

AZIM RIZAN P

21BCE448

VIT-AP

```

In [37]: ## IMPORT LIBRARIES

In [43]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from scipy import stats
```

IMPORT DATASET

```

In [38]: df=pd.read_csv("VA_Fn-UseC-HR-Employee-Attrition.csv")

In [53]: df

Out[53]:
   Age  Attrition  BusinessTravel  DailyRate  Department  DistanceFromHome  Education  EducationField  EmployeeCount  EmployeeNumber  ...  RelationshipSatisfaction  StandardHours  StockOptionLevel
0    41         Yes      Travel_Rarely      1102         Sales                1              2  Life Sciences                1                1  ...                1                80                0
1    49         No  Travel_Frequently      279         Research & Development                8              1  Life Sciences                1                2  ...                4                80                1
2    37         Yes      Travel_Rarely      1373         Research & Development                2              2      Other                1                4  ...                2                80                0
3    33         No  Travel_Frequently      1392         Research & Development                3              4  Life Sciences                1                5  ...                3                80                0
4    27         No      Travel_Rarely      591         Research & Development                2              1      Medical                1                7  ...                4                80                1
...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...
1465  36         No  Travel_Frequently      884         Research & Development                23              2      Medical                1                2061  ...                3                80                1
1466  39         No      Travel_Rarely      613         Research & Development                6              1      Medical                1                2062  ...                1                80                1
1467  27         No      Travel_Rarely      155         Research & Development                4              3  Life Sciences                1                2064  ...                2                80                1
1468  49         No  Travel_Frequently      1023         Sales                2              3      Medical                1                2065  ...                4                80                0
1469  34         No      Travel_Rarely      628         Research & Development                8              3      Medical                1                2068  ...                1                80                0
1470 rows × 35 columns
```

```

In [38]: df.head()

Out[53]:
   Age  Attrition  BusinessTravel  DailyRate  Department  DistanceFromHome  Education  EducationField  EmployeeCount  EmployeeNumber  ...  RelationshipSatisfaction  StandardHours  StockOptionLevel  To
0    41         Yes      Travel_Rarely      1102         Sales                1              2  Life Sciences                1                1  ...                1                80                0
1    49         No  Travel_Frequently      279         Research & Development                8              1  Life Sciences                1                2  ...                4                80                1
2    37         Yes      Travel_Rarely      1373         Research & Development                2              2      Other                1                4  ...                2                80                0
3    33         No  Travel_Frequently      1392         Research & Development                3              4  Life Sciences                1                5  ...                3                80                0
4    27         No      Travel_Rarely      591         Research & Development                2              1      Medical                1                7  ...                4                80                1
5 rows × 35 columns
```

```

In [38]: df.tail()

Out[53]:
   Age  Attrition  BusinessTravel  DailyRate  Department  DistanceFromHome  Education  EducationField  EmployeeCount  EmployeeNumber  ...  RelationshipSatisfaction  StandardHours  StockOptionLevel
