

```

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from scipy import stats

IMPORT DATASET

In [4]:
df=pd.read_csv("WA_Fn-UseC_HR-Employee-Attrition.csv")

In [5]:
df

Out[5]:
   Age  Attrition  BusinessTravel  DailyRate  Department  DistanceFromHome  Education  EducationField  EmployeeCount  EmployeeNumber  ...  Relationships
0  41    No       Travel_Rarely      1102      Sales              1           2      Life Sciences      1           1           ...           1
1  49    No       Travel_Frequently  279      Research & Development  8           1      Life Sciences      1           2           ...           4
2  37    Yes      Travel_Rarely      1373      Research & Development  2           2      Other            1           4           ...           2
3  31    No       Travel_Frequently  1392      Research & Development  3           4      Life Sciences      1           5           ...           3
4  27    No       Travel_Rarely      591      Research & Development  2           1      Medical           1           7           ...           4
...  ...  ...  ...
1466 36    No       Travel_Frequently  884      Research & Development  23          2      Medical           1           2061          ...           3
1468 39    No       Travel_Rarely      613      Research & Development  6           1      Medical           1           2062          ...           1
1467 27    No       Travel_Rarely      155      Research & Development  4           3      Life Sciences      1           2064          ...           2
1468 49    No       Travel_Frequently  1023      Sales              3           3      Medical           1           2065          ...           4
1469 34    No       Travel_Rarely      628      Research & Development  8           3      Medical           1           2068          ...           1

1470 rows x 35 columns

In [6]:
df.head()

Out[6]:
   Age  Attrition  BusinessTravel  DailyRate  Department  DistanceFromHome  Education  EducationField  EmployeeCount  EmployeeNumber  ...  Relationships
0  41    No       Travel_Rarely      1102      Sales              1           2      Life Sciences      1           1           ...           1
1  49    No       Travel_Frequently  279      Research & Development  8           1      Life Sciences      1           2           ...           4
2  37    Yes      Travel_Rarely      1373      Research & Development  2           2      Other            1           4           ...           2
3  31    No       Travel_Frequently  1392      Research & Development  3           4      Life Sciences      1           5           ...           3
4  27    No       Travel_Rarely      591      Research & Development  2           1      Medical           1           7           ...           4

5 rows x 35 columns

In [7]:
df.tail()

Out[7]:
   Age  Attrition  BusinessTravel  DailyRate  Department  DistanceFromHome  Education  EducationField  EmployeeCount  EmployeeNumber  ...  Relationships
1465 36    No       Travel_Frequently  884      Research & Development  23          2      Medical           1           2061          ...           3
1466 39    No       Travel_Rarely      613      Research & Development  6           1      Medical           1           2062          ...           1
1467 27    No       Travel_Rarely      155      Research & Development  4           3      Life Sciences      1           2064          ...           2
1468 49    No       Travel_Frequently  1023      Sales              3           3      Medical           1           2065          ...           4
1469 34    No       Travel_Rarely      628      Research & Development  8           3      Medical           1           2068          ...           1

5 rows x 35 columns

In [8]:
df.shape

Out[8]:
(1470, 35)

In [9]:
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   object
 4   Department           1470 non-null   object
 5   DistanceFromHome     1470 non-null   int64
 6   Education            1470 non-null   object
 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   int64
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(28), object(7)
memory usage: 492.1+ KB

In [10]:
df.describe()

Out[10]:
   Age  Attrition  BusinessTravel  DailyRate  Department  DistanceFromHome  Education  EmployeeCount  EmployeeNumber  EnvironmentSatisfaction  HourlyRate  JobInvolvement  JobLevel  ...
mean   38.923810    0.050900    1.000000    1470.000000    1470.0    1470.000000    1.470000    1470.000000    1470.000000    1470.000000    1470.000000    1470.000000    ...
std     9.353373    0.208646    0.388884    2.912325    1.0    1024.865306    2.721769    65.891456    2.712932    0.000000    0.000000    0.000000    ...
min     3.000000    0.000000    1.000000    1.000000    1.0    1.000000    1.000000    30.000000    1.000000    1.000000    1.000000    1.000000    ...
25%    38.000000    0.050900    1.000000    2.000000    1.0    451.250000    1.000000    40.000000    2.000000    1.000000    1.000000    1.000000    ...
50%    38.000000    0.050900    1.000000    2.000000    1.0    1028.500000    1.000000    66.000000    3.000000    2.000000    2.000000    2.000000    ...
75%    43.000000    1.050900    1.000000    4.000000    1.0    1555.750000    4.000000    83.750000    3.000000    3.000000    3.000000    3.000000    ...
max     66.000000    1.050900    1.000000    5.000000    1.0    2568.000000    4.000000    109.000000    4.000000    5.000000    5.000000    5.000000    ...

8 rows x 26 columns

In [11]:
df[["Age", "BusinessTravel", "DailyRate", "DistanceFromHome", "Education", "EmployeeCount", "EmployeeNumber", "EnvironmentSatisfaction", "HourlyRate", "JobInvolvement", "JobLevel", "Relationships"]].describe()

Out[11]:
   Age  BusinessTravel  DailyRate  DistanceFromHome  Education  EmployeeCount  EmployeeNumber  EnvironmentSatisfaction  HourlyRate  JobInvolvement  JobLevel  Relationships
mean   38.923810    1.000000    1470.000000    1470.000000    1470.0    1470.000000    1470.000000    1470.000000    1470.000000    1470.000000    1470.000000    1470.000000
std     9.353373    0.388884    2.912325    1.024865    2.721769    65.891456    2.712932    0.000000    0.000000    0.000000    0.000000    0.000000
min     3.000000    1.000000    1.000000    1.000000    1.000000    30.000000    1.000000    1.000000    1.000000    1.000000    1.000000    1.000000
25%    38.000000    1.000000    2.000000    451.250000    1.000000    40.000000    2.000000    1.000000    1.000000    1.000000    1.000000    1.000000
50%    38.000000    1.000000    2.000000    1028.500000    1.000000    66.000000    3.000000    2.000000    2.000000    2.000000    2.000000    2.000000
75%    43.000000    1.000000    4.000000    1555.750000    4.000000    83.750000    3.000000    3.000000    3.000000    3.000000    3.000000    3.000000
max     66.000000    1.000000    5.000000    2568.000000    4.000000    109.000000    4.000000    5.000000    5.000000    5.000000    5.000000    5.000000

8 rows x 26 columns

In [12]:
df[["Age", "BusinessTravel", "DailyRate", "DistanceFromHome", "Education", "EmployeeCount", "EmployeeNumber", "EnvironmentSatisfaction", "HourlyRate", "JobInvolvement", "JobLevel", "Relationships"]].corr()

Out[12]:
   Age  BusinessTravel  DailyRate  DistanceFromHome  Education  EmployeeCount  EmployeeNumber  EnvironmentSatisfaction  HourlyRate  JobInvolvement  JobLevel  Relationships
Age      1.000000    0.016611    -0.004686    0.208934    NaN      NaN      NaN      -0.001445    0.010148    0.024287    0.020820    0.5096
```