

```
In [1]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
```

```
In [2]: 1 filepath = 'C:/Users/HP/Documents/python_one_campus/Assignment/'
2 filename = 'WA_Fn-UseC_-HR-Employee-Attrition.csv'
3
4 df = pd.read_csv(filepath + filename)
```

## Exploratory Analysis

```
In [3]: 1 df.head() #checking the data
```

Out[3]:

|   | Age | Attrition | BusinessTravel    | DailyRate | Department             | DistanceFromHome | Education | Edinati |        |
|---|-----|-----------|-------------------|-----------|------------------------|------------------|-----------|---------|--------|
| 0 | 41  | Yes       | Travel_Rarely     | 1102      | Sales                  |                  | 1         | 2       | Life S |
| 1 | 49  | No        | Travel_Frequently | 279       | Research & Development |                  | 8         | 1       | Life S |
| 2 | 37  | Yes       | Travel_Rarely     | 1373      | Research & Development |                  | 2         | 2       |        |
| 3 | 33  | No        | Travel_Frequently | 1392      | Research & Development |                  | 3         | 4       | Life S |
| 4 | 27  | No        | Travel_Rarely     | 591       | Research & Development |                  | 2         | 1       |        |

5 rows × 35 columns



```
In [4]: 1 df.describe()
```

Out[4]:

|       | Age         | DailyRate   | DistanceFromHome | Education   | EmployeeCount | EmployeeNumb |
|-------|-------------|-------------|------------------|-------------|---------------|--------------|
| count | 1470.000000 | 1470.000000 | 1470.000000      | 1470.000000 | 1470.0        | 1470.0000    |
| mean  | 36.923810   | 802.485714  | 9.192517         | 2.912925    | 1.0           | 1024.8653    |
| std   | 9.135373    | 403.509100  | 8.106864         | 1.024165    | 0.0           | 602.0243     |
| min   | 18.000000   | 102.000000  | 1.000000         | 1.000000    | 1.0           | 1.0000       |
| 25%   | 30.000000   | 465.000000  | 2.000000         | 2.000000    | 1.0           | 491.2500     |
| 50%   | 36.000000   | 802.000000  | 7.000000         | 3.000000    | 1.0           | 1020.5000    |
| 75%   | 43.000000   | 1157.000000 | 14.000000        | 4.000000    | 1.0           | 1555.7500    |
| max   | 60.000000   | 1499.000000 | 29.000000        | 5.000000    | 1.0           | 2068.0000    |

8 rows × 26 columns



In [5]: 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  
 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
```

In [6]: 1 print(df.nunique())

```
Age                      43
Attrition                 2
BusinessTravel              3
DailyRate                  886
Department                  3
DistanceFromHome             29
Education                   5
EducationField                6
EmployeeCount                  1
EmployeeNumber                1470
EnvironmentSatisfaction            4
Gender                     2
HourlyRate                  71
JobInvolvement                 4
JobLevel                     5
JobRole                      9
JobSatisfaction                 4
MaritalStatus                  3
MonthlyIncome                 1349
MonthlyRate                  1427
NumCompaniesWorked             10
Over18                      1
OverTime                     2
PercentSalaryHike                15
PerformanceRating                 2
RelationshipSatisfaction            4
StandardHours                  1
StockOptionLevel                 4
TotalWorkingYears                 40
TrainingTimesLastYear                7
WorkLifeBalance                  4
YearsAtCompany                  37
YearsInCurrentRole                 19
YearsSinceLastPromotion                16
YearsWithCurrManager                 18
dtype: int64
```

In [8]: 1 df.columns

```
Out[8]: Index(['Age', 'Attrition', 'BusinessTravel', 'DailyRate', 'Department',
   'DistanceFromHome', 'Education', 'EducationField', 'EmployeeCount',
   'EmployeeNumber', 'EnvironmentSatisfaction', 'Gender', 'HourlyRate',
   'JobInvolvement', 'JobLevel', 'JobRole', 'JobSatisfaction',
   'MaritalStatus', 'MonthlyIncome', 'MonthlyRate', 'NumCompaniesWorked',
   'Over18', 'OverTime', 'PercentSalaryHike', 'PerformanceRating',
   'RelationshipSatisfaction', 'StandardHours', 'StockOptionLevel',
   'TotalWorkingYears', 'TrainingTimesLastYear', 'WorkLifeBalance',
   'YearsAtCompany', 'YearsInCurrentRole', 'YearsSinceLastPromotion',
   'YearsWithCurrManager'],
  dtype='object')
```

In [9]:

```
1 #removing the columns with one unique value because they are zero-variance p
2 #they do not contain any information for modeling
3 #EmployeeCount, Over18, StandardHours
4 #ea is employee attrition
5 ea = df.loc[:, ['Age', 'BusinessTravel', 'DailyRate', 'Department',
6                 'DistanceFromHome', 'Education', 'EducationField',
7                 'EmployeeNumber', 'EnvironmentSatisfaction', 'Gender', 'HourlyRate',
8                 'JobInvolvement', 'JobLevel', 'JobRole', 'JobSatisfaction',
9                 'MaritalStatus', 'MonthlyIncome', 'MonthlyRate', 'NumCompaniesWorked',
10                'OverTime', 'PercentSalaryHike', 'PerformanceRating',
11                'RelationshipSatisfaction', 'StockOptionLevel',
12                'TotalWorkingYears', 'TrainingTimesLastYear', 'WorkLifeBalance',
13                'YearsAtCompany', 'YearsInCurrentRole', 'YearsSinceLastPromotion',
14                'YearsWithCurrManager', 'Attrition']]
```

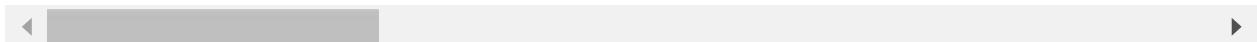
In [10]:

```
1 ea.head() #new dataframe
```

Out[10]:

|   | Age | BusinessTravel    | DailyRate | Department             | DistanceFromHome | Education | EducationField | EmployeeCount | EnvironmentSatisfaction | Gender | HourlyRate | JobInvolvement | JobLevel | JobRole | JobSatisfaction | MaritalStatus | MonthlyIncome | MonthlyRate | NumCompaniesWorked | OverTime | PercentSalaryHike | PerformanceRating | RelationshipSatisfaction | StockOptionLevel | TotalWorkingYears | TrainingTimesLastYear | WorkLifeBalance | YearsAtCompany | YearsInCurrentRole | YearsSinceLastPromotion | YearsWithCurrManager | Attrition |
|---|-----|-------------------|-----------|------------------------|------------------|-----------|----------------|---------------|-------------------------|--------|------------|----------------|----------|---------|-----------------|---------------|---------------|-------------|--------------------|----------|-------------------|-------------------|--------------------------|------------------|-------------------|-----------------------|-----------------|----------------|--------------------|-------------------------|----------------------|-----------|
| 0 | 41  | Travel_Rarely     | 1102      | Sales                  |                  | 1         | 2              | 1             | 1                       | Female | 40         | Yes            | 1        | Manager | 1               | Married       | 100000        | 100000      | 1                  | No       | 0.0               | 1                 | 1                        | 1                | 1                 | 1                     | 1               | 1              | 1                  | 1                       |                      |           |
| 1 | 49  | Travel_Frequently | 279       | Research & Development |                  | 8         | 1              | 1             | 1                       | Female | 40         | Yes            | 1        | Manager | 1               | Married       | 100000        | 100000      | 1                  | No       | 0.0               | 1                 | 1                        | 1                | 1                 | 1                     | 1               | 1              | 1                  | 1                       |                      |           |
| 2 | 37  | Travel_Rarely     | 1373      | Research & Development |                  | 2         | 2              | 1             | 1                       | Male   | 40         | Yes            | 1        | Manager | 1               | Married       | 100000        | 100000      | 1                  | No       | 0.0               | 1                 | 1                        | 1                | 1                 | 1                     | 1               | 1              | 1                  | 1                       |                      |           |
| 3 | 33  | Travel_Frequently | 1392      | Research & Development |                  | 3         | 4              | 1             | 1                       | Female | 40         | Yes            | 1        | Manager | 1               | Married       | 100000        | 100000      | 1                  | No       | 0.0               | 1                 | 1                        | 1                | 1                 | 1                     | 1               | 1              | 1                  | 1                       |                      |           |
| 4 | 27  | Travel_Rarely     | 591       | Research & Development |                  | 2         | 1              | 1             | 1                       | Male   | 40         | Yes            | 1        | Manager | 1               | Married       | 100000        | 100000      | 1                  | No       | 0.0               | 1                 | 1                        | 1                | 1                 | 1                     | 1               | 1              | 1                  | 1                       |                      |           |

5 rows × 32 columns



In [11]:

```

1 #To show data distribution
2 ea.hist()
3 plt.gcf().set_size_inches(20,20)
4 plt.show()

```



