PR(>F)
C(PercentSalaryHike)	14.0	1.805708	0.128979	0.952689	0.5006
Residual	1455.0	196.984088	0.135384	NaN	NaN

$p_val_sal = 0.5005995155055496$

Conclusion:

Since p_value 0.5005995155055496 is greater than alpha 0.05

We Accept Null Hypothesis that Average Salary Hike of Leavers is Less than or Equal to Stayers.



6c) Hypothesizing Differences in Salary Hike between Leavers and Stayers

$H_{03}: \mu_1 - \mu_2 = 0$ There is No Significant Difference in Attrition Rate Across Different Departments.
 $H_{A3}: \mu_1 - \mu_2 \neq 0$ There is Significant Difference in Attrition Rate Across Different Departments.

Contingency Table

Attrition	0	1
Department		
0	51	12
1	828	133
2	354	92

Summary Statistics

Mean Attrition for Department 0 is: 0.19047619047619047
 Mean Attrition for Department 1 is: 0.1383975026014568
 Mean Attrition for Department 2 is: 0.2062780269058296

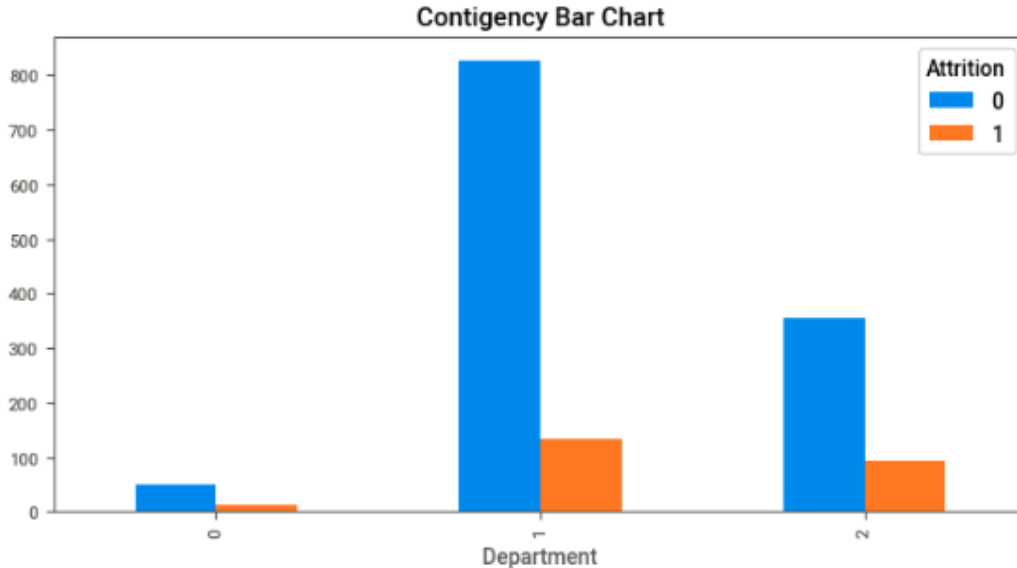
Chi-Square Statistic: 10.79600732241067
 p Value: 0.004525606574479633
 Degree of Freedom: 2
 Expected Frequencies: [[52.84285714 10.15714286]
 [806.06326531 154.93673469]
 [374.09387755 71.90612245]]

Conclusion:

Since p_value 0.004525606574479633 is less than alpha 0.05

We Reject the Null Hypothesis that there is No Significant Difference in Attrition Rate Across Different Departments.

Text(0.5, 1.0, 'Contingency Bar Chart')



7) Conducting a formal significance test for one of the hypotheses and discuss the results

7a) Significance Test for Hypothesis 1

Hypothesis 1: Hypothesizing Gender Differences in Attrition Rates.

Due to unknown standard deviation, a one-tailed t-test has been used for testing population means between female and male Genders. Since, it is a one-tailed test, at 5% significance level, 0.05 alpha (α) has been used.

$H_0: \mu_1 - \mu_2 = 0$ There is No Significant Gender Difference between Employee Attrition of Female and Male Workers.
 $H_A: \mu_1 - \mu_2 \neq 0$ There is Significant Gender Difference between Employee Attrition of Female and Male Workers.

%age Distribution of Gender

	%	Gender
1	60.0	1
0	40.0	0

Value Counts show an imbalanced Class Distribution with 60.0 % in class 1 and 40.0 % in class 0

