

# HR Analytics and Employee Attrition Prediction

## Importing Libraries

```
In [2]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
from sklearn.preprocessing import StandardScaler
from sklearn.impute import SimpleImputer
```

## Loading Data

```
In [4]: df = pd.read_csv("C:\\\\Users\\\\byash\\\\HR Analytics and Employee Attrition Prediction
```

```
In [5]: df
```

Out[5]:

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education
0	41	Yes	Travel_Rarely	1102	Sales		1
1	49	No	Travel_Frequently	279	Research & Development		8
2	37	Yes	Travel_Rarely	1373	Research & Development		2
3	33	No	Travel_Frequently	1392	Research & Development		3
4	27	No	Travel_Rarely	591	Research & Development		2
...	...	...	...	...	...	...	...
1465	36	No	Travel_Frequently	884	Research & Development		23
1466	39	No	Travel_Rarely	613	Research & Development		6
1467	27	No	Travel_Rarely	155	Research & Development		4
1468	49	No	Travel_Frequently	1023	Sales		2
1469	34	No	Travel_Rarely	628	Research & Development		8

1470 rows × 35 columns



## Data Information and Data cleaning

In [6]: `df.head()`

Out[6]:

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

5 rows × 35 columns



In [10]: `df.columns`

Out[10]:

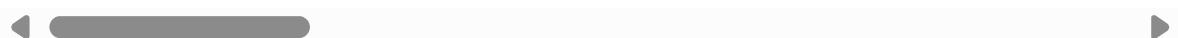
```
Index(['Age', 'Attrition', 'BusinessTravel', 'DailyRate', 'Department',
       'DistanceFromHome', 'Education', 'EducationField', 'EmployeeCount',
       'EmployeeNumber', 'EnvironmentSatisfaction', 'Gender', 'HourlyRate',
       'JobInvolvement', 'JobLevel', 'JobRole', 'JobSatisfaction',
      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

In [12]: `df.describe()`

Out[12]:

	Age	DailyRate	DistanceFromHome	Education	EmployeeCount	Employ
<b>count</b>	1470.000000	1470.000000	1470.000000	1470.000000	1470.0	1
<b>mean</b>	36.923810	802.485714	9.192517	2.912925	1.0	1
<b>std</b>	9.135373	403.509100	8.106864	1.024165	0.0	
<b>min</b>	18.000000	102.000000	1.000000	1.000000	1.0	
<b>25%</b>	30.000000	465.000000	2.000000	2.000000	1.0	
<b>50%</b>	36.000000	802.000000	7.000000	3.000000	1.0	1
<b>75%</b>	43.000000	1157.000000	14.000000	4.000000	1.0	1
<b>max</b>	60.000000	1499.000000	29.000000	5.000000	1.0	2

8 rows × 26 columns



```
In [13]: 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 [14]: df.isnull().sum()
```

```
Out[14]: Age          0
Attrition      0
BusinessTravel 0
DailyRate       0
Department      0
DistanceFromHome 0
Education        0
EducationField   0
EmployeeCount    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
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 [16]: df.duplicated().sum()
```

```
Out[16]: 0
```

## EDA(Exploratory data analysis )

Converting columns to appropriate data types and Convert categorical columns to category dtype