```
In [12]: 1 ea.isna().sum() #to check for missing values
```

```
Out[12]: Age          0  
BusinessTravel  0  
DailyRate       0  
Department      0  
DistanceFromHome 0  
Education        0  
EducationField   0  
EmployeeNumber   0  
EnvironmentSatisfaction 0  
Gender          0  
HourlyRate       0  
JobInvolvement   0  
JobLevel         0  
JobRole          0  
JobSatisfaction  0  
MaritalStatus    0  
MonthlyIncome     0  
MonthlyRate       0  
NumCompaniesWorked 0  
OverTime          0  
PercentSalaryHike 0  
PerformanceRating 0  
RelationshipSatisfaction 0  
StockOptionLevel  0  
TotalWorkingYears 0  
TrainingTimesLastYear 0  
WorkLifeBalance   0  
YearsAtCompany    0  
YearsInCurrentRole 0  
YearsSinceLastPromotion 0  
YearsWithCurrManager 0  
Attrition         0  
dtype: int64
```

In [13]:

```

1 ea.info()
2

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1470 entries, 0 to 1469
Data columns (total 32 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Age              1470 non-null    int64  
 1   BusinessTravel   1470 non-null    object  
 2   DailyRate        1470 non-null    int64  
 3   Department       1470 non-null    object  
 4   DistanceFromHome 1470 non-null    int64  
 5   Education        1470 non-null    int64  
 6   EducationField   1470 non-null    object  
 7   EmployeeNumber   1470 non-null    int64  
 8   EnvironmentSatisfaction 1470 non-null    int64  
 9   Gender            1470 non-null    object  
 10  HourlyRate       1470 non-null    int64  
 11  JobInvolvement  1470 non-null    int64  
 12  JobLevel         1470 non-null    int64  
 13  JobRole          1470 non-null    object  
 14  JobSatisfaction 1470 non-null    int64  
 15  MaritalStatus   1470 non-null    object  
 16  MonthlyIncome   1470 non-null    int64  
 17  MonthlyRate     1470 non-null    int64  
 18  NumCompaniesWorked 1470 non-null    int64  
 19  Overtime         1470 non-null    object  
 20  PercentSalaryHike 1470 non-null    int64  
 21  PerformanceRating 1470 non-null    int64  
 22  RelationshipSatisfaction 1470 non-null    int64  
 23  StockOptionLevel 1470 non-null    int64  
 24  TotalWorkingYears 1470 non-null    int64  
 25  TrainingTimesLastYear 1470 non-null    int64  
 26  WorkLifeBalance  1470 non-null    int64  
 27  YearsAtCompany  1470 non-null    int64  
 28  YearsInCurrentRole 1470 non-null    int64  
 29  YearsSinceLastPromotion 1470 non-null    int64  
 30  YearsWithCurrManager 1470 non-null    int64  
 31  Attrition        1470 non-null    object  
dtypes: int64(24), object(8)
memory usage: 367.6+ KB

```

In [14]:

```

1 #some columns have object as their datatype, they need to be changed to category
2 ea['BusinessTravel'] = ea['BusinessTravel'].astype('category')
3 ea['Department'] = ea['Department'].astype('category')
4 ea['EducationField'] = ea['EducationField'].astype('category')
5 ea['Gender'] = ea['Gender'].astype('category')
6 ea['JobRole'] = ea['JobRole'].astype('category')
7 ea['MaritalStatus'] = ea['MaritalStatus'].astype('category')
8 ea['OverTime'] = ea['OverTime'].astype('category')
9 ea['Attrition'] = ea['Attrition'].astype('category')
10
11

```

In [15]:

```

1 #categorical data encoding
2 from sklearn.preprocessing import LabelEncoder
3
4 le = LabelEncoder()
5 # Assigning numerical values and storing in another column
6
7
8 ea['BusinessTravel'] = le.fit_transform(ea['BusinessTravel'])
9 ea['Department'] = le.fit_transform(ea['Department'])
10 ea['EducationField'] = le.fit_transform(ea['EducationField'])
11 ea['Gender'] = le.fit_transform(ea['Gender'])
12 ea['JobRole'] = le.fit_transform(ea['JobRole'])
13 ea['MaritalStatus'] = le.fit_transform(ea['MaritalStatus'])
14 ea['OverTime'] = le.fit_transform(ea['OverTime'])
15 ea['Attrition'] = le.fit_transform(ea['Attrition'])

```

In [16]:

```
1 dfEncoded = ea
```

In [17]:

```

1 #putting Age column as the last column to make it easy for identification
2 #as the output variable for analysis
3 #eAge is employee age
4 eAge = dfEncoded.loc[:, ['BusinessTravel', 'DailyRate', 'Department', 'DistanceFromHome', 'Education', 'EducationField', 'EmployeeNumber', 'EnvironmentSatisfaction', 'Gender', 'JobInvolvement', 'JobLevel', 'JobRole', 'JobSatisfaction', 'MaritalStatus', 'MonthlyIncome', 'MonthlyRate', 'NumCompaniesWorked', 'OverTime', 'PerformanceRating', 'RelationshipSatisfaction', 'StockOptionLevel', 'TrainingTimesLastYear', 'WorkLifeBalance', 'YearsAtCompany', 'YearsInCurrentRole', 'YearsSinceLastPromotion', 'YearsWithCurrManager', 'Attrition', 'Age']]

```

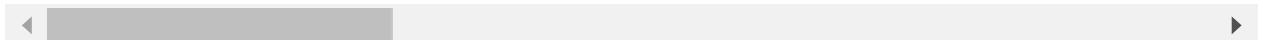
In [18]:

```
1 eAge.head()
```

Out[18]:

|   | BusinessTravel | DailyRate | Department | DistanceFromHome | Education | EducationField | EmployeeNumber |
|---|----------------|-----------|------------|------------------|-----------|----------------|----------------|
| 0 | 2              | 1102      | 2          |                  | 1         | 2              | 1              |
| 1 | 1              | 279       | 1          |                  | 8         | 1              | 1              |
| 2 | 2              | 1373      | 1          |                  | 2         | 2              | 4              |
| 3 | 1              | 1392      | 1          |                  | 3         | 4              | 1              |
| 4 | 2              | 591       | 1          |                  | 2         | 1              | 3              |

5 rows × 32 columns



In [19]: 1 eAge.columns

Out[19]: Index(['BusinessTravel', 'DailyRate', 'Department', 'DistanceFromHome',  
'Education', 'EducationField', 'EmployeeNumber',  
'EnvironmentSatisfaction', 'Gender', 'HourlyRate', 'JobInvolvement',  
'JobLevel', 'JobRole', 'JobSatisfaction', 'MaritalStatus',  
'MonthlyIncome', 'MonthlyRate', 'NumCompaniesWorked', 'OverTime',  
'PercentSalaryHike', 'PerformanceRating', 'RelationshipSatisfaction',  
'StockOptionLevel', 'TotalWorkingYears', 'TrainingTimesLastYear',  
'WorkLifeBalance', 'YearsAtCompany', 'YearsInCurrentRole',  
'YearsSinceLastPromotion', 'YearsWithCurrManager', 'Attrition', 'Age'],  
dtype='object')

In [20]:

```
1 # Rescale data (between 0 and 1)
2 from pandas import read_csv
3 from numpy import set_printoptions
4 from sklearn.preprocessing import MinMaxScaler
5
6
7 array = eAge.values
8 # separate array into input and output components
9 #X is the predictor variable
10 #Y is the outcome variable
11 X = array[:,0:31]
12 Y = array[:,31]
13 demo = MinMaxScaler(feature_range=(0, 1))
14 rescaledX = demo.fit_transform(X)
15 # summarize transformed data
16 set_printoptions(precision=3)
17
18 print(rescaledX[0:5,:])
19
```

```
[[1.000e+00 7.158e-01 1.000e+00 0.000e+00 2.500e-01 2.000e-01 0.000e+00
 3.333e-01 0.000e+00 9.143e-01 6.667e-01 2.500e-01 8.750e-01 1.000e+00
 1.000e+00 2.625e-01 6.981e-01 8.889e-01 1.000e+00 0.000e+00 0.000e+00
 0.000e+00 0.000e+00 2.000e-01 0.000e+00 0.000e+00 1.500e-01 2.222e-01
 0.000e+00 2.941e-01 1.000e+00]
[5.000e-01 1.267e-01 5.000e-01 2.500e-01 0.000e+00 2.000e-01 4.838e-04
 6.667e-01 1.000e+00 4.429e-01 3.333e-01 2.500e-01 7.500e-01 3.333e-01
 5.000e-01 2.170e-01 9.160e-01 1.111e-01 0.000e+00 8.571e-01 1.000e+00
 1.000e+00 3.333e-01 2.500e-01 5.000e-01 6.667e-01 2.500e-01 3.889e-01
 6.667e-02 4.118e-01 0.000e+00]
[1.000e+00 9.098e-01 5.000e-01 3.571e-02 2.500e-01 8.000e-01 1.451e-03
 1.000e+00 1.000e+00 8.857e-01 3.333e-01 0.000e+00 2.500e-01 6.667e-01
 1.000e+00 5.692e-02 1.213e-02 6.667e-01 1.000e+00 2.857e-01 0.000e+00
 3.333e-01 0.000e+00 1.750e-01 5.000e-01 6.667e-01 0.000e+00 0.000e+00
 0.000e+00 0.000e+00 1.000e+00]
[5.000e-01 9.234e-01 5.000e-01 7.143e-02 7.500e-01 2.000e-01 1.935e-03
 1.000e+00 0.000e+00 3.714e-01 6.667e-01 0.000e+00 7.500e-01 6.667e-01
 5.000e-01 1.001e-01 8.458e-01 1.111e-01 1.000e+00 0.000e+00 0.000e+00
 6.667e-01 0.000e+00 2.000e-01 5.000e-01 6.667e-01 2.000e-01 3.889e-01
 2.000e-01 0.000e+00 0.000e+00]
[1.000e+00 3.500e-01 5.000e-01 3.571e-02 0.000e+00 6.000e-01 2.903e-03
 0.000e+00 1.000e+00 1.429e-01 6.667e-01 0.000e+00 2.500e-01 3.333e-01
 5.000e-01 1.295e-01 5.837e-01 1.000e+00 0.000e+00 7.143e-02 0.000e+00
 1.000e+00 3.333e-01 1.500e-01 5.000e-01 6.667e-01 5.000e-02 1.111e-01
 1.333e-01 1.176e-01 0.000e+00]]
```

In [21]:

