# Week 8.2 Assignment

- Joshua Burden
- Bellevue University
- DSC550 Data Mining
- Dr. Brett Werner
- 10/23/2022

# Begin Milestone 1 with a 250-500-word narrative describing your original idea for the analysis/model building business problem.

# Clearly identify the problem you will address and the target for your model.

## Background:

During the pandemic an increase in the need for health care professionals was required. The dataset collected is a modified synthetic dataset from IBM's Watson to show a useful insight into the attrition rate for healthcare workers.

## Problem:

The data set includes information about the attrition rate for employees within the healthcare field. The meaning of employee attrition is the departure of employees from the organization for any reason whether that be voluntary or involuntary, including resignation, termination, death, or retirement. Companies to avoid attrition rates being too high is to replace those who are either leaving voluntarily or involuntary. The data set should provide insights into whether a company in the healthcare field was replacing their employees that were leaving the field, or if they continued to have a gradual but deliberate reduction in staff for any reason.

## Original Idea:

The idea behind this data set is to discover whether certain roles within the healthcare industry, hours worked, age of an employee, or any other qualifying data points stand out as to why the healthcare industry had any determining factor on whether a person was to leave their field, while also predicting whether the employee was eventually replaced.

## Dataset:

This dataset contains employee and company data useful for supervised ML, unsupervised ML, and analytics. Attrition - whether an employee left or not - is included and can be used as the

target variable. The data is synthetic and based on the IBM Watson dataset for attrition. Employee roles and departments were changed to reflect the healthcare domain. Also, known outcomes for some employees were changed to help increase the performance of ML models

## Then, do a graphical analysis creating a minimum of four graphs.

## Label your graphs appropriately and explain/analyze the information provided by each graph.

```
In [ ]: import pandas as pd
        import seaborn as sns
        from matplotlib import pyplot as plt
        import numpy as np
        import plotly.express as px
```

```
In [ ]: data_df = pd.read_csv('./DATA/watson_healthcare_modified.csv')
        data_df.head()
```

Out[ ]:

| | EmployeeID | Age | Attrition | BusinessTravel | DailyRate | Department | DistanceFromHome | Educati |
|---|---|---|---|---|---|---|---|---|
| 0 | 1313919 | 41 | No | Travel_Rarely | 1102 | Cardiology | 1 | |
| 1 | 1200302 | 49 | No | Travel_Frequently | 279 | Maternity | 8 | |
| 2 | 1060315 | 37 | Yes | Travel_Rarely | 1373 | Maternity | 2 | |
| 3 | 1272912 | 33 | No | Travel_Frequently | 1392 | Maternity | 3 | |
| 4 | 1414939 | 27 | No | Travel_Rarely | 591 | Maternity | 2 | |

5 rows × 35 columns

```
In [ ]: print("Number of duplicated data: "+str(data_df.duplicated().sum()))

        Number of duplicated data: 0
```

```
In [ ]: data_df.isnull().sum()
```