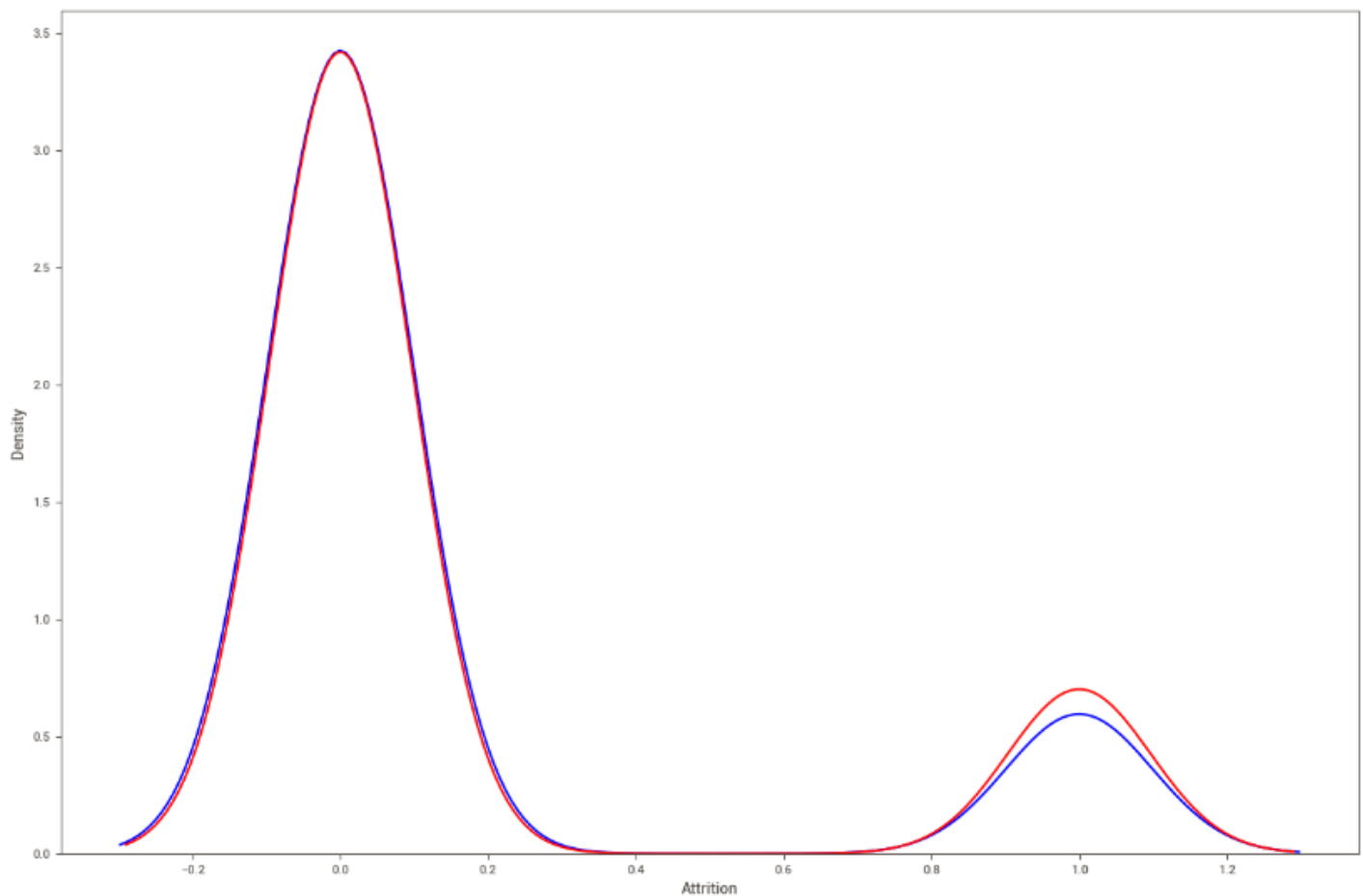
Mean Attrition for Females is: 0.14795918367346939

Mean Attrition for Males is: 0.17006802721088435

t_value1 = -1.1289761152328313

p_value1 = 0.25909236414147996

<AxesSubplot:xlabel='Attrition', ylabel='Density'>



7b) Results and Discussion

Results:

Since p_value 0.25909236414147996 is greater than alpha 0.05

We Accept the Null Hypothesis that there is No Significant Gender Difference between Attrition Rates of Female and Male Workers.

8) Suggestions for next steps in analysing this data

- EDA now should proceed towards analysing problematic data outliers for all final features to increase statistical significance of the model.
- Using Gower Distancing, further cluster analysis can be done to gain in-depth understanding of employee clusters at risk of turnover rather than utilizing one size fit all approach (see, Section 10(e)).

9) A paragraph that summarizes the quality of this data set and request for additional data if needed

Overall data quality was quite good since EDA exhibited zero null values and contained extensive variables to work on. However, while the data did contain extensive parameters, nevertheless, including other pull factors like community fit, workload, etc. may prove to be useful. For example, based on social relationship theory, pull factors like community fit, industry, etc has been found to exhibit a negative correlation with intentions to leave, especially with increased perceptions of needs fulfilment and social networking (see, Ramesh and Gelfand, 2010). Hence, including this information can shed further light on employee turnover. Then, including other mooring factors, such as, personal life involvement may also improve data quality as well as model's predictive power since these may also serve as underlying factors in employee turnover. For instance, numerous studies found higher turnover rates among employees exhibiting high family centrality and work interference with family life (see, Bagger et al., 2008; Haldorai, et al., 2019). Hence, provision of these additional parameters may render enhanced insight into employee attrition and may even change predictive outcomes.

10) Link to Other Useful Models

- a) <https://github.com/IBM/employee-attrition-aif360/blob/master/notebooks/employee-attrition.ipynb>
- b) https://github.com/JNYH/employee_attrition/blob/master/employee_attrition.ipynb
- c) <https://github.com/elastic/examples/tree/master/Machine%20Learning/Analytics%20Jupyter%20Notebooks>
- d) https://github.com/ganesh10-india/HR_Analytics-Employee_Attrition-Classification-Models/blob/main/HR_Analytics_Employee_Attrition_Classification_Models.ipynb
- e) <https://www.adam-d-mckinnon.com/posts/2020-08-04-clusteranalysis/>

11) Github Link to Assignment Notebook

<https://github.com/FATIMASP/IBM-MACHINE-LEARNING-CERTIFICATION/blob/main/Exploratory%20Data%20Analysis%20for%20Machine%20Learning/EDA%20SUPERVISED%20CLASSIFICATION%20EMPLOYEE%20ATTRITION.ipynb>