1465  36         No  Travel_Frequently      884         Research & Development                23              2      Medical                1                2061  ...                3                80                1
1466  39         No      Travel_Rarely      613         Research & Development                6              1      Medical                1                2062  ...                1                80                1
1467  27         No      Travel_Rarely      155         Research & Development                4              3  Life Sciences                1                2064  ...                2                80                1
1468  49         No  Travel_Frequently      1023         Sales                2              3      Medical                1                2065  ...                4                80                0
1469  34         No      Travel_Rarely      628         Research & Development                8              3      Medical                1                2068  ...                1                80                0
5 rows × 35 columns
```

```

In [38]: df.shape

Out[53]:
(1470, 35)
```

```

In [39]: 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    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    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   YearsCurrentRole      1470 non-null    int64
33   YearsSinceLastPromotion 1470 non-null    int64
34   YearWithCurrManager   1470 non-null    int64
dtypes: int64(26), object(9)
memory usage: 402.1 KB
```

```

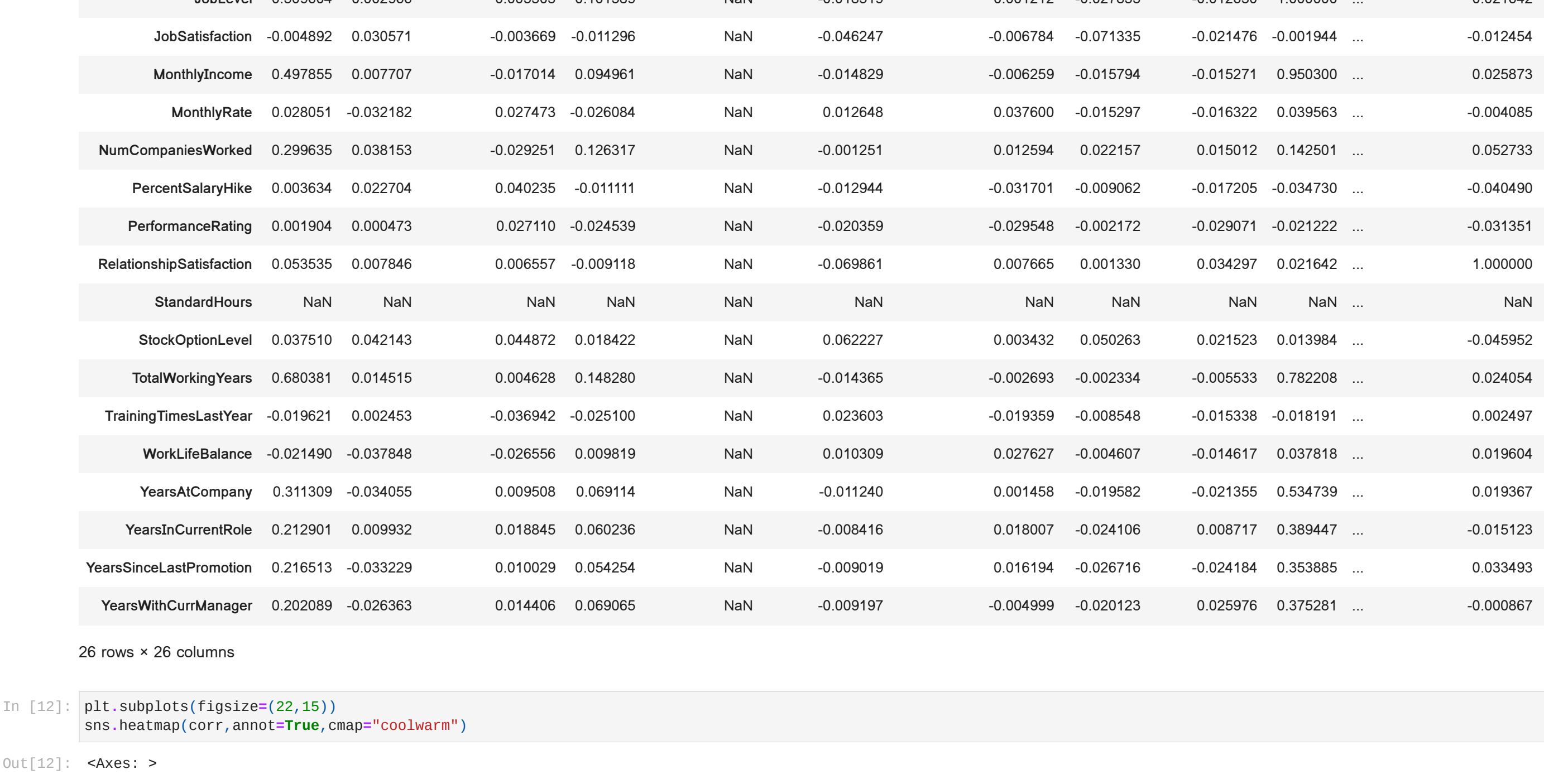
In [39]: df.describe()

Out[53]:
   count    Age      DailyRate  DistanceFromHome  Education  EmployeeCount  EmployeeNumber  EnvironmentSatisfaction  HourlyRate  JobInvolvement  JobLevel  ...  RelationshipSatisfaction  StandardHours  StockOptionLevel
mean    36.922870    802.485714    9.192517    2.912925    1.0    1024.865306    2.721769    65.891156    2.729932    2.069346    ...    2.712245
std     9.135373    403.509100    8.106864    1.024165    0.0    682.024355    1.030062    20.328428    0.711561    1.106940    ...    1.081209
min     16.000000    102.000000    1.000000    1.000000    1.0    1.000000    1.000000    30.000000    1.000000    1.000000    ...    1.000000
25%    30.000000    465.000000    2.000000    2.000000    1.0    491.260000    2.000000    48.000000    2.000000    1.000000    ...    2.000000
50%    36.000000    802.000000    7.000000    3.000000    1.0    1020.500000    3.000000    66.000000    3.000000    2.000000    ...    2.000000
75%    43.000000    1157.000000    14.000000    4.000000    1.0    1555.750000    4.000000    83.750000    3.000000    3.000000    ...    4.000000
max     60.000000    1499.000000    29.000000    5.000000    1.0    2068.000000    4.000000    100.000000    4.000000    5.000000    ...    4.000000
8 rows × 26 columns
```

```

In [41]: corr=df.corr()

C:\Users\ashan\AppData\Local\Temp\ipykernel_123216\3282248918.py:13: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.
  corr=df.corr()
```



```

In [12]: plt.subplots(figsize=(22,15),
sns.heatmap(corr,annot=True,cmap="coolwarm")

Out[12]:
<Axes: >
```



```

In [13]: df.Attrition.value_counts()

Out[13]:
Attrition
No    1223
Yes    237
Name: Attrition, dtype: int64
Checking for NULL Values
```

```

In [14]: df.isnull().any()

Out[14]:
Age                   False
Attrition             False
BusinessTravel        False
DailyRate             False
Department            False
DistanceFromHome      False
Education              False
EducationField         False
EmployeeCount         False
EmployeeNumber        False
EnvironmentSatisfaction False
Gender                False
HourlyRate            False
JobInvolvement        False
JobLevel              False
JobRole               False
JobSatisfaction        False
MaritalStatus         False
MonthlyIncome         False
MonthlyRate           False
NumCompaniesWorked    False
Over18                False
OverTime              False
PercentSalaryHike     False
PerformanceRating     False
RelationshipSatisfaction False
StandardHours         False
StockOptionLevel      False
TotalWorkingYears     False
TrainingTimesLastYear False
WorkLifeBalance       False
YearsAtCompany        False
YearsCurrentRole      False
YearsSinceLastPromotion False
YearWithCurrManager   False
dtype: bool
```