Out[ ]:
```
EmployeeID                    0
Age                           0
Attrition                     0
BusinessTravel                0
DailyRate                     0
Department                    0
DistanceFromHome              0
Education                     0
EducationField                0
EmployeeCount                 0
EnvironmentSatisfaction       0
Gender                        0
HourlyRate                    0
JobInvolvement                0
JobLevel                      0
JobRole                       0
JobSatisfaction               0
MaritalStatus                 0
MonthlyIncome                 0
MonthlyRate                   0
NumCompaniesWorked            0
Over18                        0
OverTime                      0
PercentSalaryHike             0
PerformanceRating             0
RelationshipSatisfaction      0
StandardHours                 0
Shift                         0
TotalWorkingYears             0
TrainingTimesLastYear         0
WorkLifeBalance               0
YearsAtCompany                0
YearsInCurrentRole            0
YearsSinceLastPromotion       0
YearsWithCurrManager          0
dtype: int64
```
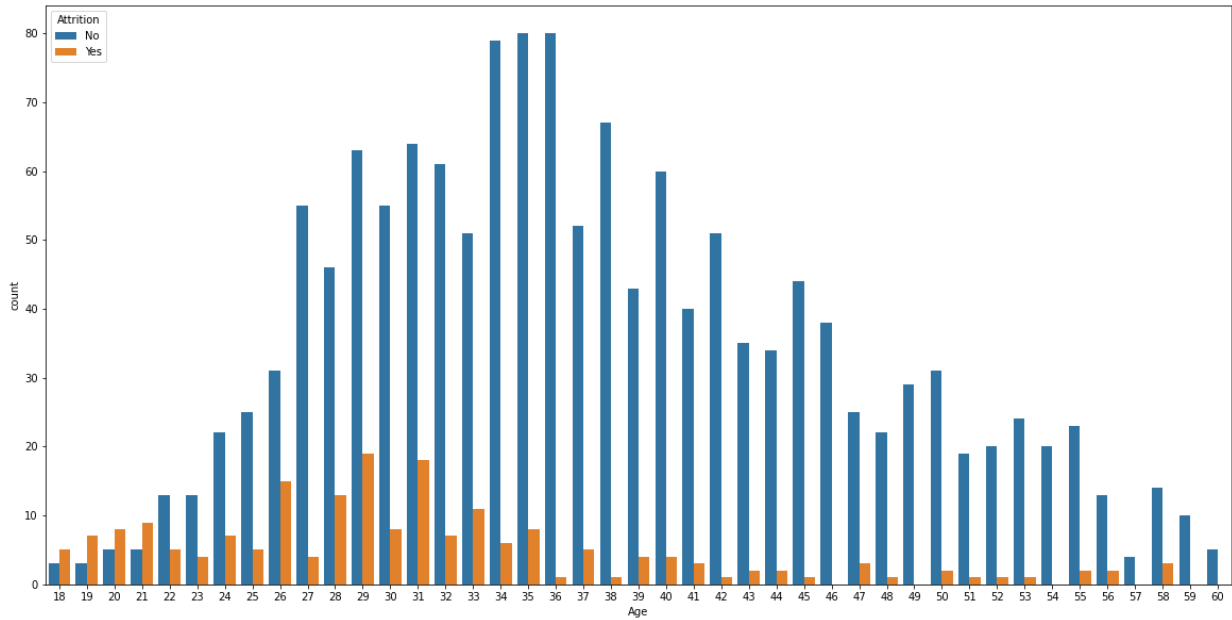
In [ ]:
```python
data_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1676 entries, 0 to 1675
Data columns (total 35 columns):
 #   Column                   Non-Null Count  Dtype
---  ------                   --------------  -----
 0   EmployeeID               1676 non-null   int64
 1   Age                      1676 non-null   int64
 2   Attrition                1676 non-null   object
 3   BusinessTravel           1676 non-null   object
 4   DailyRate                1676 non-null   int64
 5   Department               1676 non-null   object
 6   DistanceFromHome         1676 non-null   int64
 7   Education                1676 non-null   int64
 8   EducationField           1676 non-null   object
 9   EmployeeCount            1676 non-null   int64
 10  EnvironmentSatisfaction  1676 non-null   int64
 11  Gender                   1676 non-null   object
 12  HourlyRate               1676 non-null   int64
 13  JobInvolvement           1676 non-null   int64
 14  JobLevel                 1676 non-null   int64
 15  JobRole                  1676 non-null   object
 16  JobSatisfaction          1676 non-null   int64
 17  MaritalStatus            1676 non-null   object
 18  MonthlyIncome            1676 non-null   int64
 19  MonthlyRate              1676 non-null   int64
 20  NumCompaniesWorked       1676 non-null   int64
 21  Over18                   1676 non-null   object
 22  OverTime                 1676 non-null   object
 23  PercentSalaryHike        1676 non-null   int64
 24  PerformanceRating        1676 non-null   int64
 25  RelationshipSatisfaction 1676 non-null   int64
 26  StandardHours            1676 non-null   int64
 27  Shift                    1676 non-null   int64
 28  TotalWorkingYears        1676 non-null   int64
 29  TrainingTimesLastYear    1676 non-null   int64
 30  WorkLifeBalance          1676 non-null   int64
 31  YearsAtCompany           1676 non-null   int64
 32  YearsInCurrentRole       1676 non-null   int64
 33  YearsSinceLastPromotion  1676 non-null   int64
 34  YearsWithCurrManager     1676 non-null   int64
dtypes: int64(26), object(9)
memory usage: 458.4+ KB
```

## Visualization 1

```
In [ ]:  plt.figure(figsize=(20,10))
         sns.countplot(x='Age',hue='Attrition',data=data_df)
```

```
Out[ ]:  <AxesSubplot:xlabel='Age', ylabel='count'>
```

## Visualization 2

```
In [ ]:  data_df.groupby('Attrition')['MonthlyIncome'].mean().sort_values().reset_index()
```

Out[ ]:

|   | Attrition | MonthlyIncome |
|---|-----------|---------------|
| **0** | Yes | 4024.246231 |
| **1** | No | 6852.301963 |

```
In [ ]:  plt.figure(figsize=(20,10))
         sns.countplot(x='DistanceFromHome',hue='Attrition',data=data_df)
```