```

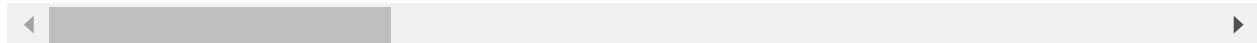
1 eAgeRescaled = pd.DataFrame(rescaledX)
2 eAgeRescaled.columns = ['BusinessTravel', 'DailyRate', 'Department', 'DistanceFromHome', 'Education', 'EducationField', 'EmployeeNumber', 'EnvironmentSatisfaction', 'Gender', 'JobInvolvement', 'JobLevel', 'JobRole', 'JobSatisfaction', 'MaritalStatus', 'MonthlyIncome', 'MonthlyRate', 'NumCompaniesWorked', 'OverTime', 'PerformanceRating', 'RelationshipSatisfaction', 'StockOptionLevel', 'TrainingTimesLastYear', 'WorkLifeBalance', 'YearsAtCompany', 'YearsInCurrentRole', 'YearsSinceLastPromotion', 'YearsWithCurrManager', 'Attrition']
3
4
5
6
7
8
9 eAgeRescaled['Age'] = eAge['Age']
10 eAgeRescaled

```

Out[21]:

|      | BusinessTravel | DailyRate | Department | DistanceFromHome | Education | EducationField | EmployeeNumber |
|------|----------------|-----------|------------|------------------|-----------|----------------|----------------|
| 0    | 1.0            | 0.715820  | 1.0        | 0.000000         | 0.25      | 0.2            | 1470           |
| 1    | 0.5            | 0.126700  | 0.5        | 0.250000         | 0.00      | 0.2            | 1470           |
| 2    | 1.0            | 0.909807  | 0.5        | 0.035714         | 0.25      | 0.8            | 1470           |
| 3    | 0.5            | 0.923407  | 0.5        | 0.071429         | 0.75      | 0.2            | 1470           |
| 4    | 1.0            | 0.350036  | 0.5        | 0.035714         | 0.00      | 0.6            | 1470           |
| ...  | ...            | ...       | ...        | ...              | ...       | ...            | ...            |
| 1465 | 0.5            | 0.559771  | 0.5        | 0.785714         | 0.25      | 0.6            | 1470           |
| 1466 | 1.0            | 0.365784  | 0.5        | 0.178571         | 0.00      | 0.6            | 1470           |
| 1467 | 1.0            | 0.037938  | 0.5        | 0.107143         | 0.50      | 0.2            | 1470           |
| 1468 | 0.5            | 0.659270  | 1.0        | 0.035714         | 0.50      | 0.6            | 1470           |
| 1469 | 1.0            | 0.376521  | 0.5        | 0.250000         | 0.50      | 0.6            | 1470           |

1470 rows × 32 columns



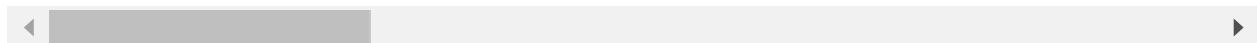
In [22]:

```
1 eAgeRescaled.describe()
```

Out[22]:

|       | BusinessTravel | DailyRate   | Department  | DistanceFromHome | Education   | EducationField |
|-------|----------------|-------------|-------------|------------------|-------------|----------------|
| count | 1470.000000    | 1470.000000 | 1470.000000 | 1470.000000      | 1470.000000 | 1470.000000    |
| mean  | 0.803741       | 0.501421    | 0.630272    | 0.292590         | 0.478231    | 0.449524       |
| std   | 0.332727       | 0.288840    | 0.263896    | 0.289531         | 0.256041    | 0.266274       |
| min   | 0.000000       | 0.000000    | 0.000000    | 0.000000         | 0.000000    | 0.000000       |
| 25%   | 0.500000       | 0.259843    | 0.500000    | 0.035714         | 0.250000    | 0.200000       |
| 50%   | 1.000000       | 0.501074    | 0.500000    | 0.214286         | 0.500000    | 0.400000       |
| 75%   | 1.000000       | 0.755190    | 1.000000    | 0.464286         | 0.750000    | 0.600000       |
| max   | 1.000000       | 1.000000    | 1.000000    | 1.000000         | 1.000000    | 1.000000       |

8 rows × 32 columns



In [23]: 1 eAge = eAgeRescaled

In [24]:

```

1 # Normalize data (length of 1)
2 from sklearn.preprocessing import Normalizer
3 from pandas import read_csv
4 from numpy import set_printoptions
5
6 # separate array into input and output components
7 array = eAge.values
8 X = array[:,0:31]
9 Y = array[:,31]
10 scaler = Normalizer().fit(X)
11 normalizedX = scaler.transform(X)
12 # summarize transformed data
13 set_printoptions(precision=5)
14 print(normalizedX[0:5,:])

```

```

[[3.10398e-01 2.22189e-01 3.10398e-01 0.00000e+00 7.75995e-02 6.20796e-02
 0.00000e+00 1.03466e-01 0.00000e+00 2.83792e-01 2.06932e-01 7.75995e-02
 2.71598e-01 3.10398e-01 3.10398e-01 8.14651e-02 2.16674e-01 2.75909e-01
 3.10398e-01 0.00000e+00 0.00000e+00 0.00000e+00 0.00000e+00 6.20796e-02
 0.00000e+00 0.00000e+00 4.65597e-02 6.89773e-02 0.00000e+00 9.12935e-02
 3.10398e-01]
[1.74127e-01 4.41238e-02 1.74127e-01 8.70635e-02 0.00000e+00 6.96508e-02
 1.68483e-04 2.32169e-01 3.48254e-01 1.54227e-01 1.16085e-01 8.70635e-02
 2.61191e-01 1.16085e-01 1.74127e-01 7.55743e-02 3.19001e-01 3.86949e-02
 0.00000e+00 2.98504e-01 3.48254e-01 3.48254e-01 1.16085e-01 8.70635e-02
 1.74127e-01 2.32169e-01 8.70635e-02 1.35432e-01 2.32169e-02 1.43399e-01
 0.00000e+00]
[3.07879e-01 2.80110e-01 1.53939e-01 1.09957e-02 7.69697e-02 2.46303e-01
 4.46849e-04 3.07879e-01 3.07879e-01 2.72693e-01 1.02626e-01 0.00000e+00
 7.69697e-02 2.05252e-01 3.07879e-01 1.75259e-02 3.73336e-03 2.05252e-01
 3.07879e-01 8.79653e-02 0.00000e+00 1.02626e-01 0.00000e+00 5.38788e-02
 1.53939e-01 2.05252e-01 0.00000e+00 0.00000e+00 0.00000e+00 0.00000e+00
 3.07879e-01]
[1.77359e-01 3.27550e-01 1.77359e-01 2.53371e-02 2.66039e-01 7.09438e-02
 6.86442e-04 3.54719e-01 0.00000e+00 1.31753e-01 2.36479e-01 0.00000e+00
 2.66039e-01 2.36479e-01 1.77359e-01 3.54906e-02 3.00026e-01 3.94132e-02
 3.54719e-01 0.00000e+00 0.00000e+00 2.36479e-01 0.00000e+00 7.09438e-02
 1.77359e-01 2.36479e-01 7.09438e-02 1.37946e-01 7.09438e-02 0.00000e+00
 0.00000e+00]
[3.81818e-01 1.33650e-01 1.90909e-01 1.36364e-02 0.00000e+00 2.29091e-01
 1.10833e-03 0.00000e+00 3.81818e-01 5.45454e-02 2.54545e-01 0.00000e+00
 9.54545e-02 1.27273e-01 1.90909e-01 4.94413e-02 2.22882e-01 3.81818e-01
 0.00000e+00 2.72727e-02 0.00000e+00 3.81818e-01 1.27273e-01 5.72727e-02
 1.90909e-01 2.54545e-01 1.90909e-02 4.24242e-02 5.09091e-02 4.49198e-02
 0.00000e+00]]

```

In [25]:

