

Assignment- HR Attrition Data

In this assignment, we need to predict whether a give employee will leave the organization or not. Your target column is Attrition We will create a model with the following steps: Import the relevant packages Download and explore the dataset Perform EDA, Apply dataset for preprocessing Predict the target columns

In [1]:

```
#IMPORTING NECESSARY PACKAGE
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

In [2]:

```
dataset = pd.read_csv('HR_Employee_Attrition-1.csv')
```

In [3]:

```
dataset.head()
```

Out [3]:

EmployeeNumber	Attrition	Age	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EmployeeCount	...	RelationshipSatisfaction	StandardHours	
0	1	Yes	41	Travel_Frequently	1102	Sales	1	2	Life Sciences	1	...	1	80
1	2	No	49	Travel_Rarely	279	Research & Development	8	1	Life Sciences	1	...	4	80
2	3	Yes	37	Travel_Rarely	1373	Research & Development	2	2	Other	1	...	2	80
3	4	No	33	Travel_Frequently	1392	Research & Development	3	4	Life Sciences	1	...	3	80
4	5	No	27	Travel_Rarely	591	Research & Development	2	1	Medical	1	...	4	80

5 rows × 35 columns

In [4]:

```
dataset.tail()
```

Out [4]:

2935	2936	No	36	Travel_Frequently	884	Research & Development	23	2	Medical	1	...	3	
2936	2937	No	39	Travel_Rarely	613	Research & Development	6	1	Medical	1	...	1	
2937	2938	No	27	Travel_Rarely	155	Research & Development	4	3	Life Sciences	1	...	2	
2938	2939	No	49	Travel_Frequently	1023	Sales	2	3	Medical	1	...	4	
2939	2940	No	34	Travel_Rarely	628	Research & Development	8	3	Medical	1	...	1	

5 rows × 35 columns

In [5]:

```
dataset.shape
```

Out [5]:

```
(2940, 35)
```

In [6]:

```
dataset.columns
```

Out [6]:

```
Index(['EmployeeNumber', 'Attrition', 'Age', 'BusinessTravel', 'DailyRate',
       'Department', 'DistanceFromHome', 'Education', 'EducationField',
       'EmployeeCount', '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 [7]:

```
dataset.describe()
```

Out [7]:

	EmployeeNumber	Age	DailyRate	DistanceFromHome	Education	EmployeeCount	EnvironmentSatisfaction	HourlyRate	JobInvolvement	JobLevel	...	RelationshipSatisfaction	StandardHours
count	2940.000000	2940.000000	2940.000000	2940.000000	2940.000000	2940.00	2940.000000	2940.000000	2940.000000	2940.000000	...	2940.000000	2940.000000
mean	1470.500000	36.823810	802.485714	9.333819	18.0	36.0	80.0	4230.428571	65.891156	2.729932	...	2.063946	80.0
std	848.849221	9.133819	403.440447	8.105485	1.023991	0.0	1.002896	20.325969	0.711440	1.106752	...	1.106752	...
min	1.000000	18.000000	102.000000	1.000000	1.000000	1.0	1.000000	30.000000	1.000000	1.000000	...	1.000000	...
25%	735.750000	30.000000	465.000000	2.000000	2.000000	1.0	2.000000	40.000000	2.000000	2.000000	...	2.000000	...
50%	1470.500000	36.000000	802.000000	7.000000	3.000000	1.0	3.000000	66.000000	3.000000	2.000000	...	3.000000	...
75%	2205.250000	43.000000	1157.000000	14.000000	4.000000	1.0	4.000000	84.000000	3.000000	3.000000	...	3.000000	...
max	2940.000000	60.000000	1499.000000	29.000000	5.000000	1.0	4.000000	100.000000	4.000000	5.000000	...	5.000000	...