Data Visualization

```

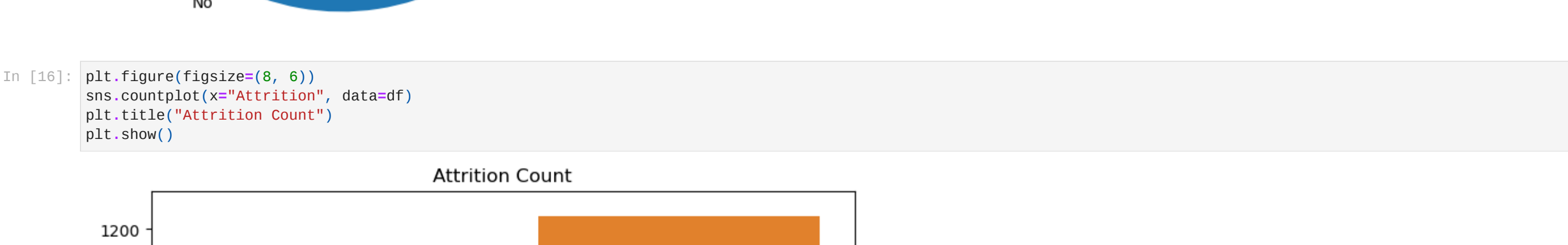
In [15]: attrition_counts = df['Attrition'].value_counts()
plt.figure(figsize=(6, 6))
plt.pie(attrition_counts, labels=attrition_counts.index, autopct='%1.1f%%', startangle=90)
plt.xlabel('Attrition Distribution')
plt.axis('equal')

plt.show()
```



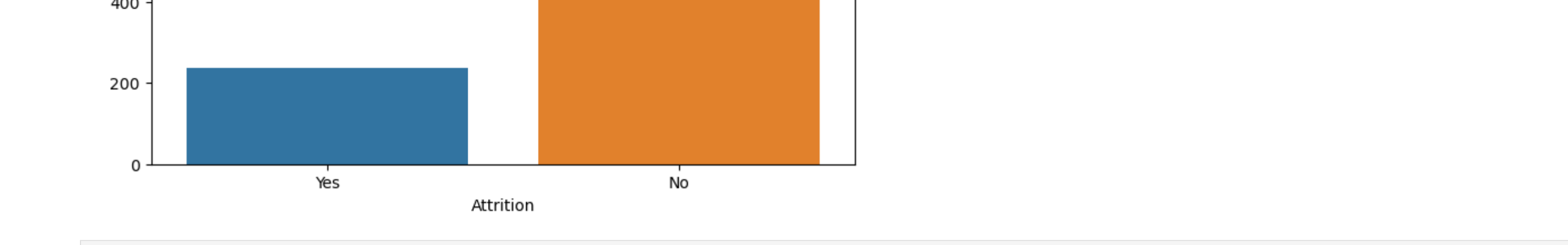
```

In [16]: plt.figure(figsize=(6, 6))
sns.countplot(x='Attrition', data=df)
plt.xlabel('Attrition Count')
plt.show()
```



```

In [17]: plt.figure(figsize=(6, 6))
sns.histplot(data=df, x='Age', kde=True)
plt.xlabel('Distribution of Age')
plt.show()
```



Outlier Detection

```

In [18]: plt.figure(figsize=(35, 8))
sns.boxplot(data=df)
plt.title('Box Plot for all the attributes')
plt.show()
```