```
In [38]: for col in df.columns:
    if col not in ['Attrition', 'BusinessTravel', 'Department', 'EducationField', 'EducationField']:
        df[col] = pd.to_numeric(df[col], errors='coerce')
categorical_cols = ['Attrition', 'BusinessTravel', 'Department', 'EducationField', 'EducationField']
for col in categorical_cols:
    df[col] = df[col].astype('category')
```

## 1. Attrition Distribution

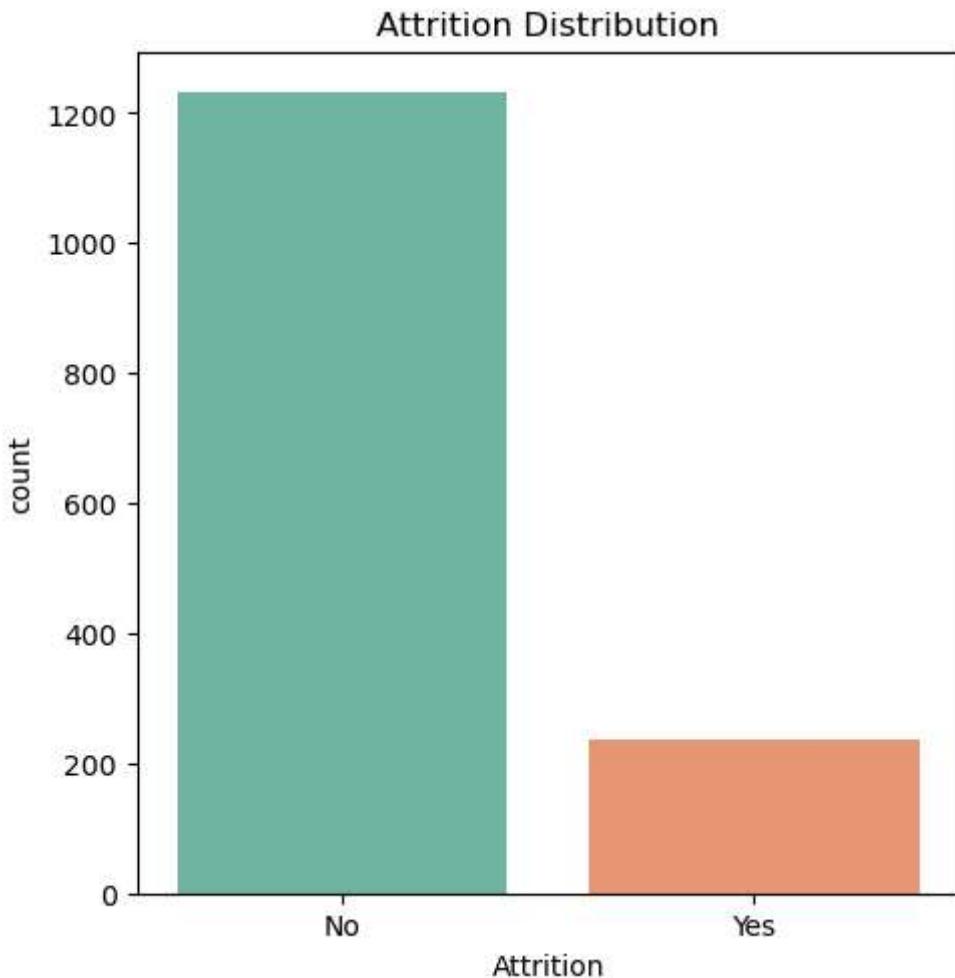
```
In [18]: plt.figure(figsize=(18, 12))
plt.subplot(2, 3, 1)
sns.countplot(data=df, x='Attrition', palette='Set2')
plt.title('Attrition Distribution')
```

C:\Users\byash\AppData\Local\Temp\ipykernel\_20496\2424425548.py:4: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

```
    sns.countplot(data=df, x='Attrition', palette='Set2')
```

Out[18]: Text(0.5, 1.0, 'Attrition Distribution')



## 2. Age Distribution by Attrition

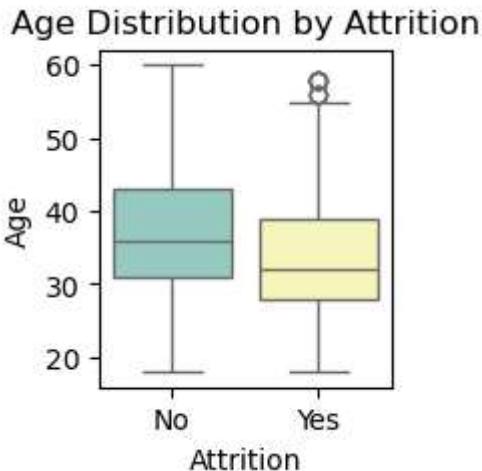
```
In [39]: plt.subplot(2, 3, 2)
sns.boxplot(data=df, x='Attrition', y='Age', palette='Set3')
plt.title('Age Distribution by Attrition')
```

```
C:\Users\byash\AppData\Local\Temp\ipykernel_20496\2314402736.py:2: FutureWarning:
```

```
Passing `palette` without assigning `hue` is deprecated and will be removed in v0.1
4.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.
```

```
sns.boxplot(data=df, x='Attrition', y='Age', palette='Set3')
```

```
Out[39]: Text(0.5, 1.0, 'Age Distribution by Attrition')
```



### 3. Monthly Income by Attrition

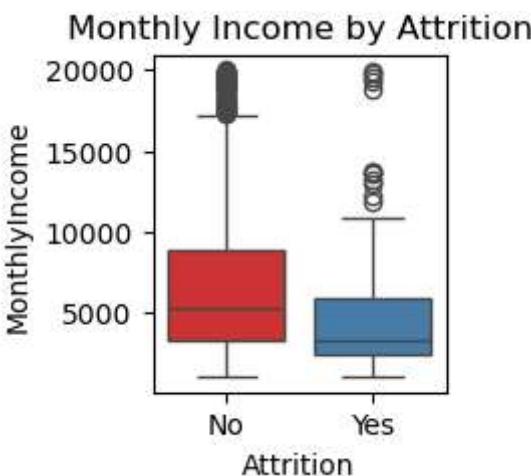
```
In [21]: plt.subplot(2, 3, 3)
sns.boxplot(data=df, x='Attrition', y='MonthlyIncome', palette='Set1')
plt.title('Monthly Income by Attrition')
```