8 rows × 20 columns

In [8]:

```
dataset.isnull().sum()
```

Out [8]:

EmployeeNumber	0
Attrition	0
Age	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
StockOptionLevel	0
TotalWorkingYears	0
TrainingTimesLastYear	0
WorkLifeBalance	0
YearsAtCompany	0
YearsInCurrentRole	0
YearsSinceLastPromotion	0
YearsWithCurrManager	0
dtype: int64	

In [9]:

```
dataset['Attrition'].value_counts()
```

Out [9]:

```
No      2466
Yes      474
Name: Attrition, dtype: int64
```

In [10]:

```
dataset.describe().T
```

Out [10]:

	count	mean	std	min	25%	50%	75%	max
EmployeeNumber	2940.0	1470.500000	848.849221	1.0	735.75	1470.5	2205.25	2940.0
Age	2940.0	36.823819	9.133819	18.0	30.00	36.0	43.00	60.00
DailyRate	2940.0	802.485714	403.440447	102.0	465.00	802.0	1157.00	1499.00
DistanceFromHome	2940.0	9.333819	8.105485	1.0	2.00	7.0	14.00	29.00
Education	2940.0	2.912905	1.023991	1.0	2.00	3.0	4.00	5.0
EmployeeCount	2940.0	1.000000	0.000000	1.0	1.00	1.0	1.00	1.0
EnvironmentSatisfaction	2940.0	2.721769	1.002896	1.0	2.00	3.0	4.00	4.0
HourlyRate	2940.0	65.891156	20.325969	30.0	40.00	66.0	84.00	100.00
JobInvolvement	2940.0	2.729932	0.711440	1.0	2.00	3.0	3.00	4.0
JobLevel	2940.0	2.063946	1.106752	1.0	1.00	2.0	3.00	5.0
JobSatisfaction	2940.0	2.728571	1.102658	1.0	2.00	3.0	4.00	4.0
MonthlyIncome	2940.0	6502.931293	4707.155770	1009.0	2911.00	4919.0	8380.00	19999.00
MonthlyRate	2940.0	14313.103491	7116.576021	2094.0	6045.00	14235.5	20462.00	26999.00
NumCompaniesWorked	2940.0	2.693197	2.497584	0.0	1.00	1.00	2.00	9.0
PercentSalaryHike	2940.0	2.502624	3.659315	1.0	12.00	14.0	18.00	25.0
PerformanceRating	2940.0	3.153741	0.360762	3.0	3.00	3.0	3.00	4.0
RelationshipSatisfaction	2940.0	2.712245	1.081025	1.0	2.00	3.0	4.00	4.0
StandardHours	2940.0	80.000000	0.000000	80.0	80.0	80.0	80.00	80.00
StockOptionLevel	2940.0	0.793878	0.851932	0.0	0.00	1.0	1.00	3.0
TotalWorkingYears	2940.0	11.279592	7.779458	0.0	6.00	10.0	15.00	40.0
TrainingTimesLastYear	2940.0	2.799320	1.289051	0.0	2.00	3.0	3.00	6.0
WorkLifeBalance	2940.0	2.761224	0.706356	1.0	2.00	3.0	3.00	4.0
YearsAtCompany	2940.0	7.008163	6.125483	0.0	3.00	5.0	9.00	40.0
YearsInCurrentRole	2940.0	4.229252	3.625251	0.0	2.00	3.0	7.00	18.0
YearsSinceLastPromotion	2940.0	2.187755	3.221822	0.0	0.00	1.0	3.00	15.0
YearsWithCurrManager	2940.0	4.123129	3.567529	0.0	2.00	3.0	7.00	17.0

In [11]:

```
dataset.nunique()
```

Out [11]:

EmployeeNumber	2940
Attrition	2
Age	43
BusinessTravel	3
DailyRate	886
Department	3
DistanceFromHome	29
Education	5
EducationField	6
EmployeeCount	1
EnvironmentSatisfaction	4
Gender	2
HourlyRate	71
JobInvolvement	4
JobLevel	5
JobRole	9
JobSatisfaction	4
MaritalStatus	3
MonthlyIncome	1349
MonthlyRate	1427
NumCompaniesWorked	19
Over18	2
OverTime	2
PercentSalaryHike	15
PerformanceRating	2
RelationshipSatisfaction	4
StandardHours	1
StockOptionLevel	4
TotalWorkingYears	49
TrainingTimesLastYear	7
WorkLifeBalance	4
YearsAtCompany	37
YearsInCurrentRole	19
YearsSinceLastPromotion	16
YearsWithCurrManager	18
dtype:	int64

In [12]:

```
dataset.dtypes
```

Out [12]:

EmployeeNumber	int64
Attrition	object
Age	int64
BusinessTravel	object
DailyRate	int64
Department	object
DistanceFromHome	object
Education	int64
EducationField	object
EmployeeCount	int64
EnvironmentSatisfaction	int64
Gender	object
HourlyRate	int64
JobInvolvement	int64
JobLevel	int64
JobRole	object
JobSatisfaction	int64
MaritalStatus	object
MonthlyIncome	int64
MonthlyRate	int64
NumCompaniesWorked	int64
Over18	object
OverTime	object
PercentSalaryHike	int64
PerformanceRating	int64
RelationshipSatisfaction	int64
StandardHours	int64
StockOptionLevel	int64
TotalWorkingYears	int64
TrainingTimesLastYear	int64
WorkLifeBalance	int64
YearsAtCompany	int64
YearsInCurrentRole	int64
YearsSinceLastPromotion	int64
YearsWithCurrManager	int64
dtype:	object

In [13]:

```
f,ax=plt.subplots(figsize=(10,10))
ax = dataset['Attrition'].value_counts().plot.pie(explode=[0,0],autopct = '%1f.1f%%',shadow=True)
ax.set_title('Attrition Probability')
```

Out [13]:

Attrition Probability

A pie chart titled 'Attrition Probability' showing the distribution of employee attrition. The 'No' category (blue) represents 82.9% of the total, and the 'Yes' category (orange) represents 17.1%.

In [14]:

```
sns.countplot(dataset['Attrition'])
```

Out [14]:

C:\Users\PC\anaconda3\Lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be "data", and passing other arguments without an explicit keyword will result in an error or misinterp...

warnings.warn(

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

A bar chart showing the count of employees for each attrition status. The x-axis is labeled 'Attrition' with categories 'Yes' and 'No'. The y-axis is labeled 'count' and ranges from 0 to 2500. The 'Yes' bar is blue and has a height of approximately 474. The 'No' bar is orange and has a height of approximately 2466.

In [15]:

```
plt.subplots(figsize=(12,4))
sns.countplot(x='YearsAtCompany',hue='Attrition',data=dataset,palette='colorblind')
```

Out [15]:

<AxesSubplot: xlabel='YearsAtCompany', ylabel='count'>

A bar chart showing the count of employees for each 'YearsAtCompany' (x-axis, 0 to 40) and 'Attrition' status (y-axis, 0 to 350). The bars are colored by 'Attrition' status: 'Yes' (blue) and 'No' (orange). The chart shows that the number of employees who attrite is generally higher in the first few years of service, particularly around year 5.

In [16]:

```
plt.subplots(figsize=(12,4))
sns.countplot(x='Age',hue='Attrition', data=dataset,palette = 'colorblind')
```

Out [16]:

<AxesSubplot: xlabel='Age', ylabel='count'>

A bar chart showing the count of employees for each 'Age' (x-axis, 18 to 60) and 'Attrition' status (y-axis, 0 to 140). The bars are colored by 'Attrition' status: 'Yes' (blue) and 'No' (orange). The chart shows that attrition is more prevalent in the 20-30 age range.

In [17]:

```
for column in dataset.columns:
    if dataset[column].dtype == object:
        print(str(column))
        print(str(dataset[column].nunique()))
        print(dataset[column].value_counts())
        print('-----')
```

Attrition: ['Yes' 'No']
No 2466
Yes 474
Name: Attrition, dtype: int64

BusinessTravel: ['Travel_Rarely' 'Travel_Frequently' 'Non-Travel']
Travel_Rarely 2089
Travel_Frequently 564
Non-Travel 386
Name: BusinessTravel, dtype: int64

Department: ['Sales' 'Research & Development' 'Human Resources']
Research & Development 1922
Sales 892
Human Resources 126
Name: Department, dtype: int64

EducationField: ['Life Sciences' 'Other' 'Medical' 'Marketing' 'Technical Degree']
Life Sciences 1212
Medical 928
Marketing 318
Technical Degree 284
Other 164
Name: EducationField, dtype: int64

Gender: ['Female' 'Male']
Male 1764
Female 2176
Name: Gender, dtype: int64

JobRole: ['Sales Executive' 'Research Scientist' 'Laboratory Technician' 'Manufacturing Director' 'Healthcare Representative' 'Manager' 'Sales Representative' 'Research Director' 'Human Resources']
Sales Executive 652
Research Scientist 584
Laboratory Technician 519
Manufacturing Director 299
Healthcare Representative 262
Manager 294
Sales Representative 166
Research Director 169
Human Resources 194
Name: JobRole, dtype: int64

MaritalStatus: ['Single' 'Married' 'Divorced']
Married 1366
Single 940
Divorced 654
Name: MaritalStatus, dtype: int64

Over18: ['Y']
Y 2948
Name: Over18, dtype: int64

OverTime: ['Yes' 'No']
No 2108
Yes 832
Name: OverTime, dtype: int64

In [18]:

```
dataset['EmployeeNumber'].unique()
```

Out [18]:

```
array([ 1,  2,  3, ..., 2938, 2939, 2940], dtype=int64)
```

In [20]:

```
dataset['StandardHours'].unique()
```

Out [20]:

```
array([80], dtype=int64)
```

In [21]:

```
dataset = dataset.drop('EmployeeCount',axis = 1)
dataset = dataset.drop('EmployeeNumber',axis = 1)
dataset = dataset.drop('StandardHours',axis = 1)
```

In [22]:

```
dataset.head()
```

Out [22]:

Attrition	Age	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EnvironmentSatisfaction	Gender	...	PerformanceRating	RelationshipSatisfaction
0	Yes	41	Travel_Frequently	1102	Sales	1	2	Life Sciences	2	Female	...	3
1	No	49	Travel_Rarely	279	Research & Development	8	1	Life Sciences	3	Male	...	4
2	Yes	37	Travel_Rarely	1373	Research & Development	2	2	Other	4	Male	...	3
3	No	33	Travel_Frequently	1392	Research & Development	3	4	Life Sciences	4	Female	...	3
4	No	27	Travel_Rarely	591	Research & Development	2	1	Medical	1	Male	...	3

5 rows × 32 columns

In [25]:

```
dataset
```

Out [25]:

Attrition	Age	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EnvironmentSatisfaction	Gender	...	PerformanceRating	RelationshipSatisfaction
0	Yes	41	Travel_Frequently	1102	Sales	1	2	Life Sciences	2	Female	...	3
1	No	49	Travel_Frequently	279	Research & Development	8	1	Life Sciences	3	Male	...	4
2	Yes	37	Travel_Rarely	1373	Research & Development	2	2	Other	4	Male	...	3
3	No	33	Travel_Frequently	1392	Research & Development	3	4	Life Sciences	4	Female	...	3
4	No	27	Travel_Rarely	591	Research & Development	2	1	Medical	1	Male	...	3
...
2935	No	36	Travel_Frequently	884	Research & Development	23	2	Medical	3	Male	...	3
2936	No	39	Travel_Rarely	613	Research & Development	6	1	Medical	4	Male	...	3
2937	No	27	Travel_Rarely	155	Research & Development	4	3	Life Sciences	2	Male	...	4
2938	No	49	Travel_Frequently	1023	Sales	2	3	Medical	4	Male	...	3
2939	No	34	Travel_Rarely	628	Research & Development	8	3	Medical	2	Male	...	3

2940 rows × 31 columns

In [26]:

```
dataset.corr()
```

Out [26]:

	Age	DailyRate	DistanceFromHome	Education	EnvironmentSatisfaction	HourlyRate	JobInvolvement	JobLevel	JobSatisfaction	MonthlyIncome	...	PerformanceRating	RelationshipSatisfaction	Perfo
Age	1.000000	0.010961	-0.001686	0.208034	0.010146	0.024287	0.029820	0.509604	-0.048922	0.497885	...	0.497885
DailyRate	0.010961	1.000000	-0.004985	0.018066	0.020875	0.022381	0.046135	0.002965	0.030571	0.007707	...	0.007707
DistanceFromHome	-0.001686	-0.004985	1.000000	0.021042	0.001840	-0.010775	0.011131	0.008783	0.005303	4.803669	...	-0.017614
Education	0.208034	0.018066	0.021042	1.000000	-0.007128	0.016775	0.042436	0.310599	-0.011296	0.094961	...	0.094961
EnvironmentSatisfaction	0.010146	0.018066	-0.001840	-0.007128	1.000000	-0.049857	0.008278	0.002112	0.006784	-0.005619	...	-0.005619
HourlyRate	0.024287	0.022381	0.011131	0.016775	-0.049857	1.000000	0.042861	0.027863	-0.021325	0.015794	...	0.015794