Out[ ]:  `<AxesSubplot:xlabel='DistanceFromHome', ylabel='count'>`



## Visualization 3

In [ ]:
```
px.histogram(data_df,x="Department",color="Attrition",barmode="group",text_auto=".2f",
            title = "Percentage of Department Type")
```

In [ ]:
```
px.histogram(data_df,x="EducationField",color="Attrition",barmode="group",text_auto=".
            title = "Percentage of EducationField Type")
```

In [ ]:
```
px.histogram(data_df,x="JobRole",color="Attrition",barmode="group",text_auto=".2f",tem
            title = "Percentage of EducationField Type")
```

## Visualization 4

In [ ]:
```
plt.figure(figsize=(10,10))
sns.countplot(x='Gender',hue='Attrition',data=data_df)
```
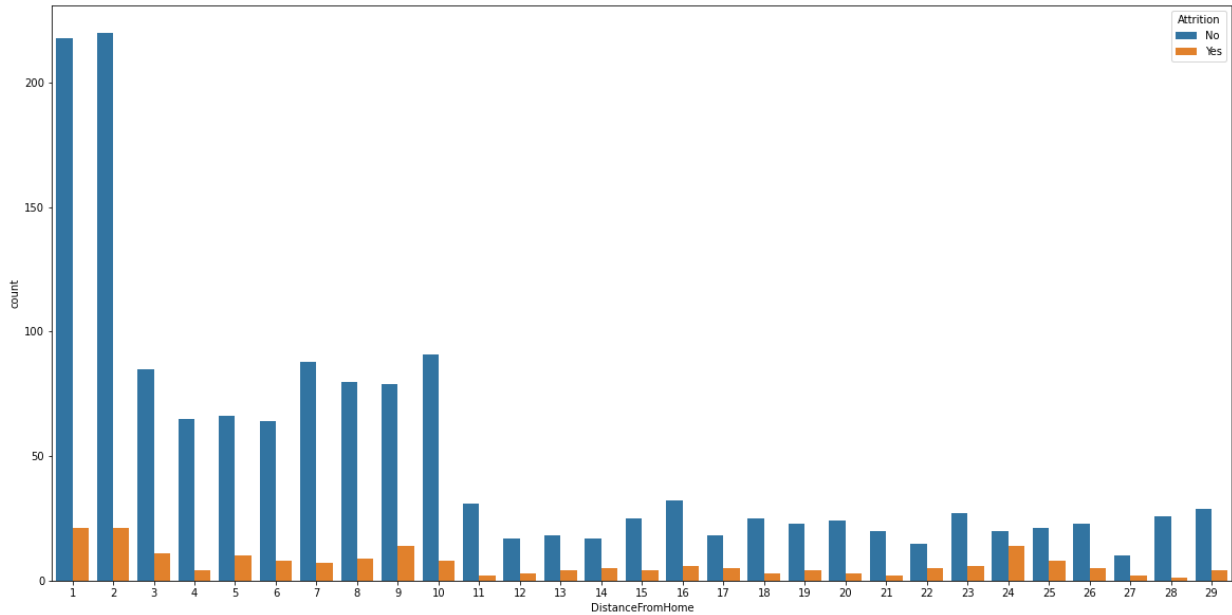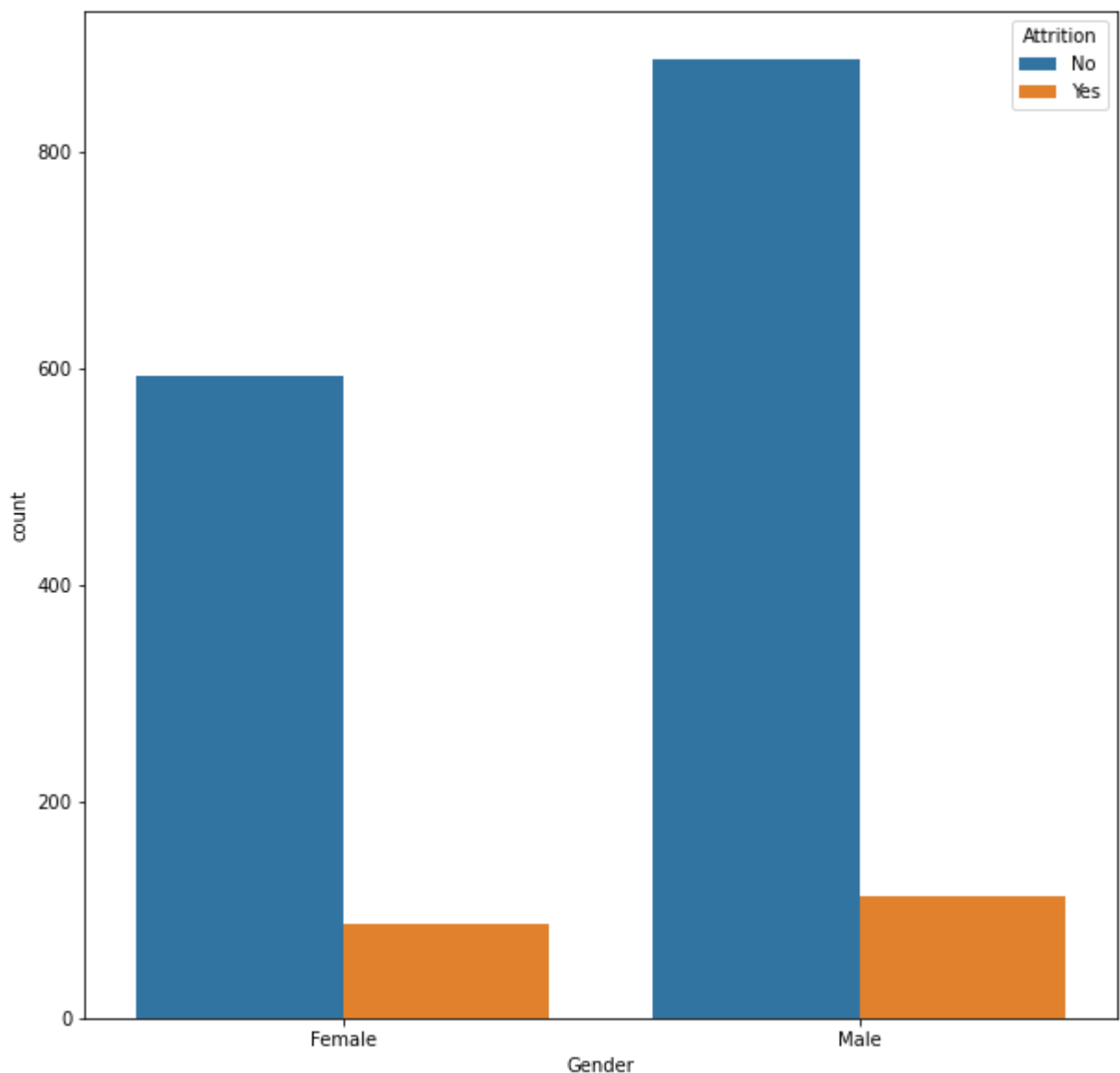
Out[ ]:
```
<AxesSubplot:xlabel='Gender', ylabel='count'>
```