```
C:\Users\byash\AppData\Local\Temp\ipykernel_20496\2874951252.py:3: FutureWarning:
```

```
Passing `palette` without assigning `hue` is deprecated and will be removed in v0.1
4.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.
```

```
sns.boxplot(data=df, x='Attrition', y='MonthlyIncome', palette='Set1')
```

```
Out[21]: Text(0.5, 1.0, 'Monthly Income by Attrition')
```



### 4. Job Satisfaction by Attrition

```
In [22]: plt.subplot(2, 3, 4)
sns.boxplot(data=df, x='Attrition', y='JobSatisfaction', palette='Set2')
plt.title('Job Satisfaction by Attrition')
```

C:\Users\byash\AppData\Local\Temp\ipykernel\_20496\3358719107.py:3: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.1 4.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

```
sns.boxplot(data=df, x='Attrition', y='JobSatisfaction', palette='Set2')
```

```
Out[22]: Text(0.5, 1.0, 'Job Satisfaction by Attrition')
```



## 5. Distance From Home by Attrition

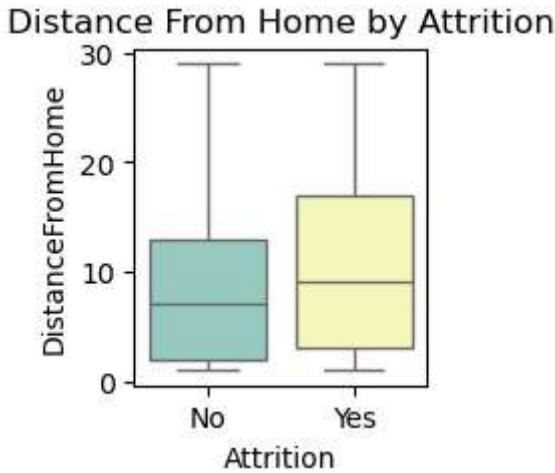
```
In [41]: plt.subplot(2, 3, 5)
sns.boxplot(data=df, x='Attrition', y='DistanceFromHome', palette='Set3')
plt.title('Distance From Home by Attrition')
```

C:\Users\byash\AppData\Local\Temp\ipykernel\_20496\2363086863.py:2: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.1 4.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

```
sns.boxplot(data=df, x='Attrition', y='DistanceFromHome', palette='Set3')
```

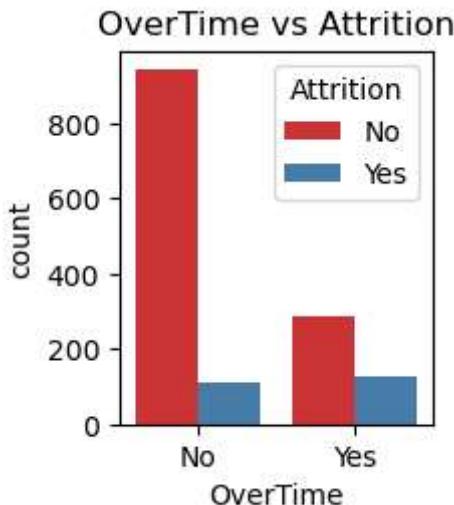
```
Out[41]: Text(0.5, 1.0, 'Distance From Home by Attrition')
```



## 6. OverTime vs Attrition

```
In [42]: plt.subplot(2, 3, 6)
sns.countplot(data=df, x='OverTime', hue='Attrition', palette='Set1')
plt.title('OverTime vs Attrition')

plt.tight_layout()
plt.show()
```



# Logistic Regression

```
In [25]: features = ['Age', 'DistanceFromHome', 'Education', 'EnvironmentSatisfaction', 'JobLevel', 'JobSatisfaction', 'MonthlyIncome', 'NumCompaniesWorked', 'PerformanceRating', 'RelationshipSatisfaction', 'StockOptionLevel', 'TrainingTimesLastYear', 'WorkLifeBalance', 'YearsAtCompany', 'YearsInCurrentRole', 'YearsSinceLastPromotion', 'YearsWithCurrManager']