```
1 eAgeNormalized = pd.DataFrame(rescaledX)
2 eAgeNormalized.columns = ['BusinessTravel', 'DailyRate', 'Department', 'DistanceFromHome', 'Education', 'EducationField', 'EmployeeNumber', 'EnvironmentSatisfaction', 'Gender', 'JobInvolvement', 'JobLevel', 'JobRole', 'JobSatisfaction', 'MaritalStatus', 'MonthlyIncome', 'MonthlyRate', 'NumCompaniesWorked', 'OverTime', 'PerformanceRating', 'RelationshipSatisfaction', 'StockOptionLevel', 'TrainingTimesLastYear', 'WorkLifeBalance', 'YearsAtCompany', 'YearsInCurrentRole', 'YearsSinceLastPromotion', 'YearsWithCurrManager', 'Attrition']
3
4
5
6
7
8
9 eAgeNormalized['Age'] = eAge['Age']
10 eAgeNormalized
```

Out[25]:

|      | BusinessTravel | DailyRate | Department | DistanceFromHome | Education | EducationField | EmployeeNumber |
|------|----------------|-----------|------------|------------------|-----------|----------------|----------------|
| 0    | 1.0            | 0.715820  | 1.0        | 0.000000         | 0.25      | 0.2            | 1470           |
| 1    | 0.5            | 0.126700  | 0.5        | 0.250000         | 0.00      | 0.2            | 1470           |
| 2    | 1.0            | 0.909807  | 0.5        | 0.035714         | 0.25      | 0.8            | 1470           |
| 3    | 0.5            | 0.923407  | 0.5        | 0.071429         | 0.75      | 0.2            | 1470           |
| 4    | 1.0            | 0.350036  | 0.5        | 0.035714         | 0.00      | 0.6            | 1470           |
| ...  | ...            | ...       | ...        | ...              | ...       | ...            | 1470           |
| 1465 | 0.5            | 0.559771  | 0.5        | 0.785714         | 0.25      | 0.6            | 1470           |
| 1466 | 1.0            | 0.365784  | 0.5        | 0.178571         | 0.00      | 0.6            | 1470           |
| 1467 | 1.0            | 0.037938  | 0.5        | 0.107143         | 0.50      | 0.2            | 1470           |
| 1468 | 0.5            | 0.659270  | 1.0        | 0.035714         | 0.50      | 0.6            | 1470           |
| 1469 | 1.0            | 0.376521  | 0.5        | 0.250000         | 0.50      | 0.6            | 1470           |

1470 rows × 32 columns



In [26]:

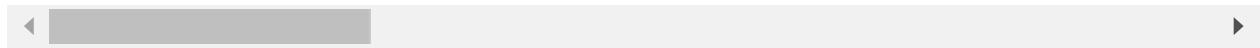
```
1 eAge = eAgeNormalized
```

In [27]: 1 eAge.describe()

Out[27]:

|              | BusinessTravel | DailyRate   | Department  | DistanceFromHome | Education   | EducationField |
|--------------|----------------|-------------|-------------|------------------|-------------|----------------|
| <b>count</b> | 1470.000000    | 1470.000000 | 1470.000000 | 1470.000000      | 1470.000000 | 1470.000000    |
| <b>mean</b>  | 0.803741       | 0.501421    | 0.630272    | 0.292590         | 0.478231    | 0.449524       |
| <b>std</b>   | 0.332727       | 0.288840    | 0.263896    | 0.289531         | 0.256041    | 0.266274       |
| <b>min</b>   | 0.000000       | 0.000000    | 0.000000    | 0.000000         | 0.000000    | 0.000000       |
| <b>25%</b>   | 0.500000       | 0.259843    | 0.500000    | 0.035714         | 0.250000    | 0.200000       |
| <b>50%</b>   | 1.000000       | 0.501074    | 0.500000    | 0.214286         | 0.500000    | 0.400000       |
| <b>75%</b>   | 1.000000       | 0.755190    | 1.000000    | 0.464286         | 0.750000    | 0.600000       |
| <b>max</b>   | 1.000000       | 1.000000    | 1.000000    | 1.000000         | 1.000000    | 1.000000       |

8 rows × 32 columns

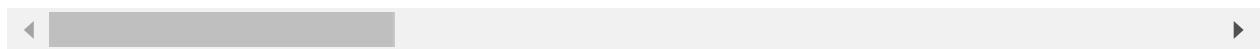


In [28]: 1 eAge.head()

Out[28]:

|          | BusinessTravel | DailyRate | Department | DistanceFromHome | Education | EducationField | Employee |
|----------|----------------|-----------|------------|------------------|-----------|----------------|----------|
| <b>0</b> | 1.0            | 0.715820  | 1.0        | 0.000000         | 0.25      |                | 0.2      |
| <b>1</b> | 0.5            | 0.126700  | 0.5        | 0.250000         | 0.00      |                | 0.2      |
| <b>2</b> | 1.0            | 0.909807  | 0.5        | 0.035714         | 0.25      |                | 0.8      |
| <b>3</b> | 0.5            | 0.923407  | 0.5        | 0.071429         | 0.75      |                | 0.2      |
| <b>4</b> | 1.0            | 0.350036  | 0.5        | 0.035714         | 0.00      |                | 0.6      |

5 rows × 32 columns



In [29]:

```

1 #to identify which attributes (and combination of attributes) contribute the most to the target variable
2 #Feature Extraction with RFE
3
4 from pandas import read_csv
5 from sklearn.feature_selection import RFE
6 from sklearn.linear_model import LogisticRegression
7 # Load data
8
9 array = eAge.values
10 X = array[:,0:31]
11 Y = array[:,31]
12 # feature extraction
13 model = LogisticRegression(max_iter = 500)
14 rfe = RFE(model, 18)
15 fit = rfe.fit(X, Y)
16 print("Num Features: {}".format(fit.n_features_))
17 print("Selected Features: {}".format(fit.support_))
18 print("Feature Ranking: {}".format(fit.ranking_))

```

C:\Users\HP\anaconda3\lib\site-packages\sklearn\utils\validation.py:70: FutureWarning: Pass n\_features\_to\_select=18 as keyword args. From version 1.0 (renaming of 0.25) passing these as positional arguments will result in an error  
`warnings.warn(f"Pass {args\_msg} as keyword args. From version "

```

Num Features: 18
Selected Features: [ True False  True False  True False  True False False False
 True  True
 True False  True  True  True False False False  True False  True
 False  True False  True  True  True]
Feature Ranking: [ 1  4  1  8  1  6  1 13  5  9  1  1  1 11  1  1  1 12 10  1
 4  1  2  1
 7  1  3  1  1  1  1]

```

In [30]:

```

1 viableCols = []
2 rankList = list(fit.ranking_)
3 columnList = ['BusinessTravel', 'DailyRate', 'Department', 'DistanceFromHome',
4               'EducationField', 'EmployeeNumber', 'EnvironmentSatisfaction', 'Gender',
5               'JobInvolvement', 'JobLevel', 'JobRole', 'JobSatisfaction', 'MaritalStatus',
6               'MonthlyIncome', 'MonthlyRate', 'NumCompaniesWorked', 'OverTime', 'PerformanceRating',
7               'RelationshipSatisfaction', 'StockOptionLevel',
8               'TrainingTimesLastYear', 'WorkLifeBalance', 'YearsAtCompany', 'YearsInCurrentRole',
9               'YearsSinceLastPromotion', 'YearsWithCurrManager', 'Attrition', 'Age']
10 for rank, col in zip(rankList, columnList):
11     idx = rank
12     colname = col
13     if idx == 1:
14         viableCols.append(colname)
15 print(viableCols)

```

```

['BusinessTravel', 'Department', 'Education', 'EmployeeNumber', 'JobInvolvement',
 'JobLevel', 'JobRole', 'MaritalStatus', 'MonthlyIncome', 'MonthlyRate', 'NumCompaniesWorked',
 'RelationshipSatisfaction', 'TotalWorkingYears', 'WorkLifeBalance',
 'YearsInCurrentRole', 'YearsSinceLastPromotion', 'YearsWithCurrManager',
 'Attrition']

```

In [31]:

```
1 eAgeModel = eAge.loc[:, viableCols] #eAgeModel is the employee age data that
2 eAgeModel
```

Out[31]:

|      | BusinessTravel | Department | Education | EmployeeNumber | JobInvolvement | JobLevel | JobRole |
|------|----------------|------------|-----------|----------------|----------------|----------|---------|
| 0    | 1.0            | 1.0        | 0.25      | 0.000000       | 0.666667       | 0.25     | 0.87    |
| 1    | 0.5            | 0.5        | 0.00      | 0.000484       | 0.333333       | 0.25     | 0.75    |
| 2    | 1.0            | 0.5        | 0.25      | 0.001451       | 0.333333       | 0.00     | 0.25    |
| 3    | 0.5            | 0.5        | 0.75      | 0.001935       | 0.666667       | 0.00     | 0.75    |
| 4    | 1.0            | 0.5        | 0.00      | 0.002903       | 0.666667       | 0.00     | 0.25    |
| ...  | ...            | ...        | ...       | ...            | ...            | ...      | ...     |
| 1465 | 0.5            | 0.5        | 0.25      | 0.996613       | 1.000000       | 0.25     | 0.25    |
| 1466 | 1.0            | 0.5        | 0.00      | 0.997097       | 0.333333       | 0.50     | 0.00    |
| 1467 | 1.0            | 0.5        | 0.50      | 0.998065       | 1.000000       | 0.25     | 0.50    |
| 1468 | 0.5            | 1.0        | 0.50      | 0.998549       | 0.333333       | 0.25     | 0.87    |
| 1469 | 1.0            | 0.5        | 0.50      | 1.000000       | 1.000000       | 0.25     | 0.25    |

1470 rows × 18 columns

|  |  |  |
|--|--|--|
|  |  |  |
|--|--|--|

In [32]:

```
1 eAgeModel['Age'] = eAge['Age']
2 eAgeModel
```

Out[32]:

|      | BusinessTravel | Department | Education | EmployeeNumber | JobInvolvement | JobLevel | JobRole |
|------|----------------|------------|-----------|----------------|----------------|----------|---------|
| 0    | 1.0            | 1.0        | 0.25      | 0.000000       | 0.666667       | 0.25     | 0.87    |
| 1    | 0.5            | 0.5        | 0.00      | 0.000484       | 0.333333       | 0.25     | 0.75    |
| 2    | 1.0            | 0.5        | 0.25      | 0.001451       | 0.333333       | 0.00     | 0.25    |
| 3    | 0.5            | 0.5        | 0.75      | 0.001935       | 0.666667       | 0.00     | 0.75    |
| 4    | 1.0            | 0.5        | 0.00      | 0.002903       | 0.666667       | 0.00     | 0.25    |
| ...  | ...            | ...        | ...       | ...            | ...            | ...      | ...     |
| 1465 | 0.5            | 0.5        | 0.25      | 0.996613       | 1.000000       | 0.25     | 0.25    |
| 1466 | 1.0            | 0.5        | 0.00      | 0.997097       | 0.333333       | 0.50     | 0.00    |
| 1467 | 1.0            | 0.5        | 0.50      | 0.998065       | 1.000000       | 0.25     | 0.50    |
| 1468 | 0.5            | 1.0        | 0.50      | 0.998549       | 0.333333       | 0.25     | 0.87    |
| 1469 | 1.0            | 0.5        | 0.50      | 1.000000       | 1.000000       | 0.25     | 0.25    |

1470 rows × 19 columns

|  |  |  |
|--|--|--|
|  |  |  |
|--|--|--|

## Regression algorithms to predict the age of

# employees leaving the firm

In [34]:

```

1 # Linear Regression
2 from pandas import read_csv
3 from sklearn.model_selection import KFold
4 from sklearn.model_selection import cross_val_score
5 from sklearn.linear_model import LinearRegression
6
7 array = eAgeModel.values
8
9 X = array[:,0:18]
10 Y = array[:,18]
11 kfold = KFold(n_splits=10, random_state=7, shuffle = True)
12 model = LinearRegression()
13 scoring = 'neg_mean_squared_error'
14 results = cross_val_score(model, X, Y, cv=kfold, scoring=scoring)
15 print(results.mean())

```

-41.646798794070435

In [35]:

```

1 # Lasso Regression
2 from pandas import read_csv
3 from sklearn.model_selection import KFold
4 from sklearn.model_selection import cross_val_score
5 from sklearn.linear_model import Lasso
6
7 array = eAgeModel.values
8 X = array[:,0:18]
9 Y = array[:,18]
10 kfold = KFold(n_splits=10, random_state=7, shuffle = True)
11 model = Lasso()
12 scoring = 'neg_mean_squared_error'
13 results = cross_val_score(model, X, Y, cv=kfold, scoring=scoring)
14 print(results.mean())

```

-71.94709308788553

In [36]:

```

1 # ElasticNet Regression
2 from pandas import read_csv
3 from sklearn.model_selection import KFold
4 from sklearn.model_selection import cross_val_score
5 from sklearn.linear_model import ElasticNet
6
7 array = eAgeModel.values
8 X = array[:,0:18]
9 Y = array[:,18]
10 kfold = KFold(n_splits=10, random_state=7, shuffle = True)
11 model = ElasticNet()
12 scoring = 'neg_mean_squared_error'
13 results = cross_val_score(model, X, Y, cv=kfold, scoring=scoring)
14 print(results.mean())

```

-75.68337458204027

In [37]:

```

1 # KNN Regression
2 from pandas import read_csv
3 from sklearn.model_selection import KFold
4 from sklearn.model_selection import cross_val_score
5 from sklearn.neighbors import KNeighborsRegressor
6
7 array = eAgeModel.values
8 X = array[:,0:18]
9 Y = array[:,18]
10 kfold = KFold(n_splits=10, random_state=7, shuffle = True)
11 model = KNeighborsRegressor()
12 scoring = 'neg_mean_squared_error'
13 results = cross_val_score(model, X, Y, cv=kfold, scoring=scoring)
14 print(results.mean())

```

-61.06560544217687

## Classification algorithms to predict the attrition of employees

In [38]:

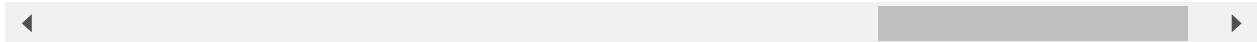
```
1 dfEncoded = ea #ea is employee attrition data
```

In [39]:

```
1 ea.describe()
```

Out[39]:

| Balance  | YearsAtCompany | YearsInCurrentRole | YearsSinceLastPromotion | YearsWithCurrManager |             |
|----------|----------------|--------------------|-------------------------|----------------------|-------------|
| 0.000000 | 1470.000000    | 1470.000000        | 1470.000000             | 1470.000000          | 1470.000000 |
| 2.761224 | 7.008163       | 4.229252           | 2.187755                | 4.123129             | (           |
| 0.706476 | 6.126525       | 3.623137           | 3.222430                | 3.568136             | (           |
| 1.000000 | 0.000000       | 0.000000           | 0.000000                | 0.000000             | (           |
| 2.000000 | 3.000000       | 2.000000           | 0.000000                | 2.000000             | (           |
| 3.000000 | 5.000000       | 3.000000           | 1.000000                | 3.000000             | (           |
| 3.000000 | 9.000000       | 7.000000           | 3.000000                | 7.000000             | (           |
| 4.000000 | 40.000000      | 18.000000          | 15.000000               | 17.000000            | (           |

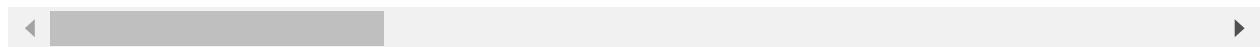


In [40]: 1 ea.head()

Out[40]:

|   | Age | BusinessTravel | DailyRate | Department | DistanceFromHome | Education | EducationField | Er |
|---|-----|----------------|-----------|------------|------------------|-----------|----------------|----|
| 0 | 41  | 2              | 1102      | 2          |                  | 1         | 2              | 1  |
| 1 | 49  | 1              | 279       | 1          |                  | 8         | 1              | 1  |
| 2 | 37  | 2              | 1373      | 1          |                  | 2         | 2              | 4  |
| 3 | 33  | 1              | 1392      | 1          |                  | 3         | 4              | 1  |
| 4 | 27  | 2              | 591       | 1          |                  | 2         | 1              | 3  |

5 rows × 32 columns



In [41]: 1 ea.columns

Out[41]: Index(['Age', 'BusinessTravel', 'DailyRate', 'Department', 'DistanceFromHome', 'Education', 'EducationField', 'EmployeeNumber', 'EnvironmentSatisfaction', 'Gender', 'HourlyRate', 'JobInvolvement', 'JobLevel', 'JobRole', 'JobSatisfaction', 'MaritalStatus', 'MonthlyIncome', 'MonthlyRate', 'NumCompaniesWorked', 'OverTime', 'PercentSalaryHike', 'PerformanceRating', 'RelationshipSatisfaction', 'StockOptionLevel', 'TotalWorkingYears', 'TrainingTimesLastYear', 'WorkLifeBalance', 'YearsAtCompany', 'YearsInCurrentRole', 'YearsSinceLastPromotion', 'YearsWithCurrManager', 'Attrition'], dtype='object')

In [42]:

```
1 # Rescale data (between 0 and 1) with attrition as Y
2 from pandas import read_csv
3 from numpy import set_printoptions
4 from sklearn.preprocessing import MinMaxScaler
5
6 array = ea.values
7 # separate array into input and output components
8 #X is the predictor variable
9 #Y is the outcome variable
10 X = array[:,0:31]
11 Y = array[:,31]
12 demo = MinMaxScaler(feature_range=(0, 1))
13 rescaledX = demo.fit_transform(X)
14 # summarize transformed data
15 set_printoptions(precision=3)
16
17 print(rescaledX[0:5,:])
18
```

```
[[5.476e-01 1.000e+00 7.158e-01 1.000e+00 0.000e+00 2.500e-01 2.000e-01
 0.000e+00 3.333e-01 0.000e+00 9.143e-01 6.667e-01 2.500e-01 8.750e-01
 1.000e+00 1.000e+00 2.625e-01 6.981e-01 8.889e-01 1.000e+00 0.000e+00
 0.000e+00 0.000e+00 0.000e+00 2.000e-01 0.000e+00 0.000e+00 1.500e-01
 2.222e-01 0.000e+00 2.941e-01]
[7.381e-01 5.000e-01 1.267e-01 5.000e-01 2.500e-01 0.000e+00 2.000e-01
 4.838e-04 6.667e-01 1.000e+00 4.429e-01 3.333e-01 2.500e-01 7.500e-01
 3.333e-01 5.000e-01 2.170e-01 9.160e-01 1.111e-01 0.000e+00 8.571e-01
 1.000e+00 1.000e+00 3.333e-01 2.500e-01 5.000e-01 6.667e-01 2.500e-01
 3.889e-01 6.667e-02 4.118e-01]
[4.524e-01 1.000e+00 9.098e-01 5.000e-01 3.571e-02 2.500e-01 8.000e-01
 1.451e-03 1.000e+00 1.000e+00 8.857e-01 3.333e-01 0.000e+00 2.500e-01
 6.667e-01 1.000e+00 5.692e-02 1.213e-02 6.667e-01 1.000e+00 2.857e-01
 0.000e+00 3.333e-01 0.000e+00 1.750e-01 5.000e-01 6.667e-01 0.000e+00
 0.000e+00 0.000e+00 0.000e+00]
[3.571e-01 5.000e-01 9.234e-01 5.000e-01 7.143e-02 7.500e-01 2.000e-01
 1.935e-03 1.000e+00 0.000e+00 3.714e-01 6.667e-01 0.000e+00 7.500e-01
 6.667e-01 5.000e-01 1.001e-01 8.458e-01 1.111e-01 1.000e+00 0.000e+00
 0.000e+00 6.667e-01 0.000e+00 2.000e-01 5.000e-01 6.667e-01 2.000e-01
 3.889e-01 2.000e-01 0.000e+00]
[2.143e-01 1.000e+00 3.500e-01 5.000e-01 3.571e-02 0.000e+00 6.000e-01
 2.903e-03 0.000e+00 1.000e+00 1.429e-01 6.667e-01 0.000e+00 2.500e-01
 3.333e-01 5.000e-01 1.295e-01 5.837e-01 1.000e+00 0.000e+00 7.143e-02
 0.000e+00 1.000e+00 3.333e-01 1.500e-01 5.000e-01 6.667e-01 5.000e-02
 1.111e-01 1.333e-01 1.176e-01]]
```

In [43]:

```

1 eaRescaled = pd.DataFrame(rescaledX)
2 eaRescaled.columns = ['Age', 'BusinessTravel', 'DailyRate', 'Department', 'D
3   'EducationField', 'EmployeeNumber', 'EnvironmentSatisfaction', 'Gender'
4   'JobInvolvement', 'JobLevel', 'JobRole', 'JobSatisfaction', 'MaritalS
5   'MonthlyIncome', 'MonthlyRate', 'NumCompaniesWorked', 'OverTime', 'Pe
6   'PerformanceRating', 'RelationshipSatisfaction', 'StockOptionLevel',
7   'TrainingTimesLastYear', 'WorkLifeBalance', 'YearsAtCompany', 'YearsI
8   'YearsSinceLastPromotion', 'YearsWithCurrManager']
9 eaRescaled['Attrition'] = ea['Attrition']
10 eaRescaled

```

Out[43]:

|      | Age      | BusinessTravel | DailyRate | Department | DistanceFromHome | Education | EducationFi |
|------|----------|----------------|-----------|------------|------------------|-----------|-------------|
| 0    | 0.547619 |                | 1.0       | 0.715820   | 1.0              | 0.000000  | 0.25        |
| 1    | 0.738095 |                | 0.5       | 0.126700   | 0.5              | 0.250000  | 0.00        |
| 2    | 0.452381 |                | 1.0       | 0.909807   | 0.5              | 0.035714  | 0.25        |
| 3    | 0.357143 |                | 0.5       | 0.923407   | 0.5              | 0.071429  | 0.75        |
| 4    | 0.214286 |                | 1.0       | 0.350036   | 0.5              | 0.035714  | 0.00        |
| ...  | ...      |                | ...       | ...        | ...              | ...       | ...         |
| 1465 | 0.428571 |                | 0.5       | 0.559771   | 0.5              | 0.785714  | 0.25        |
| 1466 | 0.500000 |                | 1.0       | 0.365784   | 0.5              | 0.178571  | 0.00        |
| 1467 | 0.214286 |                | 1.0       | 0.037938   | 0.5              | 0.107143  | 0.50        |
| 1468 | 0.738095 |                | 0.5       | 0.659270   | 1.0              | 0.035714  | 0.50        |
| 1469 | 0.380952 |                | 1.0       | 0.376521   | 0.5              | 0.250000  | 0.50        |

1470 rows × 32 columns



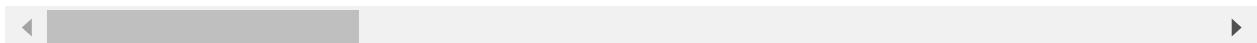
In [44]:

```
1 eaRescaled.describe()
```

Out[44]:

|       | Age         | BusinessTravel | DailyRate   | Department  | DistanceFromHome | Education   | E           |
|-------|-------------|----------------|-------------|-------------|------------------|-------------|-------------|
| count | 1470.000000 | 1470.000000    | 1470.000000 | 1470.000000 | 1470.000000      | 1470.000000 | 1470.000000 |
| mean  | 0.450567    | 0.803741       | 0.501421    | 0.630272    | 0.292590         | 0.478231    |             |
| std   | 0.217509    | 0.332727       | 0.288840    | 0.263896    | 0.289531         | 0.256041    |             |
| min   | 0.000000    | 0.000000       | 0.000000    | 0.000000    | 0.000000         | 0.000000    |             |
| 25%   | 0.285714    | 0.500000       | 0.259843    | 0.500000    | 0.035714         | 0.250000    |             |
| 50%   | 0.428571    | 1.000000       | 0.501074    | 0.500000    | 0.214286         | 0.500000    |             |
| 75%   | 0.595238    | 1.000000       | 0.755190    | 1.000000    | 0.464286         | 0.750000    |             |
| max   | 1.000000    | 1.000000       | 1.000000    | 1.000000    | 1.000000         | 1.000000    |             |

8 rows × 32 columns



In [45]: 1 ea = eaRescaled

In [46]:

```

1 # Normalize data (length of 1)
2 from sklearn.preprocessing import Normalizer
3 from pandas import read_csv
4 from numpy import set_printoptions
5
6 # separate array into input and output components
7 array = ea.values
8 X = array[:,0:31]
9 Y = array[:,31]
10 scaler = Normalizer().fit(X)
11 normalizedX = scaler.transform(X)
12 # summarize transformed data
13 set_printoptions(precision=5)
14 print(normalizedX[0:5,:])

```

```

[[1.76020e-01 3.21428e-01 2.30084e-01 3.21428e-01 0.00000e+00 8.03570e-02
 6.42856e-02 0.00000e+00 1.07143e-01 0.00000e+00 2.93877e-01 2.14285e-01
 8.03570e-02 2.81249e-01 3.21428e-01 3.21428e-01 8.43600e-02 2.24374e-01
 2.85714e-01 3.21428e-01 0.00000e+00 0.00000e+00 0.00000e+00 0.00000e+00
 6.42856e-02 0.00000e+00 0.00000e+00 4.82142e-02 7.14284e-02 0.00000e+00
 9.45376e-02]
[2.48952e-01 1.68645e-01 4.27346e-02 1.68645e-01 8.43224e-02 0.00000e+00
 6.74579e-02 1.63178e-04 2.24860e-01 3.37290e-01 1.49371e-01 1.12430e-01
 8.43224e-02 2.52967e-01 1.12430e-01 1.68645e-01 7.31949e-02 3.08958e-01
 3.74766e-02 0.00000e+00 2.89105e-01 3.37290e-01 3.37290e-01 1.12430e-01
 8.43224e-02 1.68645e-01 2.24860e-01 8.43224e-02 1.31168e-01 2.24860e-02
 1.38884e-01]
[1.44845e-01 3.20185e-01 2.91306e-01 1.60092e-01 1.14352e-02 8.00462e-02
 2.56148e-01 4.64709e-04 3.20185e-01 3.20185e-01 2.83592e-01 1.06728e-01
 0.00000e+00 8.00462e-02 2.13456e-01 3.20185e-01 1.82264e-02 3.88258e-03
 2.13456e-01 3.20185e-01 9.14813e-02 0.00000e+00 1.06728e-01 0.00000e+00
 5.60323e-02 1.60092e-01 2.13456e-01 0.00000e+00 0.00000e+00 0.00000e+00
 0.00000e+00]
[1.25681e-01 1.75953e-01 3.24953e-01 1.75953e-01 2.51362e-02 2.63930e-01
 7.03812e-02 6.80999e-04 3.51906e-01 0.00000e+00 1.30708e-01 2.34604e-01
 0.00000e+00 2.63930e-01 2.34604e-01 1.75953e-01 3.52092e-02 2.97647e-01
 3.91007e-02 3.51906e-01 0.00000e+00 0.00000e+00 2.34604e-01 0.00000e+00
 7.03812e-02 1.75953e-01 2.34604e-01 7.03812e-02 1.36852e-01 7.03812e-02
 0.00000e+00]
[8.15457e-02 3.80547e-01 1.33205e-01 1.90273e-01 1.35909e-02 0.00000e+00
 2.28328e-01 1.10463e-03 0.00000e+00 3.80547e-01 5.43638e-02 2.53698e-01
 0.00000e+00 9.51366e-02 1.26849e-01 1.90273e-01 4.92767e-02 2.22140e-01
 3.80547e-01 0.00000e+00 2.71819e-02 0.00000e+00 3.80547e-01 1.26849e-01
 5.70820e-02 1.90273e-01 2.53698e-01 1.90273e-02 4.22829e-02 5.07395e-02
 4.47702e-02]]
```

In [47]:

```
1 eaNormalized = pd.DataFrame(normalizedX)
2 eaNormalized.columns = ['Age', 'BusinessTravel', 'DailyRate', 'Department',
3                         'EducationField', 'EmployeeNumber', 'EnvironmentSatisfaction', 'Gender',
4                         'JobInvolvement', 'JobLevel', 'JobRole', 'JobSatisfaction', 'MaritalStatus',
5                         'MonthlyIncome', 'MonthlyRate', 'NumCompaniesWorked', 'OverTime', 'PerformanceRating',
6                         'RelationshipSatisfaction', 'StockOptionLevel',
7                         'TrainingTimesLastYear', 'WorkLifeBalance', 'YearsAtCompany', 'YearsInCurrentRole',
8                         'YearsSinceLastPromotion', 'YearsWithCurrManager']
9 eaNormalized['Attrition'] = ea['Attrition']
10 eaNormalized
```

Out[47]:

|      | Age      | BusinessTravel | DailyRate | Department | DistanceFromHome | Education | EducationField |
|------|----------|----------------|-----------|------------|------------------|-----------|----------------|
| 0    | 0.176020 | 0.321428       | 0.230084  | 0.321428   | 0.000000         | 0.080357  | 0.064          |
| 1    | 0.248952 | 0.168645       | 0.042735  | 0.168645   | 0.084322         | 0.000000  | 0.067          |
| 2    | 0.144845 | 0.320185       | 0.291306  | 0.160092   | 0.011435         | 0.080046  | 0.256          |
| 3    | 0.125681 | 0.175953       | 0.324953  | 0.175953   | 0.025136         | 0.263930  | 0.070          |
| 4    | 0.081546 | 0.380547       | 0.133205  | 0.190273   | 0.013591         | 0.000000  | 0.228          |
| ...  | ...      | ...            | ...       | ...        | ...              | ...       | ...            |
| 1465 | 0.143500 | 0.167417       | 0.187430  | 0.167417   | 0.263083         | 0.083708  | 0.200          |
| 1466 | 0.172268 | 0.344537       | 0.126026  | 0.172268   | 0.061524         | 0.000000  | 0.206          |
| 1467 | 0.070222 | 0.327704       | 0.012433  | 0.163852   | 0.035111         | 0.163852  | 0.065          |
| 1468 | 0.237796 | 0.161087       | 0.212400  | 0.322175   | 0.011506         | 0.161087  | 0.193          |
| 1469 | 0.132925 | 0.348929       | 0.131379  | 0.174464   | 0.087232         | 0.174464  | 0.209          |

1470 rows × 32 columns



In [48]:

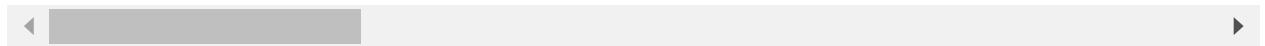
```
1 ea = eaNormalized
```

In [49]: 1 ea.describe()

Out[49]:

|              | Age         | BusinessTravel | DailyRate   | Department  | DistanceFromHome | Education   | E           |
|--------------|-------------|----------------|-------------|-------------|------------------|-------------|-------------|
| <b>count</b> | 1470.000000 | 1470.000000    | 1470.000000 | 1470.000000 | 1470.000000      | 1470.000000 | 1470.000000 |
| <b>mean</b>  | 0.149404    | 0.269174       | 0.167707    | 0.210014    | 0.097389         | 0.159918    |             |
| <b>std</b>   | 0.067989    | 0.113565       | 0.096300    | 0.084025    | 0.095689         | 0.085229    |             |
| <b>min</b>   | 0.000000    | 0.000000       | 0.000000    | 0.000000    | 0.000000         | 0.000000    |             |
| <b>25%</b>   | 0.099589    | 0.188136       | 0.089458    | 0.162288    | 0.012796         | 0.089827    |             |
| <b>50%</b>   | 0.143253    | 0.312180       | 0.165647    | 0.179955    | 0.072306         | 0.167842    |             |
| <b>75%</b>   | 0.195194    | 0.343330       | 0.251598    | 0.296860    | 0.154730         | 0.225166    |             |
| <b>max</b>   | 0.346334    | 0.479337       | 0.416890    | 0.408828    | 0.412158         | 0.447192    |             |

8 rows × 32 columns



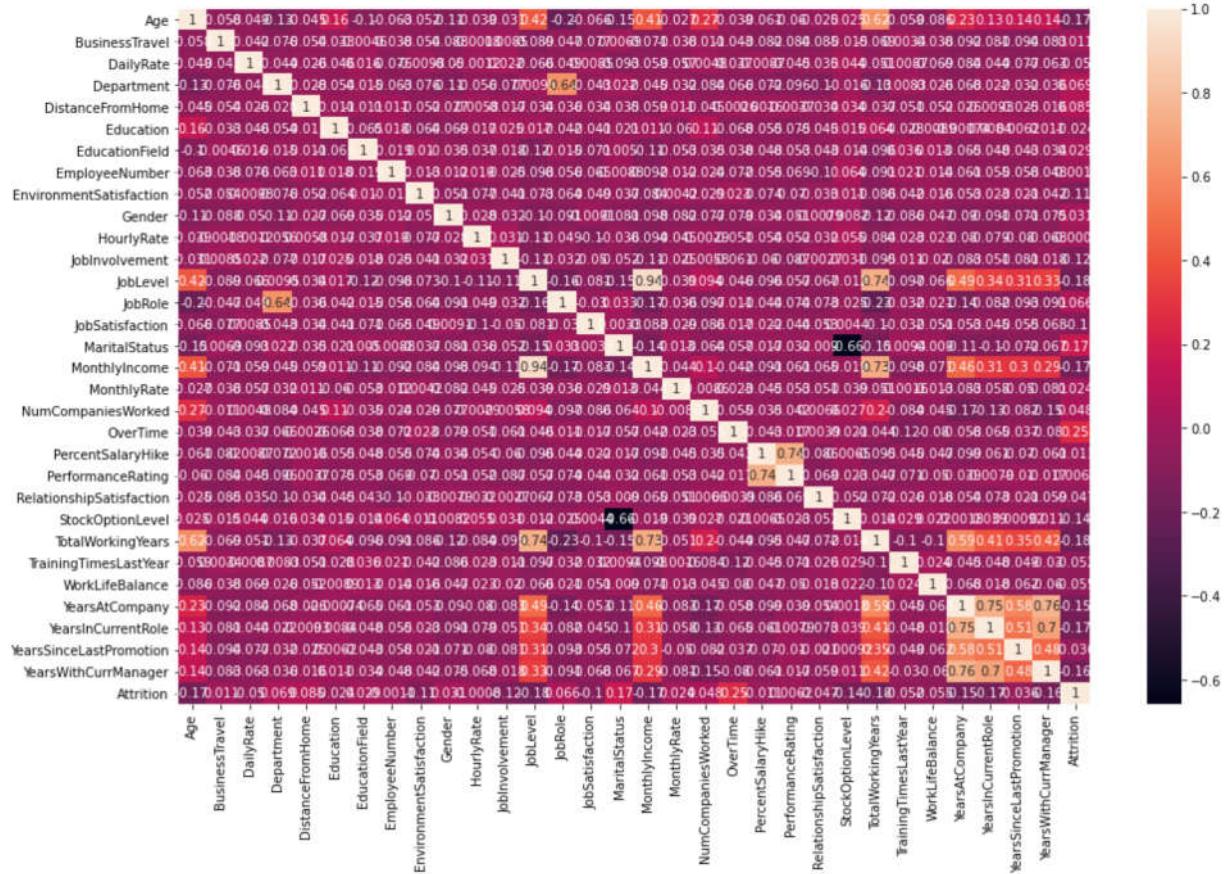
In [50]:

```

1 import seaborn as sns
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 from pandas import set_option
6 data = ea
7 set_option('display.width', 100)
8
9 plt.figure(figsize=(16,10))
10
11 sns.heatmap(data.corr(), annot = True)

```

Out[50]: &lt;AxesSubplot:&gt;



In [51]:

```

1 #to identify which attributes (and combination of attributes) contribute the most to the target variable
2 # Feature Extraction with RFE
3 from pandas import read_csv
4 from sklearn.feature_selection import RFE
5 from sklearn.linear_model import LogisticRegression
6 # Load data
7
8 array = ea.values
9 X = array[:,0:31]
10 Y = array[:,31]
11 # feature extraction
12 model = LogisticRegression()
13 rfe = RFE(model, 18)
14 fit = rfe.fit(X, Y)
15 print("Num Features: {}".format(fit.n_features_))
16 print("Selected Features: {}".format(fit.support_))
17 print("Feature Ranking: {}".format(fit.ranking_))

```

Num Features: 18  
 Selected Features: [ True False False True True False False False True False  
 False True  
 True False True True True False True True False False True True  
 True False True False True True True ]  
 Feature Ranking: [ 1 11 3 1 1 12 8 10 1 6 7 1 1 14 1 1 1 1 9 1 1  
 4 13 1 1  
 1 2 1 5 1 1 1 ]

C:\Users\HP\anaconda3\lib\site-packages\sklearn\utils\validation.py:70: FutureWarning: Pass n\_features\_to\_select=18 as keyword args. From version 1.0 (renaming of 0.25) passing these as positional arguments will result in an error  
 warnings.warn(f"Pass {args\_msg} as keyword args. From version "

In [52]:

1 ea.columns

Out[52]:

```

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

```

In [53]:

```

1 viableCols = []
2 rankList = list(fit.ranking_)
3 columnList = ['Age', 'BusinessTravel', 'DailyRate', 'Department', 'DistanceFromHome',
4               'EducationField', 'EmployeeNumber', 'EnvironmentSatisfaction', 'Gender',
5               'JobInvolvement', 'JobLevel', 'JobRole', 'JobSatisfaction', 'MaritalStatus',
6               'MonthlyIncome', 'MonthlyRate', 'NumCompaniesWorked', 'OverTime', 'PerformanceRating',
7               'RelationshipSatisfaction', 'StockOptionLevel',
8               'TrainingTimesLastYear', 'WorkLifeBalance', 'YearsAtCompany', 'YearsInCurrentRole',
9               'YearsSinceLastPromotion', 'YearsWithCurrManager', 'Attrition']
10 for rank, col in zip(rankList, columnList):
11     idx = rank
12     colname = col
13     if idx == 1:
14         viableCols.append(colname)
15 print(viableCols)

```

['Age', 'Department', 'DistanceFromHome', 'EnvironmentSatisfaction', 'JobInvolvement', 'JobLevel', 'JobSatisfaction', 'MaritalStatus', 'MonthlyIncome', 'NumCompaniesWorked', 'OverTime', 'RelationshipSatisfaction', 'StockOptionLevel', 'TotalWorkingYears', 'WorkLifeBalance', 'YearsInCurrentRole', 'YearsSinceLastPromotion', 'YearsWithCurrManager']

In [54]:

```

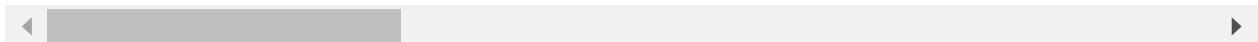
1 eaModel = ea.loc[:, viableCols] #eaModel is the employee attrition data that
2 eaModel

```

Out[54]:

|      | Age      | Department | DistanceFromHome | EnvironmentSatisfaction | JobInvolvement | JobLevel |
|------|----------|------------|------------------|-------------------------|----------------|----------|
| 0    | 0.176020 | 0.321428   | 0.000000         | 0.107143                | 0.214285       | 0.080357 |
| 1    | 0.248952 | 0.168645   | 0.084322         | 0.224860                | 0.112430       | 0.084322 |
| 2    | 0.144845 | 0.160092   | 0.011435         | 0.320185                | 0.106728       | 0.000000 |
| 3    | 0.125681 | 0.175953   | 0.025136         | 0.351906                | 0.234604       | 0.000000 |
| 4    | 0.081546 | 0.190273   | 0.013591         | 0.000000                | 0.253698       | 0.000000 |
| ...  | ...      | ...        | ...              | ...                     | ...            | ...      |
| 1465 | 0.143500 | 0.167417   | 0.263083         | 0.223222                | 0.334833       | 0.083708 |
| 1466 | 0.172268 | 0.172268   | 0.061524         | 0.344537                | 0.114846       | 0.172268 |
| 1467 | 0.070222 | 0.163852   | 0.035111         | 0.109235                | 0.327704       | 0.081926 |
| 1468 | 0.237796 | 0.322175   | 0.011506         | 0.322175                | 0.107392       | 0.080544 |
| 1469 | 0.132925 | 0.174464   | 0.087232         | 0.116310                | 0.348929       | 0.087232 |

1470 rows × 18 columns

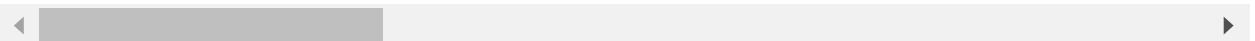


```
In [55]: 1 eaModel['Attrition'] = ea['Attrition']
          2 eaModel
```

Out[55]:

|      | Age      | Department | DistanceFromHome | EnvironmentSatisfaction | JobInvolvement | JobLeve  |
|------|----------|------------|------------------|-------------------------|----------------|----------|
| 0    | 0.176020 | 0.321428   | 0.000000         | 0.107143                | 0.214285       | 0.080351 |
| 1    | 0.248952 | 0.168645   | 0.084322         | 0.224860                | 0.112430       | 0.084321 |
| 2    | 0.144845 | 0.160092   | 0.011435         | 0.320185                | 0.106728       | 0.000000 |
| 3    | 0.125681 | 0.175953   | 0.025136         | 0.351906                | 0.234604       | 0.000000 |
| 4    | 0.081546 | 0.190273   | 0.013591         | 0.000000                | 0.253698       | 0.000000 |
| ...  | ...      | ...        | ...              | ...                     | ...            | ...      |
| 1465 | 0.143500 | 0.167417   | 0.263083         | 0.223222                | 0.334833       | 0.083708 |
| 1466 | 0.172268 | 0.172268   | 0.061524         | 0.344537                | 0.114846       | 0.172268 |
| 1467 | 0.070222 | 0.163852   | 0.035111         | 0.109235                | 0.327704       | 0.081926 |
| 1468 | 0.237796 | 0.322175   | 0.011506         | 0.322175                | 0.107392       | 0.080544 |
| 1469 | 0.132925 | 0.174464   | 0.087232         | 0.116310                | 0.348929       | 0.087232 |

1470 rows × 19 columns



## Classification Models for Attrition

In [56]:

```

1 # Logistic Regression Classification
2 from pandas import read_csv
3 from sklearn.model_selection import KFold
4 from sklearn.model_selection import cross_val_score
5 from sklearn.model_selection import cross_val_predict
6 #from sklearn.metrics import confusion_matrix
7 from sklearn.linear_model import LogisticRegression
8
9
10
11 array = eaModel.values # Converting input data to arrays
12 X = array[:,0:18] # Algorithms ONLY accept Arrays
13 Y = array[:,18]
14
15
16 num_folds = 10
17 kfold = KFold(n_splits=num_folds)#, shuffle = True, random_state=7)
18 model = LogisticRegression(max_iter = 500)
19
20 results = cross_val_score(model, X, Y, cv=kfold) # epochs
21
22 score = np.mean(results)
23 print(results)
24 print(results.mean())
25
26

```

[0.84354 0.87075 0.87075 0.88435 0.83673 0.82313 0.85714 0.87075 0.84354  
0.89796]  
0.8598639455782312

In [57]:

```

1 # LDA Classification
2 from pandas import read_csv
3 from sklearn.model_selection import KFold
4 from sklearn.model_selection import cross_val_score
5 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
6
7 array = eaModel.values
8 X = array[:,0:18]
9 Y = array[:,18]
10 num_folds = 10
11 kfold = KFold(n_splits=10,shuffle = True, random_state=7)
12 model = LinearDiscriminantAnalysis()
13 results = cross_val_score(model, X, Y, cv=kfold)
14 print(results.mean())

```

0.8727891156462585

In [58]:

```
1 #KNN Classification
2 from pandas import read_csv
3 from sklearn.model_selection import KFold
4 from sklearn.model_selection import cross_val_score
5 from sklearn.neighbors import KNeighborsClassifier
6
7 array = eaModel.values
8 X = array[:,0:18]
9 Y = array[:,18]
10 num_folds = 10
11 kfold = KFold(n_splits=10, random_state=7, shuffle = True)
12 model = KNeighborsClassifier()
13 results = cross_val_score(model, X, Y, cv=kfold)
14 print(results.mean())
```

0.8435374149659864

In [59]:

```
1 # Gaussian Naive Bayes Classification
2 from pandas import read_csv
3 from sklearn.model_selection import KFold
4 from sklearn.model_selection import cross_val_score
5 from sklearn.naive_bayes import GaussianNB
6
7 array = eaModel.values
8 X = array[:,0:18]
9 Y = array[:,18]
10 kfold = KFold(n_splits=10, random_state=7, shuffle=True)
11 model = GaussianNB()
12 results = cross_val_score(model, X, Y, cv=kfold)
13 print(results.mean())
```

0.7993197278911565

In [60]:

```
1 '''LDA was the best in predictibg the attrition rate because
2 the data had a Gaussian distribution for the numerical input variables
3 and there was no multicolinearity between the data
4 Logistic regression also performed well is prediction
5 because the value of the target variable is categorical in nature.
6 and has binary output which is either a 0 or 1.
7 KNN was also used because data is clearly labeled (0/1)
8 and there was no missing data or outliers
9 Naive Bayes didn't perform as the others because of its
10 assumption of independence among predictors
11 i.e the presence of a particular feature in a class
12 is unrelated to the presence of any other feature'''
```

Out[60]: "LDA was the best in predictibg the attrition rate because \n the data had a Gaussian distribution for the numerical input variables\n and there was no multicolinearity between the data\n Logistic regression also performed well is prediction \n because the value of the target variable is categorical in nature. \n and has binary output which is either a 0 or 1.\nKNN was also used because data is clearly labeled (0/1) \nand there was no missing data or outliers\nNaive Bayes didn't perform as the others because of its\nassumption of independence among predictors \ni.e the presence of a particular feature in a class \nis unrelated to the presence of any other feature"

In [ ]:

1

In [ ]:

1