Breakdown of all the available datapoints

```python
plt.figure(figsize=(30,50))
for index,column in enumerate(num_col):
    plt.subplot(5,5,index+1)
    sns.countplot(data=num_col,x=column)
    plt.xticks(rotation = 90)
plt.tight_layout(pad = 1.0)
plt.show()
```

```
---------------------------------------------------------------------------
NameError                                 Traceback (most recent call last)
c:\Users\Joshu\Desktop\Masters\DSC550\JoshuaBurdenAssignment8-2.ipynb Cell 25 in <cel
l line: 2>()
      <a href='vscode-notebook-cell:/c%3A/Users/Joshu/Desktop/Masters/DSC550/JoshuaBu
rdenAssignment8-2.ipynb#X33sZmlsZQ%3D%3D?line=0'>1</a> plt.figure(figsize=(30,50))
----> <a href='vscode-notebook-cell:/c%3A/Users/Joshu/Desktop/Masters/DSC550/JoshuaBu
rdenAssignment8-2.ipynb#X33sZmlsZQ%3D%3D?line=1'>2</a> for index,column in enumerate
(num_col):
      <a href='vscode-notebook-cell:/c%3A/Users/Joshu/Desktop/Masters/DSC550/JoshuaBu
rdenAssignment8-2.ipynb#X33sZmlsZQ%3D%3D?line=2'>3</a>     plt.subplot(5,5,index+1)
      <a href='vscode-notebook-cell:/c%3A/Users/Joshu/Desktop/Masters/DSC550/JoshuaBu
rdenAssignment8-2.ipynb#X33sZmlsZQ%3D%3D?line=3'>4</a>     sns.countplot(data=num_co
l,x=column)

NameError: name 'num_col' is not defined
<Figure size 2160x3600 with 0 Axes>
```

## Observations:

- Maternity departments had the highest rate of attrition followed by cardiology and neurology
- attrition rates had the highest peak at 29 years old
- 26-35 years old saw the highest range of attrition
- 42 years old and older saw the least attrition rates
- More men where likely to leave than women but Men also were more accounted for than women in the healthcare field
- Human resources were the least likely to have people quit
- Life Sciences were the Education field with the highest amount of attrition
- people that lived closer to their jobs were more likely to leave

# Milestone 2

```python
data_df.head()
```

Out[ ]:

| | EmployeeID | Age | Attrition | BusinessTravel | DailyRate | Department | DistanceFromHome | Educati |
|---|---|---|---|---|---|---|---|---|
| 0 | 1313919 | 41 | No | Travel_Rarely | 1102 | Cardiology | 1 | |
| 1 | 1200302 | 49 | No | Travel_Frequently | 279 | Maternity | 8 | |
| 2 | 1060315 | 37 | Yes | Travel_Rarely | 1373 | Maternity | 2 | |
| 3 | 1272912 | 33 | No | Travel_Frequently | 1392 | Maternity | 3 | |
| 4 | 1414939 | 27 | No | Travel_Rarely | 591 | Maternity | 2 | |

5 rows × 35 columns

In [ ]:
```python
data_df.shape
```

Out[ ]:
```
(1676, 35)
```

Dropping data columns that don't provide much value or context to the data

In [ ]:
```python
#drop some of the values
data_df1 = data_df.drop(['EmployeeID','Over18', 'EmployeeCount', 'StandardHours'], axi
data_df1.head()
```

Out[ ]:

| | Age | Attrition | BusinessTravel | DailyRate | Department | DistanceFromHome | Education | Education |
|---|---|---|---|---|---|---|---|---|
| 0 | 41 | No | Travel_Rarely | 1102 | Cardiology | 1 | 2 | Life Sci |
| 1 | 49 | No | Travel_Frequently | 279 | Maternity | 8 | 1 | Life Sci |
| 2 | 37 | Yes | Travel_Rarely | 1373 | Maternity | 2 | 2 | |
| 3 | 33 | No | Travel_Frequently | 1392 | Maternity | 3 | 4 | Life Sci |
| 4 | 27 | No | Travel_Rarely | 591 | Maternity | 2 | 1 | M |

5 rows × 31 columns

In [ ]:
```python
#change attrition rates from no/yes to 0/1
data_df1['Attrition'] = data_df1['Attrition'].str.replace('Yes', str(1))
data_df1['Attrition'] = data_df1['Attrition'].str.replace('No', str(0))
data_df1['Attrition'] = data_df1['Attrition'].astype('int')
```

converted attriations yes/no values and replaced them with 0/1 and set type to int

Look at shape and values of columns

In [ ]:
```python
data_df1.head()
```

Out[ ]:

| | Age | Attrition | BusinessTravel | DailyRate | Department | DistanceFromHome | Education | Education |
|---|---|---|---|---|---|---|---|---|
| 0 | 41 | 0 | Travel_Rarely | 1102 | Cardiology | 1 | 2 | Life Sci |
| 1 | 49 | 0 | Travel_Frequently | 279 | Maternity | 8 | 1 | Life Sci |
| 2 | 37 | 1 | Travel_Rarely | 1373 | Maternity | 2 | 2 | |
| 3 | 33 | 0 | Travel_Frequently | 1392 | Maternity | 3 | 4 | Life Sci |
| 4 | 27 | 0 | Travel_Rarely | 591 | Maternity | 2 | 1 | Me |

5 rows × 31 columns

In [ ]:
```
data_df1.describe().T
```

Out[ ]:

| | count | mean | std | min | 25% | 50% | 75% | ma |
|---|---|---|---|---|---|---|---|---|
| Age | 1676.0 | 36.866348 | 9.129126 | 18.0 | 30.00 | 36.0 | 43.00 | 60 |
| Attrition | 1676.0 | 0.118735 | 0.323573 | 0.0 | 0.00 | 0.0 | 0.00 | 1 |
| DailyRate | 1676.0 | 800.557876 | 401.594438 | 102.0 | 465.00 | 796.5 | 1157.00 | 1499 |
| DistanceFromHome | 1676.0 | 9.221957 | 8.158118 | 1.0 | 2.00 | 7.0 | 14.00 | 29 |
| Education | 1676.0 | 2.907518 | 1.025835 | 1.0 | 2.00 | 3.0 | 4.00 | 5 |
| EnvironmentSatisfaction | 1676.0 | 2.714797 | 1.097534 | 1.0 | 2.00 | 3.0 | 4.00 | 4 |
| HourlyRate | 1676.0 | 65.470167 | 20.207572 | 30.0 | 48.00 | 65.5 | 83.00 | 100 |
| JobInvolvement | 1676.0 | 2.724940 | 0.714121 | 1.0 | 2.00 | 3.0 | 3.00 | 4 |
| JobLevel | 1676.0 | 2.066826 | 1.113423 | 1.0 | 1.00 | 2.0 | 3.00 | 5 |
| JobSatisfaction | 1676.0 | 2.738663 | 1.104005 | 1.0 | 2.00 | 3.0 | 4.00 | 4 |
| MonthlyIncome | 1676.0 | 6516.512530 | 4728.456618 | 1009.0 | 2928.25 | 4899.0 | 8380.25 | 19999 |
| MonthlyRate | 1676.0 | 14287.019690 | 7138.857783 | 2094.0 | 7993.00 | 14269.5 | 20462.00 | 26999 |
| NumCompaniesWorked | 1676.0 | 2.662291 | 2.477704 | 0.0 | 1.00 | 2.0 | 4.00 | 9 |
| PercentSalaryHike | 1676.0 | 15.196897 | 3.646550 | 11.0 | 12.00 | 14.0 | 18.00 | 25 |
| PerformanceRating | 1676.0 | 3.150358 | 0.357529 | 3.0 | 3.00 | 3.0 | 3.00 | 4 |
| RelationshipSatisfaction | 1676.0 | 2.718377 | 1.078162 | 1.0 | 2.00 | 3.0 | 4.00 | 4 |
| Shift | 1676.0 | 0.806086 | 0.855527 | 0.0 | 0.00 | 1.0 | 1.00 | 3 |
| TotalWorkingYears | 1676.0 | 11.338902 | 7.834996 | 0.0 | 6.00 | 10.0 | 15.00 | 40 |
| TrainingTimesLastYear | 1676.0 | 2.805489 | 1.288431 | 0.0 | 2.00 | 3.0 | 3.00 | 6 |
| WorkLifeBalance | 1676.0 | 2.766110 | 0.702369 | 1.0 | 2.00 | 3.0 | 3.00 | 4 |
| YearsAtCompany | 1676.0 | 7.033413 | 6.098991 | 0.0 | 3.00 | 5.0 | 10.00 | 40 |
| YearsInCurrentRole | 1676.0 | 4.264916 | 3.627456 | 0.0 | 2.00 | 3.0 | 7.00 | 18 |
| YearsSinceLastPromotion | 1676.0 | 2.200477 | 3.229587 | 0.0 | 0.00 | 1.0 | 3.00 | 15 |
| YearsWithCurrManager | 1676.0 | 4.135442 | 3.559662 | 0.0 | 2.00 | 3.0 | 7.00 | 17 |

check correlation of data columns

In [ ]:  `data_df1.corr()`

Out[ ]:

| | Age | Attrition | DailyRate | DistanceFromHome | Education | Environment |
|---|---|---|---|---|---|---|
| **Age** | 1.000000 | -0.239984 | 0.001441 | -0.010079 | 0.204655 | |
| **Attrition** | -0.239984 | 1.000000 | -0.053892 | 0.105580 | -0.038843 | |
| **DailyRate** | 0.001441 | -0.053892 | 1.000000 | -0.009227 | -0.015881 | |
| **DistanceFromHome** | -0.010079 | 0.105580 | -0.009227 | 1.000000 | 0.015937 | |
| **Education** | 0.204655 | -0.038843 | -0.015881 | 0.015937 | 1.000000 | |
| **EnvironmentSatisfaction** | 0.008945 | -0.101278 | 0.010620 | -0.019730 | -0.031925 | |
| **HourlyRate** | 0.034671 | -0.036300 | 0.027128 | 0.026947 | 0.017996 | |
| **JobInvolvement** | 0.034193 | -0.166036 | 0.058864 | 0.010281 | 0.041046 | |
| **JobLevel** | 0.518333 | -0.207634 | 0.009005 | -0.023455 | 0.093227 | |
| **JobSatisfaction** | -0.015848 | -0.081881 | 0.032115 | -0.004758 | -0.003957 | |
| **MonthlyIncome** | 0.511378 | -0.193527 | 0.011030 | -0.041201 | 0.085116 | |
| **MonthlyRate** | 0.025837 | 0.045744 | -0.032211 | 0.031672 | -0.019198 | |
| **NumCompaniesWorked** | 0.296045 | 0.017279 | 0.034296 | -0.024969 | 0.126758 | |
| **PercentSalaryHike** | 0.007570 | 0.002943 | 0.019325 | 0.034172 | -0.006461 | |
| **PerformanceRating** | 0.005246 | 0.010728 | 0.003353 | 0.020482 | -0.020664 | |
| **RelationshipSatisfaction** | 0.058528 | -0.020462 | 0.014539 | 0.005482 | -0.005750 | |
| **Shift** | 0.037117 | -0.158322 | 0.054407 | 0.029180 | 0.024451 | |
| **TotalWorkingYears** | 0.692512 | -0.234182 | 0.009378 | -0.017663 | 0.143324 | |
| **TrainingTimesLastYear** | -0.015408 | -0.054836 | 0.001901 | -0.055471 | -0.014070 | |
| **WorkLifeBalance** | -0.004878 | -0.090513 | -0.028549 | -0.037821 | 0.003933 | |
| **YearsAtCompany** | 0.319012 | -0.201373 | -0.026892 | -0.007420 | 0.057461 | |
| **YearsInCurrentRole** | 0.222655 | -0.207891 | 0.019651 | 0.011448 | 0.051029 | |
| **YearsSinceLastPromotion** | 0.217212 | -0.086207 | -0.034571 | -0.000126 | 0.045785 | |
| **YearsWithCurrManager** | 0.215909 | -0.201087 | -0.025272 | 0.000403 | 0.055096 | |

24 rows × 24 columns

◄ |▬▬▬▬▬▬▬▬▬▬| ► 

dropping rows that exceed a threshold of 0.2

In [ ]:
```python
threshold = 0.2
data_df1 = data_df1.drop(data_df1.std()[data_df1.std() < threshold].index.values, axis
```

C:\Users\Joshu\AppData\Local\Temp\ipykernel_22128\2220865994.py:2: FutureWarning:

Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is de
precated; in a future version this will raise TypeError.  Select only valid columns b
efore calling the reduction.

In [ ]: `data_df1`

Out[ ]:

| | Age | Attrition | BusinessTravel | DailyRate | Department | DistanceFromHome | Education | Educa |
|---|---|---|---|---|---|---|---|---|
| **0** | 41 | 0 | Travel_Rarely | 1102 | Cardiology | 1 | 2 | Life |
| **1** | 49 | 0 | Travel_Frequently | 279 | Maternity | 8 | 1 | Life |
| **2** | 37 | 1 | Travel_Rarely | 1373 | Maternity | 2 | 2 | |
| **3** | 33 | 0 | Travel_Frequently | 1392 | Maternity | 3 | 4 | Life |
| **4** | 27 | 0 | Travel_Rarely | 591 | Maternity | 2 | 1 | |
| **...** | ... | ... | ... | ... | ... | ... | ... | |
| **1671** | 26 | 1 | Travel_Rarely | 471 | Neurology | 24 | 3 | |
| **1672** | 46 | 0 | Travel_Rarely | 1125 | Cardiology | 10 | 3 | N |
| **1673** | 20 | 0 | Travel_Rarely | 959 | Maternity | 1 | 3 | Life |
| **1674** | 39 | 0 | Travel_Rarely | 466 | Neurology | 1 | 1 | Life |
| **1675** | 27 | 0 | Travel_Rarely | 511 | Cardiology | 2 | 2 | |

1676 rows × 31 columns

In [ ]:
```python
#Get dummies
data_df2 = data_df1.copy()
data_df2 = pd.get_dummies(data_df2, drop_first=True)
data_df2
```

Out[ ]:

| | Age | Attrition | DailyRate | DistanceFromHome | Education | EnvironmentSatisfaction | HourlyRate |
|---|---|---|---|---|---|---|---|
| **0** | 41 | 0 | 1102 | 1 | 2 | 2 | 94 |
| **1** | 49 | 0 | 279 | 8 | 1 | 3 | 61 |
| **2** | 37 | 1 | 1373 | 2 | 2 | 4 | 92 |
| **3** | 33 | 0 | 1392 | 3 | 4 | 4 | 56 |
| **4** | 27 | 0 | 591 | 2 | 1 | 1 | 40 |
| **...** | ... | ... | ... | ... | ... | ... | ... |
| **1671** | 26 | 1 | 471 | 24 | 3 | 3 | 66 |
| **1672** | 46 | 0 | 1125 | 10 | 3 | 3 | 94 |
| **1673** | 20 | 0 | 959 | 1 | 3 | 4 | 83 |
| **1674** | 39 | 0 | 466 | 1 | 1 | 4 | 65 |
| **1675** | 27 | 0 | 511 | 2 | 2 | 1 | 89 |

1676 rows × 41 columns

```
In [ ]:  data_df2.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1676 entries, 0 to 1675
Data columns (total 41 columns):
 #   Column                            Non-Null Count  Dtype
---  ------                            --------------  -----
 0   Age                               1676 non-null   int64
 1   Attrition                         1676 non-null   int32
 2   DailyRate                         1676 non-null   int64
 3   DistanceFromHome                  1676 non-null   int64
 4   Education                         1676 non-null   int64
 5   EnvironmentSatisfaction           1676 non-null   int64
 6   HourlyRate                        1676 non-null   int64
 7   JobInvolvement                    1676 non-null   int64
 8   JobLevel                          1676 non-null   int64
 9   JobSatisfaction                   1676 non-null   int64
 10  MonthlyIncome                     1676 non-null   int64
 11  MonthlyRate                       1676 non-null   int64
 12  NumCompaniesWorked                1676 non-null   int64
 13  PercentSalaryHike                 1676 non-null   int64
 14  PerformanceRating                 1676 non-null   int64
 15  RelationshipSatisfaction          1676 non-null   int64
 16  Shift                             1676 non-null   int64
 17  TotalWorkingYears                 1676 non-null   int64
 18  TrainingTimesLastYear             1676 non-null   int64
 19  WorkLifeBalance                   1676 non-null   int64
 20  YearsAtCompany                    1676 non-null   int64
 21  YearsInCurrentRole                1676 non-null   int64
 22  YearsSinceLastPromotion           1676 non-null   int64
 23  YearsWithCurrManager              1676 non-null   int64
 24  BusinessTravel_Travel_Frequently  1676 non-null   uint8
 25  BusinessTravel_Travel_Rarely      1676 non-null   uint8
 26  Department_Maternity              1676 non-null   uint8
 27  Department_Neurology              1676 non-null   uint8
 28  EducationField_Life Sciences      1676 non-null   uint8
 29  EducationField_Marketing          1676 non-null   uint8
 30  EducationField_Medical            1676 non-null   uint8
 31  EducationField_Other              1676 non-null   uint8
 32  EducationField_Technical Degree   1676 non-null   uint8
 33  Gender_Male                       1676 non-null   uint8
 34  JobRole_Administrative            1676 non-null   uint8
 35  JobRole_Nurse                     1676 non-null   uint8
 36  JobRole_Other                     1676 non-null   uint8
 37  JobRole_Therapist                 1676 non-null   uint8
 38  MaritalStatus_Married             1676 non-null   uint8
 39  MaritalStatus_Single              1676 non-null   uint8
 40  OverTime_Yes                      1676 non-null   uint8
dtypes: int32(1), int64(23), uint8(17)
memory usage: 335.7 KB
```

create the final dataframe

```
In [ ]:  df_final = data_df2.copy()
         df_final
```

Out[ ]:

| | Age | Attrition | DailyRate | DistanceFromHome | Education | EnvironmentSatisfaction | HourlyRate |
|---|---|---|---|---|---|---|---|
| **0** | 41 | 0 | 1102 | 1 | 2 | 2 | 94 |
| **1** | 49 | 0 | 279 | 8 | 1 | 3 | 61 |
| **2** | 37 | 1 | 1373 | 2 | 2 | 4 | 92 |
| **3** | 33 | 0 | 1392 | 3 | 4 | 4 | 56 |
| **4** | 27 | 0 | 591 | 2 | 1 | 1 | 40 |
| **...** | ... | ... | ... | ... | ... | ... | ... |
| **1671** | 26 | 1 | 471 | 24 | 3 | 3 | 66 |
| **1672** | 46 | 0 | 1125 | 10 | 3 | 3 | 94 |
| **1673** | 20 | 0 | 959 | 1 | 3 | 4 | 83 |
| **1674** | 39 | 0 | 466 | 1 | 1 | 4 | 65 |
| **1675** | 27 | 0 | 511 | 2 | 2 | 1 | 89 |

1676 rows × 41 columns

Create a training and test model set

In [ ]:
```python
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn import tree
from sklearn.metrics import accuracy_score
from sklearn.metrics import confusion_matrix
```

In [ ]:
```python
X = df_final.drop(['Attrition'], axis='columns')
```

In [ ]:
```python
y= df_final['Attrition']
y.to_frame().head()
```

Out[ ]:

| | Attrition |
|---|---|
| **0** | 0 |
| **1** | 0 |
| **2** | 1 |
| **3** | 0 |
| **4** | 0 |

In [ ]:
```python
print("X shape",X.shape,"\n","y shape",y.shape)
```

```
X shape (1676, 40)
 y shape (1676,)
```

In [ ]:
```python
# Splitting data
X_train, X_test, y_train, y_test = train_test_split(X, y,
```

```
                                                   test_size=.25,
                                                   random_state=42,
                                                   )
```

In [ ]:
```python
from sklearn.linear_model import LogisticRegression
```

In [ ]:
```python
logreg = LogisticRegression(penalty='l2',solver='liblinear', max_iter=250)
logreg.fit(X_train, y_train)
```

Out[ ]:
```
LogisticRegression(max_iter=250, solver='liblinear')
```

Predict the model Accuracy

In [ ]:
```python
y_pred = logreg.predict(X_test)
print("Model Accuracy:",round(logreg.score(X_test, y_test) * 100 , 0),"%")
```

```
Model Accuracy: 91.0 %
```