X = df[features]
y = df['Attrition']
```

## Imputation

```
In [28]: imputer = SimpleImputer(strategy='median')
X_imputed = imputer.fit_transform(X)
```

## Split data into training and testing sets

```
X_train, X_test, y_train, y_test = train_test_split(X_imputed, y, test_size=0.3, random_state=42)
```

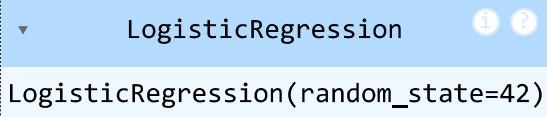
## Scale features

```
In [44]: scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

## Create and train logistic regression model

```
In [45]: model = LogisticRegression(random_state=42)
model.fit(X_train_scaled, y_train)
```

Out[45]:

A screenshot of an IPython notebook cell showing the creation of a LogisticRegression object. The cell starts with 'Out[45]:' followed by a blue box containing a dropdown menu with a downward arrow, the text 'LogisticRegression', and two small circular icons. Below the box is the code 'LogisticRegression(random\_state=42)'.

## Make predictions

```
In [32]: y_pred = model.predict(X_test_scaled)
y_pred
```

## Feature importance

```
In [48]: coefficients = model.coef_[0]
feature_importance = pd.DataFrame({'Feature': features, 'Coefficient': coefficients})
feature_importance['Absolute'] = feature_importance['Coefficient'].abs()
feature_importance = feature_importance.sort_values(by='Absolute', ascending=False)
```

## Plot feature importance

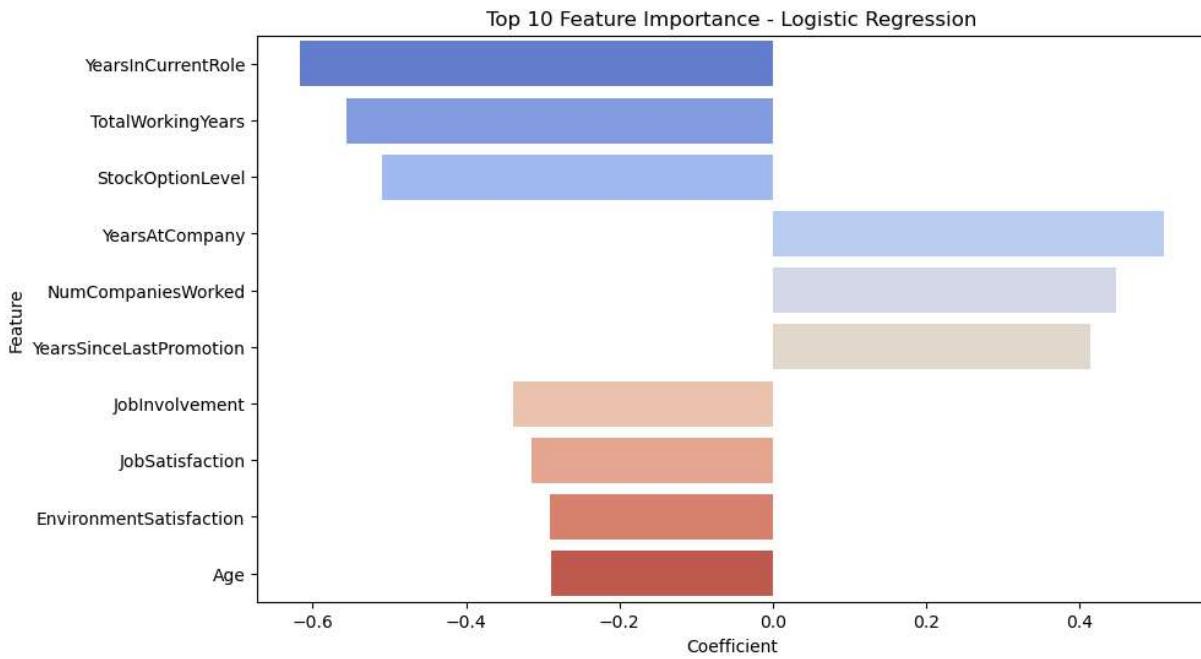
```
In [47]: plt.figure(figsize=(10, 6))
sns.barplot(x='Coefficient', y='Feature', data=feature_importance.head(10), palette
```

```
plt.title('Top 10 Feature Importance - Logistic Regression')
plt.show()
```

C:\Users\byash\AppData\Local\Temp\ipykernel\_20496\3376853221.py:2: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.1 4.0. Assign the `y` variable to `hue` and set `legend=False` for the same effect.

```
sns.barplot(x='Coefficient', y='Feature', data=feature_importance.head(10), palette='coolwarm')
```



In [ ]: